

TRABAJO FIN DE MASTER

ENHANCING THE EXFOR NUCLEAR DATA LIBRARY USING MACHINE LEARNING TECHNIQUES

TRABAJO FIN DE MASTER
PARA LA OBTENCIÓN DEL
TÍTULO DE MASTER EN
CIENCIA Y TECNOLOGÍA
NUCLEAR

SEPTIEMBRE 2023

**Juan Antonio Monleón de la
Lluvia Mazagatos**

DIRECTORES DEL TRABAJO FIN DE MASTER:

Oscar Cabellos

“All models are wrong, but some are useful.”

- George E. P. Box

ACKNOWLEDGMENTS

I would like to express my sincere gratitude to my mother, for her continuous support and encouragement throughout my life, and for always believing in me, even when I didn't believe in myself.

I am also grateful to my thesis mentor, Dr. Oscar Cabellos, for giving me the opportunity to work on this project and for his guidance and support throughout the process. His expertise and knowledge were invaluable to me.

I would also like to thank the Nuclear Department of the UPM, to the Cátedra de Seguridad Nuclear “Federico Goded” and Eduardo Gallego, for providing me with a grant that made this research possible.

Finally, I would like to thank all of the other people who have helped me along the way, including my friends, family, colleagues, and mentors. I could not have done this without you.

RESUMEN EJECUTIVO

La detección de valores atípicos es crucial en la física experimental, especialmente en el ámbito de las reacciones nucleares. Estos valores pueden sesgar análisis y generar resultados poco fiables. Esta problemática cobra especial importancia en el tratamiento de datos relativos a protones incidentes, los cuales son cada vez más relevantes en terapias médicas como el tratamiento de cáncer mediante protonterapia. En esta técnica, se usan haces de protones para irradiar células cancerígenas con alta precisión, minimizando daños a los tejidos sanos circundantes. Debido a su eficacia y precisión, la protonterapia se está consolidando como una opción de tratamiento preferente para ciertos tipos de cáncer. Sin embargo, la presencia de valores atípicos en los datos experimentales sobre protones incidentes puede comprometer la calidad y seguridad de estos tratamientos. Por tanto, es fundamental llevar a cabo un riguroso proceso de detección y eliminación de valores atípicos para asegurar resultados fiables y precisos en el ámbito de la física nuclear aplicada a la medicina y demás campos.

El principal objetivo de esta tesis de máster es evaluar la aplicabilidad de diversas técnicas de aprendizaje automático para la detección de valores atípicos en la biblioteca de datos nucleares EXFOR. No se trata solo de implementar algoritmos, sino de adaptarlos eficazmente al conjunto de datos especializado y evaluar cuáles resultan más eficaces.

Un elemento clave de la investigación es el preprocesado de datos, dado que es esencial para el éxito de cualquier algoritmo de aprendizaje automático. Dada la complejidad de la estructura de datos de EXFOR, se dedicarán esfuerzos significativos para investigar métodos eficientes de extracción, transformación y preparación de los datos. Sobre esta base, se implementarán diversos algoritmos para la detección de valores atípicos, incluidos métodos estadísticos tradicionales y técnicas más avanzadas como redes neuronales, en este caso un Autoencoder. Posteriormente, se compararán críticamente los métodos y resultados mediante análisis visuales y cuantitativos, con el fin de evaluar la fiabilidad, ventajas y limitaciones de cada enfoque.

El estudio tiene el objetivo de contextualizar estos hallazgos en la literatura científica, ofreciendo perspectivas sobre los retos y el potencial futuro de aplicar el aprendizaje automático para la detección de valores atípicos en datos nucleares.

Con el contexto y los objetivos ya establecidos, es importante profundizar en algunos de los conceptos más esenciales para esta investigación, como son la librería de datos nucleares EXFOR, la detección de valores atípicos y la naturaleza del aprendizaje automático.

EXFOR (EXchange FORmat) es una biblioteca extensa y meticulosa de datos de reacciones nucleares experimentales, desarrollada para facilitar la investigación y aplicaciones en el campo nuclear. Creada para promover el intercambio de datos entre centros de investigación, EXFOR es gestionada por la Red Internacional de Centros de Datos de Reacciones Nucleares (NRDC), bajo supervisión de la Sección de Datos Nucleares de la AIEA. Con más de 24,000 experimentos registrados hasta octubre de 2022, la plataforma busca ser una fuente completa y confiable de datos en reacciones inducidas por neutrones a baja energía, partículas cargadas y fotones. Este recurso es especialmente relevante para investigadores y profesionales en ciencias y tecnologías nucleares, incluyendo a académicos, laboratorios nacionales e industria, quienes se centran en el diseño y análisis de reactores nucleares, ciclos de combustible y otras aplicaciones nucleares.

El Aprendizaje Automático (AA) es una subdisciplina de la Inteligencia Artificial (IA) que se enfoca en diseñar sistemas que puedan aprender a tomar decisiones o hacer predicciones a partir de datos. A diferencia de programar tareas de manera manual, estos sistemas mejoran su rendimiento de forma adaptativa al ser expuestos a nuevos datos. En el mundo actual, saturado

de información, la capacidad del AA para manejar conjuntos de datos grandes y complejos ha resultado ser invaluable. Además, se pueden distinguir principalmente tres tipos de aprendizaje automático: el aprendizaje supervisado, donde el modelo se entrena con datos etiquetados; el aprendizaje no supervisado, que trabaja con datos no etiquetados para encontrar patrones; y el aprendizaje por refuerzo, donde un agente aprende a tomar decisiones mediante ensayo y error.

Estos tres tipos de aprendizaje automático, aunque distintos, comparten un objetivo común: aprender de los datos para hacer decisiones inteligentes o predicciones. La principal diferencia radica en el tipo de retroalimentación que reciben y el contexto de aprendizaje, que puede ir desde tener un mapeo claro de entradas a salidas en aprendizaje supervisado, hasta aprender de una serie de acciones y recompensas en un contexto dinámico en el caso del aprendizaje por refuerzo.

Es importante además diferenciar entre los términos relacionados pero no idénticos de Inteligencia Artificial, Aprendizaje Automático y Aprendizaje Profundo. La IA es el término más amplio que abarca cualquier técnica que permita a las máquinas imitar la inteligencia humana. El Aprendizaje Automático es una especialización dentro de la IA que usa técnicas estadísticas para mejorar con la experiencia. El Aprendizaje Profundo, por otro lado, es un subcampo del Aprendizaje Automático que se inspira en la estructura y función del cerebro humano y se aplica en tecnologías como los coches autónomos para tareas como la identificación de objetos.

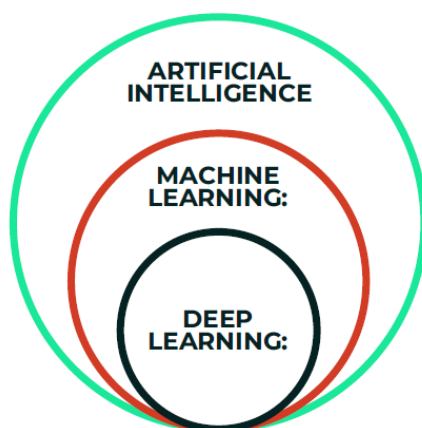


Figure 0.1: Estructura anidada de la Inteligencia Artificial [1].

Los valores atípicos son observaciones que se desvían significativamente del resto de los datos y pueden afectar los resultados del análisis de forma adversa. Su detección y gestión son cruciales en cualquier análisis de datos. Existen dos enfoques principales para detectarlos: métodos estadísticos y métodos basados en aprendizaje automático. Los métodos estadísticos, como el Z-Score y el IQR, son más adecuados para datos univariantes, mientras que las técnicas de aprendizaje automático, como los algoritmos de agrupación y detección de anomalías, son más versátiles y no requieren suposiciones sobre la distribución de los datos. La elección del método depende de las características del conjunto de datos y del objetivo del análisis. Sin embargo, es esencial tener en cuenta que no todos los valores atípicos son indeseables o errores. En algunos contextos, pueden señalar áreas para una investigación más detallada, especialmente si los datos son escasos o emergentes en esa área particular. Por lo tanto, el tratamiento de los valores atípicos debe hacerse con un entendimiento cuidadoso de su contexto.

La extracción y el preprocesamiento de datos constituyen, probablemente, la etapa más crucial en el proceso de implementación de algoritmos de aprendizaje automático (machine learning).

Contar con una base de datos limpia y adecuadamente estructurada es fundamental para el óptimo rendimiento de estos algoritmos. A continuación, se describen los aspectos clave de la metodología desarrollada para esta tarea, desde la selección de la fuente de datos hasta la generación de la base de datos final sobre la cual se aplicarán técnicas de aprendizaje automático.

- **Fuente de Datos:** El conjunto de datos principal se extrae de la biblioteca de datos nucleares EXFOR, que es una colección exhaustiva pero compleja que abarca más de 24,000 experimentos con características variadas.
- **Análisis de Datos con EXFORTABLES:** Para facilitar la manipulación y la legibilidad de los datos de EXFOR, se utiliza EXFORTABLES para descomponer la biblioteca en directorios estructurados lógicamente, estableciendo así un punto de partida efectivo para la extracción y manipulación de datos.
- **Comprensión Objetiva:** El primer paso consiste en identificar el tipo específico de información que se necesita para las aplicaciones de aprendizaje automático. El enfoque se centra en determinar la estructura en la que dicha información debería almacenarse para maximizar su utilidad.
- **Estructuración de Datos en DataFrames:** Se utilizan los dataframes de Python como la estructura de datos preferente, facilitando tanto la manipulación de datos como su almacenamiento eficiente. En estos dataframes, cada fila representa un punto de datos experimental y cada columna representa una característica.
- **Codificación de Datos:** Debido a que los algoritmos de aprendizaje automático suelen requerir datos numéricos, los valores categóricos presentes en el conjunto de datos se transforman en representaciones numéricas, a menudo mediante la técnica de codificación *One-hot*.
- **Diseño Orientado a Objetos en Python:** Se implementa una clase llamada **Experiment** para encapsular todos los atributos y comportamientos relacionados con un experimento individual. Esto permite aprovechar las ventajas inherentes a la Programación Orientada a Objetos (OOP), como la modularidad y la reutilización de código.
- **Capacidades Funcionales:** Se desarrollan diversas funciones con el objetivo de manipular y consultar el conjunto de datos de manera eficaz. Estas funciones incluyen, pero no se limitan a, el filtrado de atributos, la clasificación de experimentos y las capacidades para generar gráficos.
- **Agrupación de Datos para Análisis:** Con un enfoque especial en la detección de valores atípicos, los datos se agrupan en función de la variable que se está midiendo, lo que facilita las etapas posteriores del análisis de datos.
- **Conversión de Datos en DataFrame Final:** Finalmente, los datos organizados y clasificados se consolidan en un dataframe final, conforme a las especificaciones previamente establecidas.

Cada uno de los puntos clave descritos contribuye a asegurar que los datos se encuentren en un formato limpio, estructurado y apto para análisis, dejando el terreno preparado para la aplicación de las técnicas de aprendizaje automático. Este proceso de preparación y análisis de datos no solo se documenta detalladamente en el texto, sino que también se ha implementado de forma práctica a través de varios archivos de Python. El código fuente de estas funciones se encuentra en los anexos de este trabajo y está disponible para su descarga y uso libre en GitHub EXFOR-ProtonReactions-Analysis Repository.

Además, el código en este repositorio, en conjunción con EXFORTABLES, permite reproducir todos los análisis realizados en este trabajo. Sin embargo, la utilidad del código trasciende el mero preprocesamiento de datos para el aprendizaje automático. Una vez que los experimentos están cargados, las funciones implementadas ofrecen una serie de herramientas que facilitan tareas como consultar información de cualquier experimento, graficar múltiples experimentos simultáneamente, filtrar experimentos en base a distintos atributos y clasificar los experimentos según las necesidades específicas. Esto convierte el preprocesamiento en una herramienta versátil, útil no solo para los objetivos de este trabajo sino también para cualquier otro análisis de datos de EXFOR que se desee realizar.

El siguiente paso crucial en nuestra investigación involucra la identificación y gestión de los valores atípicos o ‘outliers’. En este sentido, hemos implementado varias técnicas de detección de outliers, tanto estadísticas como basadas en aprendizaje automático, para asegurarnos de que nuestra análisis es tanto robusto como exhaustivo. Cada uno de estos métodos ha sido cuidadosamente implementado en archivos Jupyter Notebooks, que también están disponibles para consulta y uso en el repositorio de GitHub previamente mencionado.

Método IQR (Rango Intercuartílico): Este método estadístico es especialmente eficaz para identificar *outliers* en una distribución univariante. Utiliza el rango intercuartílico, que es la diferencia entre el tercer cuartil (Q3) y el primer cuartil (Q1) del conjunto de datos. Los datos que caen fuera de los límites establecidos por $Q1 - 1.5 \times IQR$ y $Q3 + 1.5 \times IQR$ se consideran valores atípicos. Este enfoque es robusto y eficiente para distribuciones que no necesariamente siguen una normalidad.

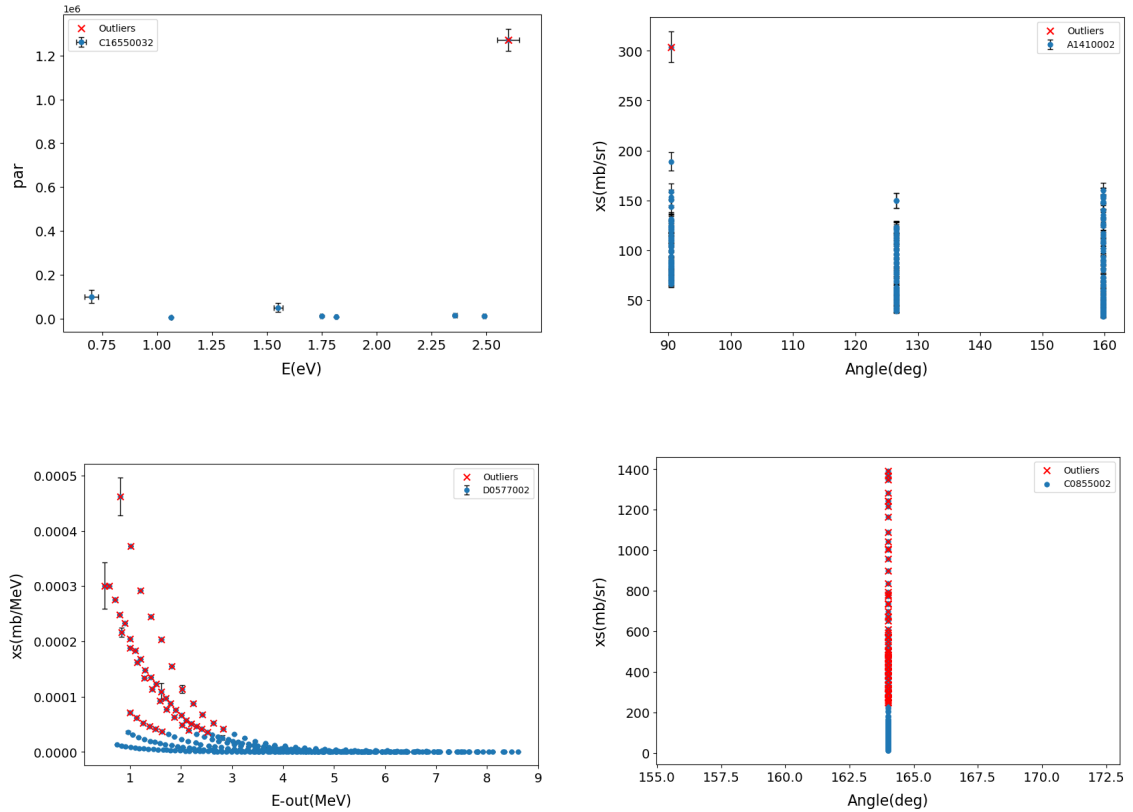


Figure 0.2: Detección de valores atípicos utilizando el método IQR.

Algunos de los resultados de aplicar este método se muestran en la Figura 0.2. En general el

método IQR es una técnica estadística que incorpora incertidumbres en sus cálculos, lo cual es beneficioso para minimizar falsos positivos. No obstante, nuestra investigación señala que el uso del método IQR se ve limitado para nuestros conjuntos de datos específicos. Este método muestra un rendimiento óptimo con datos de distribución uniforme, lo cual no corresponde al tipo de datos que estamos manejando.

En particular, cualquier pico en los datos, derivado de su distribución inherente, corre el riesgo de ser incorrectamente etiquetado como un valor atípico. Esto limita la utilidad del método IQR, ya que no puede distinguir de manera fiable entre valores atípicos auténticos y puntos de datos que forman parte natural de una distribución no uniforme.

Isolation Forest: Este método es una técnica de aprendizaje automático para la detección de anomalías que se basa en árboles de decisión. A diferencia de otros métodos, Isolation Forest no requiere un modelo previo de los datos “normales”. En lugar de eso, el algoritmo funciona aislando registros que son anómalos. Particularmente, el método selecciona al azar una característica y un valor para esa característica, y utiliza estas condiciones para dividir los datos. Este proceso se repite de manera recursiva hasta que los datos quedan completamente aislados. Cuanto más corto sea el camino para aislar un registro, más probable es que sea una anomalía. Isolation Forest es especialmente eficiente para conjuntos de datos de gran tamaño y se destaca por su capacidad para detectar outliers en un tiempo computacional relativamente bajo.

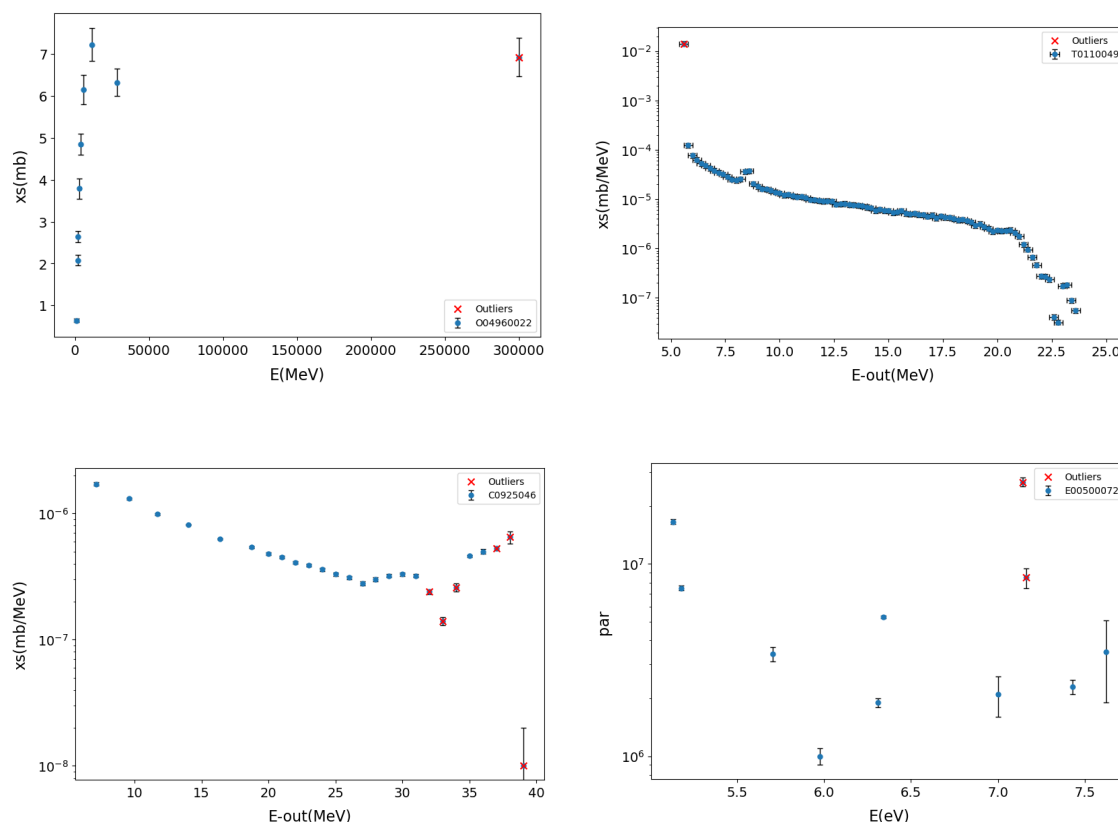


Figure 0.3: Detección de valores atípicos utilizando el método Isolation Forest.

Como se puede observar en la Figura 0.3, el método de Isolation Forest demuestra su eficacia en la detección de valores atípicos y, en general, supera a enfoques estadísticos tradicionales como el método IQR. No obstante, este algoritmo muestra una alta sensibilidad a la distribución específica del conjunto de datos evaluado. Una limitación notable, es su manejo de áreas de baja

densidad en el conjunto de datos. Puntos que están aislados debido a un muestreo insuficiente pueden ser erróneamente etiquetados como atípicos, lo que representa un desafío considerable en el contexto de nuestro estudio.

DBSCAN: Este método de aprendizaje automático es especialmente útil para identificar conglomerados en un conjunto de datos en función de la densidad de los puntos. DBSCAN agrupa puntos que están próximos entre sí según un umbral de distancia y requiere un número mínimo de puntos para formar un grupo. Los puntos que no pertenecen a ningún grupo se consideran outliers. A diferencia de otros métodos de clustering como K-means, DBSCAN no requiere especificar el número de grupos a priori y puede encontrar grupos de formas arbitrarias, lo que lo hace particularmente útil para la detección de outliers en conjuntos de datos complejos.

Según los resultados presentados en la Figura 0.4, el método DBSCAN se configura como otra herramienta útil de aprendizaje automático para la detección de valores atípicos. Este enfoque se basa en la densidad de puntos en el conjunto de datos, lo que lo convierte en una opción especialmente apta para manejar conjuntos de datos con densidades variables. Para el presente estudio, se han utilizado los hiperparámetros predeterminados de DBSCAN en una fase exploratoria inicial. Sin embargo, es importante subrayar, que el algoritmo es particularmente sensible a la selección de estos parámetros.

Es fundamental destacar que, DBSCAN puede ser computacionalmente costoso, más aún con conjuntos de datos de gran tamaño. En nuestro contexto específico, se hizo necesario subdividir alguno de los conjuntos de datos para facilitar un procesamiento computacional viable del algoritmo.

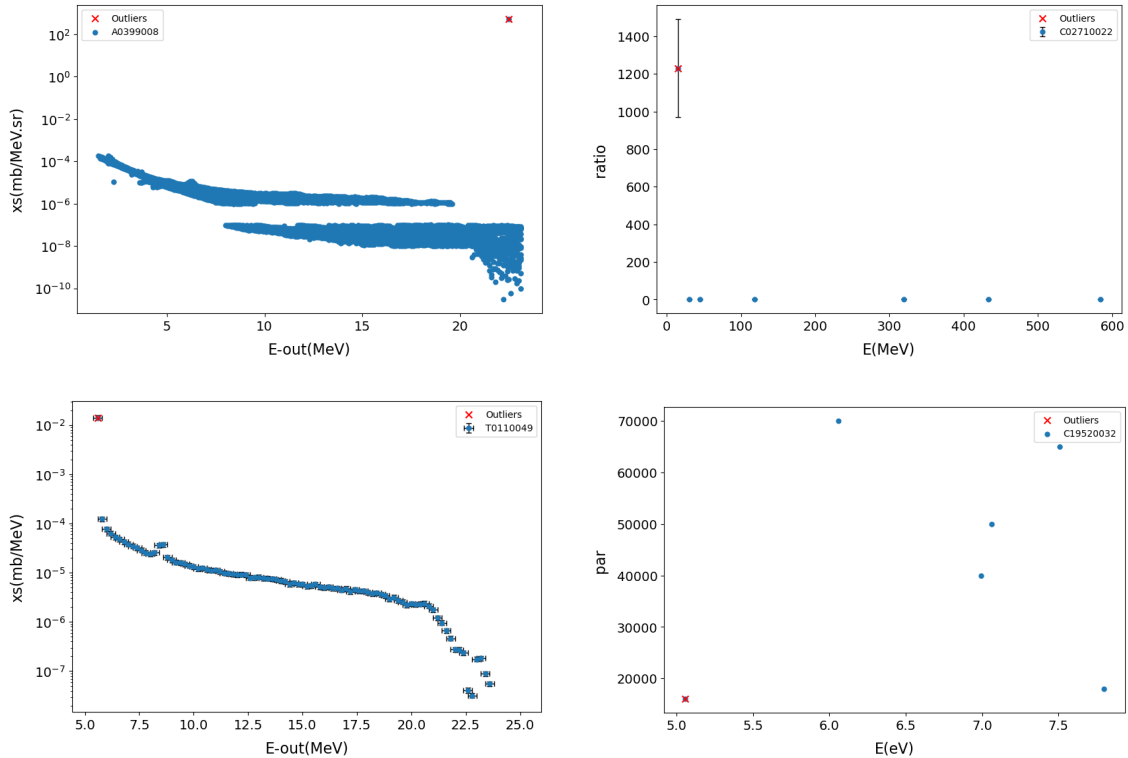


Figure 0.4: Detección de valores atípicos utilizando el método DBSCAN.

LOF (Local Outlier Factor): Este algoritmo se utiliza para identificar outliers en conjuntos de datos basándose en la densidad local de los puntos. En lugar de medir la distancia global desde un punto a todos los demás en el conjunto de datos, LOF compara la densidad de un

punto con la de sus vecinos más cercanos. Si un punto tiene una densidad significativamente menor que la de sus vecinos, se le asigna un factor de atipicidad alto, lo que lo identifica como un outlier. Este enfoque es especialmente eficaz en conjuntos de datos donde la densidad de los puntos no es uniforme, permitiendo una detección de outliers más precisa en tales escenarios.

Como se ilustra en la Figura 0.5, el método Local Outlier Factor (LOF) destaca por su versatilidad en adaptarse a diversas distribuciones de datos. A pesar de esta ventaja, el método no está exento de los retos habituales que enfrentan los algoritmos de detección de valores atípicos, tales como la sensibilidad a las particularidades del conjunto de datos. Este problema se intensifica cuando el conjunto de datos presenta huecos o falta de información en ciertos rangos, incrementando así la probabilidad de detectar falsos positivos.

La selección de los hiperparámetros es un elemento esencial para la eficacia del LOF. Este requerimiento puede convertir el proceso en una tarea pesada, ya que distintos conjuntos de datos pueden necesitar configuraciones variadas. Sin embargo, una de las fortalezas de LOF es su robustez para evitar falsos positivos en casos donde los datos presentan picos o agrupaciones pronunciadas. Este es un escenario en el que otros métodos de detección de valores atípicos suelen ser menos efectivos.

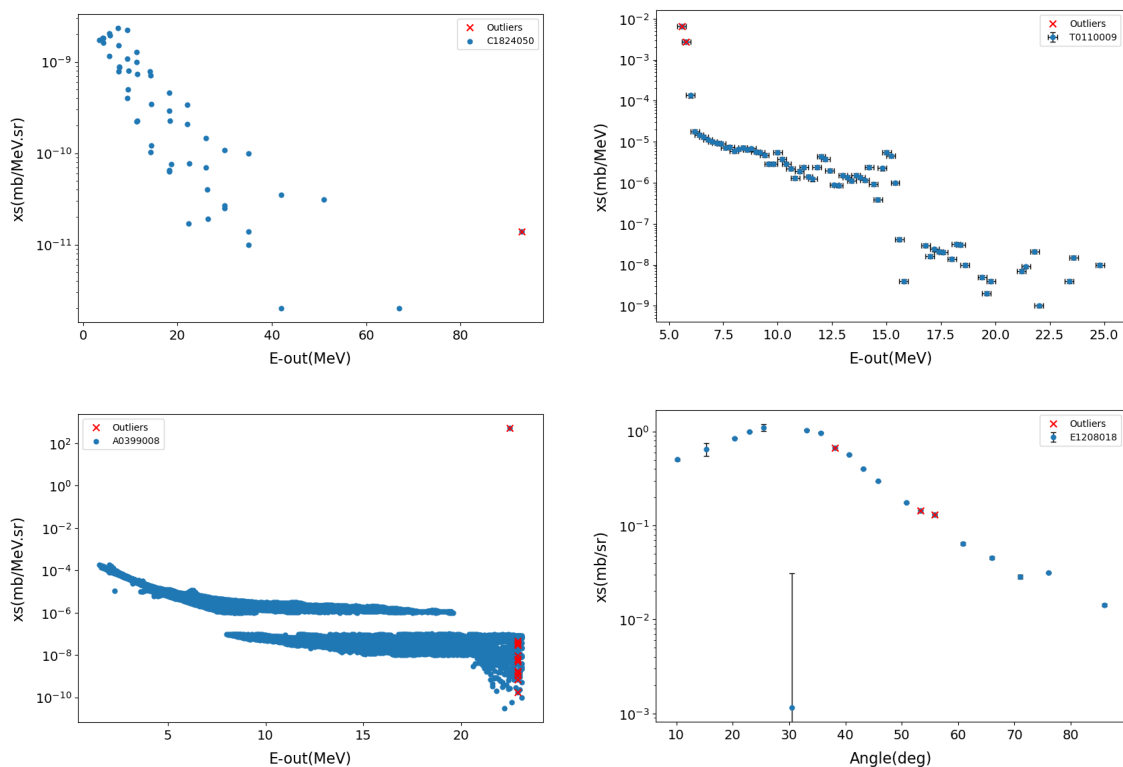


Figure 0.5: Detección de valores atípicos utilizando el método LOF.

SVM (Support Vector Machine): En el contexto de detección de outliers, se suele utilizar una variante de SVM conocida como One-class SVM. Este algoritmo intenta encontrar el hiperplano en un espacio de características de alta dimensión que mejor separa los puntos de datos "normales" de cualquier potencial outlier. La idea es maximizar el margen entre los puntos de datos y el origen en este espacio transformado. Los puntos que queden del lado incorrecto del hiperplano se consideran outliers. Este método es particularmente útil para conjuntos de datos de alta dimensión y es eficaz incluso cuando los datos son no lineales.

Conforme a los datos presentados en la Figura 0.6, el algoritmo SVM resultó ser efectivo en

la detección de valores atípicos en una variedad de conjuntos de datos. Al usar parámetros predeterminados como punto de partida, el método demostró su robustez en diferentes escenarios. Aunque estos parámetros son un buen inicio, podrían optimizarse más mediante técnicas específicas o conocimiento del dominio.

Similar a lo que ocurrió con el DBSCAN, el algoritmo SVM demanda una alta capacidad computacional, especialmente cuando se trata de grandes conjuntos de datos. Este reto fue palpable en nuestro estudio, donde fue necesario segmentar el conjunto inicial de datos en subsets más manejables para cumplir con las demandas computacionales.

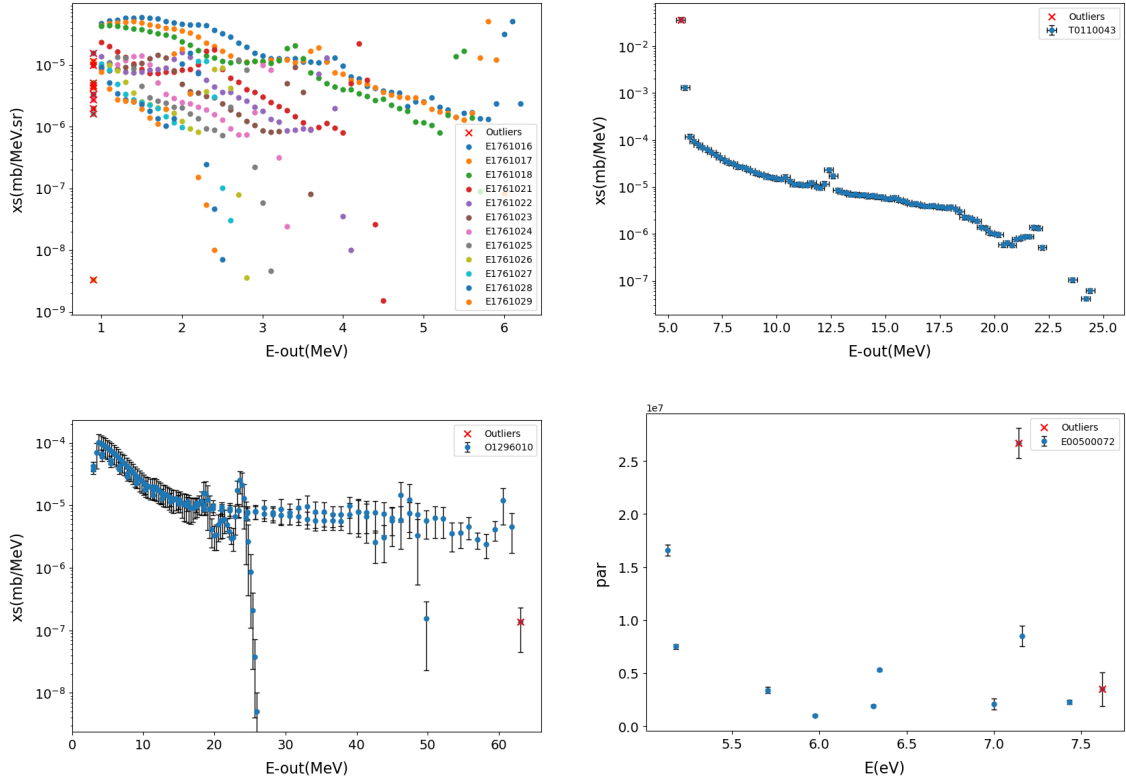


Figure 0.6: Detección de valores atípicos utilizando el método SVM.

Autoencoder: Los Autoencoders son una clase de redes neuronales utilizadas generalmente para la reducción de dimensionalidad o para el aprendizaje de representaciones. La red está compuesta típicamente por una capa de entrada, una o más capas ocultas, y una capa de salida. La red se entrena para reconstruir su entrada, forzando así a la red a aprender representaciones internas eficientes de los datos de entrada.

Arquitectura y Creación del Modelo: La red consta de dos partes principales: el codificador y el decodificador. El codificador comprime la entrada y la representa en un espacio latente de menor dimensionalidad. El decodificador toma esta representación y la reconstruye para producir una salida lo más cercana posible a la entrada original. Se entrena el modelo minimizando la diferencia entre la entrada y la salida, comúnmente usando una función de pérdida como el error cuadrático medio.

Aplicación para la Detección de Outliers: Una vez entrenado el autoencoder en datos "normales", se puede utilizar para detectar outliers al medir cuán bien el modelo puede reconstruir nuevas instancias de datos. Si un punto de datos es una observación atípica, la red tendrá dificultades para reconstruirlo correctamente, resultando en un alto valor de error de reconstrucción. Este

valor se puede utilizar como un indicador para identificar outliers en el conjunto de datos.

Según lo mostrado en la Figura 0.7, los autoencoders, que son algoritmos basados en redes neuronales, se diferencian significativamente de los métodos de aprendizaje automático tradicionales tratados en esta tesis. A pesar de requerir un nivel más profundo de entendimiento para su implementación efectiva, los beneficios potenciales de usar estas redes son considerables. Futuros trabajos podrían focalizarse en ajustar diversos parámetros, tales como el número de capas, funciones de activación y tasas de aprendizaje.

En términos de eficacia para la detección de valores atípicos, el autoencoder muestra resultados prometedores. Logró identificar diversas instancias de posibles valores atípicos con menos falsos positivos que otros métodos, aunque su rendimiento es altamente dependiente de la configuración de parámetros.

Un factor crítico para el desempeño del modelo es la cantidad de datos disponibles para entrenamiento. En ciertos casos, el tamaño limitado de los conjuntos de datos afectó negativamente el rendimiento del modelo. Una táctica alternativa podría ser el uso de un conjunto de datos más amplio, aunque esta estrategia no garantiza necesariamente una mejora en el rendimiento.

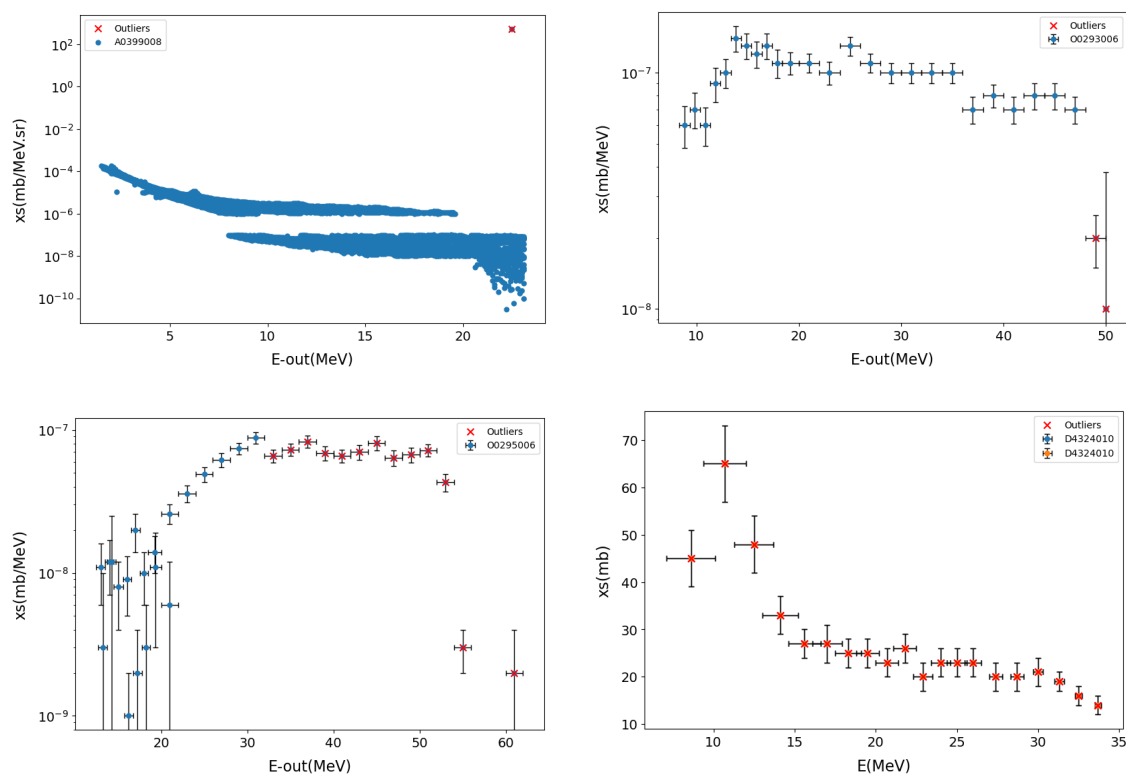


Figure 0.7: Detección de valores atípicos utilizando el método Autoencoder.

En conclusión, cada método para la detección de valores atípicos examinado en este estudio tiene sus propias ventajas y limitaciones. La efectividad de estos métodos varía según las características específicas de los conjuntos de datos analizados; algunos sobresalen en áreas donde otros pueden tener deficiencias. En el ámbito de las técnicas de aprendizaje automático, es difícil determinar un método definitivamente superior. Aunque DBSCAN mostró resultados ligeramente mejores, su eficacia depende en gran medida de la naturaleza de los datos y de los ajustes específicos elegidos para cada análisis. Los métodos basados en redes neuronales, como los autoencoders, identificaron menos valores atípicos, lo que podría sugerir una mayor

precisión. No obstante, los mayores requisitos computacionales y la complejidad en la configuración de parámetros de las redes neuronales plantean interrogantes sobre su idoneidad práctica para la tarea. Por lo tanto, cada técnica aporta su propia perspectiva para la identificación de valores atípicos, y su éxito relativo está vinculado a las complejidades de los conjuntos de datos a los que se aplican y a los objetivos específicos que buscan cumplir.

Es fundamental entender que los algoritmos de aprendizaje automático están diseñados para identificar un porcentaje fijo de valores atípicos, lo que puede llevar a falsos positivos. Por ello, la intervención de un experto en el dominio es crucial para interpretar adecuadamente los resultados. Dada la gran cantidad de datos en la biblioteca EXFOR, un análisis manual no sería práctico, resaltando la importancia de las técnicas de aprendizaje automático en este estudio. Este trabajo aporta una herramienta gratuita y fácil de usar para revisar datos, accesible para cualquiera con un ordenador con Python.

En definitiva, esta tesis de máster se ha centrado en la incorporación de técnicas de aprendizaje automático para la detección de valores atípicos en la biblioteca de datos nucleares EXFOR para reacciones inducidas por protones. El estudio ha desarrollado una metodología robusta específicamente adaptada para la complejidad de los datos de EXFOR. No solo se crearon estrategias eficientes para la extracción y preparación de datos, sino que también se estableció un marco de trabajo que sirve como modelo para futuras investigaciones en este campo.

Se exploraron diversas técnicas para la detección de valores atípicos, desde métodos estadísticos tradicionales como el rango intercuartílico (IQR) hasta algoritmos de aprendizaje automático más sofisticados como Isolation Forest, DBSCAN y el Local Outlier Factor (LOF). Además, se incorporó un modelo de red neuronal, específicamente un Autoencoder, en nuestro conjunto de algoritmos.

El estudio ha cumplido con todos los objetivos iniciales y ha establecido una metodología sólida para aplicar el aprendizaje automático en la detección de valores atípicos en la biblioteca EXFOR. La metodología desarrollada es lo suficientemente versátil como para extenderse fácilmente a otras secciones de la biblioteca EXFOR, lo que abre la puerta a mejoras integrales. Para fomentar el desarrollo y la participación comunitaria, todos los scripts y códigos generados están disponibles de forma pública en GitHub, lo que no solo cumple con el imperativo académico de compartir conocimiento, sino que también prepara el terreno para esfuerzos colaborativos futuros para refinar y ampliar el trabajo presentado.

Códigos UNESCO

- Campo:
 - 33 - Ciencias de la Tecnología
 - 12 - Ciencias de la Computación e Información
 - 22 - Física
- Disciplina:
 - 3304 - Tecnología
 - 3307 - Tecnología de los reactores nucleares
 - 3320 - Ingeniería y Tecnología Nuclear
 - 1203 - Ciencias de la Computación
 - 2202 - Física Nuclear
- Subdisciplina:
 - 330708 - Ingeniería de reactores
 - 332099 - Otras
 - 120326 - Simulación
 - 120317 - Redes Neuronales

Palabras Clave

- | | |
|------------------------------|---------------------------------|
| • EXFOR Nuclear Data Library | • Local Outlier Factor (LOF) |
| • Machine Learning | • Support Vector Machines (SVM) |
| • Proton-Induced Reactions | • Interquartile Range (IQR) |
| • Outlier Detection | • Neural Network Autoencoder |
| • Isolation Forest | • Data Preprocessing |
| • DBSCAN | • Data Mining |

ABSTRACT

The EXFOR nuclear data library is an indispensable resource for nuclear physicists and engineers. However, the database is susceptible to outliers and inconsistencies that can significantly impair its utility. This thesis focuses on the application of machine learning techniques for identifying potential outliers in proton-induced reactions within the EXFOR library. A comprehensive literature review on machine learning techniques was conducted to identify appropriate methods for outlier detection. Several algorithms, including Isolation Forest, DBSCAN, Local Outlier Factor (LOF), Support Vector Machines (SVM), along with a statistical method based on the Interquartile Range (IQR) and a neural network-based autoencoder, were rigorously tested on their capability to flag outliers effectively.

Extensive preprocessing work was carried out to prepare the data for machine learning applications. This involved creating a user-friendly tool that facilitates easy data mining of the EXFOR library, thereby serving as an invaluable resource for other researchers. Our results demonstrate that machine learning techniques can significantly improve the reliability and accuracy of the EXFOR database, with implications for a wide range of applications in nuclear science and engineering.

TABLE OF CONTENTS

ACKNOWLEDGMENTS	iv
RESUMEN EJECUTIVO	vi
ABSTRACT	xviii
1 INTRODUCTION	1
1.1 Background and Motivation	1
1.2 Objectives	1
1.3 Thesis Structure	2
2 STATE OF THE ART	3
2.1 EXFOR	4
2.2 EXFORTABLES	7
2.3 Evaluated Nuclear Data Libraries	8
2.4 Machine Learning	9
2.5 Outliers	12
2.5.1 Statistical Techniques	12
2.5.2 Machine Learning Techniques	12
3 METHODOLOGY	14
3.1 Tools and Libraries Used	14
3.2 Preprocessing and Data Mining	15
3.3 Interquartile Range (IQR)	19
3.4 Isolation Forest	20
3.5 DBSCAN	22
3.6 Local Outlier Factor (LOF)	23
3.7 Support Vector Machines (SVM)	24
3.8 Neural Network Models: Autoencoders	25

3.9	Evaluation Metrics	27
4	RESULTS	28
4.1	Data Preprocessing	28
4.2	Interquartile Range (IQR)	35
4.2.1	Implementation of IQR Method	35
4.2.2	Results and Graphical Analysis	36
4.2.3	Conclusion on IQR Method	38
4.3	Isolation Forest	39
4.3.1	Implementation of Isolation Forest Method	39
4.3.2	Results and Graphical Analysis	40
4.3.3	Conclusion on Isolation Forest Method	42
4.4	DBSCAN	42
4.4.1	Implementation of DBSCAN Method	42
4.4.2	Results and Graphical Analysis	43
4.4.3	Conclusion on DBSCAN Method	46
4.5	Local Outlier Factor (LOF)	47
4.5.1	Implementation of Local Outlier Factor (LOF) Method	47
4.5.2	Results and Graphical Analysis	48
4.5.3	Conclusion on LOF Method	51
4.6	Support Vector Machines (SVM)	51
4.6.1	Implementation of SVM Method	51
4.6.2	Results and Graphical Analysis	52
4.6.3	Conclusion on SVM Method	56
4.7	Neural Network Models: Autoencoders	56
4.7.1	Implementation of the Autoencoder	56
4.7.2	Model Evaluation	57
4.7.3	Outlier Detection Using Autoencoder	59
4.7.4	Results and Graphical Analysis	60

4.7.5	Conclusions on Autoencoder Method	62
5	CONCLUSIONS	63
5.1	Summary of Findings	63
5.2	Future Work	65
5.3	Final Conclusions	65
6	SOCIAL & PROFESSIONAL RESPONSABILITY	67
6.1	Professional Responsibility	67
6.2	Social Responsibility and Sustainable Development Goals (SDGs)	67
7	PROJECT MANAGEMENT	69
7.1	Planning	69
7.2	Budget	71
8	BIBLIOGRAPHY	73
	LIST OF TABLES	77
	LIST OF FIGURES	80
	Annexes	83
A	‘Experiment’ Class Source Code and Documentation	83
B	Utility Functions for Proton Experiment Data Analysis	89
B.1	read_experiment()	89
B.2	write_experiments_to_binary()	94
B.3	read_experiments_from_binary()	95
B.4	write_experiments_to_txt()	96
B.5	read_experiments_from_txt()	98
B.6	read_proton_experiments_from_exfortables()	101
B.7	filter_experiments()	103
B.8	get_unique_values()	105
B.9	clean_dataframe()	107

B.10	<code>classify_experiments()</code>	108
B.11	<code>classify_experiments_by_data()</code>	110
B.12	<code>plot_experiments()</code>	112
B.13	<code>plot_outliers()</code>	114
C	Interquartile Range Functions	117
C.1	<code>detect_outliers_IQR()</code>	117
C.2	<code>IQR_method()</code>	119
D	Classification Data Groups	121

1 INTRODUCTION

1.1 Background and Motivation

Outliers are observations that deviate markedly from other observations in a sample. They may indicate bad data, incorrect coding, or experimental errors. Outliers can also affect the accuracy and reliability of statistical analyses, especially when dealing with small sample sizes or non-normal distributions. Therefore, it is important to detect and treat outliers appropriately in experimental physics [2].

One area of experimental physics that requires high-quality data is nuclear reaction data, which is used for various applications such as nuclear energy, astrophysics, medical physics, and radiation protection. Nuclear reaction data are compiled in the EXFOR library [3], which contains more than 24,000 experiments involving neutrons, charged particles, and photons.

A specific type of nuclear reaction data that has gained increasing interest and necessity is incident proton data. One of the main reasons for this is that in recent years, the development of proton therapy medical treatments has gained significant interest and importance in the field of cancer treatment. This method uses proton beams to deliver radiation to cancer cells with precision, causing less damage to surrounding healthy tissues, fewer side effects, and improved outcomes than traditional radiotherapy techniques. Due to its high accuracy and effectiveness, proton therapy has become an increasingly common treatment option for certain types of cancer [4][5][6].

However, incident proton data are also subject to outliers due to various sources of uncertainty such as measurement errors, calibration errors, model uncertainties, or physical phenomena. Outliers can affect the accuracy and precision of incident proton data and consequently affect the quality and safety of proton therapy treatments. Precise knowledge of the physical properties of the proton beam and its interaction with biological tissues is necessary [7]. This knowledge is often obtained through the analysis of experimental data from proton incident particle experiments, which provide information on the behavior of the proton beam under different conditions. Therefore, the detection and removal of outliers from experimental data is an important step in obtaining reliable and accurate results.

1.2 Objectives

The primary aim of this master's thesis is to assess the implementation of various machine learning techniques for outlier detection in the EXFOR nuclear data library. In pursuit of this aim, a key objective is to develop a robust methodology that can effectively tackle this complex task. The overarching goal extends beyond mere implementation; it also includes evaluating the feasibility of applying machine learning to this specialized data set. The intent is not just to apply machine learning algorithms, but to investigate how well these algorithms can be tailored to suit the data and task at hand, and to identify which techniques emerge as most effective.

A focal point of the investigation lies in the initial stages of data mining and preprocessing, which are critical for the successful application of any machine learning algorithm. Given the complexity of EXFOR's data structures, significant efforts will be invested in researching various methods for efficiently extracting, transforming, and preparing the data for subsequent machine learning techniques.

Building upon this foundational work, the study will implement a range of outlier detection algorithms. These include traditional statistical methods like the Interquartile Range (IQR), tree-based methods such as Isolation Forest, density-based techniques like DBSCAN, and distance-based metrics like the Local Outlier Factor (LOF). Furthermore, to explore the frontiers of what machine learning can offer, a neural network model, specifically an Autoencoder, will also be deployed. This involves intricate steps from model building and parameter optimization to performance evaluation in the context of outlier detection.

The various methods and their results will then be critically compared through visual and quantitative analyses. Through this comparative lens, the study aims to gauge the reliability, advantages, and limitations inherent in each approach. Finally, the findings will be contextualized within the broader scientific literature, providing insights into the future potential and challenges of leveraging machine learning for outlier detection in the EXFOR dataset.

1.3 Thesis Structure

The structure of this thesis is organized into distinct chapters, each serving a specific role and contributing to the overall research objectives:

- **State of Art:** This opening chapter lays the groundwork by reviewing the existing literature and key concepts pivotal to this thesis. Topics like EXFOR, outlier detection, and machine learning are discussed to provide a well-rounded understanding of the research context.
- **Methodology:** This chapter is dedicated to outlining the tools, libraries, and theoretical foundations used in the research. It explains the preprocessing and data mining techniques employed, as well as the machine learning methods that are used for outlier detection.
- **Results:** The findings from the application of various methodologies for outlier detection are presented in this chapter. The results are analyzed to assess the reliability, strengths, and weaknesses of each method applied.
- **Conclusions:** This chapter serves as a repository for the key findings, implications, and future research avenues. It draws together the entire scope of the thesis to offer a consolidated view of the research carried out.
- **Social and Professional Responsibility:** This section delves into the ethical, social, and professional implications tied to the study. It highlights how these aspects have been conscientiously considered and addressed in the course of the research.
- **Project Management:** In this chapter, the managerial aspects of the research project are discussed. It covers elements such as planning, resource allocation, time management, and risk assessment.

The thesis is supplemented by two final sections: 'Glossary' and 'Annexes'. The Glossary offers explanations for key terms and concepts, while the Annexes provide additional materials and data sets that support the core research.

2 STATE OF THE ART

In the world of nuclear research, the coming together of real-world data, machine learning, and ways to spot unusual data points is a new and exciting area. At the center of this is the Exchange Format (EXFOR), a detailed collection created to help with worldwide nuclear studies and uses. But there is more to the story, with many methods, partnerships, and new technologies being part of it. In today's data-driven world, machine learning plays a big role. It helps us look at, understand, and use huge amounts of nuclear data in new ways, making research better and faster. This is not just a standalone field but something that links different areas together, allowing for new ideas and insights.

However, handling the data properly is also important. Finding and dealing with unusual or 'outlier' data points makes sure that the information is trustworthy, which means the results we get from it are more reliable. This chapter on the state of the art opens up a broad view of what's happening now, looking at:

- **EXFOR:** What it is, why it is important, and how it is part of global efforts around nuclear data.
- **Machine Learning:** A look at the latest ways it is being used, what it can do, and the challenges it faces in the nuclear field.
- **Outlier Detection:** How we find and deal with unusual data points, the tools we use, and best practices.
- **Additional Sections:** Other connected topics that add to the understanding of nuclear data, machine learning, and keeping the data of high quality.

This chapter serves as a foundational introduction to the vital elements at the intersection of nuclear research, machine learning, and data quality. It offers a starting point for those new to these fields, providing clear explanations of key concepts and terms. It also creates a connection by linking different areas together, showing how they relate to each other and why they matter. Moreover, the chapter provides insight into how these parts fit into the bigger picture of nuclear research, setting the stage for deeper exploration in subsequent chapters. Together, these elements form a comprehensive yet accessible overview, aiming to engage both experts and newcomers. By laying the groundwork, this chapter helps readers understand the core ideas, opening the door to further study and appreciation of a multifaceted and rapidly evolving discipline.

2.1 EXFOR

The EXchange FORmat, or EXFOR, is a comprehensive library of experimental nuclear reaction data, meticulously compiled to facilitate nuclear research and applications. Constructed with a clear goal, it provides a reliable and interpretive data source built around the context of each experiment [8]. The format, designed for flexibility, serves to ease the exchange of data among centers. This platform was developed to facilitate the transmission of nuclear reaction data among the members of the International Network of Nuclear Reaction Data Centres (NRDC), an international collaboration aimed at coordinating the collection, compilation, and global dissemination of nuclear data. The NRDC, overseen by the IAEA Nuclear Data Section, is responsible for maintaining and managing the EXFOR nuclear data library [9].

EXFOR, as a product of worldwide cooperation, offers a system designed for the collection, storage, exchange, and retrieval of experimental nuclear reaction data. The database, as of October 2022, holds 24,382 experiments [10]. It aims to compile a comprehensive collection of low-energy experimental neutron-induced reaction data, a less complete compilation of charged-particle-induced reaction data, and selectively compiles photon-induced, heavy-ion-induced, and high-energy neutron-induced reaction data [8].

The target audience for EXFOR encompasses researchers and professionals in nuclear science and engineering. This includes individuals in academia, national laboratories, and the industry, who are engaged in the design and analysis of nuclear reactors, nuclear fuel cycles, and other nuclear applications.

The organization of data within EXFOR follows a clear, structured format (see Figure 2.1). The basic unit, known as an “ENTRY,” corresponds to a single experiment, while divisions within an ENTRY, termed “SUBENTRIES”, contain the results of the experiment. The initial SUBENTRY typically lacks numerical data but includes information relevant to the entire ENTRY, such as bibliographic references and details about the experimental setup. The SUBENTRY is further divided into sections [10]:

- **BIB:** This section contains bibliographic, descriptive, and bookkeeping information.
- **COMMON:** This section holds data relevant to all data throughout the SUBENTRY.
- **DATA:** This section contains the experimental data tables.

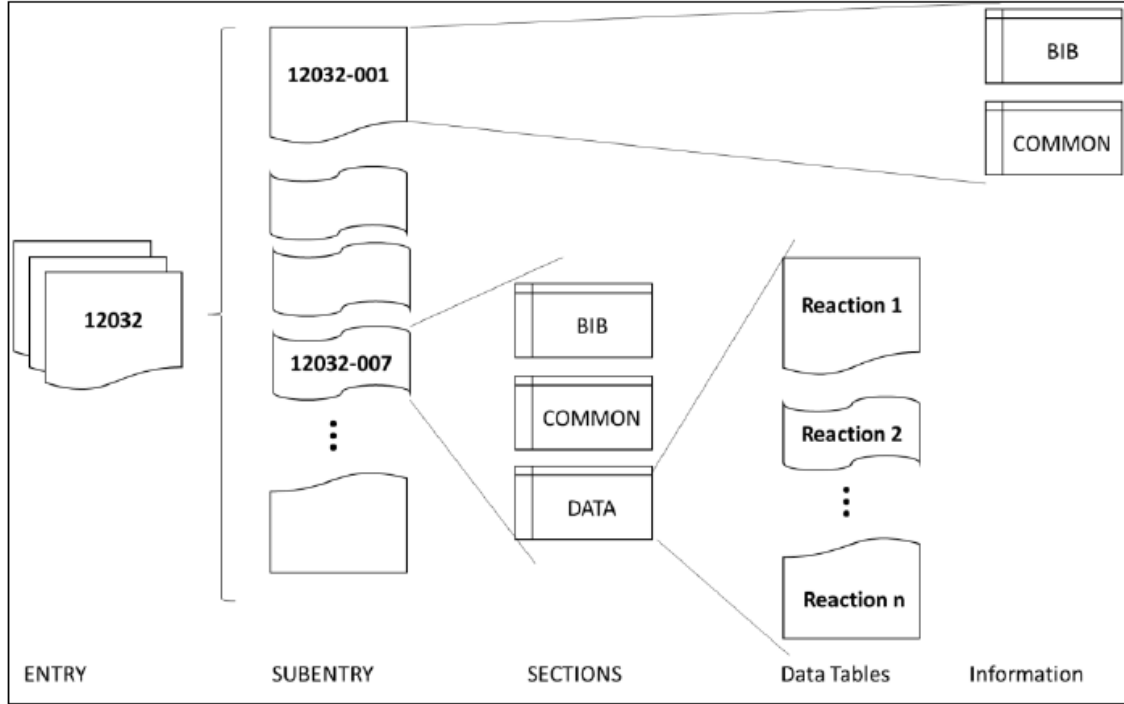


Figure 2.1: Structure of information in EXFOR [10].

Additionally, the EXFOR library incorporates “TRANS” files. These files, which are exchanged among data centers, encompass numerous ENTRIES (work). Each file is formatted as an 80 character ASCII record, wherein columns 1-66 are designated for EXFOR data. For unique identification of a record within the file, columns 67-79 are utilized [10].

EXFOR operates as a supplementary publication medium to conventional, formal publications, ensuring that any revisions made by authors post-publication are reflected in the data files. The numerical data housed within EXFOR are supplemented by explanatory text that delivers essential information about the data’s meaning and quality. An EXFOR “entry” may not necessarily correspond to the information found in a single publication. Often, a “work” is reported across multiple formal or informal publications, such as progress reports, conference papers with preliminary results, lab reports, and final articles in refereed international journals. It can also incorporate preliminary data or pre-publication data with the author’s consent. Once the final data becomes available, preliminary data are routinely replaced and any new bibliographic references describing the results are added. Importantly, EXFOR is a compilation of the author’s original published experimental data and not a collection of recommended values for each reaction. Evaluators systematically review the experimental works and recommend best evaluated cross-section data based on the experimental data from EXFOR. Their results are compiled in various evaluated data files such as ENDF/B, JEFF, JENDL, and more, most of which are also accessible from the IAEA Nuclear Data Section [9].

In the worldwide push to make the EXFOR Nuclear Data Library easier to use, a group known as the Working Party on International Nuclear Data Evaluation Co-operation Subgroup 50 (WPEC SG50), part of the Organisation for Economic Co-operation and Development’s (OECD) Nuclear Energy Agency (NEA), has been doing some interesting work. This group is building an advanced database for experimental nuclear reaction data called the Machine-readable Experimental Data User Application & Library (MEDUSAL), as explained in reference [11].

MEDUSAL is an ambitious project that aims to create a large, automatically readable database filled with experimental nuclear reaction data. It has been designed to support a wide array of users, including nuclear data evaluators, experimentalists, model builders, and those engaged in machine learning and artificial intelligence. What sets MEDUSAL apart is its ability to store not only the core data but also additional information contributed by users in areas such as model development, evaluation, and validation.

The architecture of MEDUSAL is thoughtfully planned and segmented into four levels or 'layers.' Layer 0 forms the base where data from EXFOR is translated into MEDUSAL's format, standardizing the units. Layer 1, the main access point for users, consists of information from Layer 0, augmented with machine-readable metadata and supplementary details from EXFOR and related literature. Layer 2 introduces changes to the original authors' data, such as updates to new standards and flagging inconsistent measurements. Layer 3, the most flexible layer, can hold multiple versions of the same dataset, enabling users to apply their corrections and judgments. Importantly, this layer documents precisely how a dataset was employed in an evaluation, essential for ensuring reproducibility [11].

Creating MEDUSAL has not been without its challenges, particularly in assuring the accuracy and currency of the data. However, the Lewis-led team is optimistic about MEDUSAL's future applications in nuclear data evaluations, model development, and validation. They envision expanding MEDUSAL to encompass more data types, such as covariance data, and are developing an Application Programming Interface (API) for easier database access. The potential use of MEDUSAL in machine learning and artificial intelligence amplifies the need for precise and reliable data [11].

In conjunction with the MEDUSAL project, there are ongoing developments in Python codes to parse the EXFOR nuclear data library. These codes, publicly available on the IAEA-NDS GitHub Repository [12], include tools to read the entire EXFOR data into nested dictionary structures for Python handling or .json format conversion. Since these projects are still under development and many features are pending, users should approach them with caution. A different approach, known as "EXFORTABLES", developed by Arjan Koning, will be discussed in more detail in the following subsection.

2.2 EXFORTABLES

EXFORTABLES is an innovative software tool developed by Arjan Koning and designed to translate the extensive and intricate EXFOR experimental nuclear reaction database into a more accessible and user-friendly format. This transformation enhances the ease of accessing, plotting, and utilizing nuclear data for evaluation purposes. Essentially, EXFORTABLES serves as a bridge for researchers to interact with the wealth of data within EXFOR more efficiently.

The capabilities of EXFORTABLES extend beyond mere data translation. The software’s directory structure enables users to quickly locate the required data, saving valuable time in research. Moreover, it automatically compares the data against global nuclear model calculations, identifying potential errors in the database. This dual functionality simultaneously tests nuclear models, allowing researchers to assess their accuracy and deepen their comprehension of nuclear reactions.

The open-source nature of EXFORTABLES presents an added benefit. Researchers can modify and adapt the software according to their specific needs, promoting a tailored research experience. It’s noteworthy that EXFORTABLES is assumed to be installed on a Unix/Linux operating system, and the installation process requires approximately 3 GB of free disk space.

Upon installation, the EXFORTABLES package will be stored in the `exfortables/` directory. It consists of various directories and files, each serving a distinct purpose. These include:

- `source/`: containing the Fortran code responsible for database creation and data checking.
- `files/`: housing the entire EXFOR database in different formats, alongside other related information.
- Specific directories like `'n/'`, `'p/'`, `'d/'`, etc., representing different reaction types, and containing the entire structured database.
- `doc/`: with documentation.
- `input/`: for the input file used to create the database.
- `quality`: for quality scores.
- `stat/`: with statistical information.
- `special/`: for specific data files.
- `README`: outlining the package content.

The hierarchical directory structure of EXFORTABLES is categorized by projectile type (such as neutron, proton, deuteron, etc.), element, mass, and reaction. Most users can directly utilize the database without additional installation or configuration, offering a seamless research experience.

While EXFORTABLES offers comparison with evaluated nuclear data libraries using statistical methods, it is worth acknowledging this aspect even though it falls outside the scope of this thesis, focused more on Machine Learning techniques. The underlying methodology in EXFORTABLES could also be insightful for the detection or comparison of outlier results, thus contributing to the broader field of nuclear data analysis.

2.3 Evaluated Nuclear Data Libraries

In the domain of nuclear science and technology, Evaluated Nuclear Data Libraries (ENDL) play a critical role. These libraries consist of systematic collections of evaluated nuclear reaction data, where ‘evaluated’ refers to data that has been critically assessed and validated through rigorous analysis of experimental data, theoretical models, and phenomenological trends. ENDLs are essentially a bridge between raw experimental data and practical application in nuclear technology [13].

The primary purpose of ENDLs is to serve as reliable sources of nuclear data for various applications in nuclear science and technology, which include but are not limited to nuclear energy production, medical applications, nuclear safeguards, and nuclear astrophysics. The quality and reliability of these libraries are paramount, as they directly impact the safety, efficiency, and effectiveness of these applications. In the process of obtaining the evaluated data, experimental databases such as EXFOR play an integral part. EXFOR, with its extensive collection of nuclear reaction data from a wide range of experiments conducted worldwide, serves as a primary input in the evaluation process. The raw data from EXFOR is analyzed and validated through a meticulous process of evaluation, involving critical review, data adjustment, model calculations, and uncertainty quantification. The resulting evaluated data is then compiled into the ENDLs.

To illustrate the global efforts in maintaining and developing evaluated nuclear data libraries, a range of key libraries have been listed below in Table 2.1. Each library is managed by a specific organisation or institution, demonstrating the cooperative nature of this significant scientific endeavor [13]. The libraries span across multiple continents, from the United States to Europe and Asia, each contributing a unique set of expertise and focus.

Table 2.1: Major Evaluated Nuclear Data Libraries and their respective Organisations [13]

Evaluated Nuclear Data Library	Organisation
ENDF/B-VIII (Evaluated Nuclear Data File/B Version VIII)	The National Nuclear Data Center (NNDC), Brookhaven National Laboratory, United States
JEFF (Joint Evaluated Fission and Fusion)	The Organisation for Economic Co-operation and Development (OECD) Nuclear Energy Agency (NEA) Data Bank, France
JENDL (Japanese Evaluated Nuclear Data Library)	Japan Atomic Energy Agency (JAEA)
CENDL (China Evaluated Nuclear Data Library)	China Institute of Atomic Energy (CIAE)
BROND (Russian Evaluated Nuclear Data Library)	National Nuclear Data Center of Russia, Institute for Physics and Power Engineering (IPPE), Obninsk, Russia
TENDL (TALYS Evaluated Nuclear Data Library)	Nuclear Research and Consultancy Group (NRG), Netherlands

2.4 Machine Learning

Machine Learning (ML) is a subfield of artificial intelligence (AI) that focuses on the design of systems, enabling them to learn and make decisions or predictions based on data. Rather than manually programming specific tasks, these systems improve and adapt their performance over time with exposure to new data. The roots of machine learning date back to the mid-20th century, with the earliest concepts revolving around the idea of computers or machines that can simulate aspects of human intelligence. However, the term “machine learning” was officially coined by Arthur Samuel in 1959 [14].

Machine learning is used extensively in various fields including healthcare, finance, transportation, and energy. In healthcare, machine learning algorithms can analyze large datasets to predict disease outcomes, thereby enhancing disease diagnosis and treatment [15]. In finance, machine learning aids in detecting fraudulent transactions [16]. In transportation, machine learning algorithms underpin the technologies in autonomous vehicles [17]. Similarly, in the energy sector, machine learning models can forecast energy demand and optimize energy distribution [18].

Machine learning algorithms can handle vast, complex datasets, often identifying hidden patterns that may not be immediately noticeable to human analysts. This attribute has been central to its broad application in today’s increasingly data-intensive world. However, the methodologies and approaches to machine learning can vary, with three principal types often recognized: supervised learning, unsupervised learning, and reinforcement learning [19].

- **Supervised Learning:** is a type of machine learning where the model is trained on labeled data. It learns a mapping from features to labels and applies this to accurately predict the labels for new, unseen instances. A commonly cited example of supervised learning is predicting house prices based on features such as location, size, and condition of the property. The model would be trained on historical data where the prices of sold houses (the labels) are known, and then it could predict the price of a new house on the market (an unseen instance) based on its features.
- **Unsupervised Learning:** involves training the model using unlabeled data, and the model identifies patterns and structures within this data on its own. An example of this can be seen in customer segmentation based on shopping habits, where the company does not initially know what these segments should be, so unsupervised learning is used to identify them.
- **Reinforcement Learning:** is a type of machine learning where an agent learns to make decisions by performing actions and observing the results, essentially learning through trial and error. Google’s AlphaGo, which learned to play the board game Go by playing millions of games against itself, is a notable example of reinforcement learning.

While the three types of machine learning may seem distinct, they are united by a common goal: to learn from data and make intelligent decisions or predictions. The principal difference lies in the type of feedback they are provided and the learning context, which can range from having a clear mapping of inputs to outputs (supervised learning), no specific output targets but an aim to understand data structure (unsupervised learning), or learning from a series of actions and rewards in a dynamic context (reinforcement learning) [1]. Figure 2.2 shows a diagram depicting the taxonomy of Machine Learning inside AI.

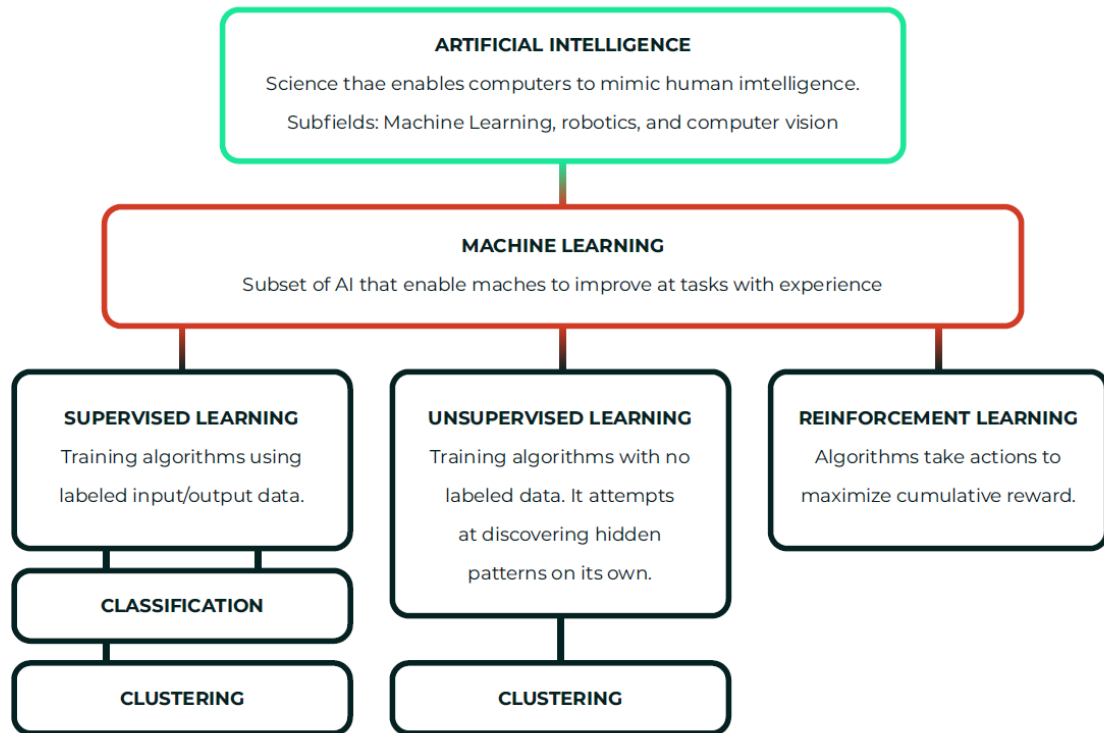


Figure 2.2: Artificial Intelligence and Machine Learning taxonomy [1].

Following this, it's essential to recognize the relationships between machine learning, deep learning, and artificial intelligence, since these are terms that are typically confused or used in an imprecise manner [1].

- **Artificial Intelligence:** is the broadest term, encompassing any technique that enables machines to mimic human intelligence, using logic, if-then rules, decision trees, and machine learning, including deep learning. For instance, chess-playing computers that examine millions of potential outcomes are using a form of AI.
- **Machine Learning:** as explained, is a subset of AI that involves the use of statistical techniques to enable machines to improve at tasks with experience. An example is email filtering systems. They learn to classify emails as 'spam' or 'not spam' based on previous examples that they have seen.
- **Deep Learning:** is a further subset of Machine Learning, inspired by the structure and function of the human brain—namely neural networks. Deep Learning (DL) models learn to perform classification tasks directly from images, text, or sound using large neural networks and a significant amount of data. For example, Deep Learning is used in driverless car technologies to identify objects, such as pedestrians and other vehicles.

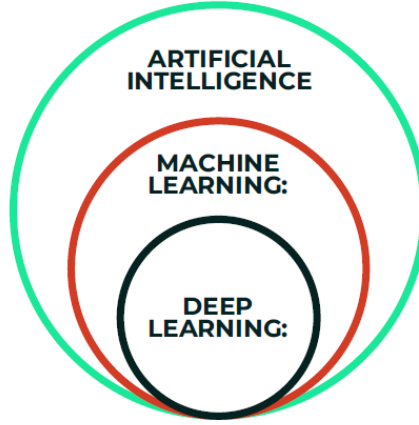


Figure 2.3: Nested structure of AI [1].

While these definitions provide a basic understanding, it's worth noting that the boundaries and interactions between AI, ML, and DL are complex and continually evolving, driven by innovations in theory and increasing computational power. In addition to the aforementioned, machine learning is generally more interpretable and easier to understand than deep learning. It allows for feature importance analysis and provides insights into the data and the underlying processes generating it. On the other hand, deep learning models, especially larger ones, often act as black boxes with little interpretability. However, they often excel at handling unstructured data and can achieve impressive performance on tasks such as image or speech recognition [1].

Python has emerged as one of the most popular programming languages for machine learning, owing to its simple syntax, extensive collection of libraries, and compatibility with a wide range of machine learning frameworks. Libraries such as scikit-learn for general machine learning tasks, TensorFlow for deep learning, and Pandas for data manipulation and analysis have made Python a preferred choice for both beginners and experts in the field. Moreover, the active online community and abundance of educational resources further contribute to its popularity, added to the fact that it is a great option for end-to-end machine learning applications due to its compatibility with diverse data sources and tools for data visualization. As a result, it acts as a complete tool for effectively developing, validating, and deploying machine learning models [1].

It is important to note that Python's role and usage in machine learning will be explored in more detail in Chapter 3 of this thesis. There, we will delve deeper into the specific libraries and techniques employed in the project, providing a better understanding of the practical application of Python in machine learning for enhancing the EXFOR Nuclear Data Library.

2.5 Outliers

Outliers are observations that fall significantly outside the norm within a dataset, typically characterized by their distance from other values in the data distribution [20]. In data analysis, these are critical to consider as they can significantly skew results, leading to inaccurate interpretations and predictions [21].

Outliers can arise for a variety of reasons, such as measurement errors, data entry errors, or genuine extreme observations. Regardless of their source, their importance in data analysis cannot be overstated. Outliers can have a significant impact on our understanding of the data: they can either provide valuable insights into anomalies and trends, or, if they are errors, they can distort the underlying structure and relationships in the data [21]. Thus, outlier detection and management form a critical step in any data analysis pipeline.

Outlier detection techniques can be categorized broadly into statistical methods and machine learning methods [22].

2.5.1 Statistical Techniques

Statistical techniques for outlier detection primarily focus on the identification of data points that deviate from a distribution's expected characteristics. These methods often assume specific statistical properties of the underlying data and are typically best suited for univariate outlier detection [23]. Some of the widely used statistical techniques include:

1. **Z-Score or Standard Score:** This method calculates the Z-score of each data point, which reflects how many standard deviations it is from the mean. Data points with a Z-score above a certain threshold (commonly 3 or -3 for a normal distribution) are considered outliers [24].
2. **Modified Z-Score Method:** This approach, based on the median absolute deviation (MAD), is more robust to outliers than the traditional Z-score method, making it suitable for datasets with heavy-tailed distributions [25].
3. **IQR Method:** This method utilizes the interquartile range (IQR), which is the range between the first quartile (25th percentile) and the third quartile (75th percentile). Any data point that falls below the first quartile minus the IQR times a factor (often 1.5) or above the third quartile plus the IQR times a factor is considered an outlier [26].

2.5.2 Machine Learning Techniques

In contrast to traditional statistical methods, machine learning techniques can handle multivariate outliers and do not require any assumptions about the distribution of the data [27]. Some popular machine learning techniques for outlier detection include:

1. **Clustering Algorithms:** Techniques such as K-Means or DBSCAN create groups (clusters) of data points based on their similarity. Data points that do not belong to any cluster, or belong to small, sparsely populated clusters, can be considered outliers [28].

2. **Anomaly Detection Algorithms:** Algorithms such as One-class SVM or Isolation Forest are specifically designed for outlier detection. They work by learning the “normal” data’s distribution and then identifying points that deviate from this norm [29].
3. **Neural Networks:** Methods such as autoencoders can be used for outlier detection. These networks are trained to reconstruct the input data, and those instances that the network has the most difficulty reconstructing are considered outliers [30].

Overall, the choice of method depends on the specific characteristics of the dataset and the task at hand. It is important to note that not all outliers are “bad” or undesirable. In some cases, outliers can represent valuable information that is worth further investigation. However, if these outliers are the result of errors or noise, it is crucial to detect and, if necessary, remove them to prevent the model from learning from these misleading instances.

3 METHODOLOGY

The Methodology chapter serves as a comprehensive guide to the procedures and tools employed in this research, detailing the process from data collection to analysis. It starts with an explanation of the programming language Python, and the specialized libraries and tools used, such as NumPy, Pandas, Scikit-learn, TensorFlow, Keras, and Matplotlib. These tools play a critical role in handling numerical computations, data manipulation, machine learning algorithms, deep learning tasks, and data visualization, essential elements of this study.

Following the discussion of tools, we address the data collection and the general approach to preprocessing the data, ensuring that it is suitable for analysis. This preprocessing is a significant step and will be explored further in chapter 4. In conjunction with preprocessing, we introduce the methods employed for outlier detection, which include traditional statistical techniques and machine learning approaches, as well as the use of a neural network model for data prediction.

Lastly, we outline the criteria by which the success and performance of the models were evaluated, focusing on the specific metrics for both outlier detection and neural networks. This chapter strives to present the methodology in a clear and understandable manner, laying out the groundwork for the work done in this research. By detailing the methods in a straightforward way, we provide a solid foundation for the implementation and evaluation stages that will be addressed in the Results chapter. The subsequent sections of this chapter will delve into these aspects in greater detail, providing readers with a clear understanding of the technical aspects of this research project.

3.1 Tools and Libraries Used

In the context of this research, Python serves as our primary language for implementing machine learning models, leveraging its powerful suite of libraries. Our methodology focuses on several key Python libraries and frameworks that are widely used in the field of machine learning:

- **NumPy:** The foundational library for numerical computing in Python, NumPy (Numerical Python), provides support for arrays, matrices, and high-level mathematical functions to operate on these. NumPy is indispensable for handling numerical data and performing mathematical operations, which are fundamental in machine learning.
- **Pandas:** This library is used for data manipulation and analysis. It provides data structures and functions needed to manipulate structured data. Features like handling missing data, merging, reshaping, selecting, replacing, data filtration, and others make it perfect for preprocessing the data for machine learning tasks.
- **Scikit-learn:** This is one of the most widely used machine learning libraries in Python. It provides a range of supervised and unsupervised learning algorithms. Scikit-learn comes with standard datasets, various machine learning algorithms (from linear regression to clustering), and also offers tools for model fitting, data preprocessing, model selection, and evaluation.
- **TensorFlow:** Originally developed by Google Brain, TensorFlow is an open-source library for numerical computation and large-scale machine learning. TensorFlow bundles together a slew of machine learning and deep learning (aka neural networking) models and algorithms and makes them useful by way of a common metaphor. It's known for its capabilities in handling deep learning tasks.

- **Keras:** Is a high-level neural networks API, capable of running on top of TensorFlow, and is used for building and training deep learning models. It's user-friendly and easy to use, making it a popular choice for beginners and experts in deep learning.
- **Matplotlib:** Data visualization is a crucial part of machine learning. Matplotlib is a popular library for creating static, animated, and interactive visualizations in Python, aiding in the analysis and interpretation of the computational results.

In terms of workflow, we first prepare the data using Pandas and NumPy, carry out preprocessing tasks, and then apply appropriate machine learning algorithms from Scikit-learn, TensorFlow, or Keras as required. The data and the performance of the models are visualized using Matplotlib.

3.2 Preprocessing and Data Mining

The preprocessing stage represents a vital phase in our methodology, dedicated to preparing the raw data for subsequent analysis. This complex and multifaceted process involves numerous tasks such as cleaning and transforming the data, selecting appropriate variables, and engineering the features that will be used in the models. These efforts set the stage for both the outlier detection and neural network modeling that follow. In this section, we will delve into the specifics of this comprehensive preprocessing effort, detailing the techniques and approaches applied to ensure that the data is in the optimal format for the analysis techniques utilized in subsequent chapters.

Preprocessing and data mining are essential components in machine learning applications. The success of algorithms relies heavily on the structure and quality of the dataset. The EXFOR nuclear data library, comprising more than 24,000 experiments, each with unique characteristics, measurements, and features, presents a significant challenge for several reasons.

First, the sheer volume of data in EXFOR is immense, demanding careful registration and management. Second, the library contains a heterogeneous mix of data, reflecting the wide array of target elements, incident particles, reactions, and quantities measured. Third, the lack of a consistent structure for storing data further complicates matters. Although there is a general structure followed within EXFOR, as mentioned in section 2.1, the earliest experiments from the 1950s and 60s lacked digitalization and common structure. When these experiments were digitized, some information might have been lost. Even in more modern experiments, not all details might be shared, resulting in a heterogeneous and complex library.

Efforts like MEDUSAL and EXFORTABLES, also referenced in 2.1, have been made over the past decade to mitigate these challenges.

EXFORTABLES, in particular, serves as a valuable starting point. It parses the entire EXFOR library, creating a logically structured directory on your computer, categorized by projectile, element, mass, and reaction. Moreover, it formats the information in a more computationally friendly way. Figure 3.1 illustrates an example of a file located at:

```
exfortables/p/Ca044/xs/004/p-Ca044-MT004-Krajewski-02136002.2013
```

As seen, the experiment-related information is neatly structured, beginning with a '#', while the data points are presented in table format without the hashtag. This approach makes the data more accessible to computer code, such as Python, compared to the original ENDF format used in EXFOR.

```
# Target Z      : 20
# Target A      : 44
# Target state:
# Projectile    : p
# Reaction      : (p,n')
# Final state   :
# Quantity      : Cross section
# Frame         : L
# MF            : 3
# MT            : 4
# X4 ID         : 02136002
# X4 code       : 20-CA-44(P,N)21-SC-44,,SIG
# Author        : Krajewski
# Year          : 2013
# Data points   : 8
#
#   E(MeV)      xs(mb)      dxs(mb)      dE(MeV)
#   -----
# 5.53000E+00   8.42000E+00  0.00000E+00  0.00000E+00
# 7.28000E+00   1.44530E+02  0.00000E+00  0.00000E+00
# 7.86000E+00   2.77330E+02  0.00000E+00  0.00000E+00
# 8.84000E+00   6.95480E+02  0.00000E+00  0.00000E+00
# 9.83000E+00   5.64270E+02  0.00000E+00  0.00000E+00
# 1.06300E+01   7.28140E+02  0.00000E+00  0.00000E+00
# 1.31400E+01   6.33630E+02  0.00000E+00  0.00000E+00
# 1.74800E+01   4.13510E+02  0.00000E+00  0.00000E+00
#
# Reference:
#S.Krajewski, I.Cydzik, K.Abbas, A.Bulgheroni, F.Simonelli, U.Holzwarth, A.Bilewicz
#Cyclotron production of 44Sc for clinical application.
#Jour. Radiochimica Acta Vol.101, p.333, 2013
```

Figure 3.1: Example entry from EXFORTABLES.

Being able to access information in a well-structured way is valuable but only the first step in the path to obtain usable data for machine learning. It becomes essential at this point to define clearly what type of information we need to apply machine learning. This involves understanding our objective and how the information should be ideally structured to use it efficiently.

Before moving on to machine learning analysis, an essential preprocessing step is the scaling of feature values. Given that most machine learning algorithms use distance-based metrics, the scale of each feature can affect the model's performance. For instance, a feature with a large range could dominate the objective function and make the estimator unable to learn from other features correctly. In this study, we utilize the `StandardScaler` provided by the `Scikit-learn` library to standardize the dataset's features by removing the mean and scaling to unit variance. The formula for standard scaling is:

$$z = \frac{(x - u)}{s} \quad (3.1)$$

Where x is the original feature vector, u is the mean of that feature vector, and s is its standard deviation. This scaling ensures that each feature contributes equally to the distance computations during machine learning model training.

In general terms, when it comes to machine learning and Python, the preferred way to structure your dataset is in a dataframe. A dataframe in Python is a two-dimensional, size-mutable, and heterogeneous tabular data structure with labeled axes (rows and columns). It's like a

spreadsheet or SQL table and is generally understood to be the most common Pandas object, allowing for the easy manipulation, analysis, and visualization of data.

But structuring the data in a dataframe is not the only requirement. Most machine learning algorithms rely on numerical values, meaning that categorical values may need to be converted into numerical ones. There are different methods to achieve this, but a common one is the One-hot encoding technique. One-hot encoding is a process that transforms categorical data variables so they can be provided to machine learning algorithms to improve predictions. It involves converting each categorical value into a new categorical column and assigns a binary value of 1 or 0. Each integer is represented as a binary vector, with all the values zero except the index, marked with a 1.

At the end of this process, we should have a dataframe, akin to a table, where each row represents a data point from an experiment, and each column corresponds to a feature that helps identify and characterize it. All the data points for the same experiment will have identical values in the columns that characterize the experiment they belong to. And the categorical values are transformed through the one-hot encoding into columns where the values are 0s or 1s, depending on whether that column represents the data point.

So first, all experiments parsed with EXFORTABLES must be read from the files and stored in a way that makes them easily accessible at any point. For this, it was decided to use Python's object-oriented coding, so a class named Experiment was defined with different attributes representing all the characteristics of the experiment. The utilization of object-oriented coding in this context offers several advantages:

1. **Encapsulation:** By bundling the data with the methods that operate on it, the implementation details are hidden from the user, ensuring consistency.
2. **Modularity:** Code is organized into objects, making it easier to maintain, modify, and understand.
3. **Reusability:** Objects and classes can be reused across different parts of the project or in different projects, promoting code efficiency.
4. **Inheritance:** It enables the creation of new classes that are built upon existing classes allowing code reuse and the addition of new features.

This approach also allows us to create different methods, such as a function for automatically plotting the results of an experiment.

Once the class is defined, we can code a function to read each experiment, parsing the information line by line from EXFORTABLES, as shown in Figure 3.1. Although saving each experiment in a Python list is not the most practical way to handle the data, our goal is to read all the experiments we are interested in (in our case, those using protons as incident particles) and save them in a file. This file can be easily read and transferred, allowing quick and straightforward access to the list of experiments without requiring each user to have the entire EXFORTABLES files on their computer.

With the data stored in this manner, we have a valuable format that is suitable for various applications such as data analysis and information retrieval. By scripting the experiments through objects and storing the information in a list, we can create functions to perform different tasks, such as:

- Getting all unique values of each attribute.
- Filtering experiments based on the value of a specific attribute (example: query all experiments in which MT=102).
- Classifying experiments by attributes (example: group experiments by MT).
- Plotting one or more experiments at once.

Indeed, one of these tasks, classifying, is crucial for the next step. While it would be feasible to create a dataframe with all the data points for every proton-induced experiment, it doesn't necessarily align with our analysis objectives. Since our main goal is to detect outliers, grouping data based on the value being measured seems more practical. Fortunately, uncertainties in the measured value or energy used are already treated, ensuring consistent column appearances. This uniformity simplifies comparison and enhances classification.

Lastly, the data within each group must be converted into a dataframe, following the specifications mentioned earlier. By meticulously organizing and preparing the data in this manner, we set the stage for accurate and efficient machine learning analysis, aligning with our primary focus on outlier detection.

To sum up this section, we could divide the preprocessing and data mining into the following key points:

1. **Utilizing EXFORTABLES:** Initial parsing of the EXFOR nuclear data library using EXFORTABLES to structure it in a computationally friendly manner.
2. **Defining Objectives:** Clarifying the information requirements and the ideal data structure for machine learning tasks.
3. **Feature Scaling:** Utilizing StandardScaler to standardize the features, thereby enhancing the algorithm's ability to learn effectively from all features.
4. **DataFrames as Data Structure:** Converting the cleaned and parsed data into Python dataframes for easy manipulation and analysis.
5. **Categorical to Numerical:** Applying One-hot encoding to transform categorical values into a format suitable for machine learning algorithms.
6. **Object-Oriented Design:** Creating a Python class named `Experiment` to encapsulate the attributes and methods associated with individual experiments.
7. **Data Reading and Storage:** Developing a function to read experiments one by one from EXFORTABLES and storing them in an easily accessible file.
8. **Data Manipulation Functions:** Implementing functionalities to filter, classify, and plot experiments based on various attributes.
9. **Grouping for Analysis:** Classifying experiments into groups based on the value being measured to facilitate outlier detection.
10. **Final DataFrame Creation:** Converting the grouped data into a final dataframe, ready for machine learning applications.

By following these steps, the data is prepared in a structured manner that is not only conducive to machine learning but also efficient for specific tasks such as outlier detection. The application and results of these preprocessing and data mining steps will be further explored in the chapter 4.

3.3 Interquartile Range (IQR)

The Interquartile Range (IQR) is a statistical measure widely used for identifying potential outliers in a dataset. Originating from the field of descriptive statistics, it provides a straightforward yet robust way to understand data variability around the median. Unlike methods that use mean and standard deviation, the IQR focuses on the middle 50% of data points, making it less sensitive to extreme values or outliers.

The IQR is calculated as the difference between the third quartile ($Q3$) and the first quartile ($Q1$) of the dataset. These quartiles divide a rank-ordered data set into four equal parts. Specifically, $Q1$ is the median of the lower half, not including the overall median if the number of elements in the dataset is odd, and $Q3$ is the median of the upper half. Mathematically, the IQR is represented as:

$$\text{IQR} = Q3 - Q1 \quad (3.2)$$

To identify outliers using the IQR, the following steps are generally followed:

- Calculate $Q1$ and $Q3$ for the dataset.
- Compute the IQR using Equation 3.2.
- Define the inner fences:

$$\text{Lower Inner Fence} = Q1 - 1.5 \times \text{IQR} \quad (3.3)$$

$$\text{Upper Inner Fence} = Q3 + 1.5 \times \text{IQR} \quad (3.4)$$

- Data points outside these inner fences are considered as outliers.

A graphical representation of how the IQR method works is shown in Figure 3.2 below.

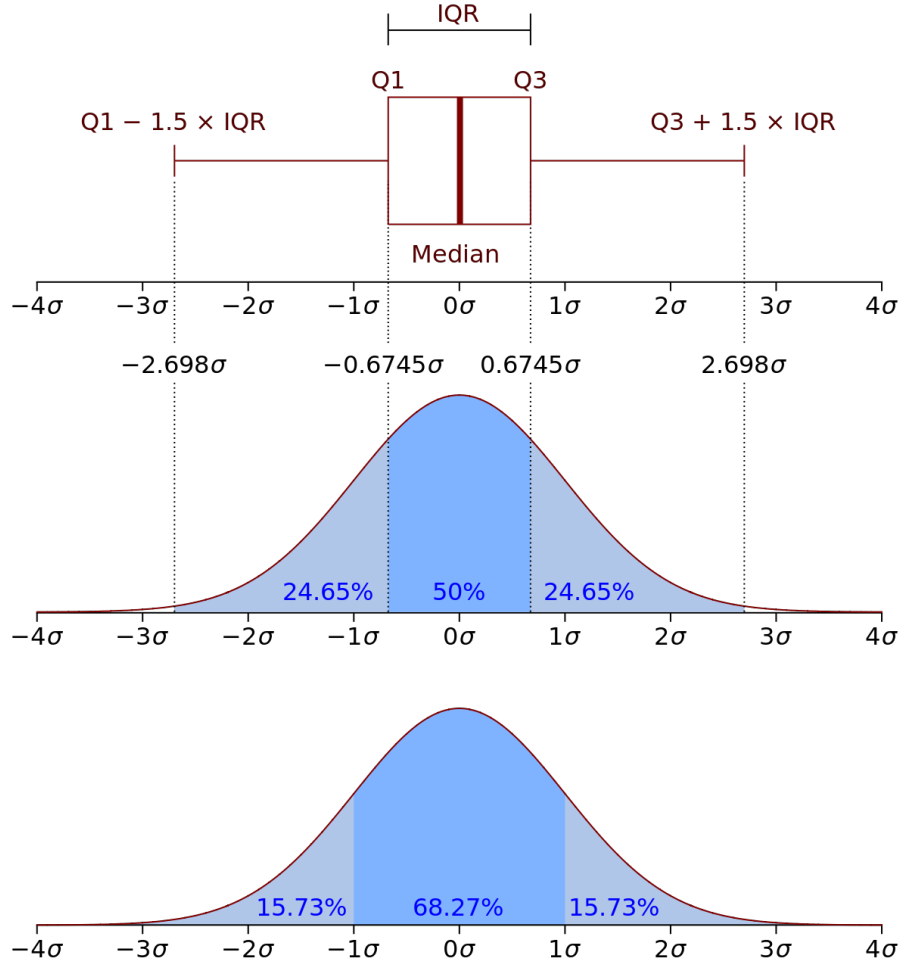


Figure 3.2: Graphical representation of the IQR method [31].

While the IQR method's simplicity and robustness are appealing, it has its limitations. For instance, the IQR may be less sensitive to outliers that are close to the inner fences. Moreover, the risk of data loss is another concern, as some genuine data points may be labeled as outliers, especially if the multiplier for calculating the inner fences is not appropriately chosen.

Additionally, the effectiveness of the IQR method is influenced by the distribution of the data points in the dataset. For skewed or heavy-tailed distributions, the IQR might not effectively identify outliers, making it crucial to select the multiplier for the inner fences carefully.

3.4 Isolation Forest

Isolation Forest is an ensemble-based machine learning algorithm for anomaly detection. Unlike traditional methods which seek to construct a profile of normal instances, Isolation Forest isolates anomalies by exploiting the inherent properties of tree-based partitioning. The basic idea is to recursively partition the feature space by randomly selecting a feature and a random value between the feature's minimum and maximum values.

The Isolation Forest algorithm has the following key steps:

- Randomly sample the data and construct an isolation tree by recursively partitioning the

feature space.

- Continue this process until a predefined number of trees are generated, forming the Isolation Forest.

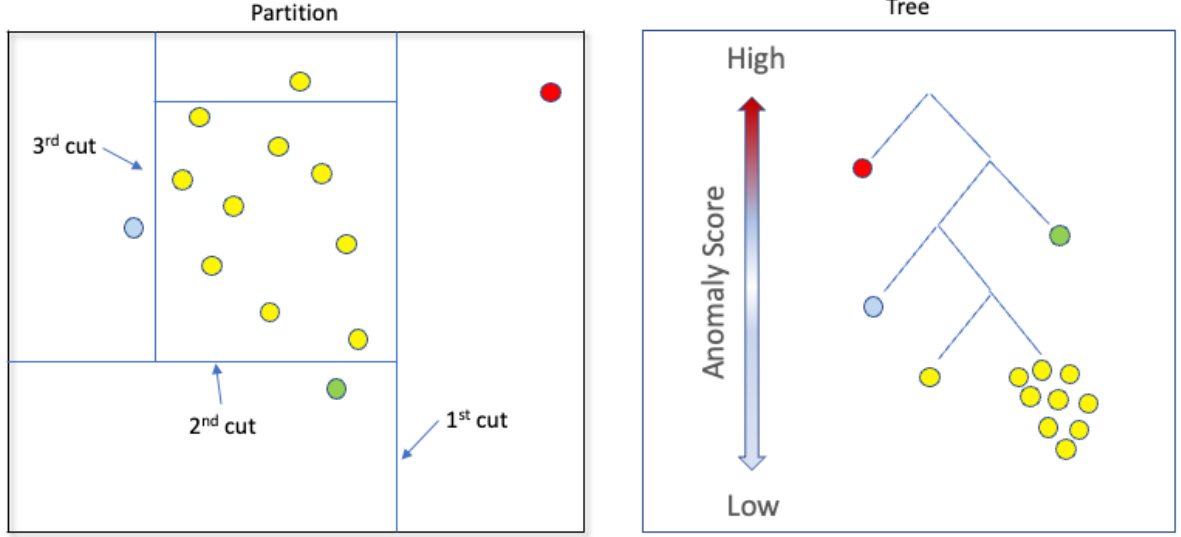


Figure 3.3: Graphical representation of the Isolation Forest method [32].

The anomaly score is calculated using the equation shown in Equation 3.5:

$$s(x, n) = 2^{-\frac{E(h(x))}{c(n)}} \quad (3.5)$$

Where $s(x, n)$ is the anomaly score of sample x , $E(h(x))$ is the average path length of x , n is the number of external nodes, and $c(n)$ is the average path length of unsuccessful search in a binary search tree.

Isolation Forest excels in its computational efficiency and its ability to handle high-dimensional datasets. It also has the advantage of not making any assumptions about the underlying distribution of the data. However, Isolation Forest may have difficulties in distinguishing between noise and actual anomalies, especially in datasets with strong overlap between the two categories. Moreover, because it uses random partitioning, the results are not always deterministic, which may be a downside in certain applications.

The distribution of points in the dataset can influence the performance of Isolation Forest. In dense clusters where anomalies are close to normal instances, the algorithm may fail to distinguish them adequately. Conversely, in well-separated data spaces, the algorithm performs remarkably well.

By understanding and implementing the Isolation Forest algorithm following these guidelines, one can effectively detect outliers in complex datasets where traditional statistical methods may fall short.

3.5 DBSCAN

DBSCAN is a density-based clustering algorithm that identifies clusters in a dataset and treats points that do not belong to any cluster as outliers. Unlike partitioning methods like K-means, DBSCAN does not require the number of clusters to be specified in advance. It works by defining a neighborhood around each data point and finding clusters based on the density of these neighborhoods.

The DBSCAN algorithm consists of the following key steps:

1. Choose two parameters: Eps , the radius of the neighborhood, and $MinPts$, the minimum number of points required to form a dense region.
2. For each unvisited point, identify its Eps -neighborhood. If it contains at least $MinPts$, a cluster is formed.
3. Expand the cluster by including points within the Eps -neighborhood of each point in the cluster.
4. Iterate over each unvisited point, repeating steps 2 and 3 until all points have been visited.

To formalize the concept of density and clustering, DBSCAN utilizes the following equation to define a neighborhood:

$$N_{Eps}(p) = \{q \in D : \text{dist}(p, q) \leq Eps\} \quad (3.6)$$

Where $N_{Eps}(p)$ is the Eps -neighborhood of point p , D is the dataset, and $\text{dist}(p, q)$ is a distance metric (usually Euclidean distance).

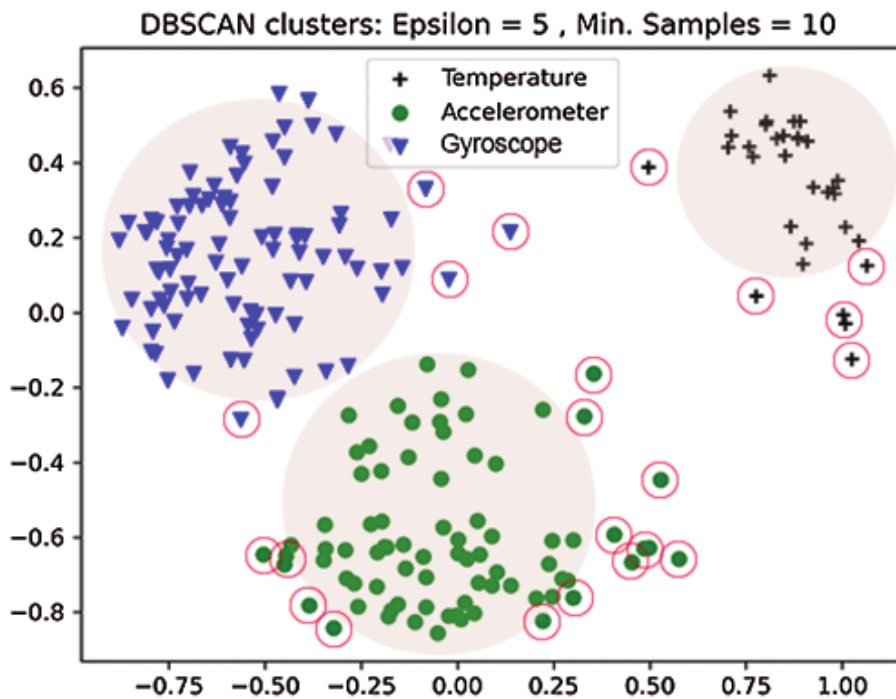


Figure 3.4: Graphical representation of the DBSCAN method [33].

DBSCAN has several strengths and weaknesses. One of its primary advantages is that it can find arbitrarily shaped clusters, unlike spherical clustering methods like K-means. Moreover, it inherently handles outliers by classifying them as points that do not belong to any cluster. On the downside, DBSCAN is sensitive to the choice of *Eps* and *MinPts*, and its performance can degrade with high-dimensional data. The algorithm can also struggle with clusters of varying densities.

The effectiveness of DBSCAN in outlier detection is directly influenced by the distribution of data points in the dataset. If points are uniformly distributed across various densities, DBSCAN may find it challenging to identify the boundaries between clusters and outliers accurately. However, in datasets where the differences in density between normal instances and outliers are distinct, DBSCAN can be incredibly effective.

3.6 Local Outlier Factor (LOF)

The Local Outlier Factor (LOF) algorithm is a method used for outlier detection that focuses on the local properties of data points rather than global characteristics. Unlike traditional methods that define outliers based on some global property of the dataset, LOF considers each data point within the context of its neighbors.

The core concept in LOF is ‘locality’, which is evaluated based on the density deviation of a data point in relation to its neighbors. The algorithm uses several steps to assess the local outlier-ness of each data point. First, for each data point, it calculates the distance to its k^{th} nearest neighbor, commonly referred to as k -distance. Mathematically, this is represented as:

$$\text{k-distance}(A) = \min_{B \in N_k(A)} d(A, B) \quad (3.7)$$

Second, the algorithm computes the Reachability Distance between a data point and its neighbors. This distance is defined as the maximum of the k -distance of the neighbor and the actual distance to the neighbor, formulated as:

$$\text{ReachabilityDistance}_k(A, B) = \max\{\text{k-distance}(B), d(A, B)\} \quad (3.8)$$

Subsequently, the Local Reachability Density (LRD) of each data point is calculated. LRD is the inverse of the average reachability distance of the k neighbors of the data point:

$$\text{LRD}_k(A) = \left(\frac{\sum_{B \in N_k(A)} \text{ReachabilityDistance}_k(A, B)}{|N_k(A)|} \right)^{-1} \quad (3.9)$$

Finally, the LOF score of each data point is computed. It is the average ratio of the LRDs between the data point and its k neighbors:

$$\text{LOF}_k(A) = \frac{\sum_{B \in N_k(A)} \frac{\text{LRD}_k(B)}{\text{LRD}_k(A)}}{|N_k(A)|} \quad (3.10)$$

LOF identifies outliers by scrutinizing the LOF score (3.10) of each data point. Specifically, a data point with a LOF score significantly greater than 1 is considered an outlier, as it is less dense than its neighbors.

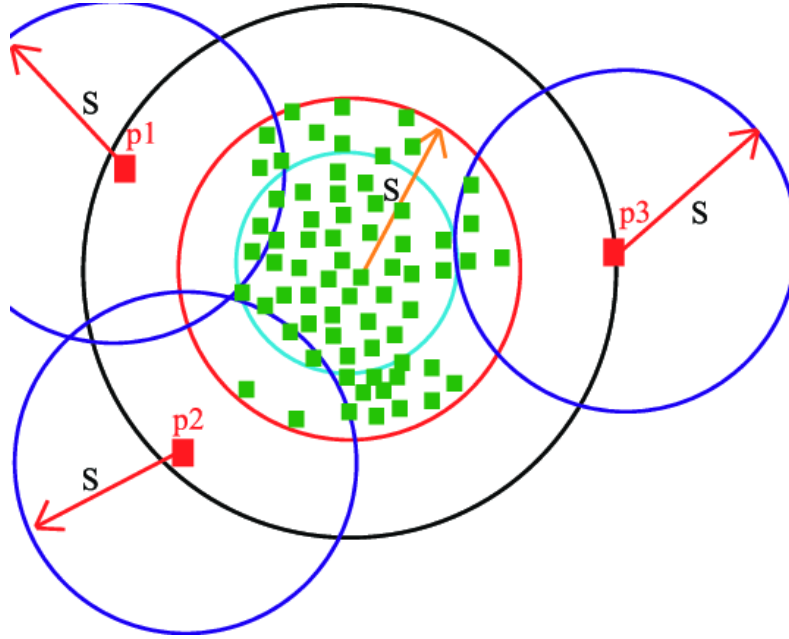


Figure 3.5: Graphical representation of the LOF method [34].

LOF has distinct advantages and disadvantages. One of its strengths is its adaptability to varying densities within a dataset, providing a nuanced method for outlier detection. It does not require prior knowledge about the number of outliers, which is an advantage for dynamically changing data. On the downside, the algorithm is computationally expensive, especially for high-dimensional datasets. The choice of the parameter k can also greatly impact the results. Most pertinent to our research, the distribution of data points within the dataset can affect the performance of LOF. Specifically, if the dataset is comprised of large regions of sparse data, LOF may have difficulties in accurately distinguishing between outliers and noise.

3.7 Support Vector Machines (SVM)

Support Vector Machines (SVM) are supervised learning models primarily used for classification and regression tasks. The main idea behind SVM is to find a hyperplane that best separates the data into different classes. In a two-dimensional space, this hyperplane is a line. For higher dimensions, it becomes a plane or a set of planes that aim to maximize the margin between different classes of data points.

The hyperplane is represented by the equation $wx - b = 0$, where w is the weight vector and b is the bias. To find the optimal hyperplane, SVM aims to maximize the margin M , which is the distance between the nearest data point from any class to the hyperplane. Mathematically, the margin is expressed as:

$$M = \frac{2}{\|w\|} \quad (3.11)$$

The algorithm operates by minimizing the following objective function with respect to w and b :

$$\min_{w,b} \frac{1}{2} \|w\|^2 \quad (3.12)$$

subject to the constraint $y_i(wx_i - b) \geq 1$ for each data point x_i and its corresponding label y_i .

For non-linearly separable data, SVM employs the kernel trick to implicitly map the input space to a higher-dimensional feature space. Common kernels include polynomial, radial basis function (RBF), and sigmoid kernels. The choice of kernel can significantly affect the performance of the model.

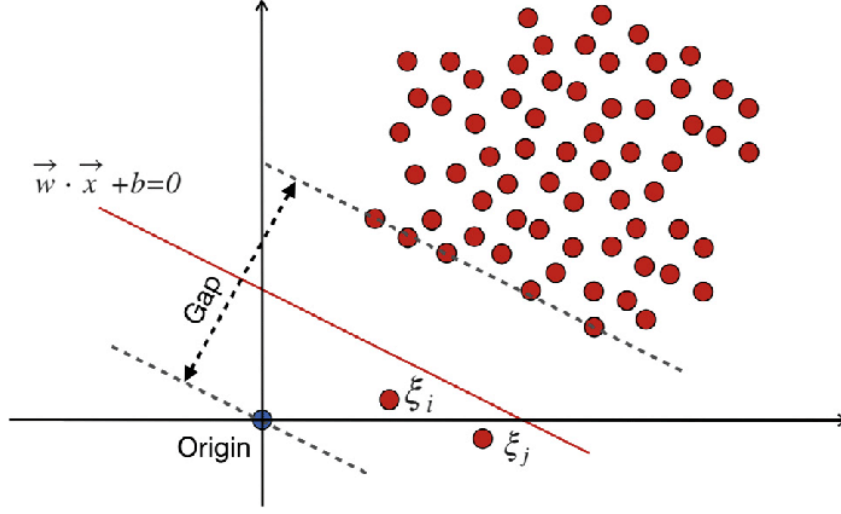


Figure 3.6: Graphical representation of the SVM method [35].

SVM offers several advantages. It is effective in high-dimensional spaces and is versatile, thanks to the various kernel functions. Furthermore, only a subset of training points, known as support vectors, are used in the decision function, making the algorithm efficient. However, SVMs have their downsides. They are sensitive to the choice of the kernel and its parameters, and they can be computationally intensive for large datasets. Additionally, SVMs are not well-suited for multi-class problems and may require complex schemes like one-vs-one or one-vs-all for such tasks.

3.8 Neural Network Models: Autoencoders

Neural networks are a subclass of machine learning algorithms modeled after the human brain. They are composed of layers of interconnected nodes or neurons. Among various types of neural network architectures, autoencoders are unsupervised neural networks employed for anomaly detection, dimensionality reduction, and feature learning, among other tasks.

An autoencoder network is typically symmetrical about its central hidden layer, known as the encoding layer. The network can be divided into two primary components: the encoder and the decoder. The encoder compresses the input into a latent-space representation and encodes the input as an internal fixed-size representation in reduced dimensionality. The decoder, on the other hand, works to reconstruct the input data from this internal representation.

Mathematically, an encoder function f maps an input x to a hidden representation h :

$$h = f(x; W) \quad (3.13)$$

where W are the weights of the encoder.

The decoder function g then maps h back to the reconstructed input \hat{x} :

$$\hat{x} = g(h; W') \quad (3.14)$$

where W' are the weights of the decoder.

The primary objective of an autoencoder is to minimize the reconstruction error, usually represented by a loss function L :

$$L(x, \hat{x}) = ||x - \hat{x}|| \quad (3.15)$$

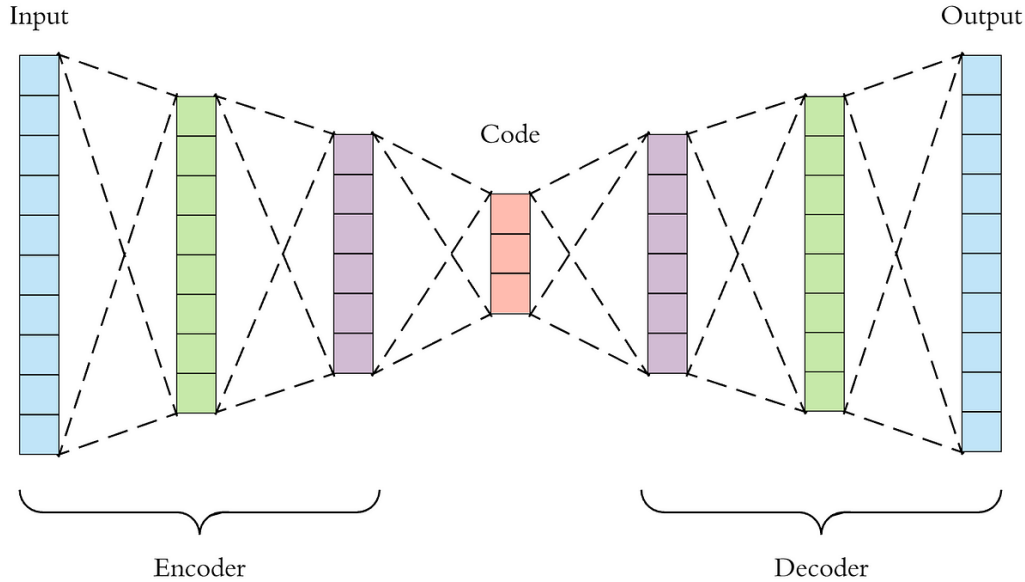


Figure 3.7: Graphical representation of an autoencoder neural network [36].

In an autoencoder network, the architecture often starts with an input layer that matches the dimensionality of the feature space. This is followed by one or more hidden layers that progressively reduce the dimensionality, leading to the central encoding layer. The decoder starts from the encoding layer and gradually increases the dimensions through its hidden layers, culminating in an output layer of the same size as the input layer.

The choice of the number of layers and neurons within those layers is often problem-dependent and influenced by the complexity of the input data. More complex data might require deeper architectures with more neurons, but care should be taken to avoid overfitting, where the model learns the training data too well but performs poorly on unseen data.

Hyperparameters are parameters whose values are set before the learning process begins. In the context of autoencoders, important hyperparameters include:

- *Learning Rate*: Controls how much to update the model parameters during training. A too-high rate can cause the model to converge too quickly and possibly overshoot the optimal values. A too-low rate could result in a very long training process.
- *Batch Size*: The number of samples processed before updating the model parameters. Smaller batch sizes often provide a regularizing effect and lower generalization error.
- *Number of Epochs*: An epoch is one complete forward and backward pass of all the training examples. The training process is often repeated for a certain number of epochs to ensure the model learns the representation effectively.
- *Activation Functions*: Functions like ReLU (Rectified Linear Unit), sigmoid, or tanh, applied to the neurons' output.

- *Regularization Techniques:* Such as dropout or L1/L2 regularization, these techniques can help prevent overfitting.

Hyperparameter optimization often involves techniques like grid search, random search, or more sophisticated methods like Bayesian optimization. The aim is to find the hyperparameter values that yield the most accurate model in terms of minimizing the loss function (3.15).

The advantage of autoencoders is their ability to learn non-linear transformations with a simple architecture, making them suitable for complex tasks such as feature learning and anomaly detection. Additionally, they can automatically learn features necessary for various tasks, thereby reducing the need for manual feature engineering. However, they have disadvantages too. Autoencoders are sensitive to the architecture and hyperparameters; a poor choice can result in overfitting or underfitting. Furthermore, training autoencoders can be computationally intensive and may require a lot of data for the model to generalize well.

In summary, autoencoders are powerful tools in the neural network arsenal that offer both sophistication and depth in solving complex machine learning problems. However, a nuanced understanding and careful selection of architecture and hyperparameters are crucial for optimal performance.

3.9 Evaluation Metrics

Evaluating the efficiency of outlier detection methods is an essential step, especially when employing a range of techniques. To aid in this evaluation, custom software scripts will be developed, serving the specific purpose of facilitating the visual inspection process. These tailored tools will enable the graphical representation of detected outliers, providing a contextual view within the dataset for more accurate assessment.

Visual inspection serves as an immediate and qualitative method to assess how well each algorithm is performing. This approach becomes particularly valuable when the dataset lacks a well-defined set of known outliers, making it challenging to establish a ground truth for comparison.

A significant point to consider is that many of the algorithms employed in this research are configured to identify a specific number of outliers. It's essential to understand that this predetermined number may not align with the actual number of outliers present in the dataset. Therefore, it's possible for these algorithms to label points as outliers even when they might not be. Given this context, thorough visual inspection becomes indispensable for distinguishing between true and false positives.

4 RESULTS

This chapter is dedicated to the comprehensive presentation and analysis of the results obtained from various outlier detection methods applied to the EXFOR Proton Reactions dataset. The purpose is to evaluate the effectiveness of these methods in identifying outliers.

The chapter is organized into several sections, each focusing on a specific outlier detection method. We begin with data preprocessing, setting the groundwork for the subsequent analyses. Then, traditional statistical and machine learning methods such as the Interquartile Range (IQR), Isolation Forest, DBSCAN, and Local Outlier Factor (LOF) are discussed. Each of these sections covers the implementation, graphical analysis, and concludes with an evaluation of the method's efficacy. This is followed by a foray into neural network models, specifically focusing on Autoencoders, which offers a more complex but potentially more effective approach for outlier detection.

Due to the large volume of data and the nature of the algorithms implemented, it is impractical to display all the results from the execution of the various scripts created. As such, only a subset of the results will be presented in the 'Results and Graphical Analysis' subsection of each method. However, all the necessary scripts for reproducing these results are included in this work and can also be accessed at GitHub EXFOR-ProtonReactions-Analysis Repository.

4.1 Data Preprocessing

Data preprocessing is a crucial and often time-consuming stage in both machine learning and various data analytics initiatives. While the core framework for this process was outlined in Section 3.2 of Chapter 3, this section provides a more in-depth look. We will focus on the results from each preprocessing step, discuss how these results are useful in machine learning and data analytics, and consider their potential applications.

To facilitate the replication or adaptation of our work, a Python Jupyter Notebook named `EXFOR.ProtonReactions.Data.Preprocessing.ipynb` has been developed. This notebook encompasses all steps required for data preprocessing. We will provide a line-by-line walkthrough of the script to clarify the logic and mechanics behind each procedure. The notebook is accessible via the GitHub repository managed by monleon96, which can be found at GitHub EXFOR-ProtonReactions-Analysis Repository. Additionally, all custom Python functions created specifically for this project are documented in the Annexes of this thesis.

Our initial assumptions are that `EXFOR` is pre-installed on your computer with its default folder structure, as outlined in reference [37]. We also presume that you have successfully imported the `Experiment` Python class from Annex A as well as the utility functions from Annex B. Both are available in the aforementioned GitHub repository under the filenames:

- `EXFOR.ProtonReactions.Experiment.Class.py`
- `EXFOR.ProtonReactions.UtilityFunctions.py`

With these prerequisites met, you can initiate a new Jupyter notebook and access the experiments from `EXFOR` by specifying the computer path to the proton-induced reactions, termed as 'p'.

```

1 path = "path_to_proton-induced_reactions"
2 experiments = read_proton_experiments_from_exfortables(path)

```

The function `read_proton_experiments_from_exfortables()` (see Annex B.6) reads each experiment's files from EXFORTABLES. It puts the important details into an `Experiment` object (see Annex A). When you run this, the variable `experiments` will hold a list of `Experiment` objects, which contain all the needed data from EXFOR's proton experiments.

To make things easier and avoid reading the EXFORTABLES files every time, we have two ways to save this data: as a binary file or as a text file (see Annex B.2 and B.4).

```

1 write_experiments_to_binary(experiments, "EXFOR_ProtonReactions_Database
  .bin")
2 write_experiments_to_txt(experiments, "EXFOR_ProtonReactions_Database.
  txt")

```

Using these functions, you can create either a binary or a text file. Later on, you only need to read one of these files to get back the saved data. This can be done with the functions `read_experiments_from_binary()` and `read_experiments_from_txt()` (see Annex B.3 and B.5).

```

1 exp_from_bin = read_experiments_from_binary("
  EXFOR_ProtonReactions_Database.bin")
2 exp_from_txt = read_experiments_from_txt("EXFOR_ProtonReactions_Database
  .txt")

```

Now, you might wonder why we offer two types of files. Binary files are quicker to load but rely on Python's pickle module. This can sometimes lead to issues with serialization—essentially, converting the data back and forth between file and program. Text files are more straightforward but might take longer to read into the program.

Once the experiments are loaded into a Python list, you can see what is inside by just printing one item from the list. This will show you all the details for that specific experiment.

```

1 print(exp_from_bin[3])

```

Output:

```

1 Title:  p-Ag000-MT002-Heiberg-004550051-ang-E0424.000.1957
2 Target Z:  47
3 Target A:  0
4 Projectile:  p
5 Reaction:  (p, el)
6 Incident energy:  424.000 MeV
7 Quantity:  Angular distribution
8 Frame:  L
9 MF:  4
10 MT:  2
11 X4 ID:  004550051
12 X4 code:  47-AG-0(P,EL)47-AG-0,,DA
13 Author:  Heiberg

```

4. RESULTS

```
14 Year: 1957
15 Data points: 15
16 Data:      Angle(deg)  xs(mb/sr)  dxs(mb/sr)  dAngle(deg)
17 0      5.69977      4710.0      500.0      NaN
18 1      6.39997      2520.0      110.0      NaN
19 2      7.59980      379.0       48.0      NaN
20 3      8.29993      425.0       28.0      NaN
21 4      9.00012      400.0       27.0      NaN
22 5      9.79996      397.0       40.0      NaN
23 6     10.40020      539.0       33.0      NaN
24 7     11.60010      263.0       14.0      NaN
25 8     12.29990      237.0       15.0      NaN
26 9     12.80010      118.0        6.0      NaN
27 10     13.90010      47.2        5.9      NaN
28 11     14.80010      51.9       11.8      NaN
29 12     15.90010      42.5        5.9      NaN
30 13     18.50010      29.5        3.5      NaN
31 14     19.20010      16.5        2.4      NaN
32 Reference:
33 E.Heiberg
34 Angular Distribution of 424-MeV Polarized Protons Elastically Scattered
35    from Various Nuclei.
    Jour. Physical Review Vol.106, p.1271, 195
```

You can also focus on just the data points by using a dataframe. This makes it easy to handle the data for further analysis.

```
1 exp_from_bin[3].data
```

Output:

```
1      Angle(deg)  xs(mb/sr)  dxs(mb/sr)  dAngle(deg)
2 0      5.69977      4710.0      500.0      NaN
3 1      6.39997      2520.0      110.0      NaN
4 2      7.59980      379.0       48.0      NaN
5 3      8.29993      425.0       28.0      NaN
6 4      9.00012      400.0       27.0      NaN
7 ..      ...      ...      ...      ...
8 10     13.90010      47.2        5.9      NaN
9 11     14.80010      51.9       11.8      NaN
10 12     15.90010      42.5        5.9      NaN
11 13     18.50010      29.5        3.5      NaN
12 14     19.20010      16.5        2.4      NaN
13
14 [15 rows x 4 columns]
```

Another helpful feature is the ability to convert an entire experiment's information into a dataframe. Here, each row represents a data point from the experiment, and the experiment's attributes become the columns of the dataframe.

```
1 exp_from_bin[3].to_dataframe()
```

Output:

```

1      Angle(deg)  xs(mb/sr)  dxs(mb/sr)  ...  MT  Author  Year
2  0      5.69977    4710.0    500.0    ...  2  Heiberg  1957
3  1      6.39997    2520.0    110.0    ...  2  Heiberg  1957
4  2      7.59980    379.0     48.0    ...  2  Heiberg  1957
5  3      8.29993    425.0     28.0    ...  2  Heiberg  1957
6  4      9.00012    400.0     27.0    ...  2  Heiberg  1957
7  ..      ...      ...      ...      ...  ..      ...      ...
8  10     13.90010    47.2      5.9     ...  2  Heiberg  1957
9  11     14.80010    51.9     11.8    ...  2  Heiberg  1957
10 12     15.90010    42.5      5.9     ...  2  Heiberg  1957
11 13     18.50010    29.5      3.5     ...  2  Heiberg  1957
12 14     19.20010    16.5      2.4     ...  2  Heiberg  1957
13
14 [15 rows x 21 columns]

```

You can also make plots directly using a method from the class. This method offers optional settings to tweak the graph's appearance, like its size or the scales of the x and y axes. If the experiment includes uncertainty data, the plot will also display error bars, giving you a full view of the data's quality.

```
1 exp_from_bin[0].plot(ylog=True)
```

Output:

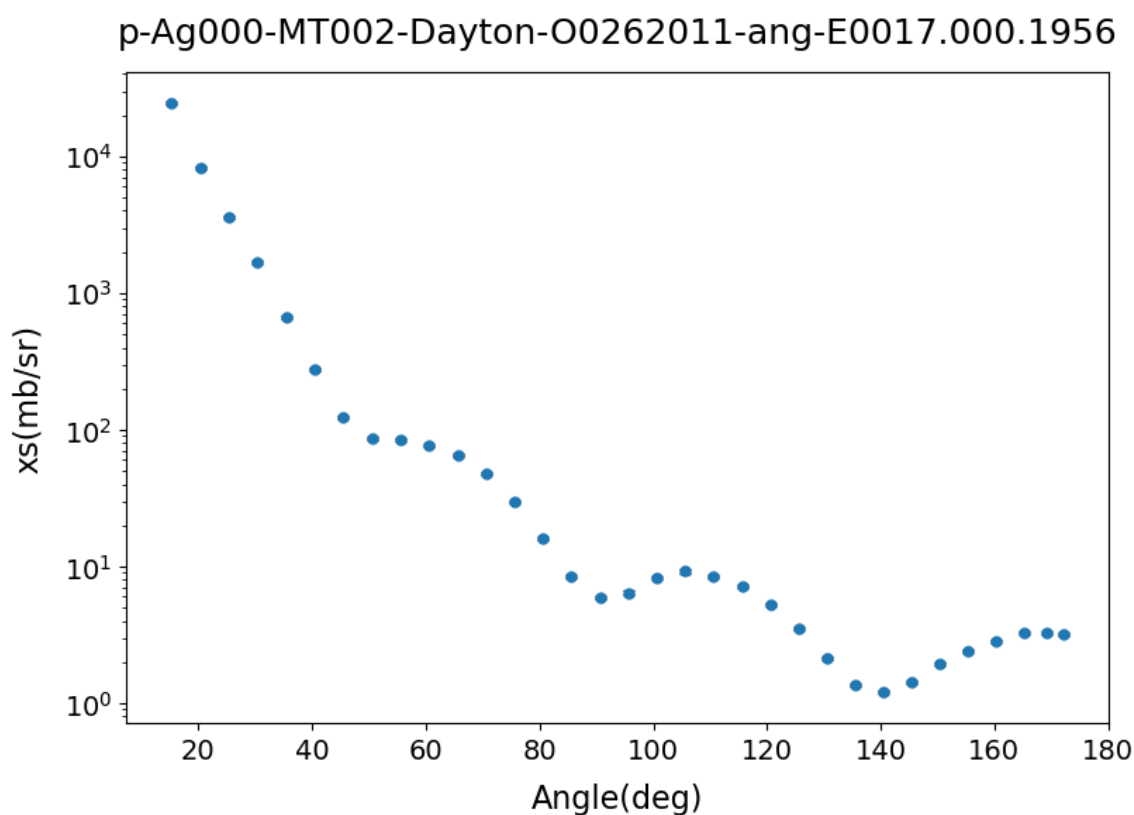


Figure 4.1: Experimental results obtained using the plot method of the `Experiment` class.

```
1 exp_from_bin[2135].plot(ylog=True)
```

Output:

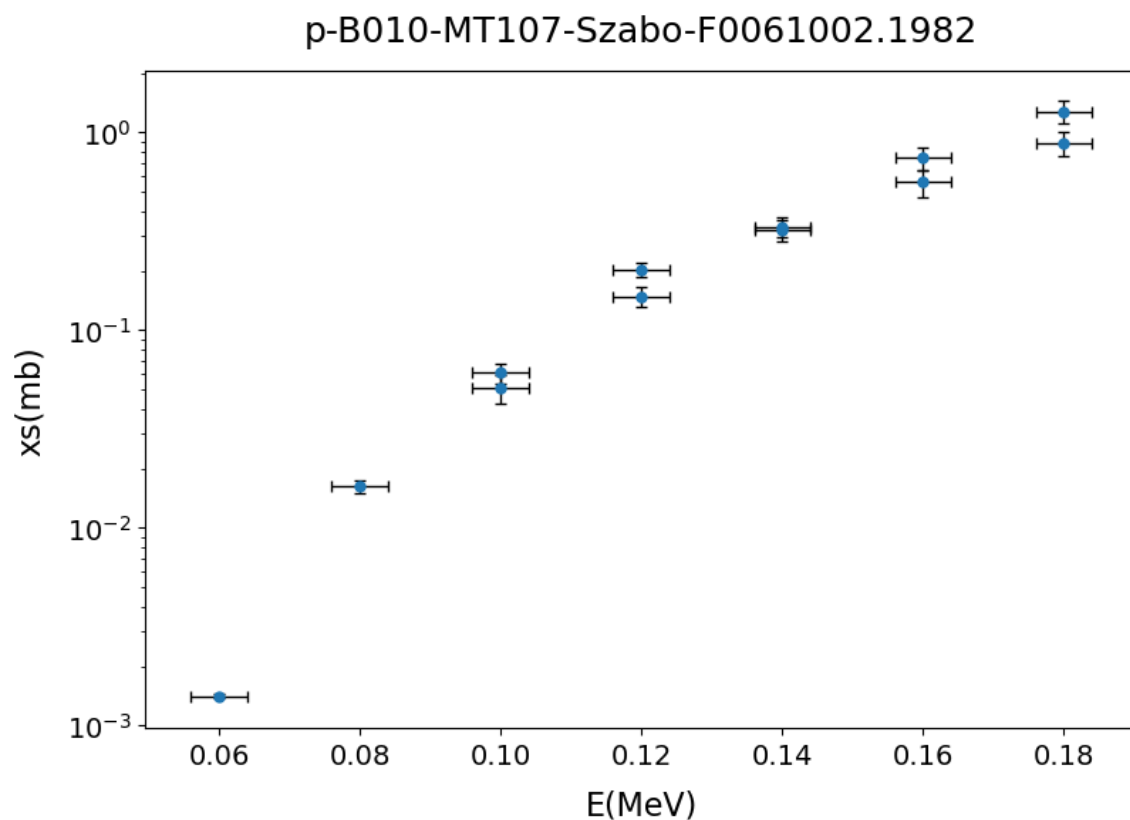


Figure 4.2: Experimental results with uncertainties obtained using the plot method of the Experiment class.

When it comes to analyzing or sorting through the experiments, there are specialized functions to make things easier. For example, the function `get_unique_values()` takes a list of experiments and an attribute name as inputs. It then returns all the unique values that appear for that attribute among the listed experiments (see Annex B.8).

```
1 MT_values = get_unique_values(exp_from_bin, "MT")
```

If you are interested in experiments that meet a specific condition, you can filter them out. The function `filter_experiments()` does this for you. For example, you can filter out experiments where $MT = 42$ from a larger list (see Annex B.7).

```
1 MT42_exp = filter_experiments(exp_from_bin, "MT", 42)
```

Output:

```
1 14 experiments with MT = 42
```

You can also utilize the `plot_experiments()` function (see Annex B.12) to display multiple similar experiments—those with the same X and Y axes—on a single graph. This is helpful

for comparing data, like in an example where 14 experiments with $MT = 42$ are showcased together.

```
1 plot_experiments(MT42_exp)
```

Output:

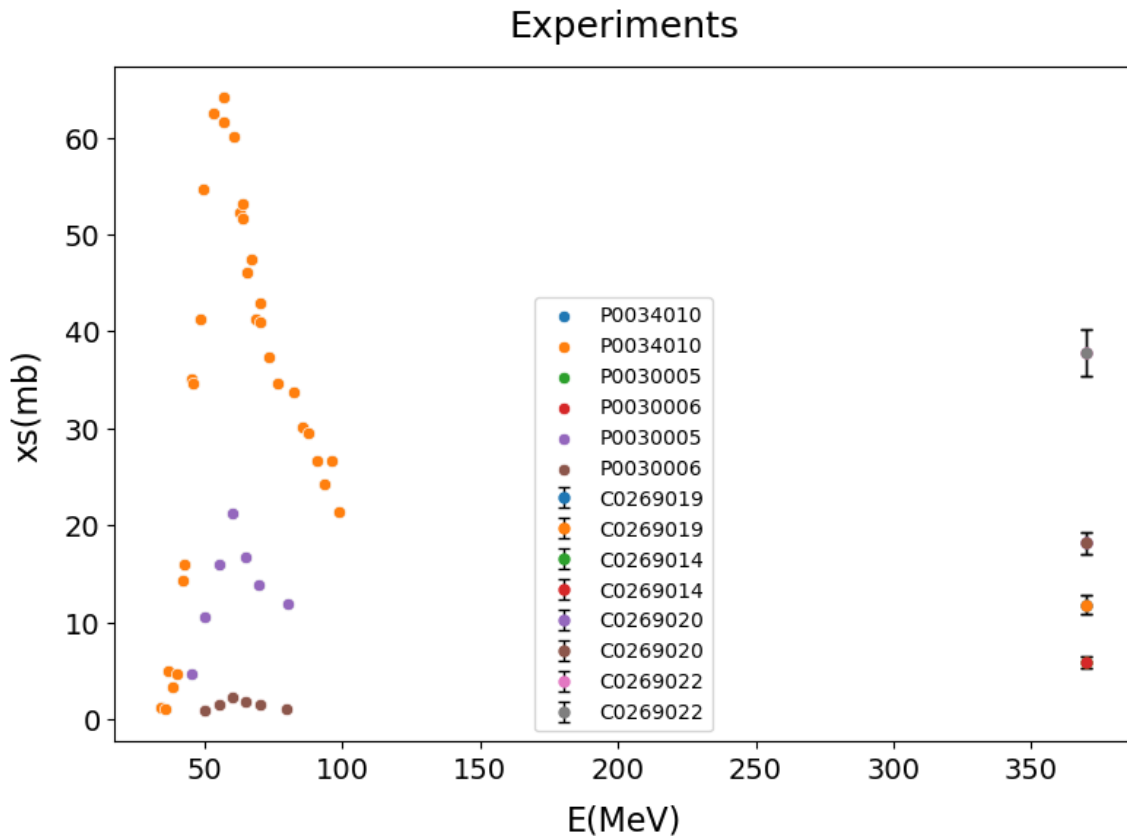


Figure 4.3: Experimental results of all proton-induced experiments with $MT = 42$.

With over 27,000 subentries for proton-induced reactions in EXFOR, not every experiment measures the same things. For machine learning applications, it is useful to sort experiments based on what is actually being measured. Essentially, this involves grouping experiments by their data. Each experiment object has a `.data` attribute that holds a table with four columns: X values, Y values, uncertainties in Y, and uncertainties in X. If there are no uncertainties, these will be marked as `NaN`. To organize the experiments in this manner, the function `classify_experiments_by_data()` is used (see Annex B.11).

```
1 df = classify_experiments_by_data(exp_from_bin)
```

This function creates a CSV file for each of the each group of similar experiments, focusing only on the details that matter for analysis. This means that certain attributes like the title or the author aren't included. Be aware that this process could take a long time since it deals with millions of data points. Once the files are created, you can easily load them into a Python dataframe using the pandas library. For additional details regarding the classification of experimental groups, please refer to Annex D.

When turning this data into CSV files, we use a technique called “one-hot encoding” to turn everything into numbers. This makes it easier to use machine learning tools on the data. Here are some additional things to consider:

1. Some columns may have the same value for every row in each group. In machine learning, these uniform columns are not useful. You can remove them using the `clean_dataframe()` function (see Annex B.9).
2. Machine learning often does not need uncertainty data. You can also remove these columns using the `clean_dataframe()` function.
3. Many machine learning tools can not handle NaN values. The `clean_dataframe()` function can usually remove these, but for Group 4, you will need an extra step. Here, the code divides Group 4 into two sub-groups to ensure there are no NaN values. The information is then saved into two different CSV files.

```
1 df = pd.read_csv("EXFOR_ProtonReactions_Classified_Group_4.csv")
2 df_nan = df[df["final_A"].isna()]
3 df_not_nan = df[df["final_A"].notna()]
4 df_nan.to_csv("EXFOR_ProtonReactions_Classified_Group_4_1.csv",
5               index=False)
6 df_not_nan.to_csv("EXFOR_ProtonReactions_Classified_Group_4_2.csv",
7                   index=False)
```

4. A last column named X4_ID is added to each row in order to identify the EXFOR experiments that corresponds to each value and be able to address it correctly. This column must be removed before applying machine learning techniques as it contains string values that can not be handle by algorithms.

There is an additional function named `classify_experiments()` (see Annex B.10). This function offers a different way to group the experiments, although it is not as well-suited for machine learning applications as the one described earlier. Specifically, it organizes the experiments based on a chosen attribute. For example, using this function with the attribute MT will create separate CSV files, where each file contains experiments that share the same MT value. Unlike the previous function, this one does not apply one-hot encoding and merely converts the information into dataframes. While it might not be the best choice for machine learning, it can still be useful for other types of data analysis.

```
1 classify_experiments(exp_from_bin, MT)
```

Finally, the function `plot_outliers()` (see Annex B.13) is designed to visualize outliers identified through various detection methods. To optimize data visualization, experiments sharing identical attributes are consolidated into a single graph. The legend serves as an intuitive guide for identifying the EXFOR experiments associated with each data point, with outliers being denoted by a red cross. Examples of this will be shown in the following sections along with the results of the outlier detection.

```
1 plot_outliers(outliers_df, exp_from_bin)
```

4.2 Interquartile Range (IQR)

4.2.1 Implementation of IQR Method

The outlier detection using the Interquartile Range (IQR) method is implemented through Python code, which applies the theoretical framework described in Section 3.3. The mathematical formulations for IQR and the lower and upper bounds used for detecting outliers are given by Equations 3.2, 3.3, and 3.4 respectively. This methodology is implemented through two functions `detect_outliers_IQR` and `IQR_method` that can be found in Annex C.

The `detect_outliers_IQR` function takes a grouped dataframe (`grouped_df`), an `Experiment` object (`experiment`), a limit factor (`limit`), and a Boolean flag to consider uncertainties (`uncertainties`).

The function calculates the first ($Q1$) and third ($Q3$) quartiles for the specified columns in the grouped dataframe:

```
1 Q1 = grouped_df[[first_column, second_column]].quantile(0.25)
2 Q3 = grouped_df[[first_column, second_column]].quantile(0.75)
```

Following this, the IQR is computed according to equation 3.2 and the conditions for determining outliers are then applied:

```
1 outlier_condition = ((grouped_df[[first_column, second_column]] < (Q1 -
    limit * IQR)) | (grouped_df[[first_column, second_column]] > (Q3 +
    limit * IQR)))
```

Finally, the function filters and returns all rows that satisfy the outlier condition.

The `IQR_method` function serves as a wrapper for `detect_outliers_IQR`. It first identifies columns by which to group the DataFrame and then applies `detect_outliers_IQR` to each valid group, ensuring the group meets the minimum number of observations (`min_observations`):

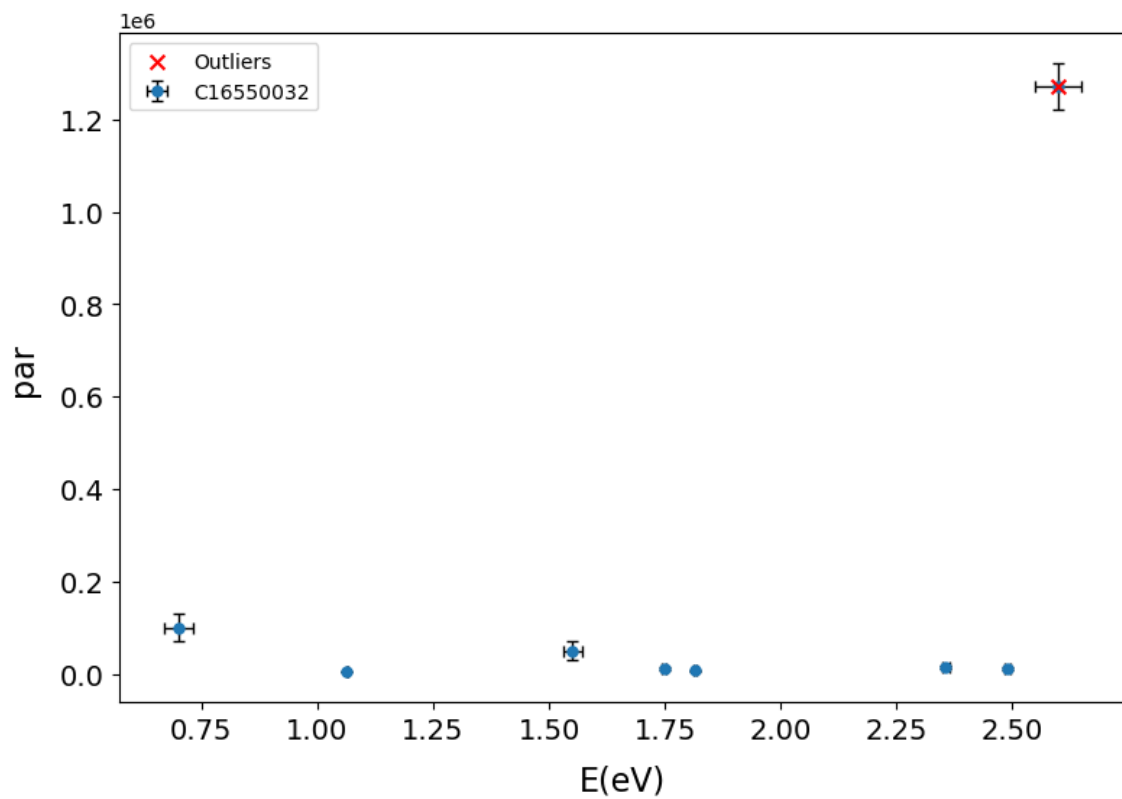
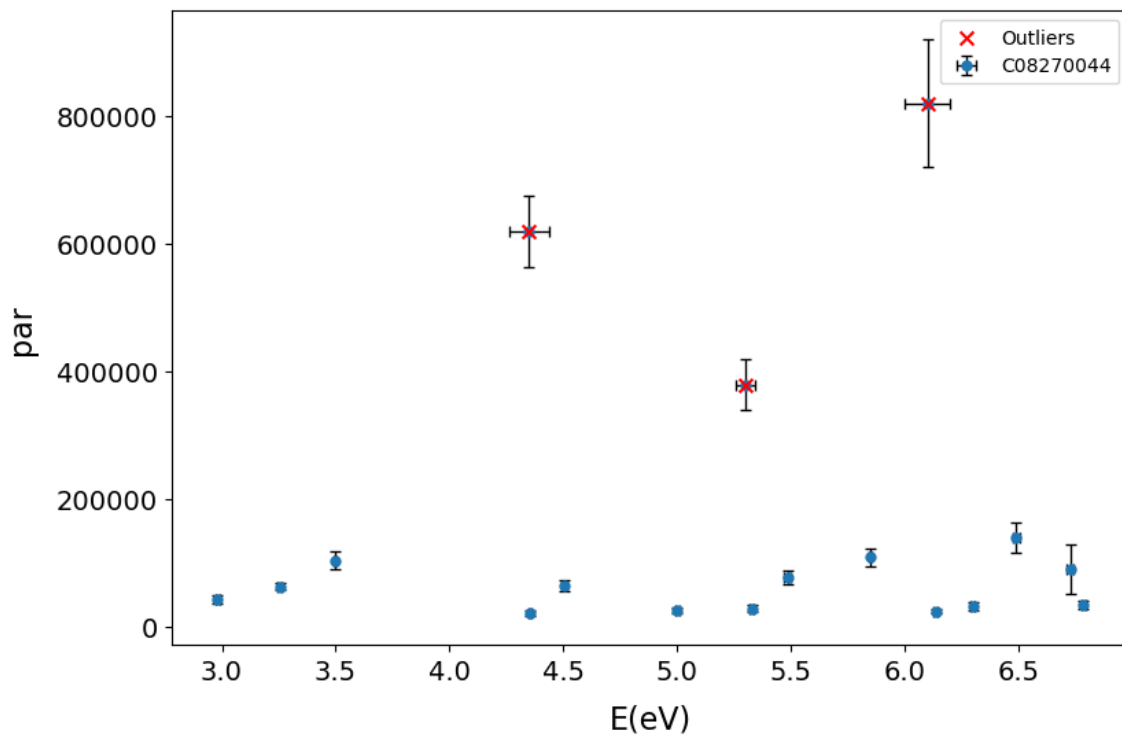
```
1 grouped = df.groupby(groupby_columns).size().reset_index(name="count")
2 valid_groups = grouped[grouped["count"] >= min_observations][
    groupby_columns].to_dict("records")
```

If uncertainties are to be considered (`uncertainties=True`), the function alters the conditions for outlier identification accordingly. The uncertainty values are used to define lower and upper values for each data point:

```
1 lower_values = grouped_df[second_column] - grouped_df[third_column]
2 upper_values = grouped_df[second_column] + grouped_df[third_column]
```

Both ends of the uncertainty range for a data point must lie outside the IQR-defined bounds for it to be considered an outlier.

4.2.2 Results and Graphical Analysis

**Figure 4.4:** Outlier Detection in Group 6 using IQR Method.**Figure 4.5:** Outlier Detection in Group 6 using IQR Method.

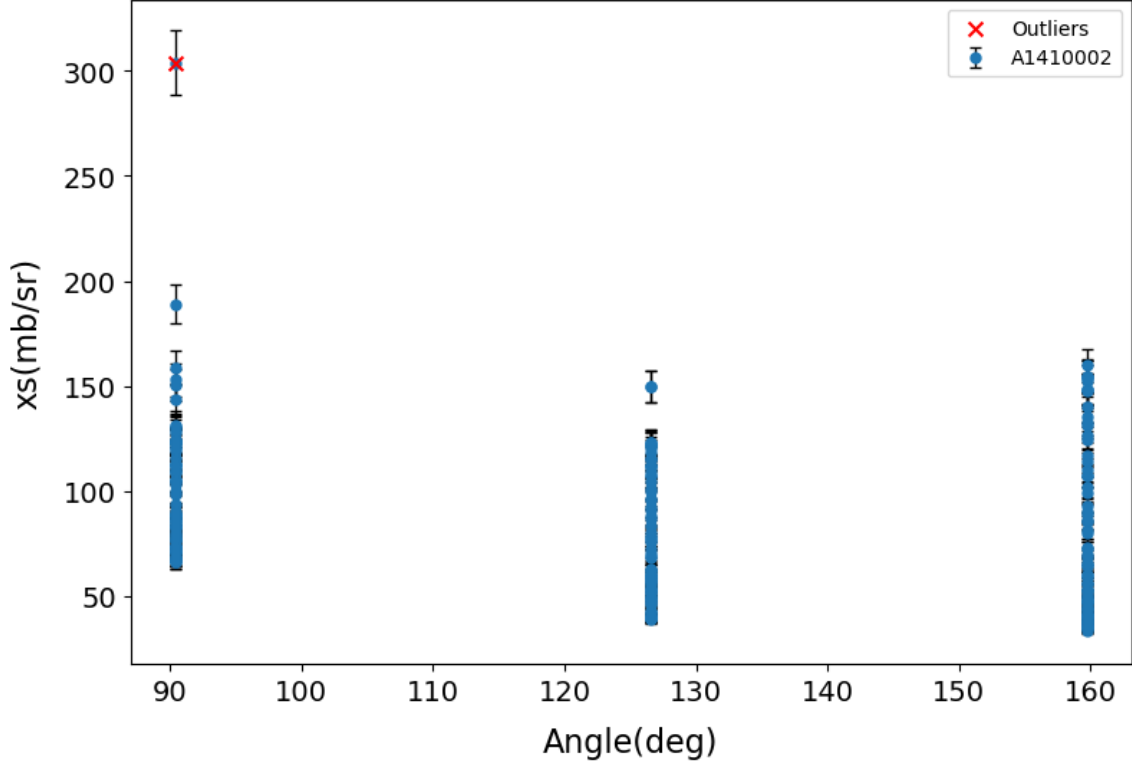


Figure 4.6: Outlier Detection in Group 1 using IQR Method.

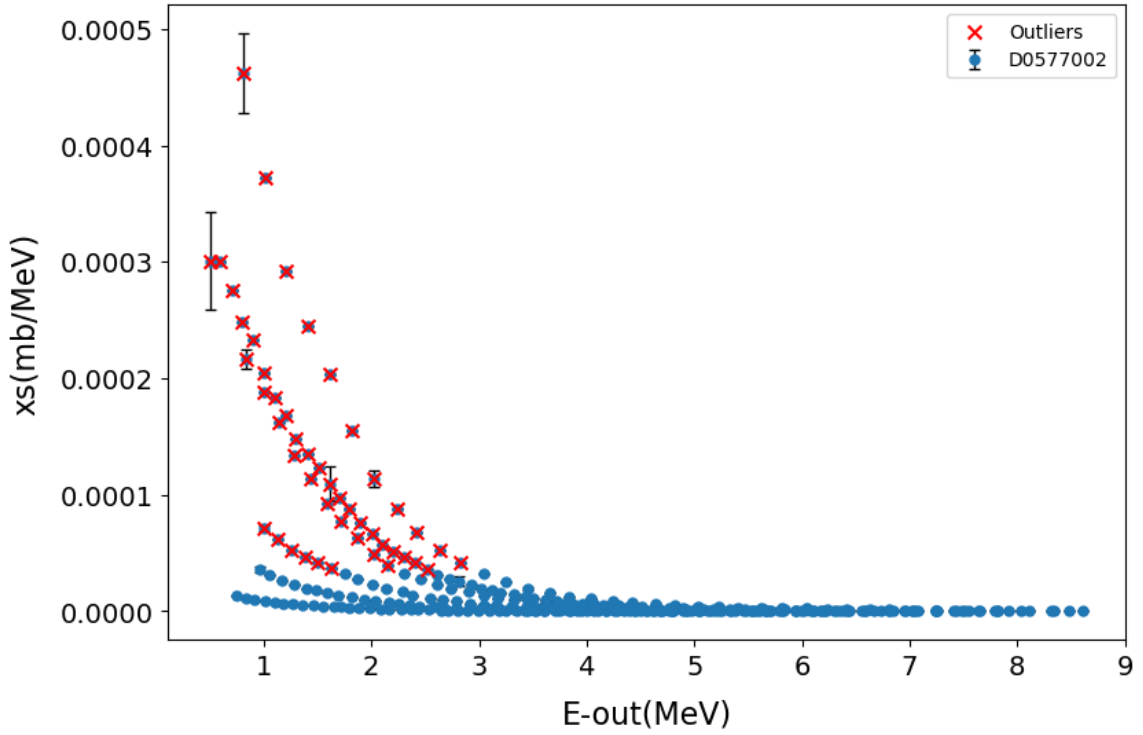


Figure 4.7: Outlier Detection in Group 5 using IQR Method.

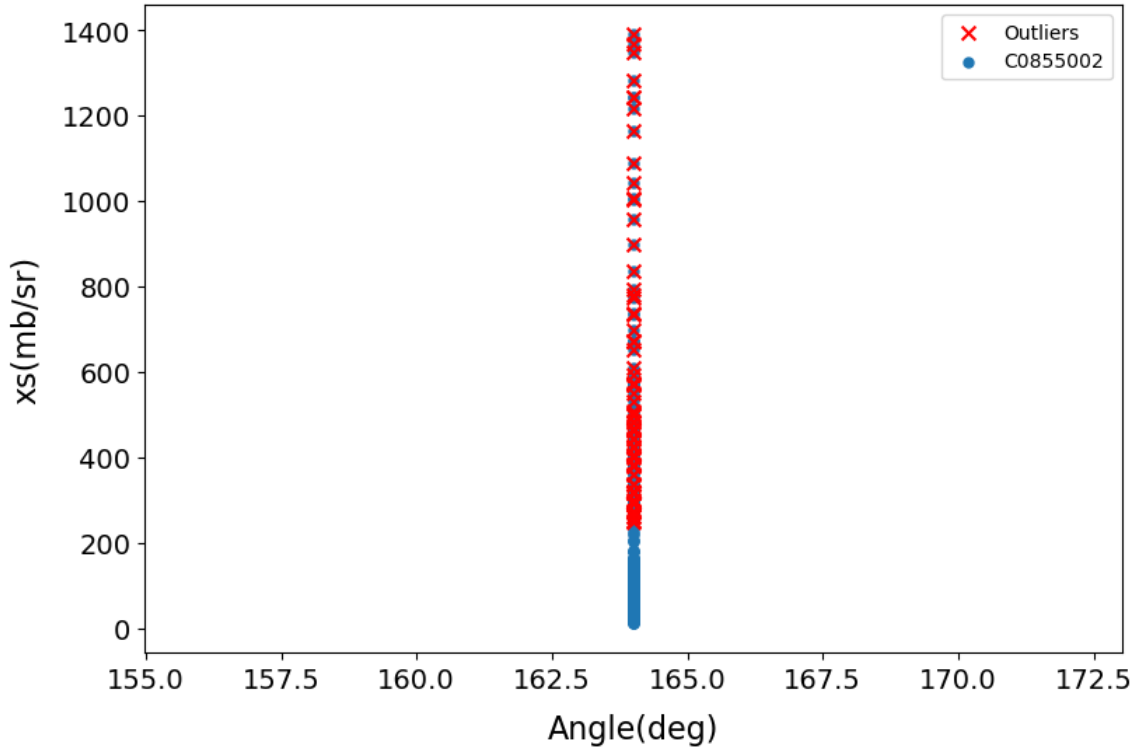


Figure 4.8: Outlier Detection in Group 1 using IQR Method.

While graphs 4.4, 4.5, and 4.6 demonstrate success in identifying potential outliers across different groups, graphs 4.7 and 4.8 highlight the limitations of the IQR method. Specifically, these latter figures exhibit a high rate of false positives. This is largely attributable to the non-uniform distributions present within the data sets, causing the IQR method to incorrectly classify a significant number of inliers as outliers.

4.2.3 Conclusion on IQR Method

The IQR method is fundamentally a statistical technique that allows for the consideration of uncertainties in its calculations. This feature can be advantageous for reducing false positives. However, our study reveals limitations in the applicability of the IQR method to our specific data sets. The nature of the data, which follows varying types of distributions, is not particularly well-suited for IQR, a method that performs better on uniformly distributed data.

In particular, any peak in the data, resulting from its underlying distribution, could be mistakenly identified as an outlier. Consequently, this makes the IQR method ill-suited for our purposes, as it cannot reliably distinguish between genuine outliers and data points that are naturally part of a non-uniform distribution.

The visualizations and analyses thus indicate that while the IQR method has its merits, it is not the most appropriate technique for outlier detection in this specific context.

4.3 Isolation Forest

4.3.1 Implementation of Isolation Forest Method

The script for this implementation can be found in the GitHub repository under the name `OutlierDetection_IsolationForest_Method.ipynb`. We utilize the `IsolationForest` class from the `sklearn.ensemble` package. Below is a snippet of the code outlining the algorithm's application to our dataset:

```
1 # Create the model
2 clf = IsolationForest(contamination=0.001) # contamination = % of
      outliers
3 # Train the model
4 clf.fit(df.iloc[:, :-1])
5 # Predict outliers
6 outliers = clf.predict(df.iloc[:, :-1])
```

After training the model, the `predict` method is used to identify the outliers in the dataset. These are then flagged and added to the dataframe:

```
1 # Add the prediction to the dataframe
2 df["is_outlier"] = np.where(outliers == -1, "Yes", "No")
3 # Get the outliers dataframe
4 outliers_df = df[df["is_outlier"] == "Yes"]
5 outliers_df = outliers_df.drop(columns=["is_outlier"])
```

The code first adds a new column to the dataframe, labeling each data point as an outlier or an inlier. Then, the proportion of detected outliers is calculated and printed.

4.3.2 Results and Graphical Analysis

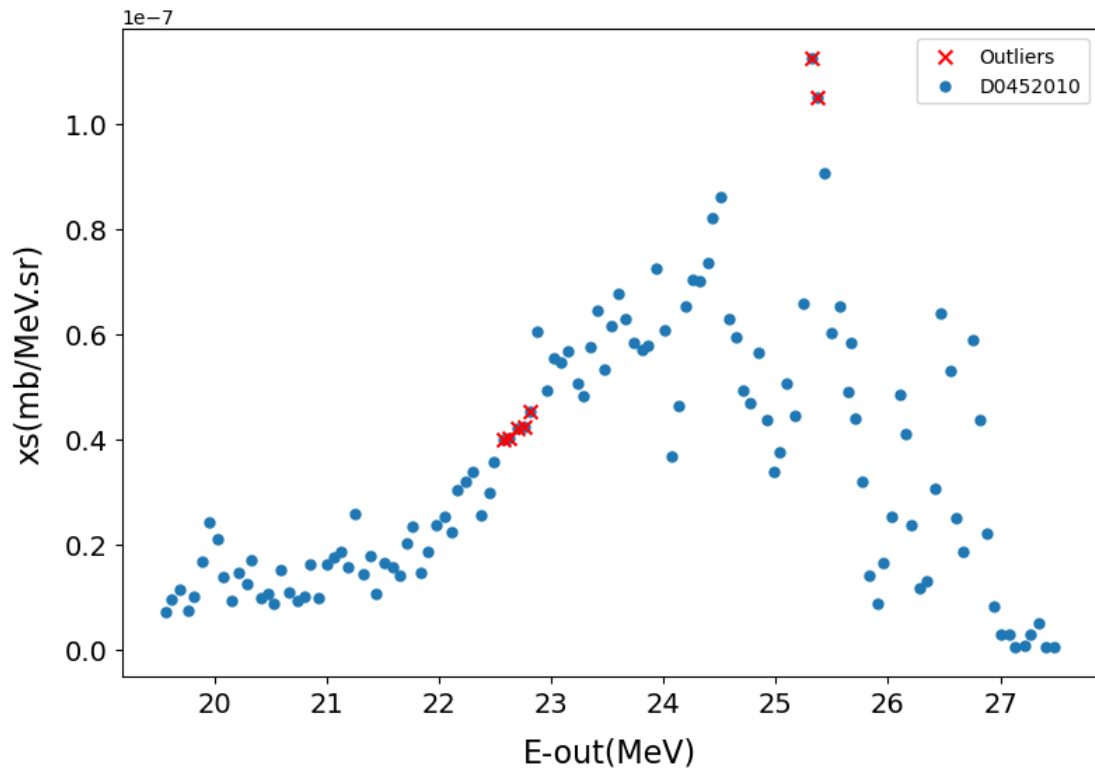


Figure 4.9: Outlier Detection in Group 2 using Isolation Forest Method.

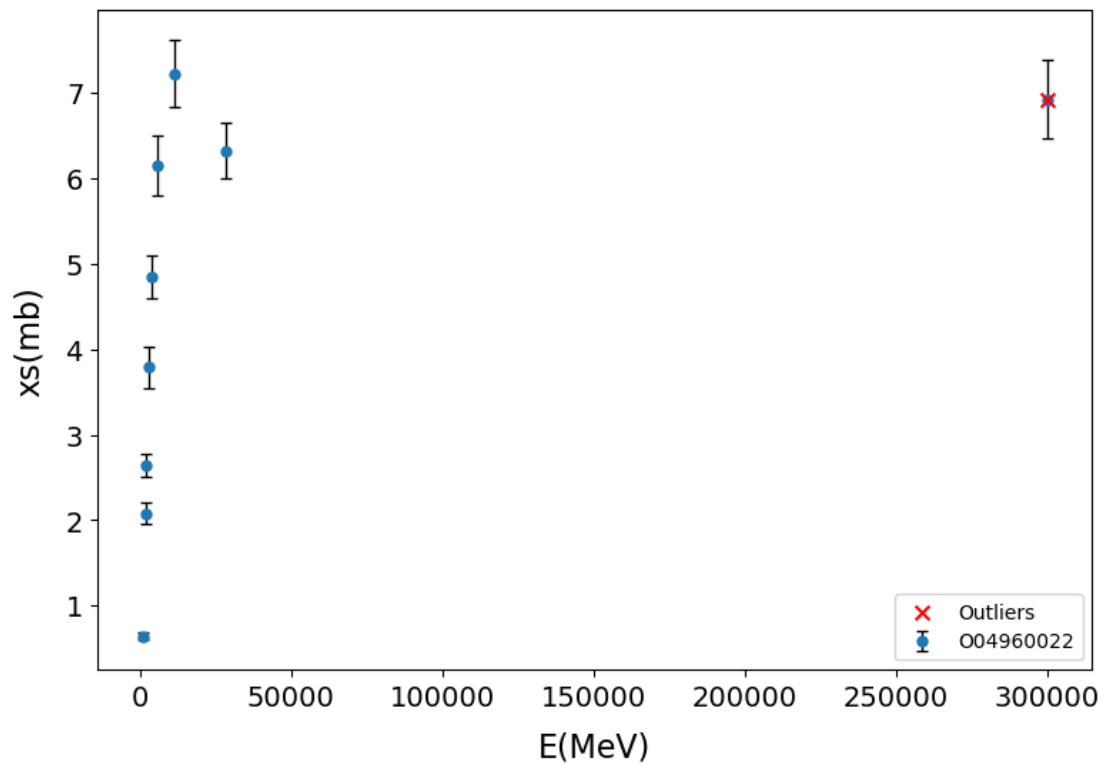


Figure 4.10: Outlier Detection in Group 4.2 using Isolation Forest Method.

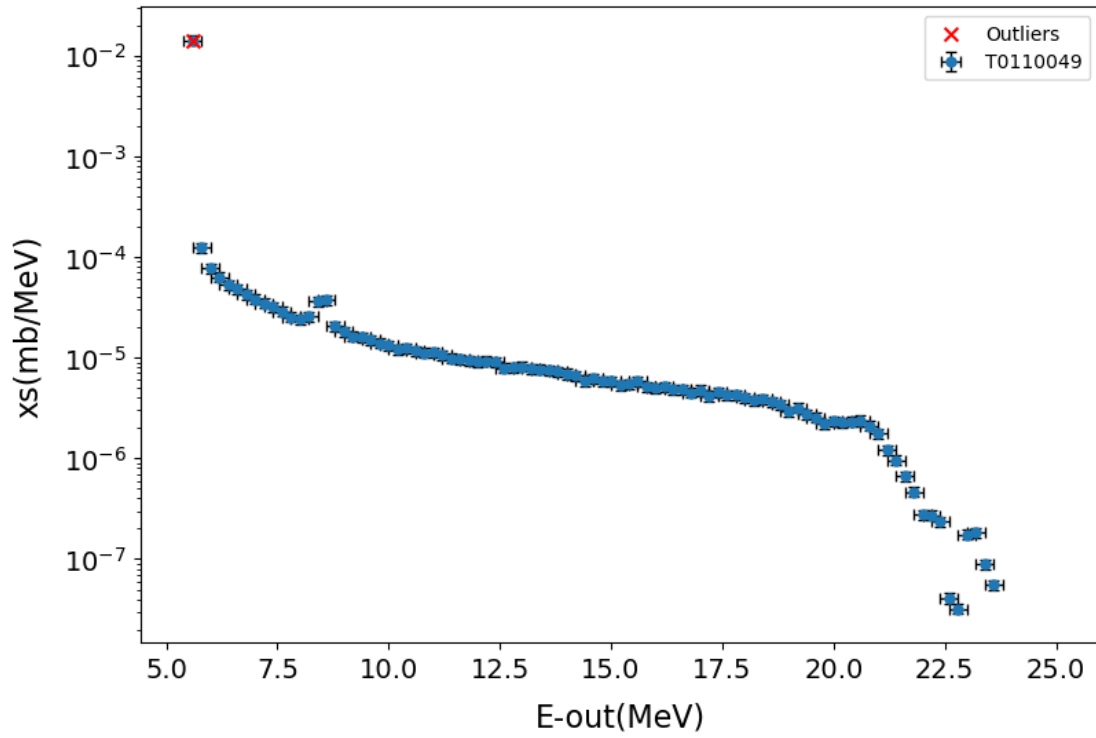


Figure 4.11: Outlier Detection in Group 5 using Isolation Forest Method.

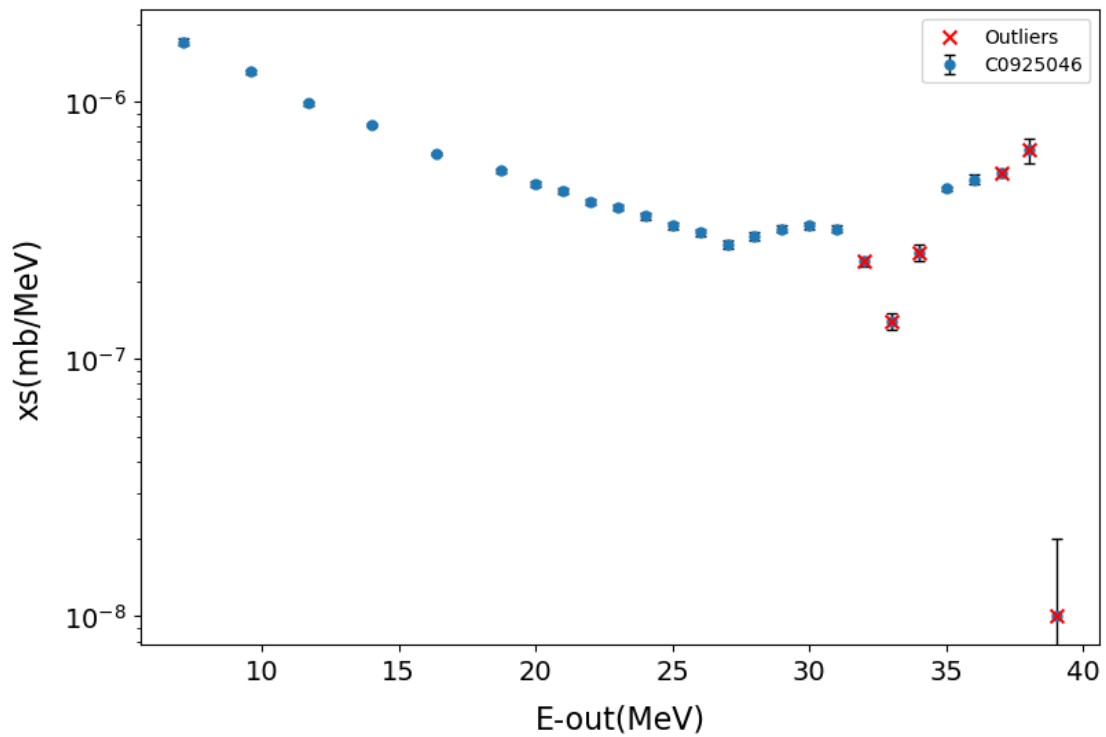


Figure 4.12: Outlier Detection in Group 5 using Isolation Forest Method.

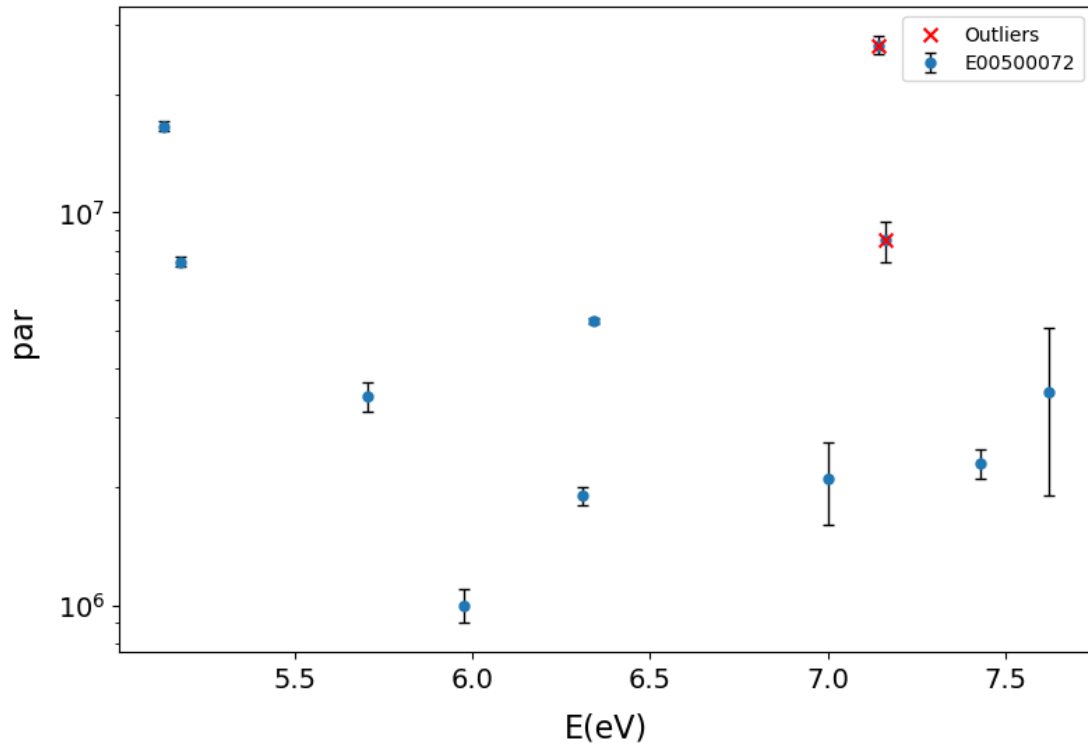


Figure 4.13: Outlier Detection in Group 6 using Isolation Forest Method.

Through graphs 4.9 to 4.13, a range of potential outliers are displayed. In some cases, the detected points are clearly outliers, supporting the efficacy of the Isolation Forest method. However, in other cases, it becomes challenging to ascertain whether the detected points are indeed outliers or simply inliers that appear isolated due to the sparsity of data points in those regions.

4.3.3 Conclusion on Isolation Forest Method

The Isolation Forest method is an effective technique for outlier detection, generally surpassing the performance of statistical methods such as the IQR method. However, the algorithm is highly sensitive to the specific distribution of the dataset. One significant limitation lies in how the method handles sparse regions of the data. When data points are lacking in a specific range—due to insufficient sampling rather than because the data points are genuine outliers—these isolated data points may be incorrectly flagged as outliers. This false positive detection represents a significant drawback in applying the Isolation Forest method in this study’s context.

4.4 DBSCAN

4.4.1 Implementation of DBSCAN Method

The script for the DBSCAN implementation is available in the GitHub repository, named as `OutlierDetection_DBSCAN_Method.ipynb`. We use the `DBSCAN` class from the `sklearn.cluster` package for this purpose. Below is a simplified version of the code to illustrate how the algorithm is applied to our dataset:

1. The data is first scaled using the `StandardScaler` from the `sklearn.preprocessing` package. This standardization is crucial as DBSCAN is sensitive to the scale of the data.

```
1 scaler = StandardScaler()
2 df_scaled = scaler.fit_transform(df_without_ids)
```

2. After scaling, the DBSCAN algorithm is applied to the data.

```
1 clustering = DBSCAN().fit(df_scaled)
```

3. Outliers are identified by finding where the clustering labels equal -1.

```
1 outlier_positions = np.where(clustering.labels_ == -1) #
   Positions of the outliers
```

4. Finally, the data is then descaled to revert it back to its original state.

```
1 df_descaled = pd.DataFrame(scaler.inverse_transform(df_scaled),
   columns=df_without_ids.columns)
```

4.4.2 Results and Graphical Analysis

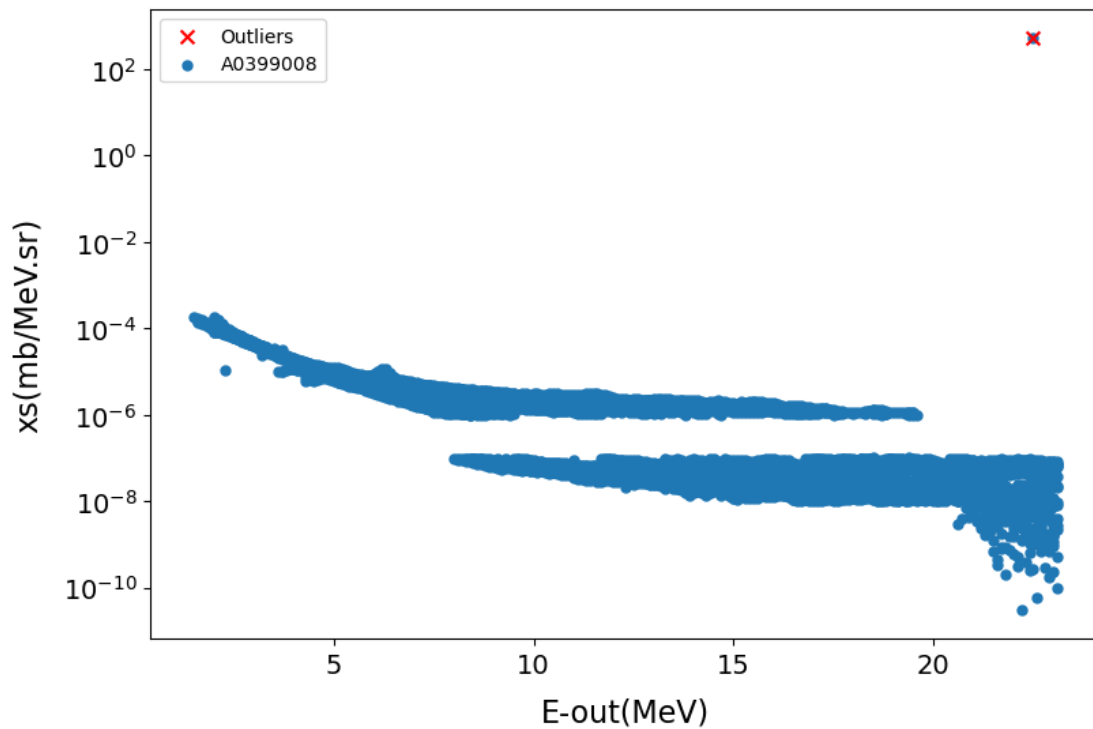


Figure 4.14: Outlier Detection in Group 2 using DBSCAN Method.

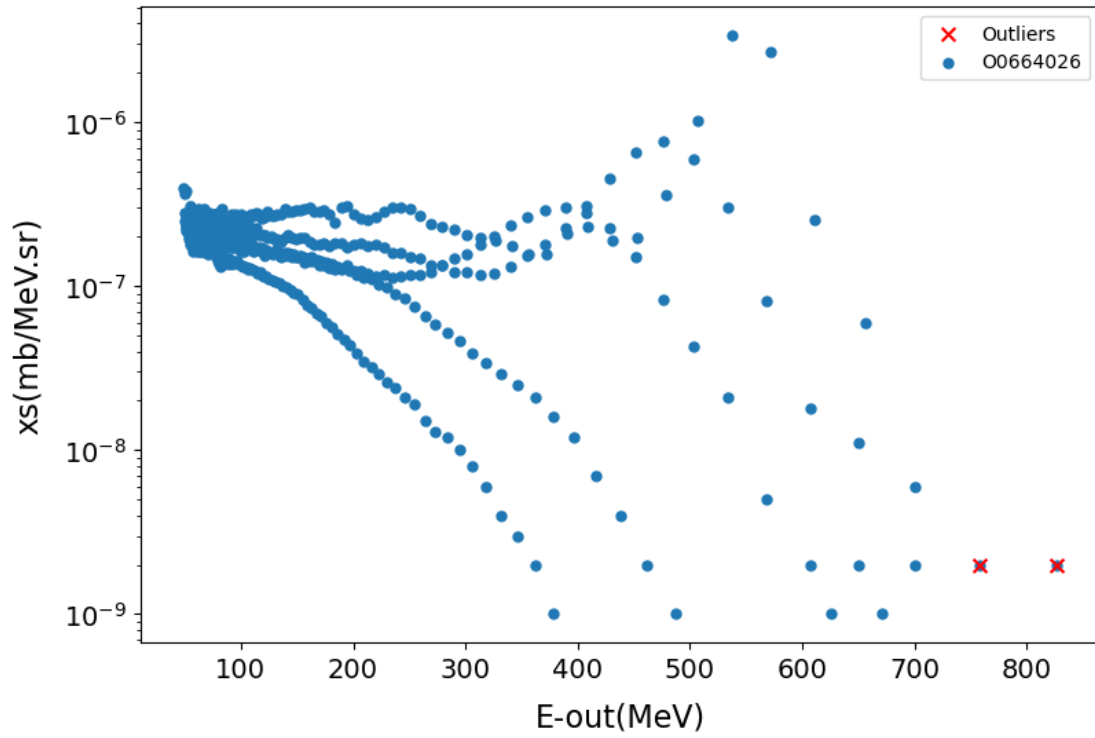


Figure 4.15: Outlier Detection in Group 2 using DBSCAN Method.

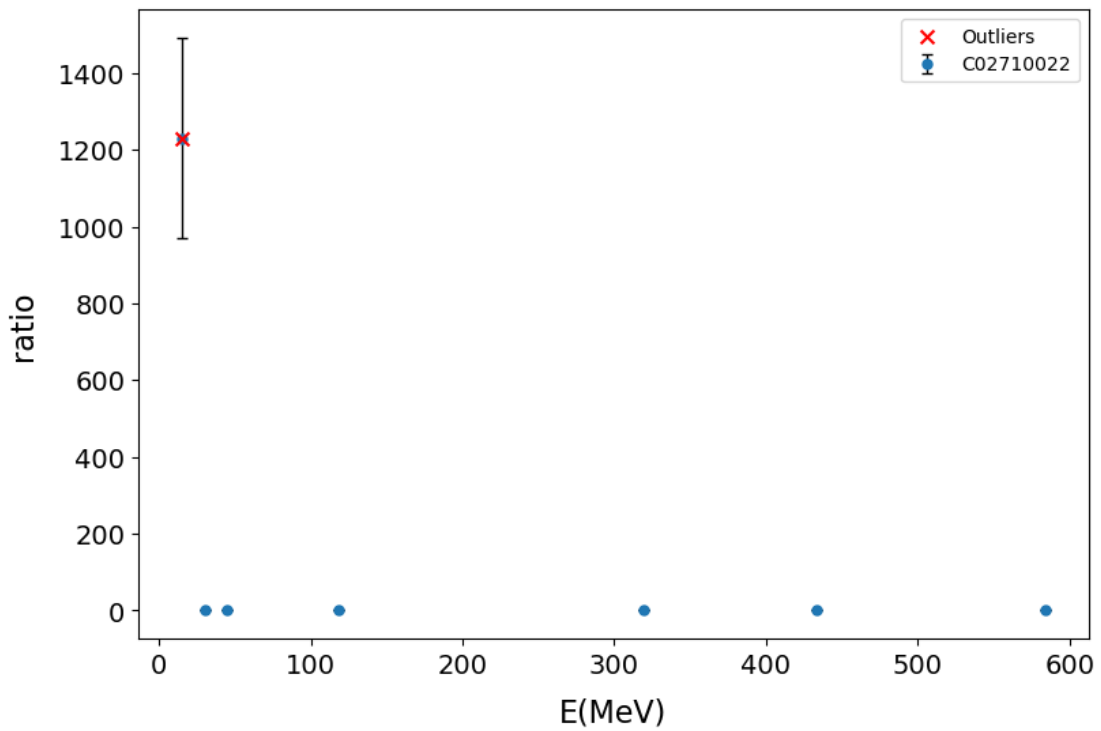


Figure 4.16: Outlier Detection in Group 3 using DBSCAN Method.

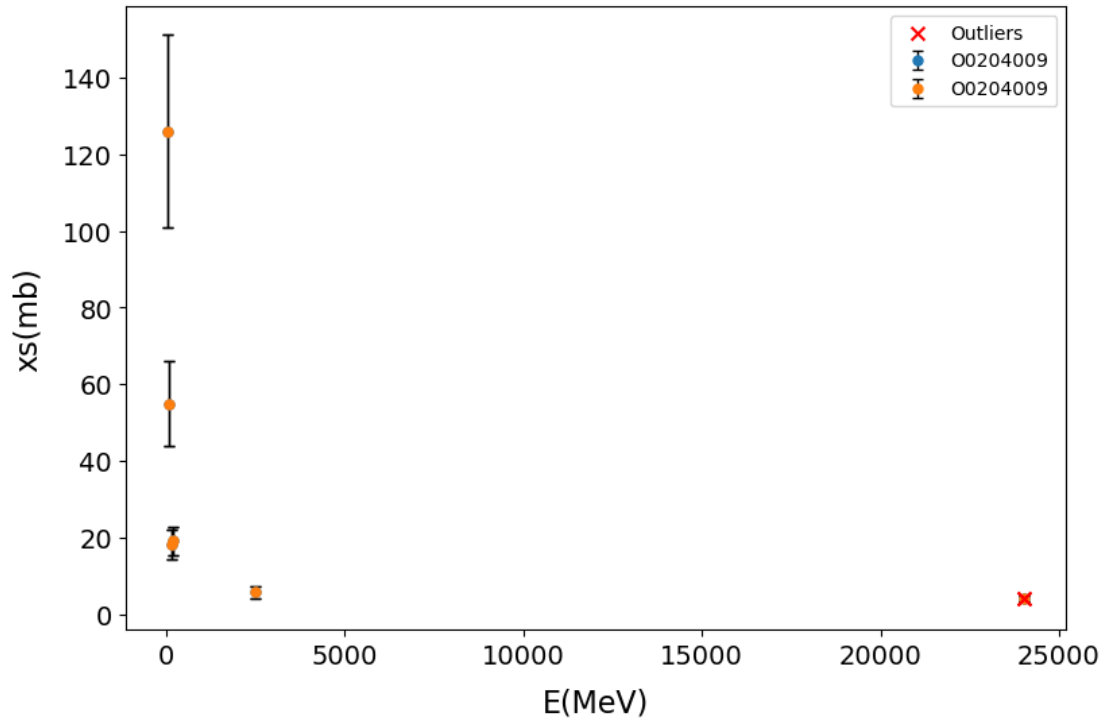


Figure 4.17: Outlier Detection in Group 4.2 using DBSCAN Method.

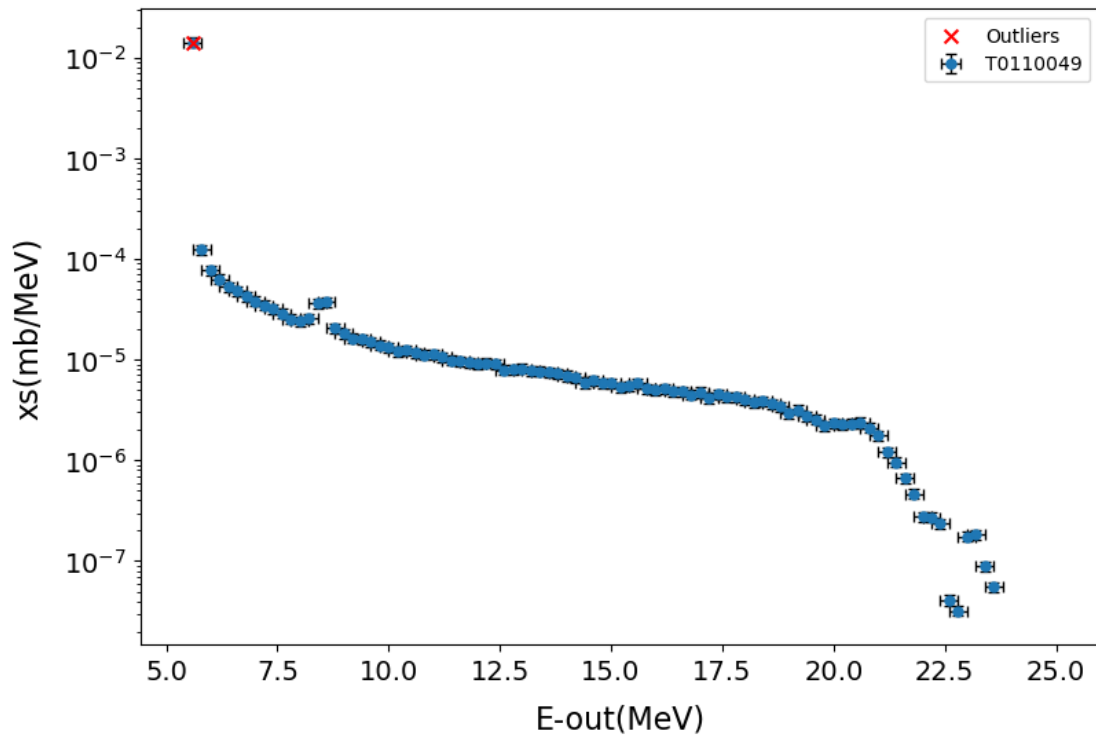


Figure 4.18: Outlier Detection in Group 5 using DBSCAN Method.

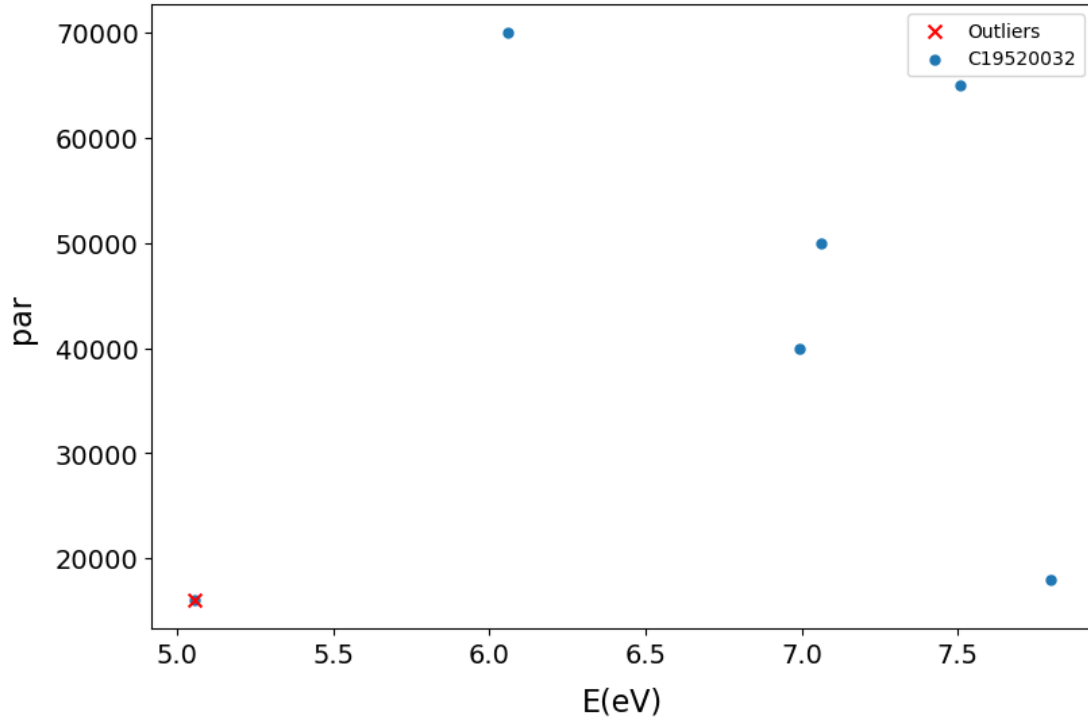


Figure 4.19: Outlier Detection in Group 6 using DBSCAN Method.

Figures 4.14 to 4.19 illustrate the outcomes of applying the DBSCAN algorithm to various data groups. Generally speaking, the results appear promising, although some instances of false positives can be observed, along with some outliers that are debatable. Nevertheless, a portion of the detected points clearly display characteristics consistent with being potential outliers.

4.4.3 Conclusion on DBSCAN Method

DBSCAN is another machine learning-based method for outlier detection that offers a unique approach. Unlike methods that are primarily statistical, DBSCAN considers the density of points in a data space to identify outliers. This can be advantageous as it might better handle data with varying densities.

For the purpose of this study, the default hyperparameters are used in the DBSCAN implementation. While the method is known to be highly sensitive to its hyperparameters, such as the radius of the neighborhood and the minimum number of points needed to form a dense region, these parameters are kept at their default settings as an initial exploration. Detailed explanations for these hyperparameters can be found in Section 3.5. Future work may involve hyperparameter tuning to optimize the method for the specific characteristics of the dataset. However, even with default settings, DBSCAN provides a distinct and insightful approach to outlier detection compared to purely statistical methods.

Moreover, it is important to note that DBSCAN can be computationally intensive, particularly when dealing with large datasets. In this study, for example, the data for Group 1 had to be subdivided into smaller datasets to make the algorithm computationally feasible. This is a limitation to consider, especially when aiming to apply the method on a larger scale.

4.5 Local Outlier Factor (LOF)

4.5.1 Implementation of Local Outlier Factor (LOF) Method

The script for the LOF implementation is accessible in the GitHub repository, named as `OutlierDetection_LOF`. We make use of the `LocalOutlierFactor` class available in the `sklearn.neighbors` package. Below is a segment of code demonstrating how the algorithm is applied:

```
1 # Scaling the Data
2 scaler = StandardScaler()
3 df_scaled = scaler.fit_transform(df_without_id)
4
5 # Applying the LOF algorithm
6 lof = LocalOutlierFactor(n_neighbors=20, contamination=0.01)
7 outliers = lof.fit_predict(df_scaled)
8
9 # Identifying Outliers
10 is_outlier = outliers == -1
11
12 # Reverting the Scaling
13 df_descaled = pd.DataFrame(scaler.inverse_transform(df_scaled), columns=
    df_without_id.columns)
14
15 # Adding the IDs and extracting the outliers
16 df_descaled["X4_ID"] = X4_ID
17 df_descaled["is_outlier"] = is_outlier
```

In the implementation of the LOF (Local Outlier Factor) method, the first step involves standardizing the dataset using the `StandardScaler` class from the `sklearn.preprocessing` package. This ensures that each feature has a mean of zero and a standard deviation of one. After scaling, the outlier detection is carried out by utilizing the `LocalOutlierFactor` class. The algorithm is parameterized with `n_neighbors=20`, denoting the number of neighbors to consider for each data point, and a `contamination` factor, which is suitably selected to indicate the expected proportion of outliers. The `fit_predict` method is subsequently invoked to classify each data point as an outlier or inlier, with outliers receiving a label of `-1`. The data is then reverted to its original scale for interpretability. Lastly, identification labels (`X4_ID`) and outlier statuses (`is_outlier`) are appended to the dataframe for visual verification.

4.5.2 Results and Graphical Analysis

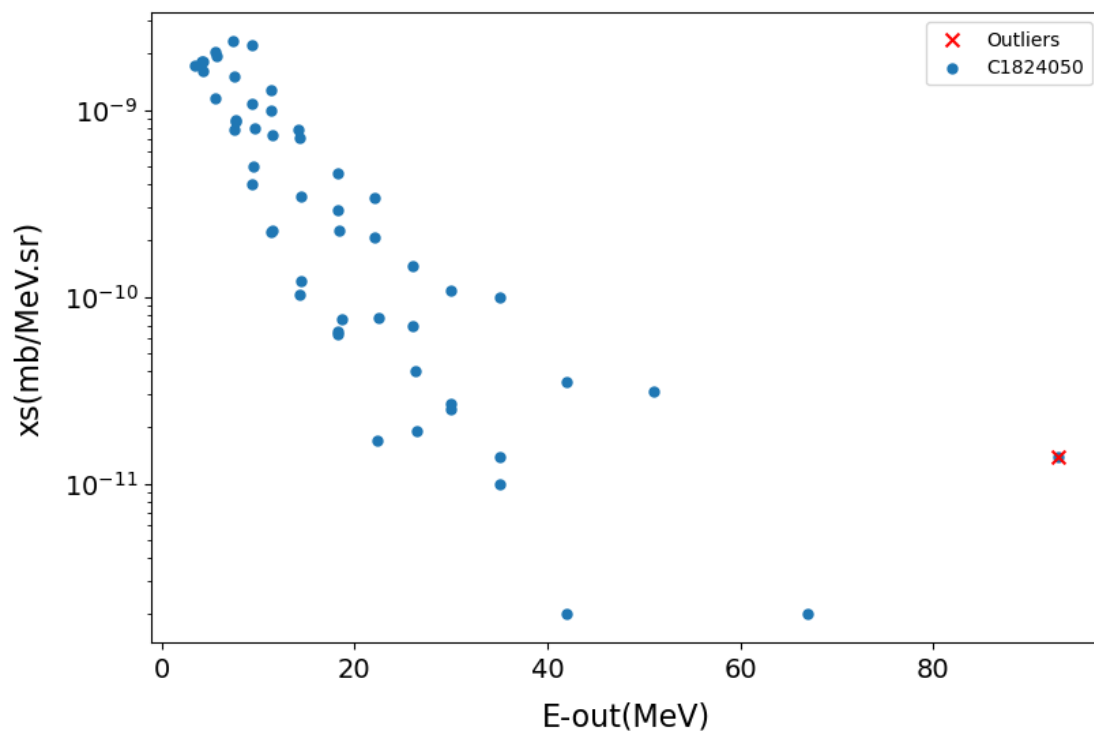


Figure 4.20: Outlier Detection in Group 2 using LOF Method.

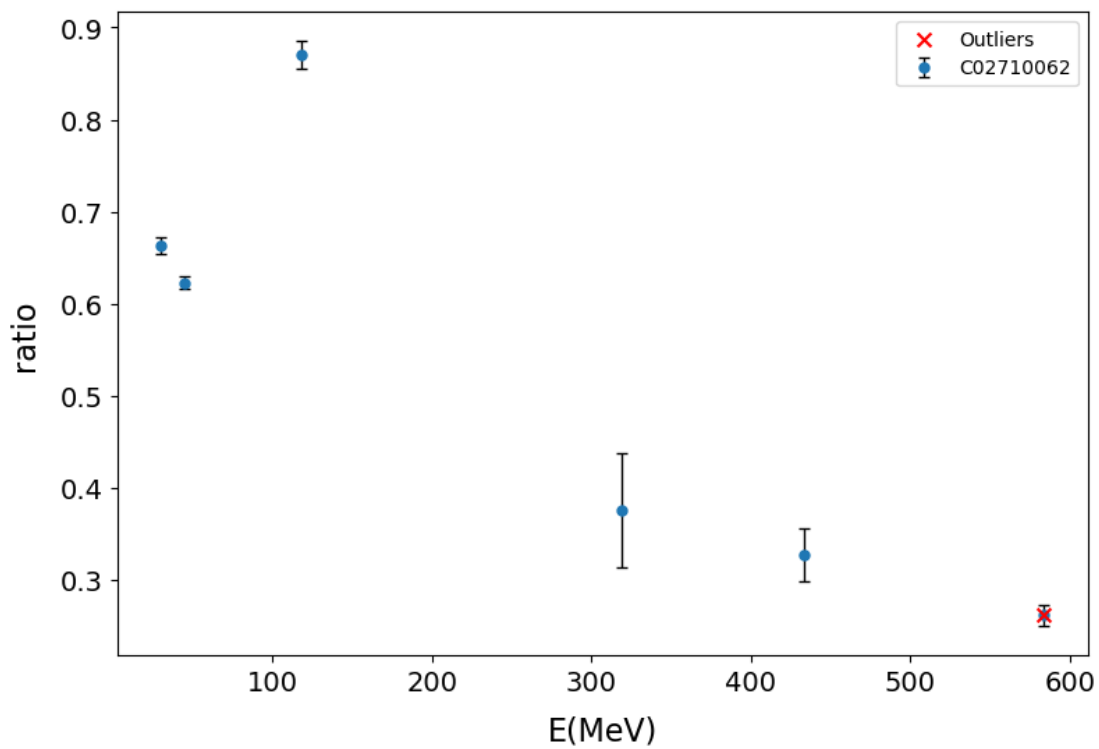


Figure 4.21: Outlier Detection in Group 3 using LOF Method.

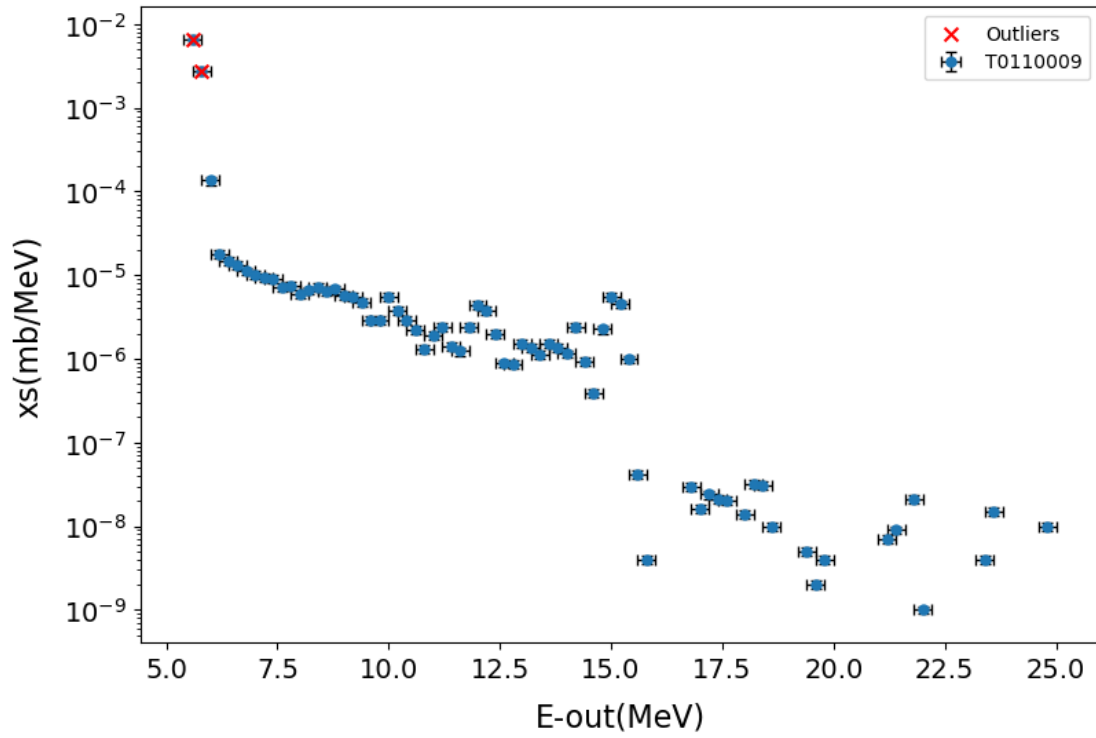


Figure 4.22: Outlier Detection in Group 5 using LOF Method.

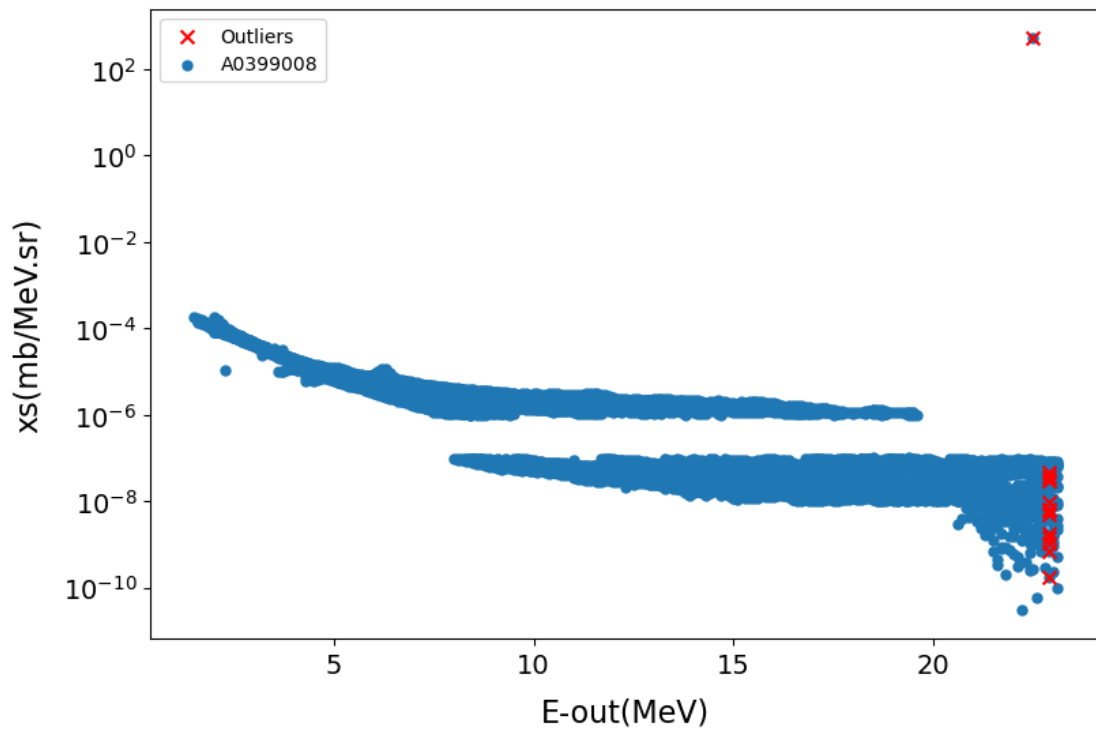


Figure 4.23: Outlier Detection in Group 2 using LOF Method.

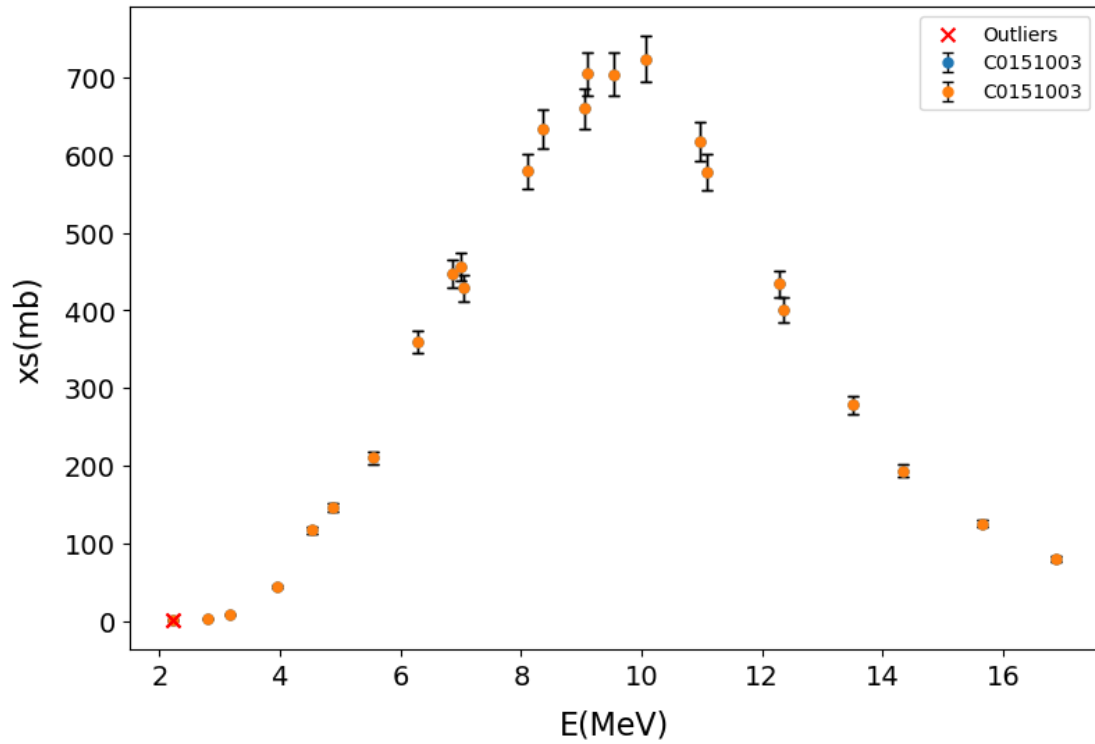


Figure 4.24: Outlier Detection in Group 4.2 using LOF Method.

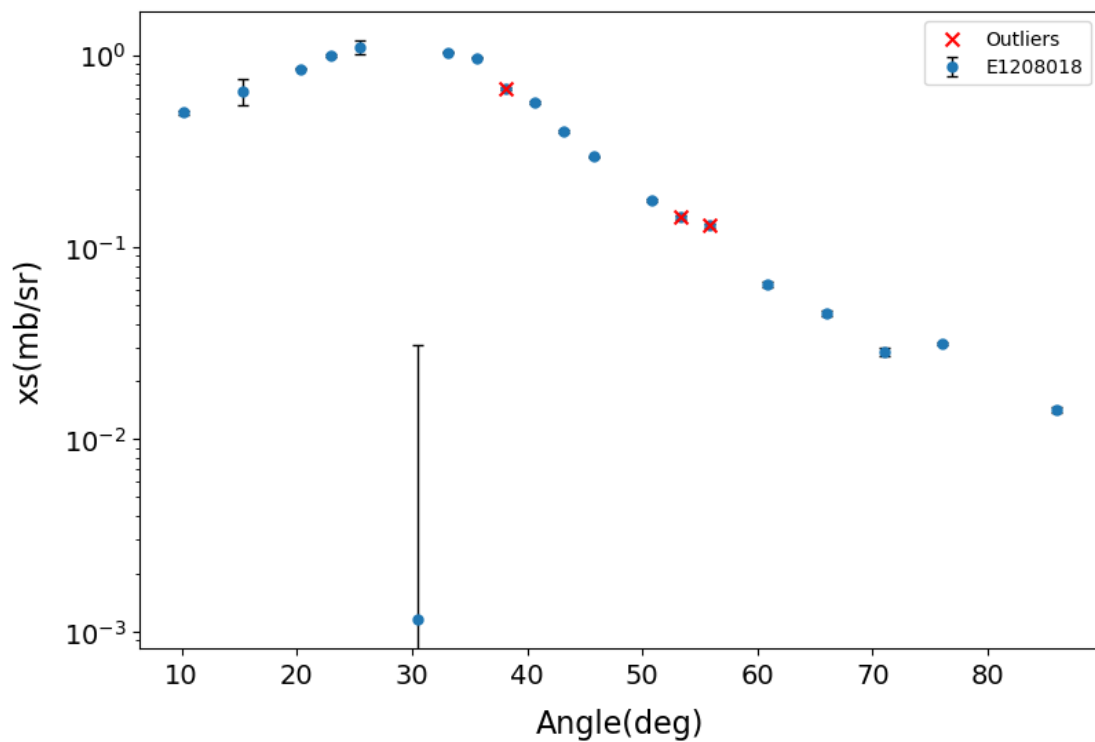


Figure 4.25: Outlier Detection in Group 1 using LOF Method.

The performance of the LOF method exhibits notable variability across different scenarios. In certain instances, such as demonstrated in Figures 4.20 to 4.22, the algorithm adeptly pinpoints data points that could be considered legitimate outliers. However, there are instances where

the detection is less reliable. For example, Figure 4.23 highlights a case where the algorithm successfully identifies a potential outlier but also erroneously flags other points that are more likely to be inliers. In more extreme situations, such as depicted in Figures 4.24 and 4.25, the algorithm falters entirely, failing to correctly identify obvious outliers within the dataset.

4.5.3 Conclusion on LOF Method

The Local Outlier Factor (LOF) method is lauded for its adaptability in handling various types of data distributions. However, it also grapples with challenges common to many outlier detection algorithms. One such issue is its sensitivity to the specific characteristics of the dataset. This becomes particularly evident when the dataset lacks adequate data points in certain ranges, thereby increasing the likelihood of identifying false positives.

Moreover, the necessity for fine-tuning hyperparameters cannot be understated. The effectiveness of LOF in outlier detection is significantly influenced by the choice of hyperparameters, such as the number of neighbors and the contamination factor. Unfortunately, these settings may require different configurations for each group of data, making the process somewhat time-consuming to arrive at satisfactory results.

Despite these challenges, LOF exhibits strengths where other methods may falter. It is particularly robust in avoiding false positives in datasets that feature sharp peaks or clusters. Many algorithms struggle in these areas and tend to misclassify inliers as outliers, whereas LOF demonstrates superior performance in such conditions.

4.6 Support Vector Machines (SVM)

4.6.1 Implementation of SVM Method

The script for the SVM-based outlier detection is available in the GitHub repository, named as `OutlierDetection_SVM_Method.ipynb`. For this approach, we leverage the `OneClassSVM` class from the `sklearn.svm` package. The corresponding segment of the code detailing the algorithm's application is presented below:

```
1 # Scaling the Data
2 scaler = StandardScaler()
3 scaled_df = scaler.fit_transform(df_without_id)
4
5 # Applying OneClassSVM
6 ocsvm = OneClassSVM(kernel="rbf", nu=0.01)
7 ocsvm.fit(scaled_df)
8
9 # Identifying Outliers
10 pred = ocsvm.predict(scaled_df)
11 outliers = (pred == -1)
12
13 # Reverting the Scaling
14 descaled_df = scaler.inverse_transform(scaled_df)
15
16 # Adding the IDs and extracting the outliers
17 result_df = pd.DataFrame(descaled_df, columns=df_without_id.columns)
18 result_df["X4_ID"] = x4_id_column
```

```

19 result_df["outliers"] = outliers
20 outliers_df = result_df[result_df["outliers"] == 1].iloc[:, :-1]
21 print("Percentage of outliers: {:.2f}%".format(len(outliers_df)/len(df)
22       *100))
outliers_df

```

In the SVM-based outlier detection methodology, the initial phase entails standardizing the dataset using the `StandardScaler` class from the `sklearn.preprocessing` package. By doing so, every feature is normalized to have zero mean and unit variance. The `OneClassSVM` class is then applied, with the kernel set to `rbf` (Radial Basis Function). The `nu` parameter is set to its default value, which serves as both an upper bound on the fraction of margin errors and a lower bound on the fraction of support vectors. Once the SVM is trained, predictions are made on the scaled dataset to identify outliers, which are labeled as `-1`. After identifying outliers, the dataset is transformed back to its original scale. Lastly, ID labels (`X4_ID`) and outlier indicators are appended to the dataframe, which is then filtered to extract the detected outliers for further analysis.

4.6.2 Results and Graphical Analysis

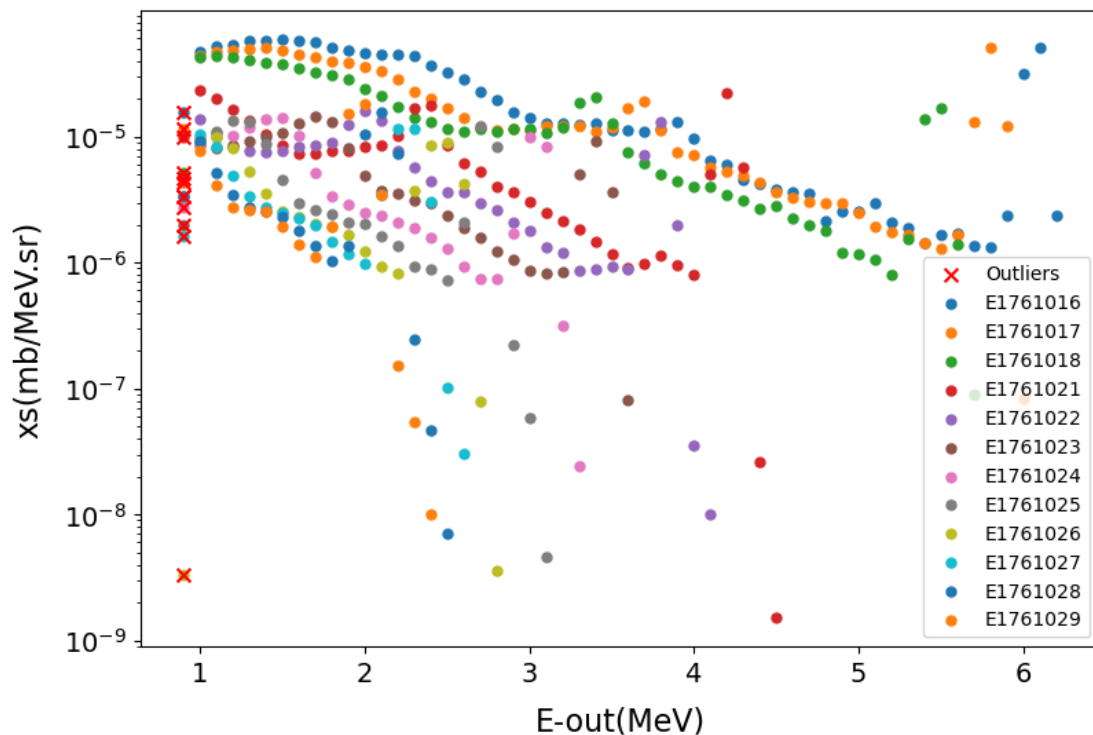


Figure 4.26: Outlier Detection in Group 2 using SVM Method.

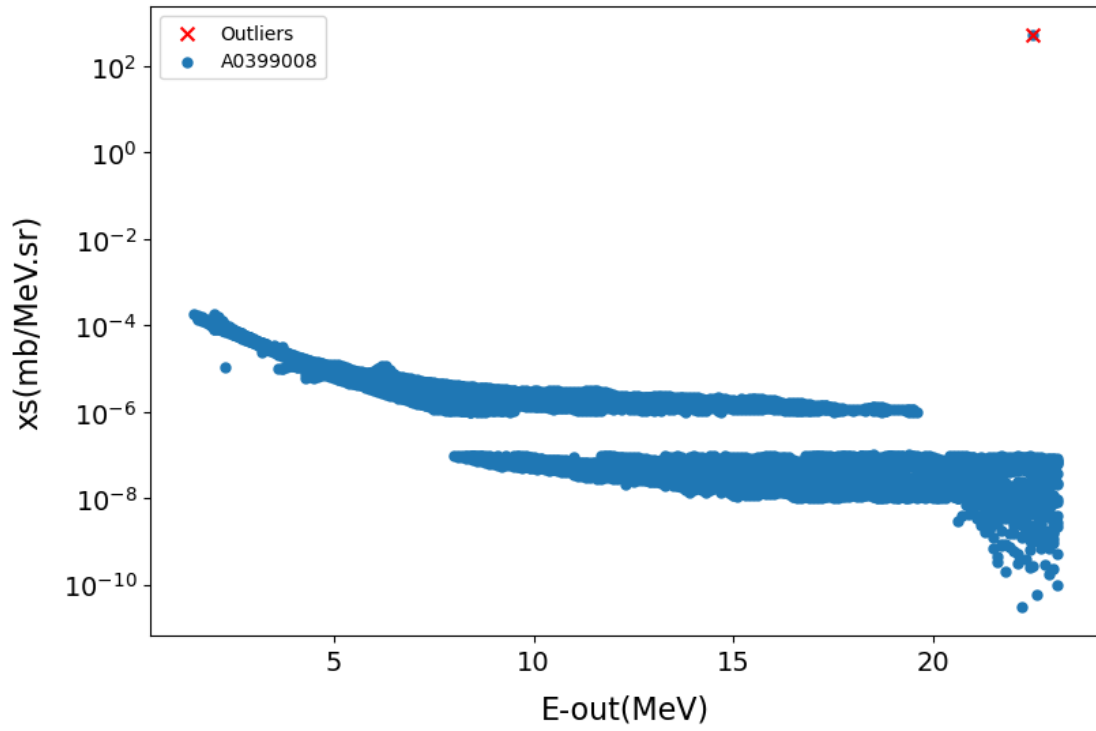


Figure 4.27: Outlier Detection in Group 2 using SVM Method.

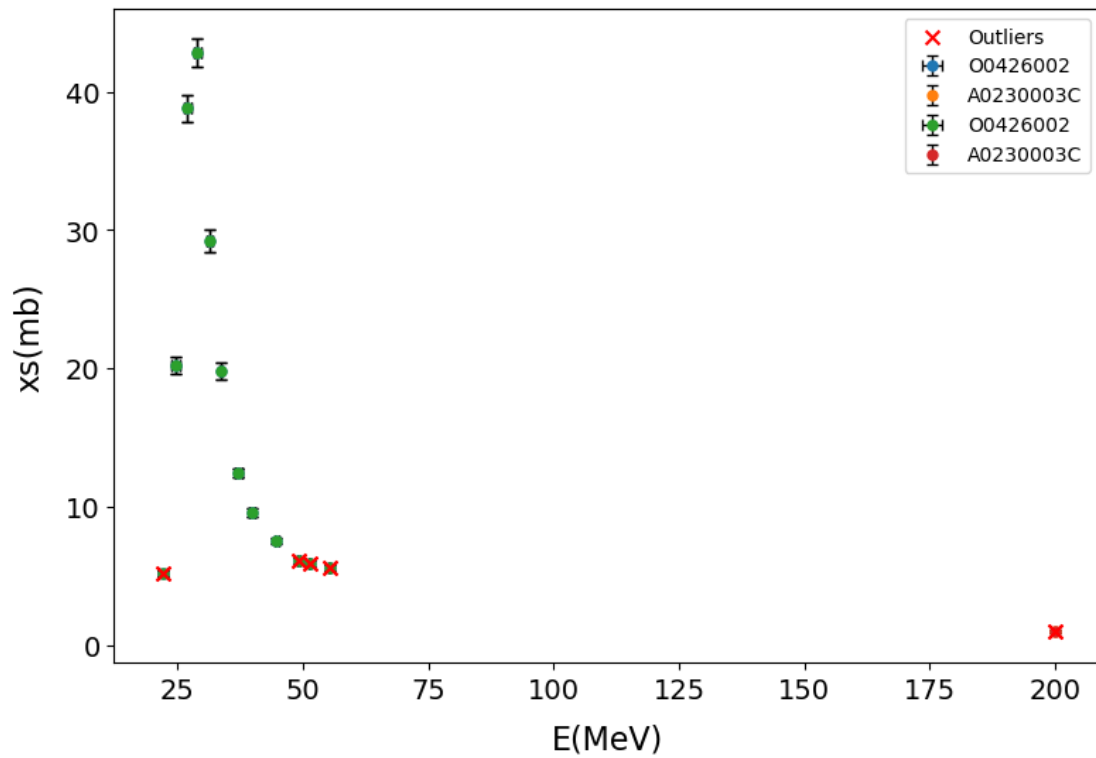


Figure 4.28: Outlier Detection in Group 4.1 using SVM Method.

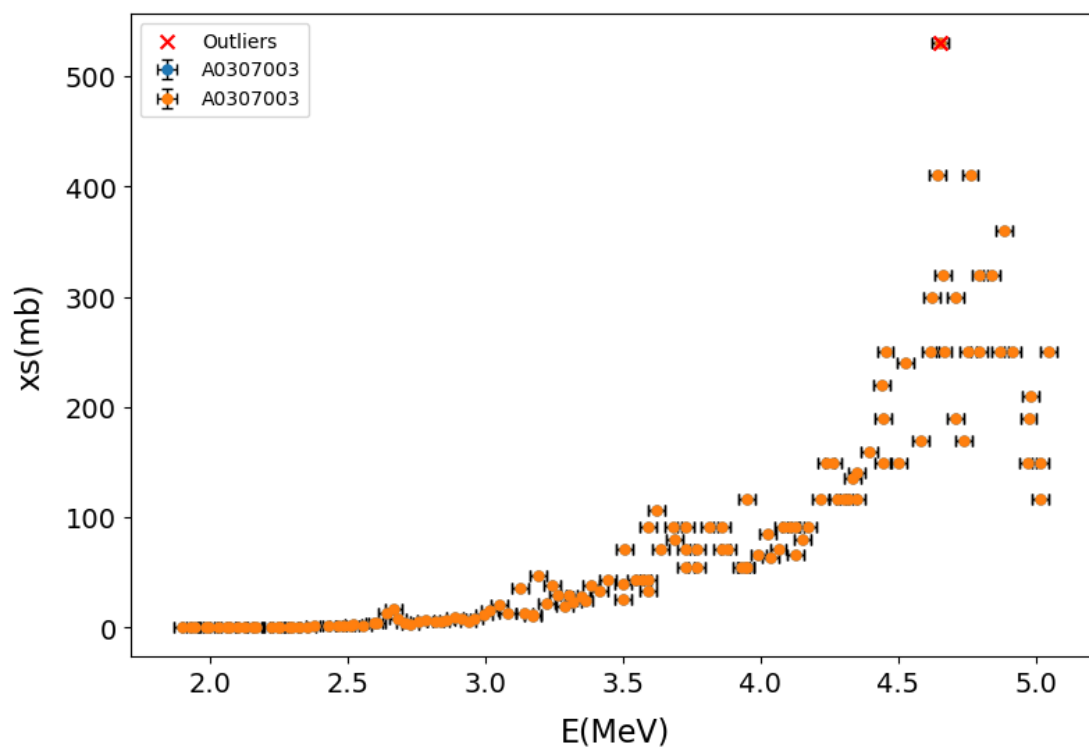


Figure 4.29: Outlier Detection in Group 4.2 using SVM Method.

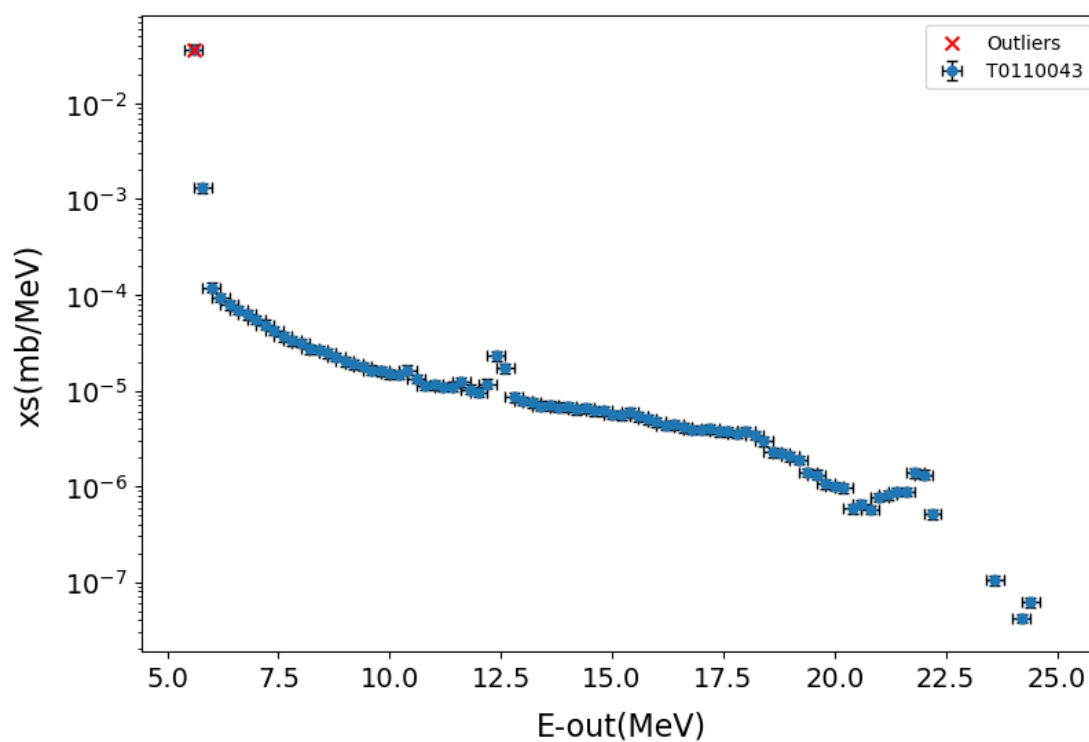


Figure 4.30: Outlier Detection in Group 5 using SVM Method.

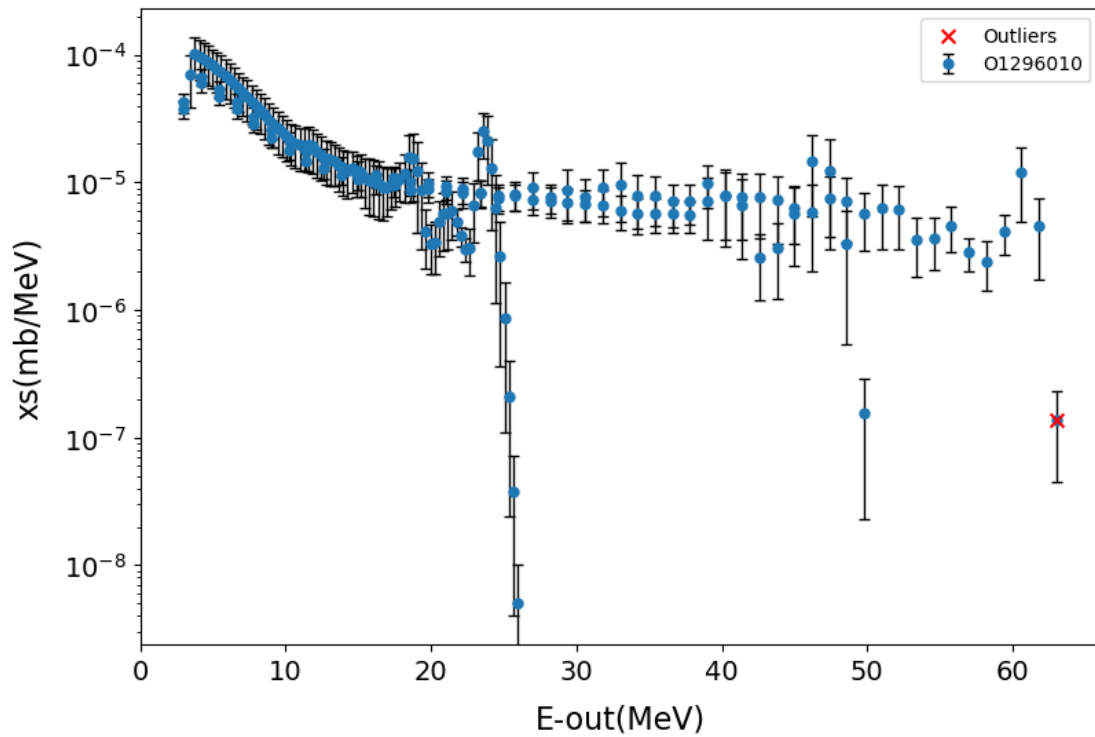


Figure 4.31: Outlier Detection in Group 5 using SVM Method.

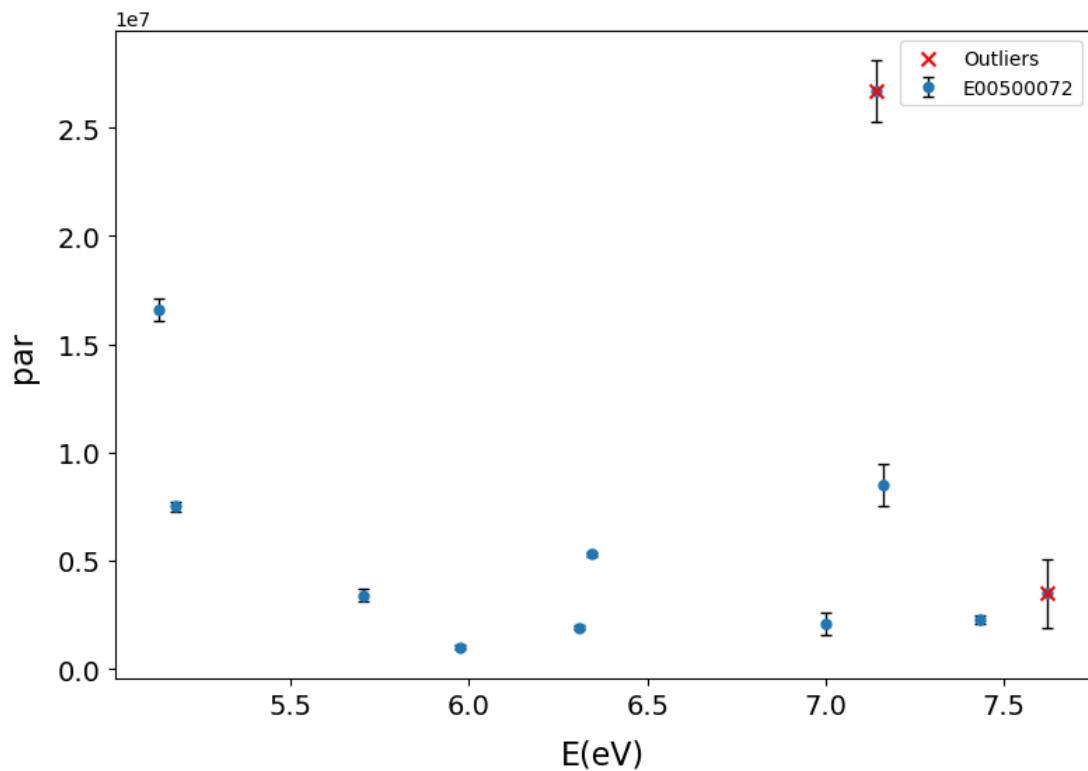


Figure 4.32: Outlier Detection in Group 6 using SVM Method.

Figures 4.26 to 4.32 show the results of using the SVM method for outlier detection. Overall, the method performs well across all data groups. While some false positives do appear, which

is common for these types of algorithms, SVM does a good job of identifying likely outliers in each dataset.

4.6.3 Conclusion on SVM Method

The SVM algorithm proved effective in detecting outliers across all our data groups. The use of default parameters serves as a reasonable starting point, demonstrating the algorithm's robustness across diverse scenarios. While these default settings are well-calibrated for a broad range of problems, they can be further fine-tuned either through domain-specific knowledge or hyperparameter optimization techniques.

Despite its overall effectiveness, the algorithm is computationally intensive, particularly for large datasets. This limitation was evident in our study where we had to divide the first data group into smaller subsets to manage the computational requirements. Thus, while SVM is powerful, its scalability can pose challenges for larger datasets.

4.7 Neural Network Models: Autoencoders

While previous subsections focused on machine learning algorithms, we now turn to a neural network-based methodology—specifically, autoencoders implemented using the Keras library. Due to the inherent complexities associated with neural networks, this section is subdivided to cover implementation details, evaluation strategies, and a specialized use-case for outlier detection. This structured approach aims to provide a comprehensive overview from architectural design to practical applications.

4.7.1 Implementation of the Autoencoder

Data Preprocessing:

```
1 # Data Splitting (Train/Test)
2 X_train, X_test = train_test_split(df, test_size=0.2, random_state=42)
3 # Data Scaling
4 scaler = MinMaxScaler()
5 X_train = scaler.fit_transform(X_train)
6 X_test = scaler.transform(X_test)
```

Firstly, the dataset is split into a training and test set using scikit-learn's `train_test_split` function, reserving 20% of the data for testing. The `random_state=42` is set for reproducibility. Then, the `MinMaxScaler` is applied to scale the features between 0 and 1, aiding in model convergence.

Model Architecture:

```
1 # Model Architecture
2 input_dim = X_train.shape[1]
3 encoding_dim = int(input_dim / 2)
4 input_layer = Input(shape=(input_dim,))
5 encoder = Dense(encoding_dim, activation="relu")(input_layer)
6 decoder = Dense(input_dim, activation="sigmoid")(encoder)
```


The architecture of the autoencoder is simple but effective, consisting of an input layer, an encoder layer, and a decoder layer. The encoder uses a **ReLU** activation function, and the decoder uses a **sigmoid** activation function. The encoding dimension is set to half of the input dimension, a parameter that can be adjusted depending on specific needs.

Model Compilation and Training:

```
1 # Model Compilation
2 autoencoder.compile(optimizer=Adam(learning_rate=0.001), loss="
   mean_squared_error")
3 # Model Training
4 history = autoencoder.fit(X_train, X_train, epochs=50, batch_size=256,
   validation_data=(X_test, X_test))
```

The model is compiled using the **Adam** optimizer with a learning rate of 0.001 and Mean Squared Error as the loss function. Finally, the model is trained for 50 epochs with a batch size of 256, using the test data for validation during training.

4.7.2 Model Evaluation

For the evaluation of our autoencoder model, we rely on a combination of visual inspection of the training and validation loss curves, as well as quantitative metrics, including the Mean Squared Error (MSE).

Visual Inspection of Loss:

```
1 # Loss Plot
2 plt.plot(history.history["loss"], label="Training Loss")
3 plt.plot(history.history["val_loss"], label="Validation Loss")
4 plt.legend()
5 plt.yscale("log")
6 plt.show()
```

The loss curves for both the training and validation sets are illustrated in Figure 4.33. A logarithmic scale is employed to enhance visualization clarity. A downward trend across epochs typically signifies effective model learning. For data group 1, as shown in the figure, the loss curve stabilizes and flattens, indicating that the chosen number of epochs is sufficient for model convergence.

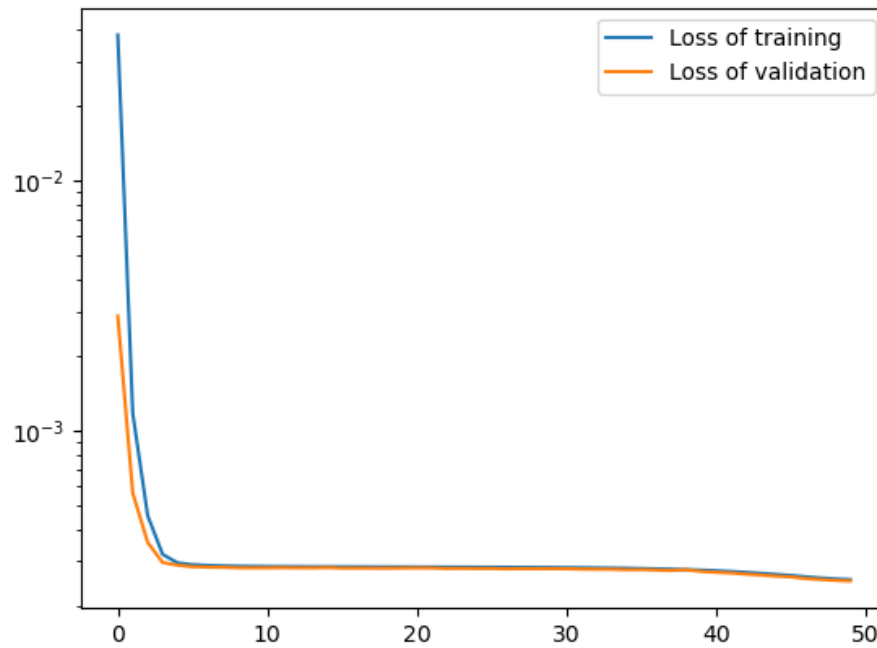


Figure 4.33: Training and Validation Loss Curves for Group 1 data.

Immediately following the visualization, we examine several performance metrics:

```
1 # Performance Metrics
2 predictions = autoencoder.predict(X_test)
3 mse = np.mean(np.power(X_test - predictions, 2), axis=1)
4 print(f"MSE: {np.mean(mse)}")
5
6 mean_value = np.mean(X_train, axis=0)
7 mse_naive = np.mean((X_test - mean_value) ** 2)
8 print(f"MSE of naive model: {mse_naive}")
9
10 std_value = np.std(X_train, axis=0)
11 print(f"Standard deviation: {std_value}")
```

Interpreting Metrics: To better comprehend the MSE value, consider these approaches:

- *Comparison to Naive Model:* A low MSE value could be a good sign, but it's informative to compare it with a naive model that always predicts the mean of the training set. A lower MSE for the autoencoder compared to the naive model would highlight its relative effectiveness.
- *Comparison to Data Variability:* Furthermore, comparing the MSE to the standard deviation of the features can offer insights. If the MSE is significantly lower than the standard deviation, the model is likely doing a good job in capturing the underlying patterns in the data.

By examining both the visual loss curves and these quantitative metrics, a comprehensive assessment of the model's performance can be attained.

4.7.3 Outlier Detection Using Autoencoder

After evaluating the model's performance, we utilize it for outlier detection. The primary objective is to identify data points that deviate significantly from the normal data distribution.

Data Preparation:

```

1 # Preparing the DataFrame for Outlier Detection
2 df = pd.read_csv(path)
3 df = clean_dataframe(df)
4 X4_ID_column = df["X4_ID"].copy()
5 df_without_X4_ID = df.drop("X4_ID", axis=1)

```

First, we prepare the dataframe by reading the data and applying any preprocessing steps such as cleaning. We preserve the X4_ID column for later use and remove it before scaling and making predictions.

Scaling and Predicting:

```

1 # Scaling and Predictions
2 df_scaled_complete = scaler.transform(df_without_X4_ID)
3 predictions = autoencoder.predict(df_scaled_complete)

```

The model is then applied to the scaled data to generate the reconstructed values.

Identifying Anomalies:

```

1 # Detecting Outliers
2 reconstruction_error = np.mean(np.power(df_scaled_complete - predictions
3     , 2), axis=1)
4 threshold = np.percentile(reconstruction_error, 99.0)
5 anomalies_col2 = reconstruction_error > threshold

```

The reconstruction error, calculated as the mean squared error between the original and reconstructed data, serves as the metric for identifying outliers. We then select a threshold based on the 99th percentile of the reconstruction error to flag anomalies.

Post-Processing:

```

1 # Post-processing
2 df_original_values = pd.DataFrame(scaler.inverse_transform(
3     df_scaled_complete), columns=df_without_X4_ID.columns, index=df.index
4 )
5 df_original_values["X4_ID"] = X4_ID_column
6 df_original_values["Outliers"] = anomalies_col2
7
8 # Extracting the Outliers
9 outliers_df = df_original_values[df_original_values["Outliers"] == True
10     ].drop("Outliers", axis=1)

```

Finally, the dataframe is post-processed to extract the outliers based on the flag set earlier. We also include the X4_ID column and calculate the percentage of detected outliers.

The autoencoder's architecture inherently learns to compress the most important features of the data. When it attempts to reconstruct an outlier, the reconstruction error tends to be high, making the model an effective tool for outlier detection.

4.7.4 Results and Graphical Analysis

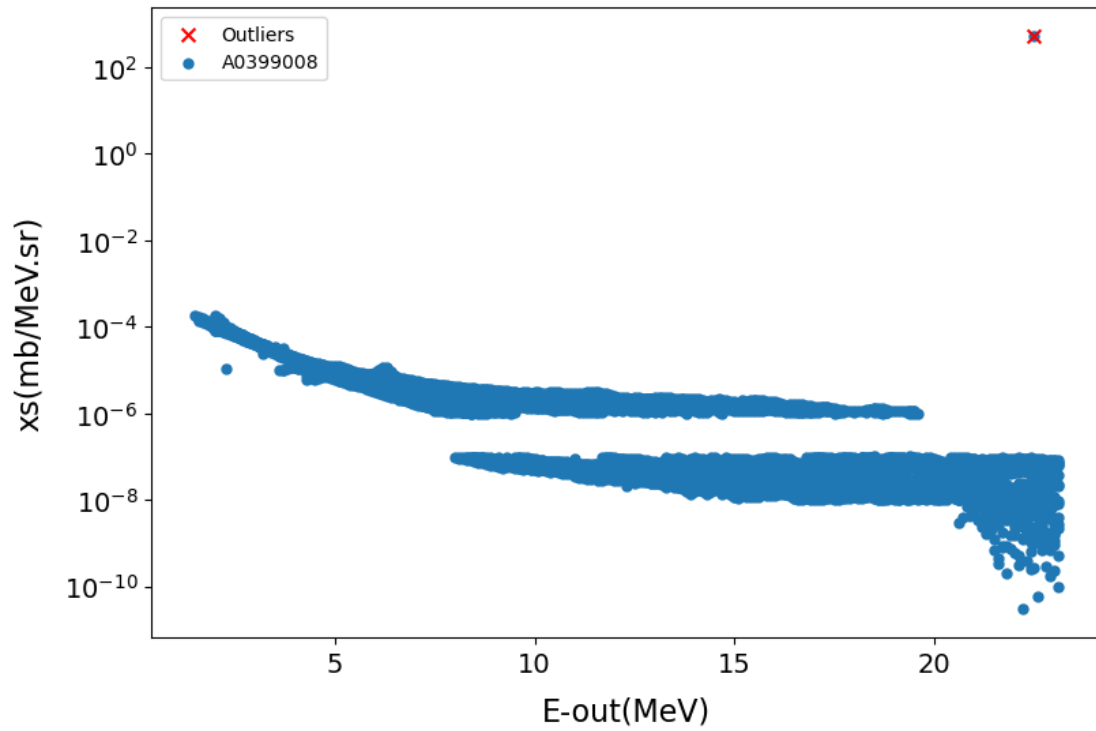


Figure 4.34: Outlier Detection in Group 2 using Autoencoder Method.

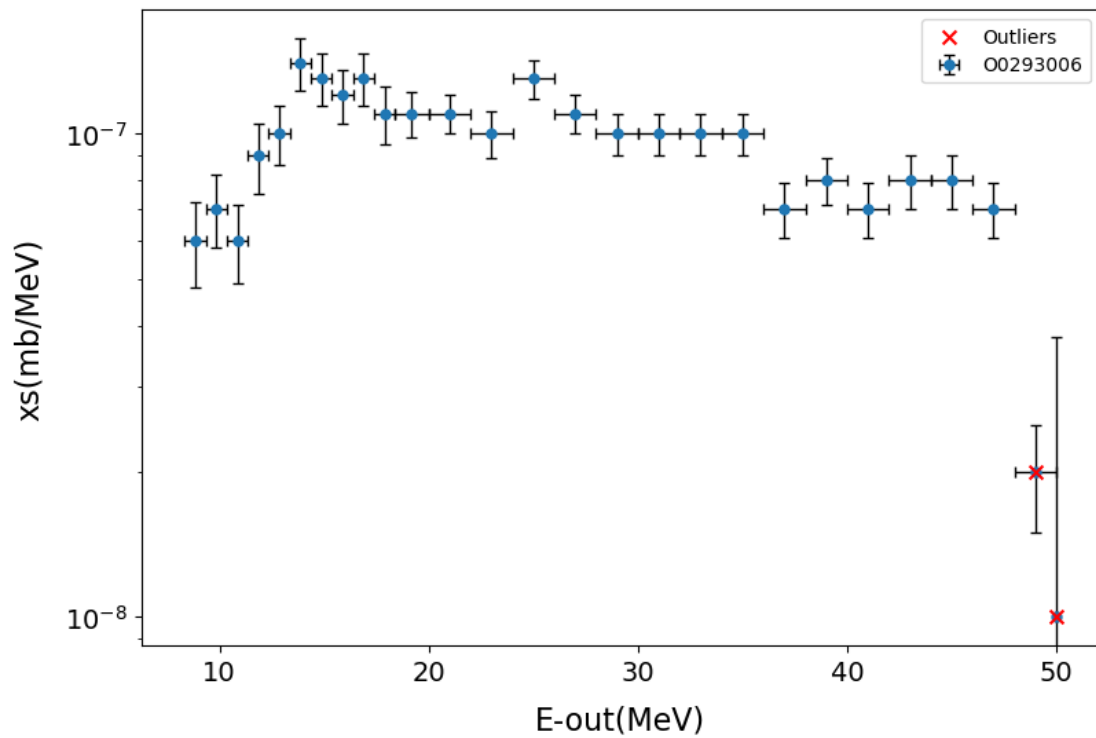


Figure 4.35: Outlier Detection in Group 5 using Autoencoder Method.

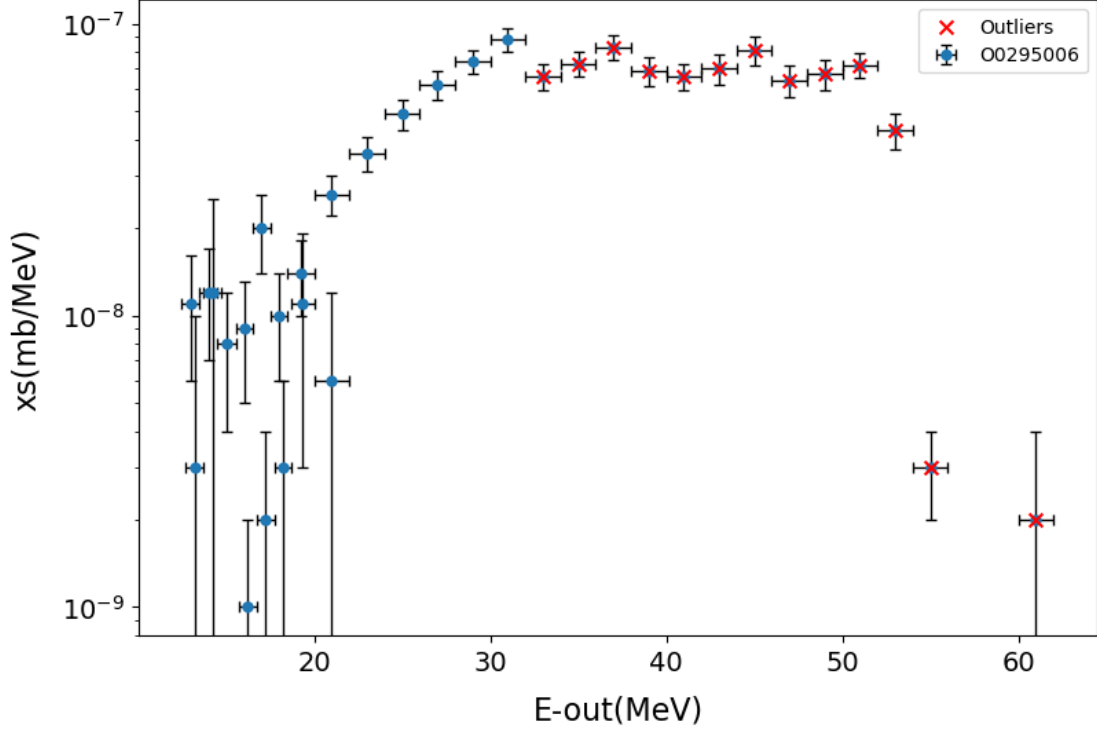


Figure 4.36: Outlier Detection in Group 5 using Autoencoder Method.

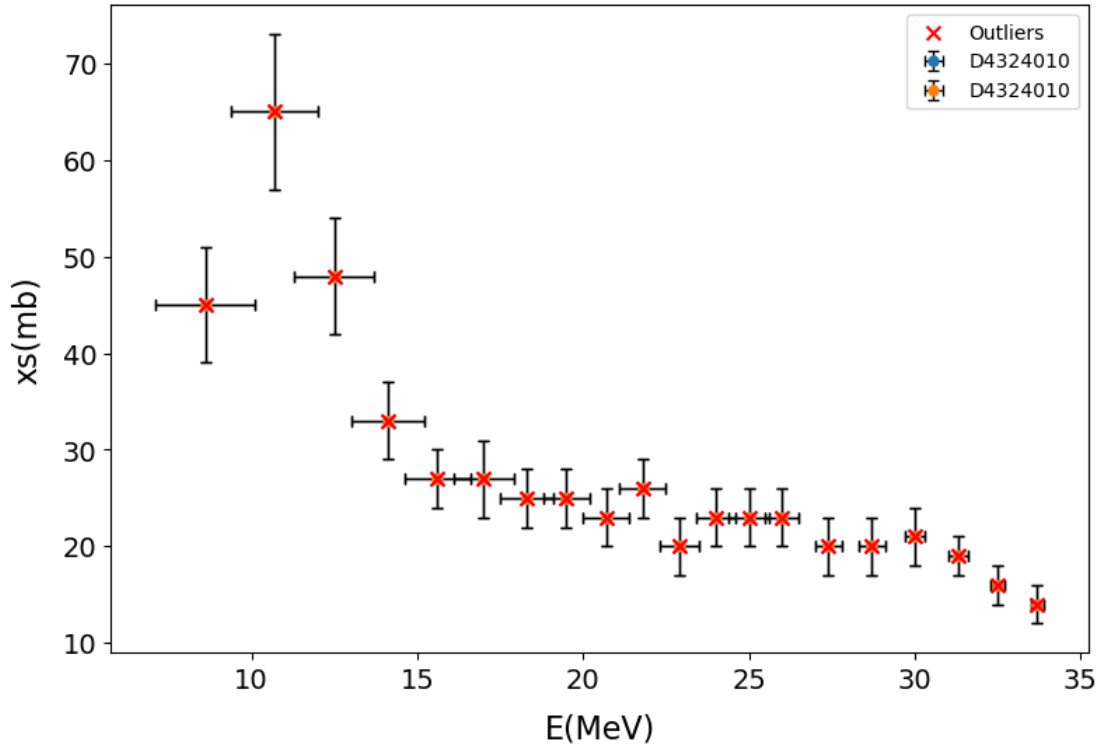


Figure 4.37: Outlier Detection in Group 4.2 using Autoencoder Method.

The results of the outlier detection present varied outcomes as evidenced by the figures. Some instances appear to be potential outliers, as shown in Figures 4.34 and 4.35. However, in other cases, the detection includes false positives, identifying inliers as outliers, such as in Figure 4.36.

Additionally, certain experiments are entirely flagged as outliers despite having a substantial number of data points, as illustrated in Figure 4.37.

4.7.5 Conclusions on Autoencoder Method

In contrast to traditional machine learning methodologies explored earlier in this thesis, the utilization of autoencoders, a neural network-based approach, necessitates a deeper level of understanding for effective application. Despite this complexity, the potential benefits of employing neural networks are significant. Future work could extend to optimizing various parameters such as the number of layers, activation functions, and learning rates.

Regarding the performance of the autoencoder in outlier detection, the results are promising but not without challenges. While the model successfully identified several potential outliers, it also produced a lower rate of false positives compared to other methods. However, it is important to note that the efficacy of the model is heavily reliant on the selected parameters.

One critical consideration is the volume of data available for training. The data groups used for training the autoencoder were identical to those used in traditional machine learning techniques, and in some cases, these groups had limited data points. This scarcity can adversely affect the model's performance. An alternative strategy could be to use the entire proton reaction database without specific categorization into groups. While this approach could potentially eliminate the need for data classification, thereby saving time, there is no guarantee of improved performance. Since classification is required for the machine learning methods, maintaining a consistent approach across methodologies serves as a more effective basis for comparison.

5 CONCLUSIONS

In this chapter, we consolidate the conclusions that arise from our thorough research, scripting, and data analysis efforts. The insights gained from the analyses in Chapter 4 are critical in shaping the conclusions we reach. This chapter is organized into three primary sections: Summary of Findings, Future Work, and Final Conclusions.

Each section delves into the specific conclusions reached from the utilization of various outlier detection techniques applied to the EXFOR nuclear data library. The methods employed in this thesis range from traditional statistical approaches such as the Interquartile Range (IQR) to more advanced machine learning algorithms, including Isolation Forest, DBSCAN, Local Outlier Factor (LOF), Support Vector Machine (SVM), and Autoencoders. Each methodology contributes valuable insights into the intricacies of outlier detection in complex scientific datasets.

Concluding this chapter is a section titled ‘Final Conclusions.’ Within this section, overarching conclusions about the entirety of the research conducted in this thesis are articulated. This final section provides a comprehensive summation of the work and signifies the completion of this thesis.

5.1 Summary of Findings

The preprocessing step is vital for the effective application of machine learning techniques. In this study, we opted to classify EXFOR proton reaction data into eight distinct groups based on the logical criteria of data result similarity. It’s important to note that these groups vary in the number of data points they contain. Generally, machine learning algorithms perform better with larger datasets, which posed a limitation for groups like 3, 6, 7, 8, and to some extent, 5, as they contain fewer data points. Alternate classification schemes that allow for larger, more robust data sets could potentially yield improved results, especially when deploying neural networks, which are better suited for handling complex data.

Machine learning algorithms typically partition the dataset into a training set and a test set. This partitioning is often performed randomly, resulting in variable outcomes with each script execution. Although the impact of this randomization should ideally be minimal, it is strongly influenced by the type and quantity of data at hand. To reproduce the results presented in this thesis, the `random_state` hyperparameter in the machine learning models should be set to 42. For varied results, which could be better, worse, or comparable, one can either remove the `random_state` parameter or modify its value.

It is crucial to understand that machine learning algorithms are typically designed to identify a predetermined percentage of outliers in a dataset, which does not necessarily correspond to the actual number of outliers present. This leads to an inevitable occurrence of false positives. Hence, the expertise of a domain specialist is indispensable for the thorough interpretation of algorithmic results. While machine learning can flag potential outliers, a human expert ultimately needs to evaluate these flagged data points to distinguish between legitimate outliers, false positives, and genuine but anomalous data points. Given the sheer volume of data in the EXFOR library, manual analysis would be unfeasible, highlighting the invaluable role of machine learning techniques in this study. This work contributes significantly to the field by providing a freely accessible, user-friendly tool for data review, available to anyone with a computer equipped with Python.

To finish this section, we summarize the key findings and limitations of each outlier detection method applied in this study to provide a consolidated view of their performance and applicability.

- **Interquartile Range Method:** The IQR method stands out for its statistical foundation and its capability to incorporate uncertainties, potentially reducing false positives. However, it demonstrated limitations when applied to our complex data sets. Specifically, the IQR method is more suited for uniformly distributed data, which contrasts with the varied distributions present in the EXFOR library. The method struggled with identifying peaks in the data as part of the distribution rather than as outliers. Hence, while IQR has its merits, it is not the most compatible technique for outlier detection in the specialized context of this study.
- **Isolation Forest Method:** This method generally outperformed traditional statistical methods like the IQR, highlighting its effectiveness in detecting outliers. However, the algorithm's sensitivity to the distribution of the dataset revealed some drawbacks. Specifically, the method struggles with sparse regions in the data, often misidentifying isolated data points as outliers when they are simply under-sampled regions. This tendency towards false positive detection limits the method's applicability for our specific use-case.
- **DBSCAN Method:** The DBSCAN stands out for its density-based approach to outlier detection, offering a unique angle compared to traditional statistical methods. While the algorithm is known to be sensitive to hyperparameter settings, even its default configuration yielded insightful results in this study, making it a valuable addition to our outlier detection toolbox. However, DBSCAN does come with computational challenges. Specifically, the method proved to be resource-intensive for larger datasets, requiring us to subdivide our data to make the algorithm computationally feasible. This poses a scalability issue that needs to be addressed for broader applications.
- **Local Outlier Factor Method:** The LOF method shines in its adaptability, handling a variety of data distributions more robustly than some other methods. Particularly, it performs well in datasets with sharp peaks or clusters, avoiding the pitfall of misclassifying inliers as outliers. However, LOF is not without its challenges. The method is sensitive to dataset characteristics and demands careful hyperparameter tuning. In specific cases, different groupings of data may necessitate individualized hyperparameter settings, which adds to the time and effort required for effective implementation.
- **Support Vector Machine Method:** The SVM algorithm stood out for its effectiveness in outlier detection across various data groups, even when using default parameters. This speaks to its robustness and potential for fine-tuning through either domain-specific expertise or hyperparameter optimization techniques. However, the algorithm's computational intensity is a notable limitation, especially when dealing with large datasets. This was demonstrated in our study, as the first data group had to be split into smaller subsets to meet computational demands, highlighting challenges in scalability.
- **Autoencoder Method:** Utilizing autoencoders, a form of neural networks, presented a unique set of advantages and challenges. Despite the complexity associated with neural network-based approaches, the autoencoder's performance was promising, yielding a lower rate of false positives compared to traditional techniques. Nevertheless, the model's efficacy was sensitive to parameter selection. The study also revealed that data scarcity in certain groups negatively impacted the model's performance. Future work may consider a unified approach, training the model on the entire proton reaction database, although it remains uncertain whether this would guarantee better results while sacrificing data classification consistency.

In conclusion, each method for detecting outliers that was examined in this study has distinct advantages and limitations. Their effectiveness in identifying potential outliers varies depending on the specific characteristics of the datasets being analyzed; some methods excel in areas where others may fall short. In the realm of machine learning techniques, it is difficult to single out a definitive best performer. Although DBSCAN demonstrated marginally better outcomes, this efficacy is largely dependent on the nature of the data and the specific settings chosen for each analysis. Neural network methods, such as autoencoders, identified fewer outliers, which suggests they may be more accurate. However, the greater computational requirements and more complex parameter settings for neural networks raise questions about their practical suitability for the task at hand. Therefore, each technique offers its own perspective for identifying outliers, and their relative success is tied to the complexities of the datasets they are applied to and the specific objectives they aim to fulfill.

5.2 Future Work

The foundations established by this research pave the way for further study in the realm of outlier detection for specialized scientific databases like EXFOR. Our methodology has proven its utility in various respects but leaves room for advancements and refinements.

Future research could delve deeper into the fine-tuning of hyperparameters for different machine learning algorithms used in this study. While the default settings employed here provided a reasonable starting point, more meticulous hyperparameter optimization could potentially enhance the algorithms' performance. However, this fine-tuning would demand a substantial investment of time, and the benefits must be carefully weighed, especially as the results will still require expert evaluation and visual verification for final validation.

Another noteworthy consideration for future work centers around Group 7, which presents unique challenges due to its peculiar data structure. Modifying the preprocessing steps specifically for this group could make the existing methodologies more applicable and accurate, warranting further investigation.

One ambitious avenue that future research might explore is the comparison of EXFOR data with evaluated nuclear data libraries. These complex libraries consider various factors beyond just experimental results, offering an additional layer of validation for outlier detection. Yet, this approach is not without its complications, notably the need for rigorous preprocessing to make the data from both sources comparable. Ensuring that data are formatted equivalently will be crucial for meaningful comparative analyses.

By tackling these issues, subsequent research can build upon the groundwork of this study, potentially yielding more robust and refined techniques for outlier detection in scientific databases.

5.3 Final Conclusions

This master's thesis set out with a multifaceted objective aimed at integrating machine learning techniques into the realm of outlier detection within the proton induced reactions of EXFOR nuclear data library. The journey from ideation to realization was encapsulated through a series of structured steps, fulfilling each set aim in a systematic manner.

The study successfully developed a robust methodology tailored specifically for the complexity of the EXFOR data. Through comprehensive research and implementation, we not only devised

efficient strategies for data extraction, transformation, and preparation but also established a framework that serves as a blueprint for future research in applying machine learning to specialized scientific datasets. This foundation ensures that the machine learning algorithms implemented could operate effectively, fulfilling the thesis’s objectives and setting the stage for continued exploration in this field.

Our exploration spanned various outlier detection techniques, from traditional statistical methods like the Interquartile Range (IQR) to more sophisticated machine learning algorithms such as Isolation Forest, DBSCAN, and the Local Outlier Factor (LOF). Each method was implemented, and its results were analyzed to assess its reliability, strengths, and weaknesses. Moreover, we took an ambitious step by integrating a neural network model, specifically an Autoencoder, into our repertoire of algorithms.

In the course of completing this work, we have successfully met all initial objectives and established a robust methodology for applying machine learning to outlier detection in the EXFOR library for proton-induced reactions. Perhaps even more significantly, the developed methodology is versatile enough to be effortlessly extended to other sections of the EXFOR library, offering the potential for comprehensive improvements. To foster further development and community engagement, all scripts and code generated through this research have been made publicly available on the GitHub EXFOR-ProtonReactions-Analysis Repository. This open-source approach not only fulfills the academic imperative of sharing knowledge but also sets the stage for ongoing collaborative efforts to refine and expand upon the work presented here.

6 SOCIAL & PROFESSIONAL RESPONSABILITY

This chapter examines the crucial aspects of social and professional responsibility within the context of enhancing the EXFOR Nuclear Data Library using machine learning techniques. The inherent ethical, professional, and social responsibilities that accompany the development and deployment of these technologies are explored. The chapter will also consider the alignment of the project with the United Nations' Sustainable Development Goals (SDGs). It is imperative that, as researchers and practitioners in the field of nuclear science and machine learning, we recognise the profound impact our work can have, not just scientifically, but also socially, ethically, and professionally.

6.1 Professional Responsibility

When we talk about enhancing the EXFOR Nuclear Data Library using machine learning techniques, we are discussing a process that is far from simple. It requires a deep commitment to doing our work accurately and respectfully, understanding that we are building on the work of many others in our field.

Machine learning is an area that's always evolving. It's a bit like a race, where the finish line keeps moving further away just as you think you're getting close. So, it's our job as professionals in this field to keep learning, always staying up to speed with the latest techniques and ethical considerations. This is incredibly important because any inaccuracies or errors could potentially impact the evaluation of nuclear data libraries. In other words, we need to make sure our machine learning models are reliable and high-quality.

As professionals, we also have to share our work with our peers and contribute to the evaluation of others' work. This means giving constructive feedback and always seeking ways to improve our own work. We need to be mindful of others' confidentiality and intellectual property rights, and we need to follow professional guidelines or rules of conduct. Also, when we share our findings, we need to do so in a way that is clear and understandable, not only to other scientists but to the general public as well.

In essence, our task is not just to improve the EXFOR Nuclear Data Library. We have a responsibility to uphold the highest standards of scientific work. We must do so with honesty, transparency, and a deep respect for our field, always striving to push the boundaries of what we know and what we can achieve.

6.2 Social Responsibility and Sustainable Development Goals (SDGs)

Improving the EXFOR Nuclear Data Library through machine learning has extensive implications for sustainable development. Here, we specifically discuss the alignment of this project with several United Nations' Sustainable Development Goals (SDGs):

- **SDG 3: Good Health and Well-being:** Enhanced nuclear data libraries can improve the production of medical isotopes, which are critical for various diagnostic and therapeutic procedures in the healthcare sector. The accurate detection of outliers in the data could increase the efficiency and reliability of these isotopes, potentially enabling better disease diagnosis and treatment strategies.

- **SDG 7: Affordable and Clean Energy:** The improvements in the EXFOR Nuclear Data Library can lead to more efficient, safer nuclear energy production. Accurate and robust nuclear data are essential for the operation and safety assessment of nuclear power plants. Enhanced data accuracy could potentially reduce operational costs, thereby contributing to the affordability aspect of this goal.
- **SDG 9: Industry, Innovation and Infrastructure:** The application of machine learning techniques to the nuclear data library represents a significant innovation in the field. This could inspire further technological advancements and infrastructural improvements in the nuclear industry.
- **SDG 13: Climate Action:** Nuclear energy, being a low-carbon source of power, is pivotal in our attempts to mitigate climate change. Better nuclear data can result in more efficient and safer nuclear power plants, thereby helping to reduce greenhouse gas emissions.

By aligning the enhancement of the EXFOR Nuclear Data Library with these SDGs, we aim to ensure that our project contributes to sustainable development and the betterment of society as a whole.

In conclusion, the integration of machine learning techniques into the enhancement of the EXFOR Nuclear Data Library brings with it considerable professional and social responsibilities. As professionals, we have a duty to ensure the reliability, validity, and robustness of our research and development activities. As members of society, we are compelled to make certain that our work aligns with the global commitment to sustainable development and peaceful usage of nuclear technology. By fulfilling these responsibilities, we can ensure that our work not only advances our field but also contributes positively to society and our planet's future.

7 PROJECT MANAGEMENT

This chapter delves into the crucial aspects of project management involved in enhancing the EXFOR Nuclear Data Library using machine learning techniques. The facets covered in this chapter include both planning and budgeting processes integral to the successful execution of the project. The planning section offers a detailed account of the task distribution and project timeline, presenting a chronological breakdown of the project's structure and execution. We will also discuss some of the challenges faced during the project, including strategies and solutions employed to overcome them. Furthermore, the budgeting aspect of the project, an integral element that will be addressed in a subsequent section, plays a pivotal role in dictating resource allocation and influencing project timelines.

7.1 Planning

The planning phase was meticulously executed to ensure an optimal distribution of tasks, thereby facilitating the smooth progress of the project. Table 7.1 and Figure 7.1 both represent the timeline and distribution of tasks for this master's thesis project. The project was initiated on 1st June and concluded on 31st August, spanning over three months. The tasks were divided into two levels: Task Level 1, representing major phases of the project, and Task Level 2, indicating specific actions within these phases.

The project started with preparatory tasks, such as kick-off and data collection. Subsequent steps included data preprocessing, a crucial task that was allocated more time due to its complexity and importance to the overall project. Machine learning model training, outlier detection, and final model training followed this. Each phase was broken down into subtasks such as feature extraction, data cleaning, model optimization, and cross-validation to ensure comprehensive coverage of each aspect. The final stages of the project involved thesis writing and review, a process started soon enough to ensure adequate time for drafting, review, and revision.

Table 7.1: Project Timeline and Task Distribution

Task	Duration (days)	Begin Date	End Date
Project Kick-off and Preparation	9	1st June	9th June
- Data Collection	5	4th June	9th June
Data Preprocessing	31	9th June	10th July
- Feature Extraction	10	9th June	19th June
- Data Cleaning	21	19th June	10th July
Initial Machine Learning Model Training	16	10th July	26th July
- Model Optimization	9	10th July	19th July
- Cross-validation	7	19th July	26th July
Outlier Detection	15	27th July	10th Aug
- Outlier Analysis and Verification	15	27th July	10th Aug
Final Machine Learning Model Training	12	11th Aug	23th Aug
- Model Testing	12	11th Aug	23th Aug
Thesis Writing & Review	58	4th July	31st Aug
- Drafting	41	4th July	14th Aug
- Review and Revision	16	15th Aug	31st Aug

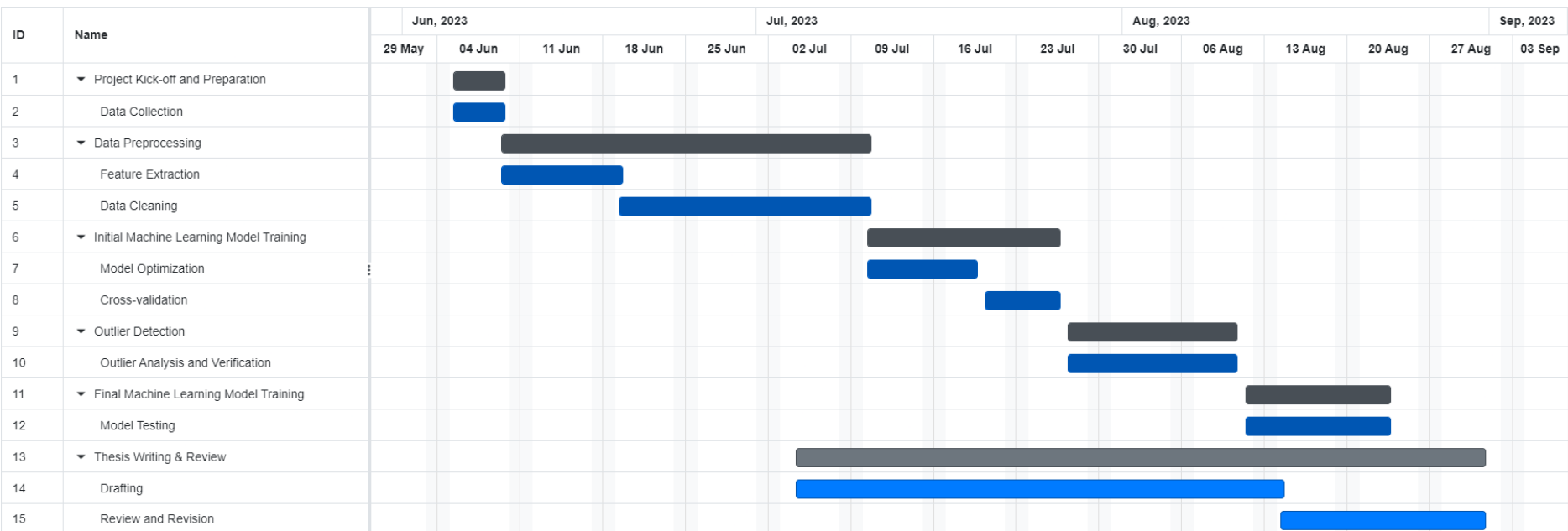


Figure 7.1: Gantt chart.

One of the most significant challenges encountered during the project was managing the tight timeline, particularly concerning thesis drafting. Due to the short duration of the project, drafting started concurrently with the later stages of the project. This parallel approach enabled progress on the thesis while still fine-tuning the project’s technical aspects. Data preprocessing also posed a considerable challenge, given its critical role in the project. We approached this by allocating more time to this task, ensuring that the data was appropriately prepared for machine learning model training. This involved careful scheduling to avoid rushing the process, and iterative refinement as necessary.

In terms of the machine learning aspect, model optimization and outlier detection were challenging due to the complexity of the nuclear data and the potential impact of outliers on the results. We tackled this through a rigorous approach involving multiple iterations of model training, testing, optimization, and validation. This meticulous process ensured the reliable performance of the model and accurate outlier detection. By proactively addressing these challenges, we were able to effectively manage the project timeline and deliver the desired outcomes. This experience underscores the importance of robust project planning and adaptability in the face of unforeseen challenges.

7.2 Budget

In project management, budgeting is a critical component, serving as a financial map that aids the successful execution of a project. For this project, workforce expenses form the primary budgetary aspect, meticulously recorded in Table 7.2.

The student worked a total of 360 hours, corresponding to the 12-credit requirement for the thesis. The student was awarded a grant of 500€ per month for the project’s execution, amounting to a total of 1500€ for the three-month duration of the project. This grant effectively means an hourly rate of 4.17€ for the student’s work. On the other hand, the mentor dedicated approximately 30 hours to this project. With a rate of 35€ per hour, the cost associated with the mentor’s contribution amounts to 1050€.

Table 7.2: Workforce Expenses

Concept	Hours	Cost per hour (€)	Subtotal (€)
Student	360	4.17	1500.00
Mentor	30	35.00	1050.00
Total			2550.00

In terms of resource expenses, a typical project would consider the cost of software licenses, peripherals, accommodation, utilities, and other operational costs. However, for this project, a decision was made to regard these resource expenses as being effectively amortized. This is largely because the project leveraged my personal resources. For instance, my personal computer was used for all tasks, negating the need for additional hardware investments.

Furthermore, the software tools employed for the project, including those for machine learning and document preparation like Overleaf, are open-source or freely available. This decision not only promotes the use of accessible resources but also significantly brings down the cost associated with software procurement. Regarding accommodation and utility expenses, these have been absorbed into my personal living costs as I carried out the project from my family

home. Thus, no separate accommodation expenses are factored into the project budget.

When considering the tax implications of the project budget, it's important to note that in Spain, grants are tax-exempt. Therefore, the grant amount of 1500€ for the student does not factor into the tax calculation. The taxable amount is solely the mentor's fee of 1050€. Applying the 21% tax rate to the mentor's fee, we obtain a tax amount of 220.50€. Therefore, the final budget for the project, including tax, becomes 1500 (student grant) + 1050 (mentor's fee) + 220.50 (tax) = 2770.50€.

8 BIBLIOGRAPHY

References

- [1] S. Team and M. Santos, “Machine learning de a a la z: R y python para data science.” [urlhttps://www.udemy.com/course/machinelearning-es/](https://www.udemy.com/course/machinelearning-es/), 2021. Udemy.
- [2] N. I. of Standards and Technology, “1.3.5.17. detection of outliers.” <https://www.itl.nist.gov/div898/handbook/eda/section3/eda35h.htm>.
- [3] IAEA, “Experimental nuclear reaction data — iaea.” <https://www.iaea.org/resources/databases/experimental-nuclear-reaction-data>.
- [4] ASCO.org, “Proton therapy — cancer.net.” <https://www.cancer.net/navigating-cancer-care/how-cancer-treated/radiation-therapy/proton-therapy>.
- [5] J. Hopkins, “Proton therapy — johns hopkins medicine.” <https://www.hopkinsmedicine.org/health/treatment-tests-and-therapies/proton-therapy>.
- [6] M. Clinic, “Proton therapy - mayo clinic.” <https://www.mayoclinic.org/tests-procedures/proton-therapy/about/pac-20384758>.
- [7] M. Ofuya, L. McParland, L. Murray, S. Brown, D. Sebag-Montefiore, and E. Hall, “Systematic review of methodology used in clinical studies evaluating the benefits of proton beam therapy,” *Clinical and Translational Radiation Oncology*, vol. 19, pp. 17–26, 11 2019.
- [8] A. e. a. Author, *EXFOR Formats Manual*, 1979. In J. Nucl. Phys., Vol. 12, p. 345. EXFOR A0123.012.
- [9] I. A. E. Agency, “Exfor formats description for users.” https://www-nds.iaea.org/exfor/exfor_formats.pdf, 2008. [Accessed: 05/07/2023].
- [10] O. Cabellos and N. Otsuka, “Handbook on nuclear data great-pioneer eu project,” tech. rep., Great-Pioneer, 2022.
- [11] A. Lewis, D. Neudecker, A. Koning, D. Barry, J. Brown, and G. Schnabel, “Wpec sg50: Developing an automatically readable, comprehensive and curated experimental nuclear reaction database,” *EPJ Web of Conferences*, vol. 284, p. 18003, 2023. This is an open access article distributed under the terms of the Creative Commons Attribution License 4.0.
- [12] I. NDS, “Github repository iaea nds.” <https://github.com/IAEA-NDS>. Accessed: 15/07/23.
- [13] “Oecd nuclear energy agency.” <https://www.oecd-neo.org/>. Accessed: 13/07/2023.
- [14] A. Samuel, “Some studies in machine learning using the game of checkers,” *IBM Journal of research and development*, vol. 3, no. 3, pp. 210–229, 1959.
- [15] K. Kourou, T. P. Exarchos, K. P. Exarchos, M. V. Karamouzis, and D. I. Fotiadis, “Machine learning applications in cancer prognosis and prediction,” *Computational and structural biotechnology journal*, vol. 13, pp. 8–17, 2015.
- [16] R. J. Bolton and D. J. Hand, “Statistical fraud detection: A review,” *Statistical Science*, pp. 235–249, 2002.

- [17] A. K. Sikder and A. R. Ganguly, “An overview of computing approaches for predicting short term renewable energy,” *Current Sustainable/Renewable Energy Reports*, vol. 5, no. 4, pp. 135–146, 2018.
- [18] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [19] S. Russell and P. Norvig, *Artificial intelligence: a modern approach*. Malaysia: Pearson Education Limited, 2016.
- [20] F. E. Grubbs, “Procedures for detecting outlying observations in samples,” *Technometrics*, vol. 11, no. 1, pp. 1–21, 1969.
- [21] D. M. Hawkins, *Identification of outliers*, vol. 11. Springer, 1980.
- [22] V. J. Hodge and J. Austin, “A survey of outlier detection methodologies,” in *Artificial intelligence review*, vol. 22, pp. 85–126, Springer, 2004.
- [23] V. Barnett and T. Lewis, *Outliers in statistical data*. Wiley, 1994.
- [24] H. W. Ives, “Evaluation of tests for normality,” *Statistical Distributions in Scientific Work*, vol. 4, pp. 37–52, 1986.
- [25] C. Leys, C. Ley, O. Klein, P. Bernard, and L. Licata, “Detecting outliers: Do not use standard deviation around the mean, use absolute deviation around the median,” *Journal of Experimental Social Psychology*, vol. 49, no. 4, pp. 764–766, 2013.
- [26] D. C. Hoaglin, F. Mosteller, and J. W. Tukey, *Exploring data tables, trends, and shapes*. Wiley, 1986.
- [27] V. Chandola, A. Banerjee, and V. Kumar, “Anomaly detection: A survey,” *ACM computing surveys (CSUR)*, vol. 41, no. 3, pp. 1–58, 2009.
- [28] A. K. Jain, M. Murty, and P. Flynn, “Data clustering: a review,” *ACM computing surveys (CSUR)*, vol. 31, no. 3, pp. 264–323, 1999.
- [29] B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson, “Estimating the support of a high-dimensional distribution,” in *Neural computation*, vol. 13, pp. 1443–1471, MIT Press, 2001.
- [30] N. Chauhan and L. F. Vargas, “Anomaly detection using autoencoders,” *Journal of Physics: Conference Series*, vol. 1222, no. 1, p. 012018, 2019.
- [31] Wikipedia contributors, “Interquartile range — Wikipedia, the free encyclopedia.” https://en.wikipedia.org/w/index.php?title=Interquartile_range&oldid=1171498340, 2023. [Accessed: 25/08/2023].
- [32] Chris Kuo and Dr. Dataman, “Handbook of anomaly detection with python outlier detection.” <https://towardsdatascience.com/use-the-isolated-forest-with-pyod-3818eea68f08>, 2019. [Accessed: 25/08/2023].
- [33] O. Ali, M. K. Ishak, and M. K. Bhatti, “A machine learning approach for early covid-19 symptoms identification,” *Computers, Materials and Continua*, vol. 70, pp. 3803–3820, 09 2021.
- [34] A. Elmogy, H. Rizk, and A. Sarhan, “Ofcod: On the fly clustering based outlier detection framework,” *Data*, vol. 6, p. 1, 12 2020.

- [35] B. Koo and B. Shin, “Using geometry based anomaly detection to check the integrity of ifc classifications in bim models,” *Journal of KIBIM*, vol. 7, pp. 18–27, 03 2017.
- [36] Arden Dertat, “Applied deep learning - part 3: Autoencoders.” <https://towardsdatascience.com/applied-deep-learning-part-3-autoencoders-1c083af4d798>, 2017. [Accessed: 25/08/2023].
- [37] A. Koning, “EXFORTABLES-1.0 User Manual,” 2019.

LIST OF TABLES

2.1 Major Evaluated Nuclear Data Libraries and their respective Organisations [13] . 8

7.1 Project Timeline and Task Distribution 69

7.2 Workforce Expenses 71

LIST OF FIGURES

2.1	Structure of information in EXFOR [10].	5
2.2	Artificial Intelligence and Machine Learning taxonomy [1].	10
2.3	Nested structure of AI [1].	11
3.1	Example entry from EXFORTABLES.	16
3.2	Graphical representation of the IQR method [31].	20
3.3	Graphical representation of the Isolation Forest method [32].	21
3.4	Graphical representation of the DBSCAN method [33].	22
3.5	Graphical representation of the LOF method [34].	24
3.6	Graphical representation of the SVM method [35].	25
3.7	Graphical representation of an autoencoder neural network [36].	26
4.1	Experimental results obtained using the plot method of the Experiment class. . .	31
4.2	Experimental results with uncertainties obtained using the plot method of the Experiment class.	32
4.3	Experimental results of all proton-induced experiments with $MT = 42$	33
4.4	Outlier Detection in Group 6 using IQR Method.	36
4.5	Outlier Detection in Group 6 using IQR Method.	36
4.6	Outlier Detection in Group 1 using IQR Method.	37
4.7	Outlier Detection in Group 5 using IQR Method.	37
4.8	Outlier Detection in Group 1 using IQR Method.	38
4.9	Outlier Detection in Group 2 using Isolation Forest Method.	40
4.10	Outlier Detection in Group 4.2 using Isolation Forest Method.	40
4.11	Outlier Detection in Group 5 using Isolation Forest Method.	41
4.12	Outlier Detection in Group 5 using Isolation Forest Method.	41
4.13	Outlier Detection in Group 6 using Isolation Forest Method.	42
4.14	Outlier Detection in Group 2 using DBSCAN Method.	43
4.15	Outlier Detection in Group 2 using DBSCAN Method.	44
4.16	Outlier Detection in Group 3 using DBSCAN Method.	44
4.17	Outlier Detection in Group 4.2 using DBSCAN Method.	45

4.18	Outlier Detection in Group 5 using DBSCAN Method.	45
4.19	Outlier Detection in Group 6 using DBSCAN Method.	46
4.20	Outlier Detection in Group 2 using LOF Method.	48
4.21	Outlier Detection in Group 3 using LOF Method.	48
4.22	Outlier Detection in Group 5 using LOF Method.	49
4.23	Outlier Detection in Group 2 using LOF Method.	49
4.24	Outlier Detection in Group 4.2 using LOF Method.	50
4.25	Outlier Detection in Group 1 using LOF Method.	50
4.26	Outlier Detection in Group 2 using SVM Method.	52
4.27	Outlier Detection in Group 2 using SVM Method.	53
4.28	Outlier Detection in Group 4.1 using SVM Method.	53
4.29	Outlier Detection in Group 4.2 using SVM Method.	54
4.30	Outlier Detection in Group 5 using SVM Method.	54
4.31	Outlier Detection in Group 5 using SVM Method.	55
4.32	Outlier Detection in Group 6 using SVM Method.	55
4.33	Training and Validation Loss Curves for Group 1 data.	58
4.34	Outlier Detection in Group 2 using Autoencoder Method.	60
4.35	Outlier Detection in Group 5 using Autoencoder Method.	60
4.36	Outlier Detection in Group 5 using Autoencoder Method.	61
4.37	Outlier Detection in Group 4.2 using Autoencoder Method.	61
7.1	Gantt chart.	70

GLOSSARY

AI Artificial Intelligence

DBSCAN Density-based clustering non-parametric algorithm

DL Deep Learning

ENDL Evaluated Nuclear Data Libraries

EXFOR EXchange FORmat

IQR Interquartile Rrange

LOF Local Outlier Factor

MEDUSAL Machine-readable Experimental Data User Application & Library

ML Machine Learning

NEA Nuclear Energy Agency

OECD Organisation for Economic Co-operation and Development's

SDG Sustainable Development Goals

SVM Support Vector Machines

WPEC Working Party on International Nuclear Data Evaluation Co-operation

ANNEXES

A ‘Experiment’ Class Source Code and Documentation

```

1  class Experiment:
2      """
3      This class models an experimental setup for scientific data
4      collection.
5      Attributes include various experimental parameters and data in the
6      form of a Pandas DataFrame.
7      """
8      def __init__(self):
9          """
10         Initializes an Experiment object with default values set to None
11         and an empty Pandas DataFrame for data.
12         """
13         self.title = None
14         self.target_Z = None
15         self.target_A = None
16         self.target_state = None
17         self.projectile = None
18         self.reaction = None
19         self.E_inc = None
20         self.final_Z = None
21         self.final_A = None
22         self.final_state = None
23         self.MTrat = None
24         self.Ratio_isomer = None
25         self.quantity = None
26         self.frame = None
27         self.MF = None
28         self.MT = None
29         self.X4_ID = None
30         self.X4_code = None
31         self.author = None
32         self.year = None
33         self.data_points = None
34         self.data = pd.DataFrame()
35         self.reference = None
36
37     def __str__(self):
38         """
39         Overloaded string method to print non-empty attributes of an
40         Experiment object.
41         """
42         print("Title: ", self.title)
43         if self.target_Z != None: print("Target Z: ", self.target_Z)
44         if self.target_A != None: print("Target A: ", self.target_A)
45         if self.target_state != None: print("Target state: ", self.
46         target_state)
47         if self.projectile != None: print("Projectile: ", self.
48         projectile)
49         if self.reaction != None: print("Reaction: ", self.reaction)
50         if self.E_inc != None: print("Incident energy: ", self.E_inc)
51         if self.final_Z != None: print("Final Z: ", self.final_Z)
52         if self.final_A != None: print("Final A: ", self.final_A)

```

```

48     if self.final_state != None: print("Final state: ", self.
        final_state)
49     if self.MTrat != None: print("MT ratio: ", self.MTrat)
50     if self.Ratio_isomer != None: print("Ratio isomer: ", self.
        Ratio_isomer)
51     if self.quantity != None: print("Quantity: ", self.quantity)
52     if self.frame != None: print("Frame: ", self.frame)
53     if self.MF != None: print("MF: ", self.MF)
54     if self.MT != None: print("MT: ", self.MT)
55     if self.X4_ID != None: print("X4 ID: ", self.X4_ID)
56     if self.X4_code != None: print("X4 code: ", self.X4_code)
57     if self.author != None: print("Author: ", self.author)
58     if self.year != None: print("Year: ", self.year)
59     if self.data_points != None: print("Data points: ", self.
        data_points)
60     if not self.data.empty: print("Data: ", self.data)
61     if self.reference != None: print("Reference:\n", self.reference)
62     return ""
63
64
65     def to_dataframe(self):
66         """
67         Converts the attributes of the object into a single DataFrame
        where each attribute is a column.
68         Returns a DataFrame containing the enriched data.
69         """
70         # Starting with the data in self.data, we are going to add
        # information from the other attributes
71         # The information will be added as a new column in the dataframe
        # for all the rows in the dataframe
72
73         # List with all the attributes to add to the dataframe
74         attributes = ["Experiment", "Target Z", "Target A", "Target
            state", "Reaction", "Incident energy",
75                       "Final Z", "Final A", "Final state", "MT ratio", "
            Ratio isomer", "Quantity", "Frame",
76                       "MF", "MT", "Author", "Year"]
77
78         # Create a copy of the data
79         data = self.data.copy()
80         # Add the attributes to the dataframe
81         for attribute in attributes:
82             # Create a list with the value of the attribute for all the
            # rows in the dataframe
83             if attribute == "Experiment": value = [self.title]*len(data)
84             elif attribute == "Target Z": value = [self.target_Z]*len(
            data)
85             elif attribute == "Target A": value = [self.target_A]*len(
            data)
86             elif attribute == "Target state": value = [self.target_state
            ]*len(data)
87             elif attribute == "Reaction": value = [self.reaction]*len(
            data)
88             elif attribute == "Incident energy": value = [self.E_inc]*
            len(data)
89             elif attribute == "Final Z": value = [self.final_Z]*len(data)
            )

```

```

90         elif attribute == "Final A": value = [self.final_A]*len(data
91         )
92         elif attribute == "Final state": value = [self.final_state]*
93         len(data)
94         elif attribute == "MT ratio": value = [self.MTrat]*len(data)
95         elif attribute == "Ratio isomer": value = [self.Ratio_isomer
96         ]*len(data)
97         elif attribute == "Quantity": value = [self.quantity]*len(
98         data)
99         elif attribute == "Frame": value = [self.frame]*len(data)
100        elif attribute == "MF": value = [self.MF]*len(data)
101        elif attribute == "MT": value = [self.MT]*len(data)
102        elif attribute == "Author": value = [self.author]*len(data)
103        elif attribute == "Year": value = [self.year]*len(data)
104        # Add the column to the dataframe
105        data[attribute] = value
106    # Return the dataframe
107    return data
108
109
110    def add_numeric_attributes(self):
111        """
112        Adds numeric attributes to the DataFrame by converting suitable
113        string values to their respective numeric forms.
114        """
115        attributes = ["E_inc", "MF", "MT", "MTrat", "Ratio_isomer", "
116        final_A", "final_Z", "target_A", "target_Z"]
117
118        for attr in attributes:
119            if attr == "E_inc":
120                value = getattr(self, attr)
121                if value is not None:
122                    first_element = value.split()[0]
123                    self.data[attr] = float(first_element)
124                else:
125                    self.data[attr] = None
126            elif attr in ["MTrat", "final_A", "final_Z"]:
127                value = getattr(self, attr)
128                if value is not None:
129                    self.data[attr] = int(value)
130                else:
131                    self.data[attr] = None
132            else:
133                self.data[attr] = getattr(self, attr)
134
135
136    def encode_categorical_attributes(self):
137        """
138        Transforms categorical attributes to a format suitable for model
139        training by using one-hot encoding.
140        """
141        categorical_attributes = {
142            "projectile": ["projectile_p"],
143            "final_state": ["final_state_+", "final_state_1", "
144            final_state_2", "final_state_G", "final_state_M"],
145            "frame": ["frame_C", "frame_L"],
146            "quantity": ["qty_Angular distribution", "qty_Cross section"
147            , "qty_Cross section ratio", "qty_Delayed nubar", "

```

```

139         qty_Differential cross section", "qty_Fission yields", "
            qty_Prompt nubar", "qty_Resonance Parameters", "qty_Total
            nubar"],
140     "reaction": ["reaction_(n, )", "reaction_(p, el)", "
            reaction_(p, f)", "reaction_(p, x)", "reaction_(p, xa)",
            "reaction_(p, xd)", "reaction_(p, xg)", "reaction_(p, xh)
            ", "reaction_(p, xn)", "reaction_(p, xp)", "reaction_(p,
            xt)", "reaction_(p,2a)", "reaction_(p,2n)", "reaction_(p
            ,2n)g", "reaction_(p,2n)m", "reaction_(p,2na)", "
            reaction_(p,2np)", "reaction_(p,2p)", "reaction_(p,2p)g",
            "reaction_(p,2p)m", "reaction_(p,3a)", "reaction_(p,3n)
            ", "reaction_(p,3n)g", "reaction_(p,3n)m", "reaction_(p,3n
            )n", "reaction_(p,3na)", "reaction_(p,3np)", "reaction_(p
            ,3np)g", "reaction_(p,3np)m", "reaction_(p,4n)", "
            reaction_(p,4n)g", "reaction_(p,4n)m", "reaction_(p,a)",
            "reaction_(p,a)g", "reaction_(p,a)m", "reaction_(p,d)", "
            reaction_(p,d2a)", "reaction_(p,da)", "reaction_(p,f)", "
            reaction_(p,f)g", "reaction_(p,f)m", "reaction_(p,f)n", "
            reaction_(p,g)", "reaction_(p,g)g", "reaction_(p,g)m", "
            reaction_(p,h)", "reaction_(p,h)g", "reaction_(p,n)", "
            reaction_(p,n)g", "reaction_(p,n)m", "reaction_(p,n)n"
            , "reaction_(p,n)_01", "reaction_(p,n)_40", "reaction_(
            p,n2a)", "reaction_(p,n2p)", "reaction_(p,n3a)", "
            reaction_(p,na)", "reaction_(p,na)g", "reaction_(p,na)m",
            "reaction_(p,non)", "reaction_(p,np)", "reaction_(p,np)g
            ", "reaction_(p,np)m", "reaction_(p,npa)", "reaction_(p,p
            )", "reaction_(p,p)m", "reaction_(p,pa)", "reaction_(p,pd
            )", "reaction_(p,pt)", "reaction_(p,t)", "reaction_(p,xa)
            ", "reaction_(p,xd)", "reaction_(p,xg)", "reaction_(p,xh)
            ", "reaction_(p,xn)", "reaction_(p,xp)", "reaction_(p,xt)
            ", "reaction_Exchange_scattering", "
            reaction_Inelastic_scattering", "reaction_ratio"],
141     "target_state": ["target_state_m"]
142 }
143 new_columns_data = {}
144
145 for attribute, possible_values in categorical_attributes.items():
146     # Get the current attribute value of this instance
147     current_attribute_value = getattr(self, attribute)
148
149     for value in possible_values:
150         # Get the value after the last underscore
151         comparison_value = value.split("_")[-1]
152         # Check if the attribute value of the current instance
153         # matches the comparison value
154         new_columns_data[value] = int(current_attribute_value ==
155                                     comparison_value)
156
157     # Convert the dictionary into a DataFrame and concatenate it
158     # with the original DataFrame
159     new_columns_df = pd.DataFrame(new_columns_data, index=self.data.
160                                  index)
161     self.data = pd.concat([self.data, new_columns_df], axis=1)
162
163 def prepare_data(self):

```

```

161     """
162     Combines the functionalities of add_numeric_attributes and
        encode_categorical_attributes to prepare the data for
        analysis.
163     Returns the prepared DataFrame.
164     """
165     self.add_numeric_attributes()
166     self.encode_categorical_attributes()
167     self.data["X4_ID"] = str(getattr(self, "X4_ID"))
168     return self.data
169
170
171 def plot(self, xlog=False, ylog=False, fig_size=(9,6)):
172     """
173     Creates a plot of the stored data with optional logarithmic
        scaling for x and/or y axes.
174     Parameters:
175     - xlog (bool): Whether to use log scale on the x-axis.
176     - ylog (bool): Whether to use log scale on the y-axis.
177     - fig_size (tuple): Tuple specifying the dimensions of the plot.
178     """
179     if self.data.empty:
180         print("No data to plot")
181     else:
182         # check if the data in the third column are all NaN or 0
183         y_err = self.data.iloc[:,2].isnull().values.all() or self.
            data.iloc[:,2].eq(0).all()
184         # check if the data in the fourth column are all NaN or 0
185         x_err = self.data.iloc[:,3].isnull().values.all() or self.
            data.iloc[:,3].eq(0).all()
186
187         # Get the headers of the data
188         headers = list(self.data.columns.values)
189
190         # Set the size of the plot
191         plt.figure(figsize=fig_size)
192         # Plot the data with a scatter plot
193         # If y_err and x_err are False, plot the data with error
            bars
194         if y_err == False and x_err == False:
195             plt.errorbar(x=self.data[headers[0]], y=self.data[
                headers[1]],
196                         xerr=self.data[headers[3]], yerr=self.data[
                headers[2]],
197                         fmt="o", ecolor="black", capsize=3,
                elinewidth=1, markersize=5)
198         # If y_err is False, plot the data with error bars on the y-
            axis
199         elif y_err == False:
200             plt.errorbar(x=self.data[headers[0]], y=self.data[
                headers[1]], yerr=self.data[headers[2]],
201                         fmt="o", ecolor="black", capsize=3,
                elinewidth=1, markersize=5,)
202         # If x_err is False, plot the data with error bars on the x-
            axis
203         elif x_err == False:
204             plt.errorbar(x=self.data[headers[0]], y=self.data[
                headers[1]], xerr=self.data[headers[3]],

```

```
205         fmt="o", ecol="black", capsize=3,
206             elinewidth=1, markersize=5,)
207     # If y_err and x_err are False, plot the data without error
208     # bars
209     else:
210         sns.scatterplot(x=headers[0], y=headers[1], data=self.
211             data)
212
213     # Set the scale of the plot
214     if xlog == True: plt.xscale("log")
215     if ylog == True: plt.yscale("log")
216
217     # Set the title of the plot
218     plt.title(self.title, fontsize=18, pad=15)
219     # Set the labels of the plot
220     plt.xlabel(headers[0], fontsize=16, labelpad=10)
221     plt.ylabel(headers[1], fontsize=16, labelpad=10)
222     # Set the ticks of the plot
223     plt.xticks(fontsize=14)
224     plt.yticks(fontsize=14)
225     # Show the plot
226     plt.show()
```


B Utility Functions for Proton Experiment Data Analysis

B.1 read_experiment()

```

1  def read_experiment(filename):
2      """
3      This function reads a file specified by "filename", processes its
4      content line by line, and populates the properties
5      of an Experiment object according to specific keywords present in
6      the file. It finally returns this populated
7      Experiment object.
8
9      File Format:
10     The file should contain lines starting with specific keywords,
11     followed by a colon and the value (EXFORFORMATS format).
12     (e.g., "# Target Z: 12")
13
14     Properties of Returned Experiment Object:
15     - title (str): Extracted from the filename.
16     - target_Z (int or None): Atomic number of the target.
17     - target_A (int or None): Mass number of the target.
18     - target_state (int/str): the state of the target.
19     - projectile (str): the projectile used in the experiment.
20     - reaction (str): the reaction that occurred during the experiment.
21     - E_inc (str): the incident energy of the projectile.
22     - final_Z (int/str): the atomic number of the final state.
23     - final_A (int/str): the mass number of the final state.
24     - final_state (int/str): the state of the final state.
25     - MTrat (float/str): the ratio of the number of gamma rays to the
26       number of neutrons.
27     - Ratio_isomer (float/str): the ratio of the number of isomers to
28       the number of ground state.
29     - Quantity (str): the quantity of the experiment.
30     - Frame (str): the frame of the experiment.
31     - MF (int/str): the MF of the experiment.
32     - MT (int/str): the MT of the experiment.
33     - X4_ID (str): the EXFOR ID of the experiment.
34     - X4_code (str): the EXFOR code of the experiment.
35     - author (str): the author of the experiment.
36     - year (int/str): the year of the experiment.
37     - data_points (int/str): the number of data points in the experiment
38       .
39     - data (pd.DataFrame): the data of the experiment.
40     - Reference (str): the reference of the experiment.
41
42     If a property value is not found in the file or cannot be converted
43     to its respective type, the value is set to None.
44
45     Dependencies:
46     - Pandas: Required for storing experimental data in a DataFrame.
47     - NumPy: Required for handling "nan" values.
48
49     Parameters:
50         filename (str): Name or path of the file to read.
51
52     Returns:
53         experiment (Experiment): Populated Experiment object.

```

```

47
48     Raises:
49         FileNotFoundError: If the file specified by "filename" does not
50             exist.
51
52     """
53
54     # Create the experiment
55     experiment = Experiment()
56     experiment.title = filename.split("\\")[-1]
57     print("Reading experiment: ", experiment.title)
58     read_header = False          # Flag to read the header
59     read_ref = False             # Flag to read the reference
60     temp_line = ""              # Temporary line to store the reference
61
62     # Open the file
63     with open(filename, "r") as f:
64         # Read the first line
65         lines = f.readlines()
66         for line in lines:
67
68             if line.startswith("#") and not read_ref and not read_header
69                 :
70                 if line.startswith("# Target Z"):
71                     try:
72                         experiment.target_Z = int(line.split(":")[1].
73                             strip())
74                     except ValueError:
75                         experiment.target_Z = line.split(":")[1].strip()
76                         if line.split(":")[1].strip() != "" else
77                             None
78                 elif line.startswith("# Target A"):
79                     try:
80                         experiment.target_A = int(line.split(":")[1].
81                             strip())
82                     except ValueError:
83                         experiment.target_A = int(line.split(":")[1].
84                             strip()) if line.split(":")[1].strip() != ""
85                             else None
86                 elif line.startswith("# Target state"):
87                     try:
88                         experiment.target_state = int(line.split(":")
89                             [1].strip())
90                     except ValueError:
91                         experiment.target_state = line.split(":")[1].
92                             strip() if line.split(":")[1].strip() != ""
93                             else None
94                 elif line.startswith("# Projectile"):
95                     experiment.projectile = line.split(":")[1].strip()
96                     if line.split(":")[1].strip() != "" else None
97                 elif line.startswith("# Reaction      :"):
98                     #
99                     This includes the ":" because there is another line
100                     with the same beginning in some files
101                     experiment.reaction = line.split(":")[1].strip() if
102                         line.split(":")[1].strip() != "" else None
103                 elif line.startswith("# E-inc"):
104                     experiment.E_inc = line.split(":")[1].strip() if
105                         line.split(":")[1].strip() != "" else None
106                 elif line.startswith("# Final Z"):

```

```

89         try:
90             experiment.final_Z = int(line.split(":")[1].
91                                     strip())
92         except ValueError:
93             experiment.final_Z = line.split(":")[1].strip()
94             if line.split(":")[1].strip() != "" else None
95     elif line.startswith("# Final A"):
96         try:
97             experiment.final_A = int(line.split(":")[1].
98                                     strip())
99         except ValueError:
100             experiment.final_A = line.split(":")[1].strip()
101             if line.split(":")[1].strip() != "" else None
102     elif line.startswith("# Final state"):
103         try:
104             experiment.final_state = int(line.split(":")[1].
105                                         strip())
106         except ValueError:
107             experiment.final_state = line.split(":")[1].
108                                     strip() if line.split(":")[1].strip() != ""
109                                     else None
110     elif line.startswith("# MTrat"):
111         try:
112             experiment.MTrat = float(line.split(":")[1].
113                                     strip())
114         except ValueError:
115             experiment.MTrat = float(line.split(":")[1].
116                                     strip()) if line.split(":")[1].strip() != ""
117                                     else None
118     elif line.startswith("# Ratio isomer"):
119         try:
120             experiment.Ratio_isomer = float(line.split(":")
121                                             [1].strip())
122         except ValueError:
123             experiment.Ratio_isomer = float(line.split(":")
124                                             [1].strip()) if line.split(":")[1].strip() !=
125                                             "" else None
126     elif line.startswith("# Quantity"):
127         experiment.quantity = line.split(":")[1].strip() if
128         line.split(":")[1].strip() != "" else None
129     elif line.startswith("# Frame"):
130         experiment.frame = line.split(":")[1].strip() if
131         line.split(":")[1].strip() != "" else None
132     elif line.startswith("# MF"):
133         try:
134             experiment.MF = int(line.split(":")[1].strip())
135         except ValueError:
136             experiment.MF = int(line.split(":")[1].strip())
137             if line.split(":")[1].strip() != "" else None
138     elif line.startswith("# MT"):
139         try:
140             experiment.MT = int(line.split(":")[1].strip())
141         except ValueError:
142             experiment.MT = int(line.split(":")[1].strip())
143             if line.split(":")[1].strip() != "" else None
144     elif line.startswith("# X4 ID"):
145         experiment.X4_ID = line.split(":")[1].strip() if
146         line.split(":")[1].strip() != "" else None

```

```

129         elif line.startswith("# X4 code"):
130             experiment.X4_code = line.split(":")[1].strip() if
                line.split(":")[1].strip() != "" else None
131         elif line.startswith("# Author"):
132             experiment.author = line.split(":")[1].strip() if
                line.split(":")[1].strip() != "" else None
133         elif line.startswith("# Year"):
134             try:
135                 experiment.year = int(line.split(":")[1].strip()
                    )
136             except ValueError:
137                 experiment.year = int(line.split(":")[1].strip()
                    ) if line.split(":")[1].strip() != "" else
                    None
138         elif line.startswith("# Data points"):
139             try:
140                 experiment.data_points = int(line.split(":")[1].
                    strip())
141             except ValueError:
142                 experiment.data_points = int(line.split(":")[1].
                    strip()) if line.split(":")[1].strip() != ""
                    else None
143             read_header = True
144         elif line.startswith("# Reference"):
145             read_ref = True
146
147     elif read_header:
148         header = line.split()[1:]
149         data_list = [[] for i in range(len(header))]
150         read_header = False
151
152     elif read_ref:
153         # If line only contains "#\n", then it is the end of the
            reference
154         if line == "#\n":
155             # add the reference to the experiment except for the
                last two characters (which are "# ")
156             experiment.reference = temp_line[:-2]
157             read_ref = False           # Reset the flag
158             temp_line = ""            # Reset the temporary line
159             continue
160         temp_line += line[1:]
161
162     else:
163         data = line.split()
164         # Get number of splits
165         n = len(data)
166         data_list[0].append(float(data[0]))
167         data_list[1].append(float(data[1]))
168         # If there are 2 splits, then there are no values for dE
            and dxs
169         if n == 2:
170             data_list[2].append(np.nan)
171             data_list[3].append(np.nan)
172         # If there are 3 splits, then there are no values for dE
173         elif n == 3:
174             data_list[2].append(float(data[2]))
175             data_list[3].append(np.nan)

```

```
176         # If there are 4 splits, then there are values for dE
           and dxs
177     elif n == 4:
178         data_list[2].append(float(data[2]))
179         data_list[3].append(float(data[3]))
180
181     if read_ref:          # This being True means that the file ended
                           directly after the reference
182         experiment.reference = temp_line[:-2]
183
184     # Close the file
185     f.close()
186
187     # Create a data frame using "header" and "E", "xs", "dxs" and "dE"
       lists
188     experiment.data = pd.DataFrame(data_list, index=header).T
189
190     # Return the experiment
191     return experiment
```

B.2 write_experiments_to_binary()

```
1 def write_experiments_to_binary(experiments, filename):
2     """
3     This function writes a list of Experiment objects to a binary file.
4     It uses Python's pickle library for object serialization.
5
6     Parameters:
7         experiments (list): A list of Experiment objects to be written
8         to the file.
9         filename (str): The name of the file where the experiments will
10        be written.
11
12    Returns:
13        None
14
15    Example:
16        experiments = [experiment1, experiment2, experiment3]
17        filename = "experiments.bin"
18        write_experiments_to_binary(experiments, filename)
19
20    Note:
21        The function uses Python's pickle library, so it's essential to
22        ensure that the Experiment objects are picklable.
23        Be aware of the potential security risks if you're unpickling
24        data from an untrusted source.
25    """
26
27    with open(filename, "wb") as f:
28        for experiment in experiments:
29            pickle.dump(experiment, f)
```

B.3 read_experiments_from_binary()

```
1 def read_experiments_from_binary(filename):
2     """
3     This function reads a list of Experiment objects from a binary file
4     and returns them as a list.
5     It uses Python's pickle library for object deserialization.
6
7     Parameters:
8         filename (str): The name of the binary file to read the
9         experiments from.
10
11     Returns:
12         loaded_experiments (list): A list of experiments that were read
13         from the binary file.
14
15     Note:
16         The function uses Python's pickle library, so be cautious of
17         potential security risks if you're unpickling data from an
18         untrusted source.
19     """
20
21     loaded_experiments = []
22
23     # Use a with statement to ensure the file is properly closed after
24     # reading
25     with open(filename, "rb") as f:
26         while True:
27             try:
28                 loaded_experiments.append(pickle.load(f))
29             except EOFError:
30                 break
31             except pickle.PickleError:
32                 print("Error in deserializing object. Skipping...")
33                 break
34
35     return loaded_experiments
```

B.4 write_experiments_to_txt()

```

1  def write_experiments_to_txt(experiments, filename):
2      """
3      write_experiments_to_txt(experiments, filename)
4
5      Serializes a list of experiment objects into a text file, capturing
6      various attributes of each experiment.
7      The function handles empty string attributes by writing them as such
8      in the output file. Once all experiment details are written, the
9      file is closed.
10
11     Parameters:
12         experiments (list): A list of experiment objects.
13
14         filename (str): The name of the output text file where the
15         experiment details will be serialized.
16
17     Returns:
18         None
19
20     Example:
21         experiments = [experiment1, experiment2]
22         filename = "experiments.txt"
23         write_experiments_to_txt(experiments, filename)
24
25     Note:
26         Ensure that each experiment object in the list has all the
27         mentioned attributes. Missing attributes will lead to an
28         error.
29     """
30
31     with open(filename, "w") as f:
32         for experiment in experiments:
33             # Write the header
34             f.write("# Title      : " + experiment.title + "\n")
35             # Write the information. If the information is "", then
36             # write an empty string
37             f.write("# Reaction   : " + experiment.reaction + "\n" if
38                 experiment.reaction is not None else "# Reaction   : \n"
39                 )
40             f.write("# Ratio isomer: " + str(experiment.Ratio_isomer) +
41                 "\n" if experiment.Ratio_isomer is not None else "# Ratio
42                 isomer: \n")
43             f.write("# Quantity   : " + experiment.quantity + "\n" if
44                 experiment.quantity is not None else "# Quantity   : \n"
45                 )
46             f.write("# Frame      : " + experiment.frame + "\n" if
47                 experiment.frame is not None else "# Frame      : \n")
48             f.write("# MF        : " + str(experiment.MF) + "\n" if
49                 experiment.MF is not None else "# MF        : \n")
50             f.write("# MT        : " + str(experiment.MT) + "\n" if
51                 experiment.MT is not None else "# MT        : \n")
52             f.write("# X4 ID     : " + experiment.X4_ID + "\n" if
53                 experiment.X4_ID is not None else "# X4 ID     : \n")
54             f.write("# X4 code    : " + experiment.X4_code + "\n" if
55                 experiment.X4_code is not None else "# X4 code    : \n")

```



```
38     f.write("# Author      : " + experiment.author + "\n" if
        experiment.author is not None else "# Author      : \n")
39     f.write("# Year        : " + str(experiment.year) + "\n" if
        experiment.year is not None else "# Year        : \n")
40     f.write("# Data points : " + str(experiment.data_points) + "
        \n" if experiment.data_points is not None else "# Data
        points : \n")
41     # Write the data
42     f.write(experiment.data.to_string(header=True, index=False))
43     # Write the reference
44     # If the reference is "", then write an empty string
45     f.write("\n# Reference  : \n" + experiment.reference + "\n"
        if experiment.reference is not None else "# Reference
        : \n")
46     # Write a line to separate the reference from the next
        experiment
47     f.write("# END\n")
48
49     f.write("# END OF FILE")
```

B.5 read_experiments_from_txt()

```
1 def read_experiments_from_txt(filename):
2     """
3     Reads a text file containing a structured list of experiments and
4     returns a list of Experiment objects.
5
6     The function expects the file to be structured as follows:
7     - Each experiment block starts with a line for the title ("# Title
8       :") and ends with "# END".
9     - Between these lines, each line should follow the format: "# [
10      field_name] : [field_value]".
11     - After "# END", the data of the experiment is present with the
12       first line as the header.
13     - References start with "# Reference :" and end with "# END".
14     - The file should end with a line containing "# END OF FILE".
15
16     If a field is missing, the corresponding attribute in the Experiment
17     object is set to an empty string ("").
18
19     Parameters:
20     -----
21     filename : str
22         The name of the text file containing the experiments" details.
23
24     Returns:
25     -----
26     experiments : list
27         A list of Experiment objects populated with the details read
28         from the text file.
29
30     Raises:
31     -----
32     ValueError:
33         If the data in the file is not structured as expected.
34
35     Example:
36     -----
37     experiments = read_experiments_from_txt("experiment_data.txt")
38
39     Note:
40     ----
41     The function expects that each Experiment object has attributes
42     corresponding to the fields in the text file.
43     """
44
45     # Create a list of experiments
46     experiments = []
47     # Open the file
48     f = open(filename, "r")
49     # Read the first line
50     line = f.readline()
51     # While you dont read "# END OF FILE" read lines
52     while True:
53         # If the line starts with "# END OF FILE", then you have reached
54         # the end of the file
55         if line.startswith("# END OF FILE"):
56             # Close the file
```

```

49     f.close()
50     # Return the list of experiments
51     return experiments
52 # While you dont read "# END" read lines
53 while not line.startswith("# END"):
54     # Create a new experiment
55     experiment = Experiment()
56     # Read the title exluding any posible "\n" at the end
57     experiment.title = line.split(":")[1][:-1]
58     # Read the rest of the information. If the information is
        empty, then set it to an empty string. Exclude any
        posible "\n" at the end
59     line = f.readline()
60     experiment.reaction = line.split(":")[1][:-1] if line.split(
        ":")[1] != "" else ""
61     line = f.readline()
62     try:
63         experiment.Ratio_isomer = float(line.split(":")[1])
64     except ValueError:
65         experiment.Ratio_isomer = ""
66     line = f.readline()
67     experiment.quantity = line.split(":")[1][:-1] if line.split(
        ":")[1] != "" else ""
68     line = f.readline()
69     experiment.frame = line.split(":")[1][:-1] if line.split(":")
        [1] != "" else ""
70     line = f.readline()
71     try:
72         experiment.MF = float(line.split(":")[1])
73     except ValueError:
74         experiment.MF = ""
75     line = f.readline()
76     try:
77         experiment.MT = float(line.split(":")[1])
78     except ValueError:
79         experiment.MT = ""
80     line = f.readline()
81     experiment.X4_ID = line.split(":")[1][:-1] if line.split(":")
        [1] != "" else ""
82     line = f.readline()
83     experiment.X4_code = line.split(":")[1][:-1] if line.split(":")
        [1] != "" else ""
84     line = f.readline()
85     experiment.author = line.split(":")[1][:-1] if line.split(":")
        [1] != "" else ""
86     line = f.readline()
87     try:
88         experiment.year = int(line.split(":")[1])
89     except ValueError:
90         experiment.year = ""
91     line = f.readline()
92     try:
93         experiment.data_points = int(line.split(":")[1])
94     except ValueError:
95         experiment.data_points = ""
96     # Read the data. All data starts without "#". The first line
        is the header.
97     # Read the header

```

```
98     line = f.readline()
99     header = line.split()
100     # Read the data
101     data = []
102     line = f.readline()
103     while not line.startswith("#"):
104         data.append(line.split())
105         line = f.readline()           # In the last iteration,
                                     # the line is "# Reference :\\n"
106     # Create a dataframe with the data
107     experiment.data = pd.DataFrame(data, columns=header)
108     # Read the references. References start with "# Reference
                                     # :\\n" and end with "# END"
109     line = f.readline()
110     experiment.reference = ""
111     while not line.startswith("# END"):
112         experiment.reference += line
113         line = f.readline()           # In the last iteration,
                                     # the line is "# END"
114     # If the reference is empty, then set it to "". If not, then
                                     # remove the last "\\n"
115     if experiment.reference != "": experiment.reference[:-1]
116     # Add the experiment to the list of experiments
117     experiments.append(experiment)
118     # Read the line "# END OF FILE" or next experiment
119     line = f.readline()
```

B.6 read_proton_experiments_from_exfortables()

```

1 def read_proton_experiments_from_exfortables(path):
2     """
3     Reads all proton experiment files located in a given directory and
4     its subdirectories.
5
6     The function recursively traverses through all directories and
7     subdirectories under the specified path.
8     It reads files that do not have names ending with "list" or "ruth"
9     and returns a list of Experiment objects.
10
11     Parameters:
12     -----
13     path : str
14         The absolute or relative path to the root directory containing
15         proton experiment files.
16
17     Returns:
18     -----
19     experiments : list
20         A list of Experiment objects containing data read from the files
21         under the specified path.
22
23     Raises:
24     -----
25     FileNotFoundError:
26         If the specified path does not exist.
27     ValueError:
28         If the data in the files is not structured as expected.
29
30     Example:
31     -----
32     experiments = read_proton_experiments_from_exfortables("./
33         proton_experiment_data/")
34
35     Notes:
36     -----
37     - The function assumes that each Experiment object has attributes
38       corresponding to the fields in the text files.
39     - The function relies on 'read_experiment()' for reading individual
40       experiment files.
41     """
42
43     # Get name of all directories in the path
44     dirs = [d for d in os.listdir(path) if os.path.isdir(os.path.join(
45         path, d))]
46
47     experiments = []
48
49     # Go through all directories
50     for d in dirs:
51         # Get name of all subdirectories in the path that dont end with
52         # "list"
53         subdirs = [sd for sd in os.listdir(os.path.join(path, d)) if os.
54             path.isdir(os.path.join(path, d, sd)) and not sd.endswith("
55             list")]
56         # Go through all subdirectories

```

```
45     for sd in subdirs:
46         # Check if there is any subsubdirectory inside the
            subdirectory
47         contains_subsubdir = False
48         for item in os.scandir(os.path.join(path, d, sd)):
49             if item.is_dir():
50                 contains_subsubdir = True
51                 break
52         # If there is no subsubdirectory, then the files are in the
            subdirectory
53         if contains_subsubdir:
54             # Get name of all subdirectories in the path
55             subsubdirs = [ssd for ssd in os.listdir(os.path.join(
                path, d, sd)) if os.path.isdir(os.path.join(path, d,
                    sd, ssd))]
56             # Go through all subdirectories
57             for ssd in subsubdirs:
58                 # Get the name of all the files which dont end with
                    "list" or "ruth"
59                 files = [f for f in os.listdir(os.path.join(path, d,
                    sd, ssd)) if os.path.isfile(os.path.join(path, d
                        , sd, ssd, f)) and not f.endswith("list") and not
                            f.endswith("ruth")]
60                 # Go through all files
61                 for f in files:
62                     # Read the experiment
63                     experiment = read_experiment(os.path.join(path,
                        d, sd, ssd, f))
64                     # Add the experiment to the list of experiments
65                     experiments.append(experiment)
66         else:
67             # Get the name of all the files which dont end with "
                list" or "ruth"
68             files = [f for f in os.listdir(os.path.join(path, d, sd)
                ) if os.path.isfile(os.path.join(path, d, sd, f)) and
                    not f.endswith("list") and not f.endswith("ruth")]
69             # Go through all files
70             for f in files:
71                 # Read the experiment
72                 experiment = read_experiment(os.path.join(path, d,
                    sd, f))
73                 # Add the experiment to the list of experiments
74                 experiments.append(experiment)
75     return experiments
```

B.7 filter_experiments()

```

1 def filter_experiments(experiments, attribute, value):
2     """
3     Filters a list of Experiment objects based on the specified
4         attribute and its value.
5
6     The function evaluates the attribute for each Experiment object in
7     the list, keeping only
8     those objects where the attribute matches the given value. If the
9     attribute does not exist,
10    it will return None and print a list of available attributes.
11
12    Parameters:
13    -----
14    experiments : list
15        A list of Experiment objects to filter.
16    attribute : str
17        The attribute name based on which the filtering is to be done.
18    value : str | int | float
19        The value of the attribute for filtering the list of Experiment
20        objects.
21
22    Returns:
23    -----
24    filtered_list : list | None
25        A list of Experiment objects that have the specified attribute
26        and value.
27        Returns None if the attribute does not exist or no experiments
28        meet the condition.
29
30    Raises:
31    -----
32    AttributeError:
33        If the attribute does not exist in the Experiment objects.
34
35    Example:
36    -----
37    filtered_exp = filter_experiments(experiment_list, "year", 2021)
38
39    Notes:
40    -----
41    - Make sure the attribute exists in the Experiment object.
42    - This function is case-sensitive for string attributes.
43    """
44
45    # Check if the attribute exists
46    if not hasattr(experiments[0], attribute):
47        print("Attribute {} does not exist\n".format(attribute))
48        # Print the list of available attributes
49        print("Available attributes:")
50        print([attr for attr in dir(experiments[0]) if not attr.
51              startswith("__")])
52        return None
53
54    # Filter the list of experiments to get only the ones with the given
55    # attribute and value
56    filtered_list = [obj for obj in experiments if getattr(obj,
57        attribute) == value]

```

```
48
49     if len(filtered_list) == 0:
50         print("No experiments with {} = {}".format(attribute, value))
51         # Print the list of available values for the given attribute
52         print("Available values for {}".format(attribute))
53         print(list(np.unique([getattr(obj, attribute) for obj in
54                               experiments])))
54         return None
55
56     else:
57         print("{} experiments with {} = {}".format(len(filtered_list),
58                                                     attribute, value))
59
60     # Return the filtered list
61     return filtered_list
```


B.8 get_unique_values()

```

1 def get_unique_values(experiments, attribute):
2     """
3     Retrieves the distinct values for a specified attribute from a list
4     of Experiment objects.
5
6     The function examines each Experiment object in the list to identify
7     all the unique values
8     of the specified attribute. If the attribute is not present, it
9     returns None and displays
10    a list of available attributes for the Experiment objects.
11
12    Parameters:
13    -----
14    experiments : list
15        A list of Experiment objects from which to extract unique values
16        of the attribute.
17    attribute : str
18        The name of the attribute whose unique values are to be
19        determined.
20
21    Returns:
22    -----
23    unique_values : list | None
24        A list of unique values for the given attribute.
25        Returns None if the attribute does not exist within the
26        Experiment objects.
27
28    Raises:
29    -----
30    AttributeError:
31        If the attribute does not exist in the Experiment objects.
32
33    Example:
34    -----
35    unique_years = get_unique_values(experiment_list, "year")
36
37    Notes:
38    -----
39    - Ensure the attribute exists in the Experiment objects for accurate
40      results.
41    - The function employs the numpy "unique" function for identifying
42      unique values.
43    """
44
45    # Check if the attribute exists
46    if not hasattr(experiments[0], attribute):
47        print("Attribute {} does not exist\n".format(attribute))
48        # Print the list of available attributes
49        print("Available attributes:")
50        print([attr for attr in dir(experiments[0]) if not attr.
51              startswith("__")])
52        return None
53
54    # Get the unique values of the attribute
55    unique_values = list(np.unique([getattr(obj, attribute) for obj in
56                                  experiments]))
57    # Return the unique values

```

47

```
return unique_values
```

B.9 clean_dataframe()

```

1 def clean_dataframe(df, uncertainties=False):
2     """
3     Cleans a DataFrame by removing redundant columns for analysis.
4
5     This function eliminates all columns with only one unique value and,
6     if the 'uncertainties' flag is False, all columns starting with "d".
7
8     Parameters:
9     -----
10    df : pd.DataFrame
11        The DataFrame to be cleaned.
12    uncertainties : bool, optional
13        Whether to retain columns that start with "d". Default is False.
14
15    Returns:
16    -----
17    pd.DataFrame
18        A DataFrame with the redundant columns removed.
19
20    Example:
21    -----
22    cleaned_df = clean_dataframe(original_df)
23
24    Notes:
25    -----
26    - The function relies on pandas for DataFrame operations.
27    """
28    # Drop the columns which heading starts with "d"
29    if not uncertainties: df = df.drop([col for col in df.columns if col
30        .startswith("d")], axis=1)
31
32    # Get the number of unique values per column
33    unique_values = df.apply(lambda x: len(x.unique()))
34    # Get the columns with only one unique value
35    columns_to_drop = unique_values[unique_values == 1].index
36    # Drop the columns
37    df = df.drop(columns_to_drop, axis=1)
38    # Return the cleaned dataframe
39    return df

```

B.10 `classify_experiments()`

```
1 def classify_experiments(experiments, attribute):
2     """
3     Classifies a list of Experiment objects based on a specified
4     attribute.
5
6     The function generates individual DataFrames for each unique value
7     of the given attribute
8     and writes these as CSV files. Finally, it returns a dictionary
9     containing all these
10    DataFrames.
11
12    Parameters:
13    -----
14    experiments : list
15        A list of Experiment objects to be classified.
16    attribute : str
17        The attribute based on which the classification is performed.
18
19    Returns:
20    -----
21    dfs : dict of pd.DataFrame
22        A dictionary with keys as the unique values of the attribute and
23        the corresponding
24        DataFrames as values.
25
26    Example:
27    -----
28    classified_dfs = classify_experiments(experiment_list, "year")
29
30    Notes:
31    -----
32    - The function employs the 'get_unique_values' function to identify
33      unique attribute values.
34    - The output DataFrames are written as CSV files in a directory
35      corresponding to the attribute name.
36    """
37    # Get the unique values of the attribute
38    unique_values = get_unique_values(experiments, attribute)
39
40    # Create an empty dataframe per unique value of the attribute
41    dfs = {}
42    for value in unique_values:
43        dfs[value] = pd.DataFrame()
44
45    # Go through all experiments
46    for experiment in experiments:
47        # Get the value of the attribute
48        value = getattr(experiment, attribute)
49        # Add the experiment to the dataframe with the corresponding
50        # value of the attribute
51        dfs[value] = pd.concat([dfs[value], experiment.to_dataframe()]).
52            reset_index(drop=True)
53
54    # Create a new csv file for each dataframe
55    for value in unique_values:
```

```
48     dfs[value].to_csv("EXFOR_ProtonReactions_Classssified_by_{}/{ }_
49                        {}.csv".format(attribute, attribute, value))
50
51     # Return the dataframe
52     return dfs
```

B.11 `classify_experiments_by_data()`

```
1 def classify_experiments_by_data(experiments):
2     """
3     Classifies a list of Experiment objects based on their column
4     headers after data preparation.
5
6     The function processes each Experiment object, extracts its
7     dataframe representation using
8     the 'prepare_data' method, and then groups these dataframes by their
9     column headers.
10    Each grouped dataframe is then saved as a CSV file. Finally, it
11    prints the summary of the
12    classification process.
13
14    Parameters:
15    -----
16    experiments : list
17        A list of Experiment objects to be classified.
18
19    Returns:
20    -----
21    None
22
23    Example:
24    -----
25    classify_experiments_by_data(experiment_list)
26
27    Notes:
28    -----
29    - Assumes each Experiment object has a 'prepare_data' method that
30      returns a DataFrame.
31    - The output grouped DataFrames are written as CSV files with a
32      specific naming convention.
33    """
34    # Dictionary to store dataframes grouped by their column headers
35    grouped_dataframes = {}
36
37    total_experiments = len(experiments)
38    print(f"Processing {total_experiments} experiments...")
39
40    for idx, experiment in enumerate(experiments, start=1):
41        print(f"Processing experiment {idx} out of {total_experiments}
42              {...}")
43
44        # Prepare the data
45        df = experiment.prepare_data()
46
47        # Create a key based on the column headers
48        key = tuple(df.columns)
49
50        # If the key already exists in the dictionary, concatenate the
51        # dataframe. Otherwise, create a new entry.
52        if key in grouped_dataframes:
53            grouped_dataframes[key] = pd.concat([grouped_dataframes[key],
54                                                  df], axis=0)
55        else:
56            grouped_dataframes[key] = df
```

```
48
49 # Display the detected groups of headers on the screen
50 group_number = 1
51 for key in grouped_dataframes.keys():
52     print(f"\nGroup {group_number}:")
53     print(", ".join(key))
54     print("="*50)
55     group_number += 1
56
57 # Save each grouped dataframe to a CSV file
58 total_groups = len(grouped_dataframes)
59 print(f"\nSaving {total_groups} grouped dataframes to CSV files...")
60 for idx, (key, df) in enumerate(grouped_dataframes.items(), start=1)
61 :
62     filename = f"EXFOR_ProtonReactions_Classified_Group_{idx}.csv"
63     df.to_csv(filename, index=False)
64     print(f"Group {idx}'s dataframe saved as {filename}")
65
66 print(f"\nFinished! {len(grouped_dataframes)} groups of experiments
67 found based on column headers.")
```

B.12 plot_experiments()

```
1 def plot_experiments(experiments, xlog=False, ylog=False, fig_size=(9,6))
2     :
3     """
4     Plots a list of experiments using matplotlib and seaborn libraries.
5
6     The function accepts a list of experiment objects and plots them on
7     a 2D graph based on their data attributes. It can plot with or without error bars along
8     the x and y axes.
9     Additionally, it supports logarithmic scaling for both axes and
10    allows figure size customization.
11
12    Parameters:
13    -----
14    experiments : list
15        A list of experiment objects to be plotted. Each object should
16        have a "data" attribute
17        which is a pandas DataFrame containing the relevant data.
18    xlog : bool, optional
19        A flag to specify if the x-axis should be logarithmic. Default
20        is False.
21    ylog : bool, optional
22        A flag to specify if the y-axis should be logarithmic. Default
23        is False.
24    fig_size : tuple, optional
25        The dimensions of the figure to be plotted. Default is (9, 6).
26
27    Returns:
28    -----
29    None :
30        The function does not return anything; it generates and displays
31        the plot.
32
33    Example:
34    -----
35    plot_experiments([experiment1, experiment2], xlog=True, ylog=False,
36                      fig_size=(10,8))
37
38    Notes:
39    -----
40    - The function first validates that all experiments have the same
41      DataFrame columns.
42    - The plotting accounts for possible error attributes in the x and y
43      axes.
44    - The resulting plot can be customized through optional parameters
45      like 'xlog', 'ylog', and 'fig_size'.
46    """
47
48    # First, check if there are experiments to plot
49    if not experiments:
50        print("No experiments to plot")
51        return
52
53    # Then, verify that all dataframes have the same headers
54    headers = list(experiments[0].data.columns.values)
55    for experiment in experiments[1:]:
```



```

45     if list(experiment.data.columns.values) != headers:
46         print("Mismatch in headers between experiments.")
47         return
48
49     # Now proceed to set up and plot each experiment
50     plt.figure(figsize=fig_size)
51     for experiment in experiments:
52         data = experiment.data
53         y_err = data.iloc[:,2].isnull().values.all() or data.iloc[:,2].
54             eq(0).all()
55         x_err = data.iloc[:,3].isnull().values.all() or data.iloc[:,3].
56             eq(0).all()
57
58         # Depending on y_err and x_err, we decide how to plot
59         if y_err == False and x_err == False:
60             plt.errorbar(x=data[headers[0]], y=data[headers[1]],
61                         xerr=data[headers[3]], yerr=data[headers[2]],
62                         fmt="o", ecolor="black", capsize=3, elinewidth
63                             =1, markersize=5, label=experiment.X4_ID)
64         elif y_err == False:
65             plt.errorbar(x=data[headers[0]], y=data[headers[1]], yerr=
66                 data[headers[2]],
67                 fmt="o", ecolor="black", capsize=3, elinewidth
68                     =1, markersize=5, label=experiment.X4_ID)
69         elif x_err == False:
70             plt.errorbar(x=data[headers[0]], y=data[headers[1]], xerr=
71                 data[headers[3]],
72                 fmt="o", ecolor="black", capsize=3, elinewidth
73                     =1, markersize=5, label=experiment.X4_ID)
74         else:
75             sns.scatterplot(x=headers[0], y=headers[1], data=data, label
76                             =experiment.X4_ID)
77
78     # Set scale, labels, title, and display the plot
79     if xlog: plt.xscale("log")
80     if ylog: plt.yscale("log")
81
82     plt.xlabel(headers[0], fontsize=16, labelpad=10)
83     plt.ylabel(headers[1], fontsize=16, labelpad=10)
84     plt.xticks(fontsize=14)
85     plt.yticks(fontsize=14)
86     plt.legend()
87     plt.show()

```

B.13 plot_outliers()

```
1 def plot_outliers(outliers_df, experiments, xlog=False, ylog=False,
2   fig_size=(9,6)):
3     """
4     Plots outliers along with their corresponding experiments using
5     matplotlib and seaborn libraries.
6
7     The function accepts a DataFrame containing outlier data and a list
8     of experiment objects.
9     It plots these outliers and the normal data on a 2D graph. Both
10    outlier and normal data points
11    can be visualized with or without error bars along the x and y axes.
12    Logarithmic scaling for both axes
13    and figure size customization are also supported.
14
15    Parameters:
16    -----
17    outliers_df : pd.DataFrame
18        A DataFrame containing the outlier data points. Should have an "
19        X4_ID" column which matches the experiment IDs.
20    experiments : list
21        A list of experiment objects to be plotted. Each object should
22        have a "data" attribute which is a pandas DataFrame.
23    xlog : bool, optional
24        A flag to specify if the x-axis should be logarithmic. Default
25        is False.
26    ylog : bool, optional
27        A flag to specify if the y-axis should be logarithmic. Default
28        is False.
29    fig_size : tuple, optional
30        The dimensions of the figure to be plotted. Default is (9, 6).
31
32    Returns:
33    -----
34    None :
35        The function does not return anything; it generates and displays
36        the plot.
37
38    Example:
39    -----
40    plot_outliers(outliers_dataframe, [experiment1, experiment2], xlog=
41        True, ylog=False, fig_size=(10,8))
42
43    Notes:
44    -----
45    - The function first identifies the DataFrame columns that are
46      specific to outliers, not present in normal data.
47    - It then groups outliers based on these specific columns before
48      plotting.
49    - Normal data and outliers from the same experiment are plotted
50      together for better visualization.
51    """
52
53    # Get the column names from the "data" attribute of a corresponding
54    # experiment
55    example_experiment = next(experiment for experiment in experiments
56                              if experiment.X4_ID in outliers_df["X4_ID"].tolist())
```

```

41 data_columns = example_experiment.data.columns.values.tolist()
42
43 # Get columns from outliers_df that are not in "data" and are also
   not "X4_ID"
44 groupby_columns = [col for col in outliers_df.columns if col not in
   data_columns and col != "X4_ID"]
45
46 grouped_outliers = outliers_df.groupby(groupby_columns)
47
48 print("Grouped outliers according to columns: {}".format(
   groupby_columns))
49
50 for _, group in grouped_outliers:
51     plt.figure(figsize=fig_size)
52     x4_ids = group["X4_ID"].tolist()
53
54     # Add a legend entry for the outliers without actually plotting
   them
55     plt.scatter([], [], color="red", s=50, label="Outliers", marker=
   "x")
56
57     for experiment in experiments:
58         if experiment.X4_ID in x4_ids:
59             data = experiment.data
60             headers = list(data.columns.values)
61
62             y_err = not (data.iloc[:,2].isnull().values.all() or
   data.iloc[:,2].eq(0).all())
63             x_err = not (data.iloc[:,3].isnull().values.all() or
   data.iloc[:,3].eq(0).all())
64
65             # Plot normal data
66             if y_err and x_err:
67                 plt.errorbar(x=data[headers[0]], y=data[headers[1]],
   xerr=data[headers[3]], yerr=data[headers
68                 [2]],
69                 fmt="o", ecolor="black", capsize=3,
   elinewidth=1, markersize=5, label=
   experiment.X4_ID)
70             elif y_err:
71                 plt.errorbar(x=data[headers[0]], y=data[headers[1]],
   yerr=data[headers[2]],
72                 fmt="o", ecolor="black", capsize=3,
   elinewidth=1, markersize=5, label=
   experiment.X4_ID)
73             elif x_err:
74                 plt.errorbar(x=data[headers[0]], y=data[headers[1]],
   xerr=data[headers[3]],
75                 fmt="o", ecolor="black", capsize=3,
   elinewidth=1, markersize=5, label=
   experiment.X4_ID)
76             else:
77                 plt.scatter(x=data[headers[0]], y=data[headers[1]],
   s=25, label=experiment.X4_ID)
78
79     # Actually plot the outliers
80     for experiment in experiments:
81         if experiment.X4_ID in x4_ids:

```

```
82         specific_outliers = outliers_df[outliers_df["X4_ID"] ==  
83             experiment.X4_ID]  
84         plt.scatter(specific_outliers[headers[0]],  
85             specific_outliers[headers[1]], color="red", s=50,  
86             zorder=3, marker="x")  
87  
88     if xlog: plt.xscale("log")  
89     if ylog: plt.yscale("log")  
90     plt.xlabel(headers[0], fontsize=16, labelpad=10)  
91     plt.ylabel(headers[1], fontsize=16, labelpad=10)  
92     plt.xticks(fontsize=14)  
93     plt.yticks(fontsize=14)  
94     plt.legend()  
95     plt.show()
```

C Interquartile Range Functions

C.1 detect_outliers_IQR()

```

1  def detect_outliers_IQR(grouped_df, experiment, limit=1.5, uncertainties
    =False):
2      """
3      Identifies outliers in the data of an Experiment object using the
        Interquartile Range (IQR) method.
4
5      The function takes a grouped dataframe and an Experiment object,
        then calculates the first (Q1) and third (Q3) quartiles
6      along with the interquartile range (IQR) for specified columns.
        Outliers are determined based on these statistics.
7
8      Parameters:
9      -----
10     grouped_df : DataFrame
11         The grouped dataframe containing the data.
12     experiment : object
13         An Experiment object containing the data and metadata.
14     limit : float, optional
15         The limit factor to multiply with the IQR to determine the range
        for outliers (default is 1.5).
16     uncertainties : bool, optional
17         Whether to consider uncertainties in the data (default is False)
        .
18
19     Returns:
20     -----
21     DataFrame
22         A dataframe containing only the outliers.
23
24     Example:
25     -----
26     detect_outliers_IQR(grouped_df, experiment_object)
27
28     Notes:
29     -----
30     - Assumes the Experiment object contains a dataframe representation
        in its 'data' attribute.
31     - If uncertainties are considered, it assumes the Experiment object"
        s dataframe includes columns for uncertainties.
32     """
33
34     if not uncertainties:
35         # Get the names of the first two columns
36         if experiment.data.columns[0] != "Z":
37             first_column = experiment.data.columns[0]
38             second_column = experiment.data.columns[1]
39         else:
40             first_column = experiment.data.columns[2]
41             second_column = experiment.data.columns[3]
42
43     Q1 = grouped_df[[first_column, second_column]].quantile(0.25)
44     Q3 = grouped_df[[first_column, second_column]].quantile(0.75)
45     IQR = Q3 - Q1

```

```

46     # Here is where we define the condition for outliers
47     outlier_condition = ((grouped_df[[first_column, second_column]]
48         < (Q1 - limit * IQR)) |
49         (grouped_df[[first_column, second_column]] >
50             (Q3 + limit * IQR)))
51
52     return grouped_df[outlier_condition.any(axis=1)]
53
54 else:
55     if experiment.data.columns[0] != "Z":
56         first_column = experiment.data.columns[0]
57         second_column = experiment.data.columns[1]
58         if experiment.data[experiment.data.columns[2]].isnull().
59             values.any():
60             Q1 = grouped_df[second_column].quantile(0.25)
61             Q3 = grouped_df[second_column].quantile(0.75)
62             IQR = Q3 - Q1
63
64             # Here is where we define the condition for outliers
65             outlier_condition = ((grouped_df[second_column] < (Q1 -
66                 limit * IQR)) |
67                 (grouped_df[second_column] > (Q3 +
68                     limit * IQR)))
69
70             return grouped_df[outlier_condition]
71         else:
72             third_column = experiment.data.columns[2]
73
74     else:
75         first_column = experiment.data.columns[2]
76         second_column = experiment.data.columns[3]
77         third_column = experiment.data.columns[4]
78
79     # Calculate the bounds using IQR
80     Q1 = grouped_df[second_column].quantile(0.25)
81     Q3 = grouped_df[second_column].quantile(0.75)
82     IQR = Q3 - Q1
83
84     # Convert lower_bound and upper_bound to Series
85     lower_bound = (Q1 - limit * IQR)
86     upper_bound = (Q3 + limit * IQR)
87
88     # Check if both ends of the uncertainty range are outside the bounds
89     lower_values = grouped_df[second_column] - grouped_df[third_column]
90     upper_values = grouped_df[second_column] + grouped_df[third_column]
91
92     # Check if either end of the uncertainty range crosses the bounds
93     outliers_condition = ((lower_values < lower_bound) & (upper_values <
94         lower_bound)) | ((upper_values > upper_bound) & (lower_values >
95         upper_bound))
96
97     # Return the points that meet the condition of being outliers
98     return grouped_df[outliers_condition]

```

C.2 IQR_method()

```

1  def IQR_method(df, experiments, min_observations=20, uncertainties=False
2  ):
3      """
4      Applies the IQR method to identify outliers from a DataFrame using
5      data from a list of Experiment objects.
6
7      The function first identifies columns by which to group the data. It
8      then filters the groups based on the
9      minimum number of observations specified. For each valid group, it
10     applies the 'detect_outliers_IQR' function
11     to identify outliers.
12
13     Parameters:
14     -----
15     df : DataFrame
16         The input DataFrame containing the data to be checked for
17         outliers.
18     experiments : list
19         A list of Experiment objects containing relevant data and
20         metadata.
21     min_observations : int, optional
22         The minimum number of observations required to consider a group
23         for outlier detection (default is 20).
24     uncertainties : bool, optional
25         Whether to consider uncertainties in the data (default is False)
26         .
27
28     Returns:
29     -----
30     DataFrame
31         A DataFrame containing only the outliers.
32
33     Example:
34     -----
35     IQR_method(df, experiment_list)
36
37     Notes:
38     -----
39     - Assumes the Experiment object contains a dataframe representation
40       in its 'data' attribute.
41     - If uncertainties are considered, it assumes the Experiment object"
42       s dataframe includes columns for uncertainties.
43     """
44
45     # Get the name of the "data" columns from a corresponding experiment
46     example_exp = next(experiment for experiment in experiments if
47                        experiment.X4_ID in df["X4_ID"].tolist())
48     data_columns = example_exp.data.columns.values.tolist()
49
50     # Get columns from outliers_df that are not in data and also not
51     X4_ID
52     groupby_columns = [col for col in df.columns if col not in
53                       data_columns and col != "X4_ID"]
54
55     # Group by columns and perform initial size count

```

```
43 grouped = df.groupby(groupby_columns).size().reset_index(name="count")
44
45 # Filter groups by minimum number of observations
46 valid_groups = grouped[grouped["count"] >= min_observations][
47     groupby_columns].to_dict("records")
48 filtered_df = df[df.set_index(groupby_columns).index.isin([tuple(d.
49     values()) for d in valid_groups])]
50
51 print("Number of groups with at least {} observations: {}".format(
52     min_observations, len(valid_groups)))
53
54 print("Calculating outliers {} uncertainties...\n".format("with" if
55     uncertainties else "without"))
56
57 # Detect outliers
58 outliers_df = (filtered_df.groupby(groupby_columns, group_keys=False)
59     )
60     .apply(lambda x: detect_outliers_IQR(x, example_exp,
61         uncertainties=uncertainties))
62     .reset_index(drop=True))
63
64 print("Percentage of outliers: {:.2f}%".format(len(outliers_df) /
65     len(filtered_df) * 100))
66
67 return outliers_df
```


D Classification Data Groups

In the development of a machine learning application for analyzing proton-induced reactions from the EXFOR library, data is classified into eight distinct groups. Each group is characterized by specific measured parameters that serve as features for the machine learning model. The classification is as follows:

- **Group 1: Differential Cross-Section by Angle**
 - **Angle(deg)**: The angle in degrees at which the reaction occurs.
 - **xs(mb/sr)**: The differential cross-section in millibarns per steradian.
 - **dxs(mb/sr)**: The uncertainty in the differential cross-section.
 - **dAngle(deg)**: The uncertainty in the angle.
- **Group 2: Energy Outgoing and Cross-Section**
 - **E-out(MeV)**: The outgoing energy in mega electronvolts.
 - **xs(mb/MeV.sr)**: The differential cross-section in millibarns per MeV per steradian.
 - **dxs(mb/MeV.sr)**: The uncertainty in the differential cross-section.
 - **dE(MeV)**: The uncertainty in the outgoing energy.
- **Group 3: Energy and Ratio**
 - **E(MeV)**: The energy in mega electronvolts.
 - **ratio**: The ratio of the observed value to the expected value.
 - **dratio**: The uncertainty in the ratio.
 - **dE(MeV)**: The uncertainty in the energy.
- **Group 4: Energy and Total Cross-Section**
 - **E(MeV)**: The energy in mega electronvolts.
 - **xs(mb)**: The total cross-section in millibarns.
 - **dxs(mb)**: The uncertainty in the total cross-section.
 - **dE(MeV)**: The uncertainty in the energy.
- **Group 5: Energy Outgoing and Cross-Section per MeV**
 - **E-out(MeV)**: The outgoing energy in mega electronvolts.
 - **xs(mb/MeV)**: The differential cross-section in millibarns per MeV.
 - **dxs(mb/MeV)**: The uncertainty in the differential cross-section.
 - **dE(MeV)**: The uncertainty in the outgoing energy.
- **Group 6: Energy in eV and Parity**
 - **E(eV)**: The energy in electronvolts.
 - **par**: Refers to a parameter, such as "resonance width" (WID).
 - **dpar**: The uncertainty in the parameter.
 - **dE(eV)**: The uncertainty in the energy.
- **Group 7: Atomic and Mass Numbers with Yield**

- **Z**: The atomic number.
- **A**: The mass number.
- **Yield**: The yield of the reaction.
- **dYield**: The uncertainty in the yield.

- **Group 8: Energy and Average Number of Neutrons Emitted**

- **E(MeV)**: The energy in mega electronvolts.
- **nubar**: The average number of neutrons emitted.
- **dnubar**: The uncertainty in the average number of neutrons emitted.
- **dE(MeV)**: The uncertainty in the energy.

Each of these groups serves a specific purpose in the analysis and provides a unique set of features for the machine learning model to learn from. By classifying the data into these groups, the model can be more effectively trained to make accurate predictions for various types of proton-induced reactions.



POLITÉCNICA

**ESCUELA TÉCNICA SUPERIOR DE INGENIEROS INDUSTRIALES
UNIVERSIDAD POLITÉCNICA DE MADRID**

José Gutiérrez Abascal, 2. 28006 Madrid
Tel.: 91 336 3060
info.industriales@upm.es

www.industriales.upm.es