



COLEGIO DE CIENCIAS E INGENIERÍAS
INGENIERÍA EN CIENCIAS DE LA COMPUTACIÓN

Entregable 1 del Proyecto Integrador

Tutor: Ricardo Flores Moyano

Autor: John Ochoa Abad

Quito – Ecuador
2026



1. Título del Proyecto:

Arquitectura Mixture of Experts para Clasificación de Ataques en Datasets con Desbalance Extremo de Clases

2. Resumen de actividades realizadas:

1. Revisión del estado del arte: Se realizó una revisión de las investigaciones principales relacionadas al proyecto de forma general a manera de estar al tanto de la información concerniente al estado del arte. A continuación, se listan los principales documentos revisados:
 - “Fusion of statistical importance for feature selection in deep neural network-based intrusion detection system” por A. Thakkar y R. Lohiya. doi: 10.1016/j.inffus.2022.09.026.
 - “A cognitive security framework for detecting intrusions in IoT and 5G utilizing deep learning” por U. K. Lilhore, S. Dalal, and S. Simaiya. doi: 10.1016/j.cose.2023.103560.
 - “Early network intrusion detection enabled by attention mechanisms and RNNs” por T. E. T. Djaidja, B. Brik, S. M. Senouci, A. Boualouache, and Y. Ghamri-Doudane. doi: 10.1109/TIFS.2024.3441862
 - “Convolutional neural networks and mixture of experts for intrusion detection in 5G networks and beyond” por L. Ilias, G. Doukas, V. Lamprou, C. Ntanos, and D. Askounis. doi: 10.3389/frai.2025.1708953
2. Lectura de documentación respectiva a manejo de datasets desbalanceados: Con el objetivo de conocer más respecto al procesamiento de datasets desbalanceados se efectuó una lectura de documentos e investigaciones recientes para saber cómo llevar a cabo un procesamiento correcto y justificado del dataset en intrusiones cibernéticas. Dichos documentos son:
 - “Potential Anchoring for imbalanced data classification” por M. Koziarski. doi: 10.1016/j.patcog.2021.108114
 - “Imbalanced data problem in machine learning: A review” por M. Altalhan, A. Algarni, and M. Turki-Hadj Alouane. doi: 10.1109/ACCESS.2025.3531662
 - “Machine Learning for Imbalanced Data: Tackle imbalanced datasets using machine learning and deep learning techniques” por Kumar Abhishek y Mounir Abdelaziz. Packet Publishing, 2023
 - “A Review of Unlabeled and Imbalanced Data Challenges in Machine Learning: Strategies and Solutions” por N. MS. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 2025
3. Obtención de dataset de intrusiones: Se descargó el dataset de intrusiones “Intrusion Detection Evaluation Dataset (CIC-IDS2017)” para el proyecto A la par, se llevó a cabo una lectura de la información presentada por la University of New Brunswick respecto al dataset, con el fin de comprender el proceso de captura del tráfico en un entorno controlado, los escenarios de ataque simulados y las características extraídas mediante CICFlowMeter. A la par, se realizó una investigación para tener más información respecto a los tipos de intrusiones de los labels. Se seleccionó el CIC-IDS2017 porque contiene tráfico benigno y múltiples tipos de ataques etiquetados a nivel de flujo, incluyendo alrededor de 80 features estadísticas y temporales, lo cual



permitió establecer una base sólida para el análisis exploratorio y el posterior procesamiento del dataset. La fuente de este es: <https://www.unb.ca/cic/datasets/ids-2017.html>

4. Generación de repositorio para proyecto: Se generó el siguiente repositorio en Github para almacenar el código y documentos pertinentes al desarrollo del proyecto con el fin de tener un historial de los avances de este. El enlace del repositorio de Github que he creado es el siguiente: <https://github.com/johnnyredwood/ArquitecturaMoEClasificacionIntrusiones>
5. Caracterización del dataset: Se procedió a cargar el dataset y se efectuó un análisis estructural inicial basado en identificar dimensiones, tipos de datos, estadísticas descriptivas y distribución de clases con el objetivo de comprender la naturaleza del problema y el formato de los datos. A su vez, se estudiaron las features con el objetivo de entender la definición de cada característica provista de modo que se pueda reconocer la importancia de cada una de ellas respecto al label de clasificación de intrusiones.
6. Gestión de accesos para servidor de IA de la Universidad: Se gestionó el acceso a la VPN de la Universidad San Francisco de Quito y a la par se dio de alta un contenedor y un usuario personal en el servidor H200 para el desarrollo de mi proyecto con el objetivo de poder aprovechar los recursos de dicho servidor de IA. De igual forma se generó la conexión con dicho servidor mediante SSH desde VS Code a manera de vincular el IDE con el servidor.
7. Análisis de desbalance de clases: Se evaluó el grado de desbalance mediante el cálculo del Imbalance Ratio y la entropía normalizada de clases, confirmando la presencia de una distribución no uniforme entre tráfico benigno y ataques. Se generaron visualizaciones para analizar la proporción de cada clase.
8. Limpieza y depuración de datos: Se identificaron valores negativos ilógicos en ciertas features relacionadas con métricas de red, eliminando columnas problemáticas y tratando valores inconsistentes según su impacto por clase. Asimismo, se eliminaron muestras con valores duplicados y negativos en la clase mayoritaria para reducir redundancia sin afectar la representatividad de las clases minoritarias.
9. Tratamiento de valores faltantes e infinitos: Se detectaron valores infinitos derivados de cálculos con divisores pequeños, los cuales fueron convertidos a datos tipo NaN. Posteriormente, se aplicó imputación por clase utilizando Iterative Imputer, comparando dos enfoques: Bayesian Ridge y Random Forest Regressor. Se seleccionó Bayesian Ridge por su estabilidad matemática y comportamiento consistente en las distribuciones imputadas.
10. Separación estratificada, codificación y escalamiento: Se realizó una división estratificada en conjuntos de entrenamiento, validación y prueba, preservando la distribución original de clases. Posteriormente, se codificó la variable objetivo mediante Label Encoding para su uso en modelos supervisados. Por último, se aplicó el Standard Scaler para escalar los datos a manera de apoyar con esto a la convergencia de los modelos.
11. Estrategias de oversampling: Se implementaron técnicas avanzadas de oversampling como ADASYN y SMOTE, incluyendo variantes como KMeansSMOTE, incorporando análisis de artefactos mediante Local Outlier Factor y evaluación de distancias a vecinos más cercanos. A su vez, se definió una estrategia automática de aplicación de SMOTE basada en el IR por clase, evitando su uso en clases extremadamente raras para reducir riesgo de sobreajuste y contaminación sintética.
12. Feature Selection Supervisada: Se definió e implementó el uso del método de feature selection supervisado Mutual Information debido a que se considera estándar para datos desequilibrados porque no asume una relación lineal midiendo cuánta información aporta una característica



13. Definición de estructura de pipelines para entrenamiento de modelos baseline: Se definió el siguiente pipeline para el proceso de entrenamiento de los modelos: Imputación por clase, Escalado, Oversampling, Feature Selection Supervisado, Modelo de clasificación
14. Generación inicial de código para entrenamiento y evaluación de modelos baseline: Se generó el código para entrenar los modelos Random Forest, SVM y XGBoost sobre el dataset procesado incluyendo búsqueda de hiperparámetros, valoración en set de validación, función para obtención de métricas de clasificación orientadas a clases minoritarias e inclusión de pesos por clase en modelos para ajustar las funciones de pérdida de manera inversamente proporcional a la frecuencia de cada clase con base en distribución del dataset. Con esto se efectuaron primeros entrenamientos.
15. Redacción de documento final: Los avances relacionados al proyecto indicados en los puntos anteriores fueron incluidos en el documento final acorde a lo indicado en la siguiente sección

3. Secciones o capítulos del documento final desarrollados

1. Portada
2. Derechos de autor
3. Tabla de contenidos
4. Introducción
5. Estado del arte
6. Desarrollo del Proyecto: Parte 1. Data Engineering

4. Revisión y firma del tutor del proyecto

Yo, Ricardo Flores Moyano, profesor de la carrera de Ingeniería en Ciencias de la Computación, hago constar que he revisado y, por lo tanto, apruebo las actividades realizadas durante este período de trabajo. Por otra parte, considero que el avance del proyecto integrador es adecuado y se corresponde con el cronograma definido en el documento de planificación.

Fdo: Ricardo Flores Moyano

Quito, 27 de Febrero de 2026