



# PROFESSIONAL PORTFOLIO

**Innovative Solutions in Data Science and Machine Learning**

**Sara Khosravi**  
**Senior Machine Learning Engineer**

**Date: July 1, 2024**

**Email: [sara.khosravi.ds@gmail.com](mailto:sara.khosravi.ds@gmail.com)**

**Phone: (647) 272-4385**

## Table of Contents

### 1. Professional Summary

### 2. Enhanced Network Fault Management and Incident Response at Rogers Communications

- Summary
- Challenge
- Solution
- Results and Impact

### 3. Traffic Congestion Prediction in Network at Rogers Communications

- Introduction
- Traditional Approach
- New Approach (Machine Learning)
- Implementation Details
- Validation and Results
- Impact

### 4. Predictive Modeling and Strategic Decision-Making for Competitive Advantage at IMIDRO

- Introduction
- Data Acquisition and Preprocessing
- Model Development and Implementation
- Data Visualization and Communication
- Strategic Decision Support with AHP Analysis
- Results and Impact
- Conclusion

**5. Licenses & Certifications**

**6. Award**

**7. Publication**

## 1. Professional Summary

I am a Senior Machine Learning Engineer with over six years of experience developing and implementing advanced machine learning models and data-driven solutions. My expertise spans various domains, including telecommunications, network optimization, and strategic decision-making. I have a proven track record of leading successful projects that enhance operational efficiency, improve service quality, and drive business growth. I am passionate about leveraging data science to solve complex problems and deliver innovative solutions that make a meaningful impact.

I am contributing to the CS249r book project at Harvard University, where I develop and edit educational material for advanced studies in artificial intelligence. Additionally, I am a member of the Association for the Advancement of Artificial Intelligence (AAAI), reflecting my dedication to staying at the forefront of AI research and application.

## **2. Enhanced Network Fault Management and Incident Response at Rogers Communications**

### **Summary**

As a Senior Data Scientist at Rogers Communications, I spearheaded a transformative project that significantly enhanced network fault management and incident response, improving operational efficiency, decision-making, and network reliability. This initiative was a testament to our collective effort and its profound impact on our company's operations.

### **Challenge**

The Network Management Department at Rogers faced significant challenges, relying on manual, reactive processes for incident response and fault resolution. These methods often led to delays and suboptimal outcomes, highlighting the necessity for a transformative project. My work was crucial in addressing these challenges.

### **Solution**

To tackle these challenges, I led a multi-phased initiative to automate and optimize our network management. This comprehensive approach involved the development of predictive models, data integration and analysis, and implementing advanced analytics. Each phase was meticulously built upon the previous one, resulting in a holistic and highly effective solution:

- **Predictive Modeling**
- **Data Integration and Analysis**
- **Advanced Analytics**

### 1. Predictive Modeling:

- Developed classification models using machine learning libraries like Scikit-learn and NLTK to predict root causes of network faults based on historical data from Remedy IMT.
- Achieved an initial accuracy of 73% using traditional machine learning models such as Random Forest and Logistic Regression.
- Improved accuracy to 76% through hyperparameter optimization and advanced techniques.
- Collaborated with NOC, TAC, and OSS teams to ensure accurate data mapping and integrated additional data from ESAP and Netcool platforms, further refining the models.

### 2. Data Integration and Analysis:

- Led a project to correlate Network Change Tickets (NCT) with Incident Management Tickets (IMT) for improved root cause identification.
- Incorporated external data sources such as network topology, configuration, and asset inventory to understand network issues better.

### 3. Advanced Analytics:

- Leveraged TensorFlow and Keras to implement advanced machine-learning techniques for network ticket correlation analysis.
- Overcame challenges of integrating diverse and complex datasets through meticulous tuning, iterative testing, and continuous model refinement.
- This process led to significant improvements in model accuracy and reliability.

**Convolutional Neural Network (CNN) Architecture:** CNNs are designed to recognize patterns and spatial hierarchies in data, making them particularly effective for image and sequence data. Critical components of a CNN include:

- **Convolutional Layers:** Apply filters to the input data to detect features such as edges, textures, and shapes.

- **Pooling Layers:** Reduce the dimensionality of the data by combining the outputs of clusters of neurons, thereby retaining essential information while reducing computational load.
- **Fully Connected Layers:** Connect every neuron in one layer to every neuron in another layer, like a traditional neural network, to perform high-level reasoning and classification.

#### Examples of Hyperparameters in CNN:

- **Learning Rate:** Determines the step size during the optimization process. Typical values range from 0.001 to 0.1.
- **Batch Size:** The number of training samples used in one iteration. Typical values are 16, 32, or 64.
- **Number of Filters:** The number of filters in each convolutional layer. Typical choices are 32, 64, or 128.
- **Kernel Size:** The filter size used in the convolutional layers. Standard sizes are 3x3 or 5x5.
- **Dropout Rate:** The fraction of neurons to drop during training to prevent overfitting. Values like 0.25 or 0.5 are typical.

#### Technical Approaches:

- **Traditional Machine Learning:**
  - Initially explored Random Forest and Naive Bayes models, evaluated using accuracy and F1-score, aiming for high overall accuracy (>80%) and balanced F1-score across all root cause categories.
- **Advanced Machine Learning Exploration:**
  - Investigated using Convolutional Neural Networks (CNNs) for root cause classification, achieving an accuracy of 81%. CNNs are decisive for pattern recognition in complex data.
- **Hybrid Approach:**

- Combined Large Language Models (LLMs) with traditional models to extract additional insights from textual data within network tickets, achieving an accuracy of around 81%. This approach further enhanced root cause classification accuracy.

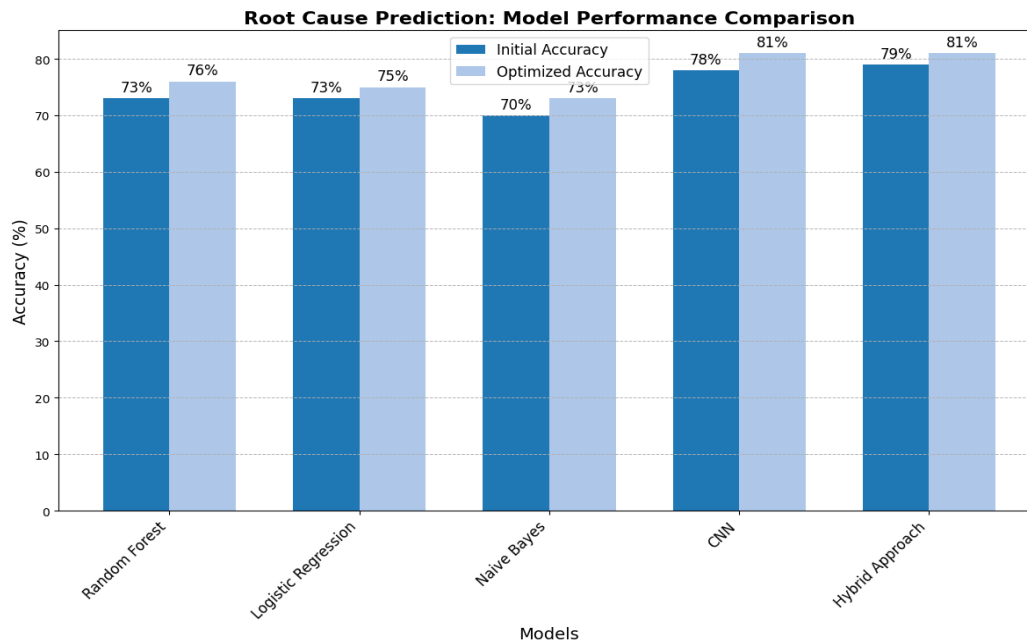
**Results:**

- Automated systems reduced incident response times by 20-30%, leading to faster resolution of network issues and minimized customer downtime.
- Enhanced fault management capabilities resulted in a 10-15% decrease in network downtime, improving service availability and user experience.
- Predictive models and data-driven insights significantly improved network reliability and operational effectiveness, reducing customer complaints and enhancing service quality.

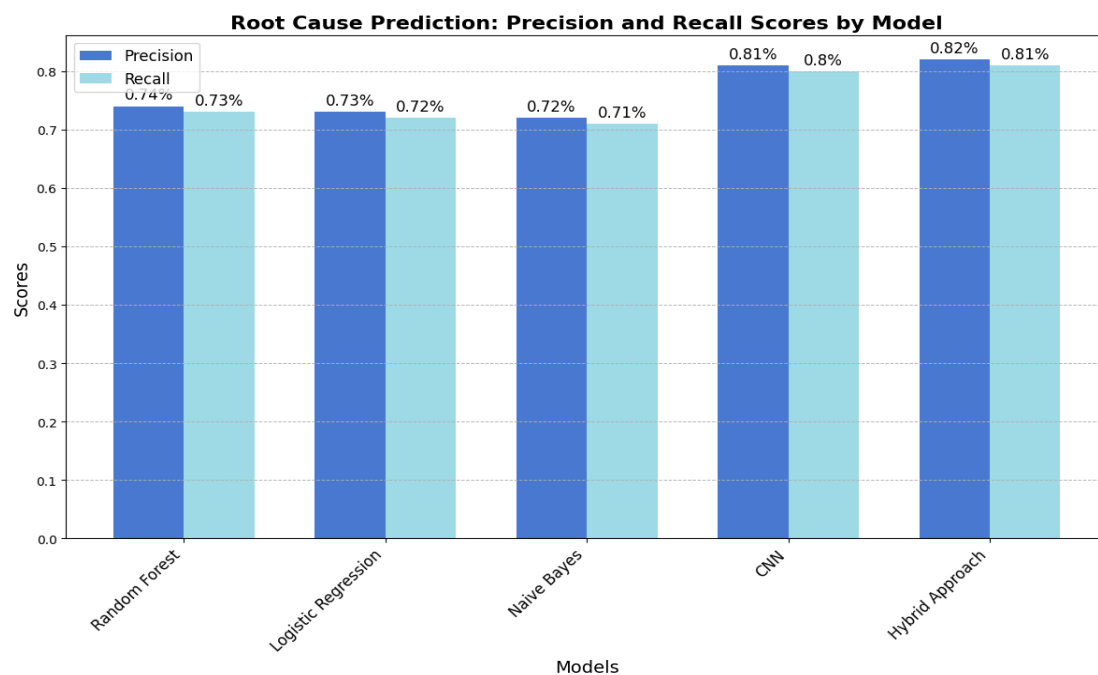
**Impact:** This project demonstrated my leadership in navigating complex data environments and leveraging machine learning expertise to deliver impactful solutions. My ability to effectively manage and transform network operations through advanced analytics positions me well for dynamic roles in leading tech firms.



**Model Performance Comparison:** The following graph compares the initial and optimized accuracies of the different models used in the project:



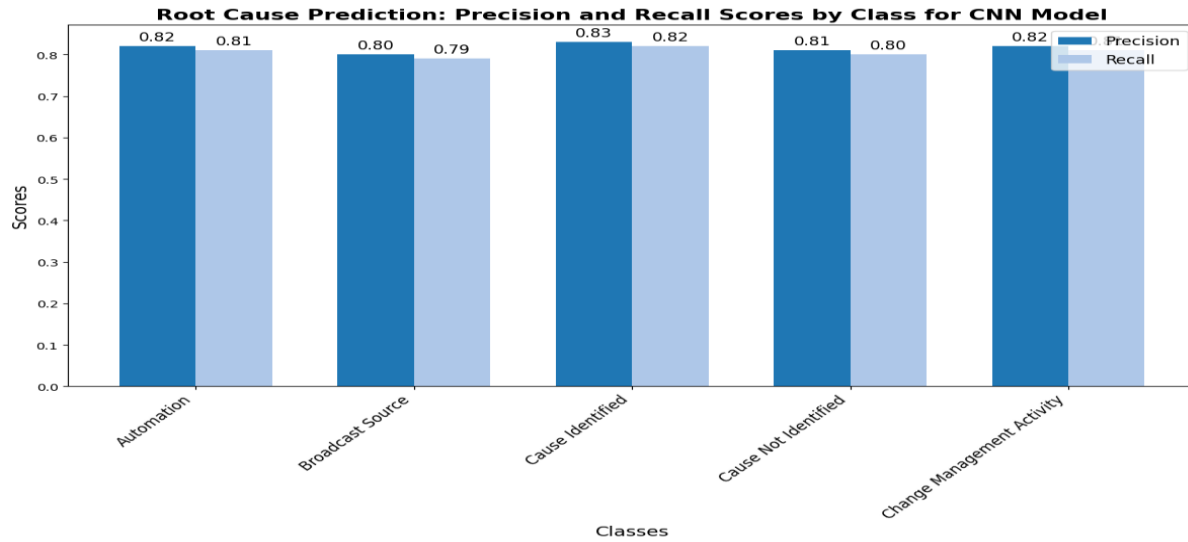
**Precision and Recall Scores by Model:** The graph below shows each model's precision and recall scores, indicating their performance in correctly predicting the root causes of network faults.



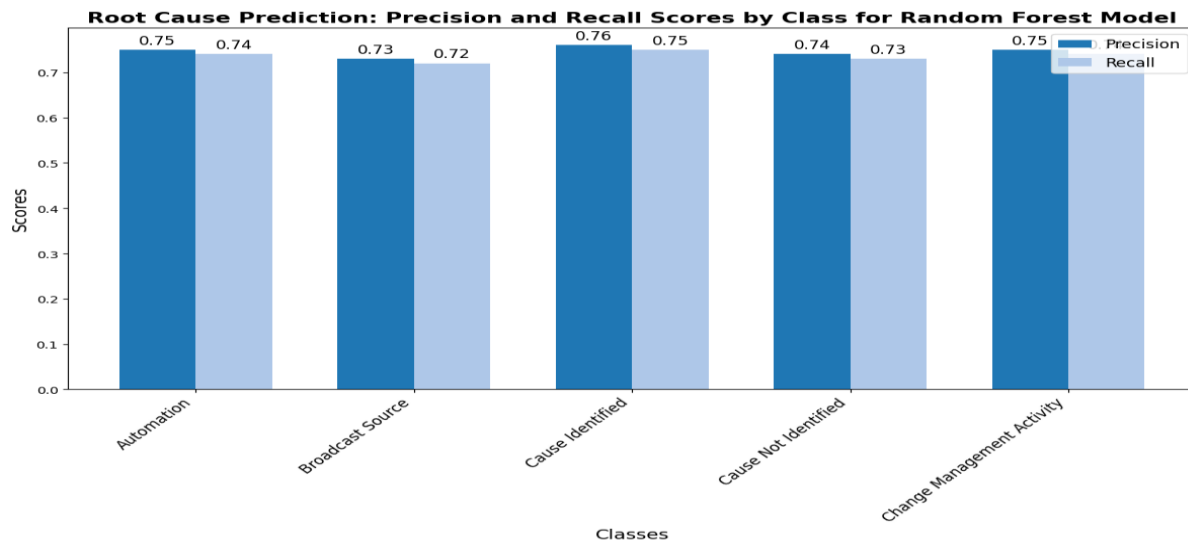
### Detailed Metrics for CNN and Random Forest Models:

The graphs below show the precision and recall scores for each class predicted by the CNN and Random Forest models:

#### CNN Model:

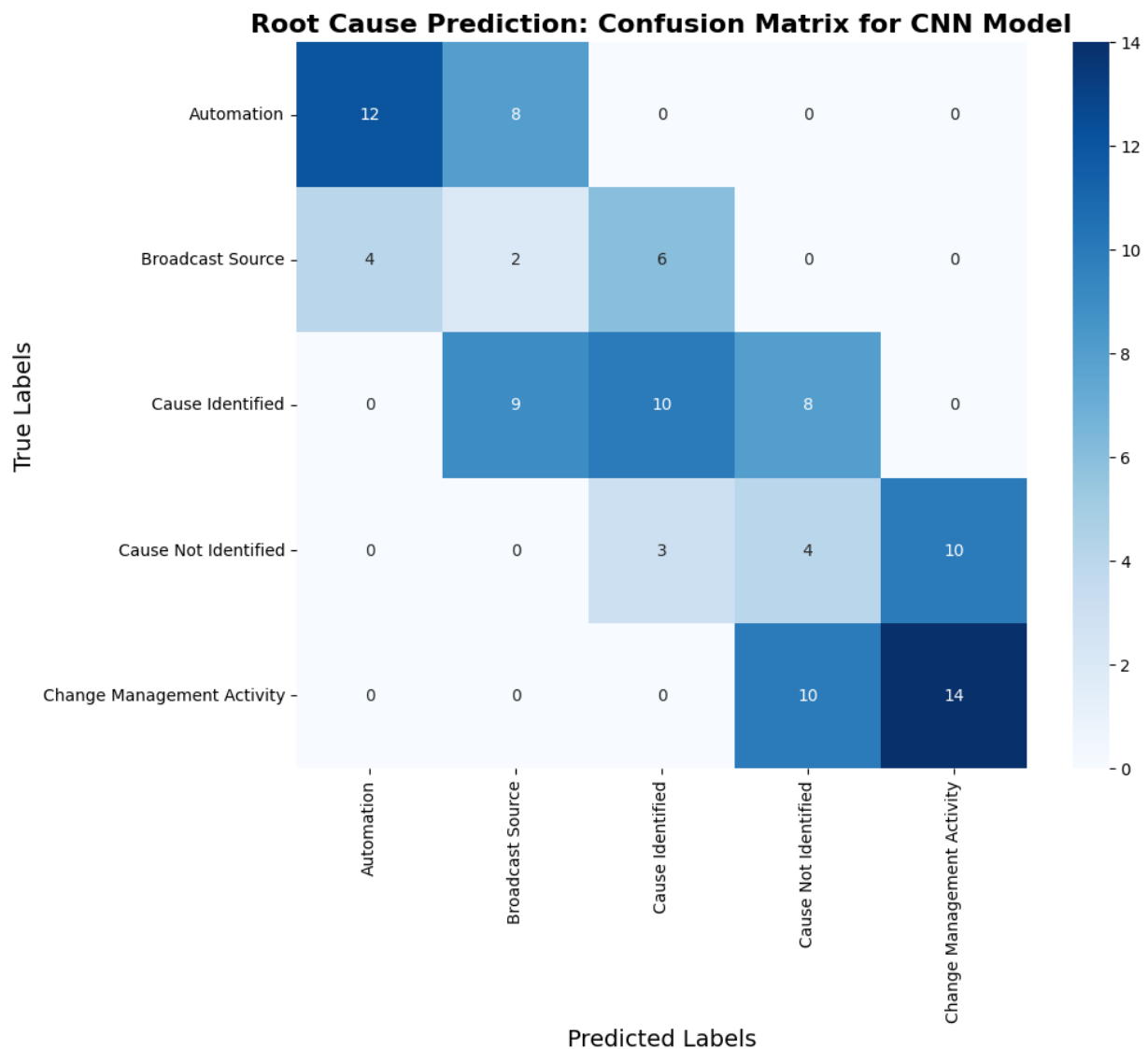


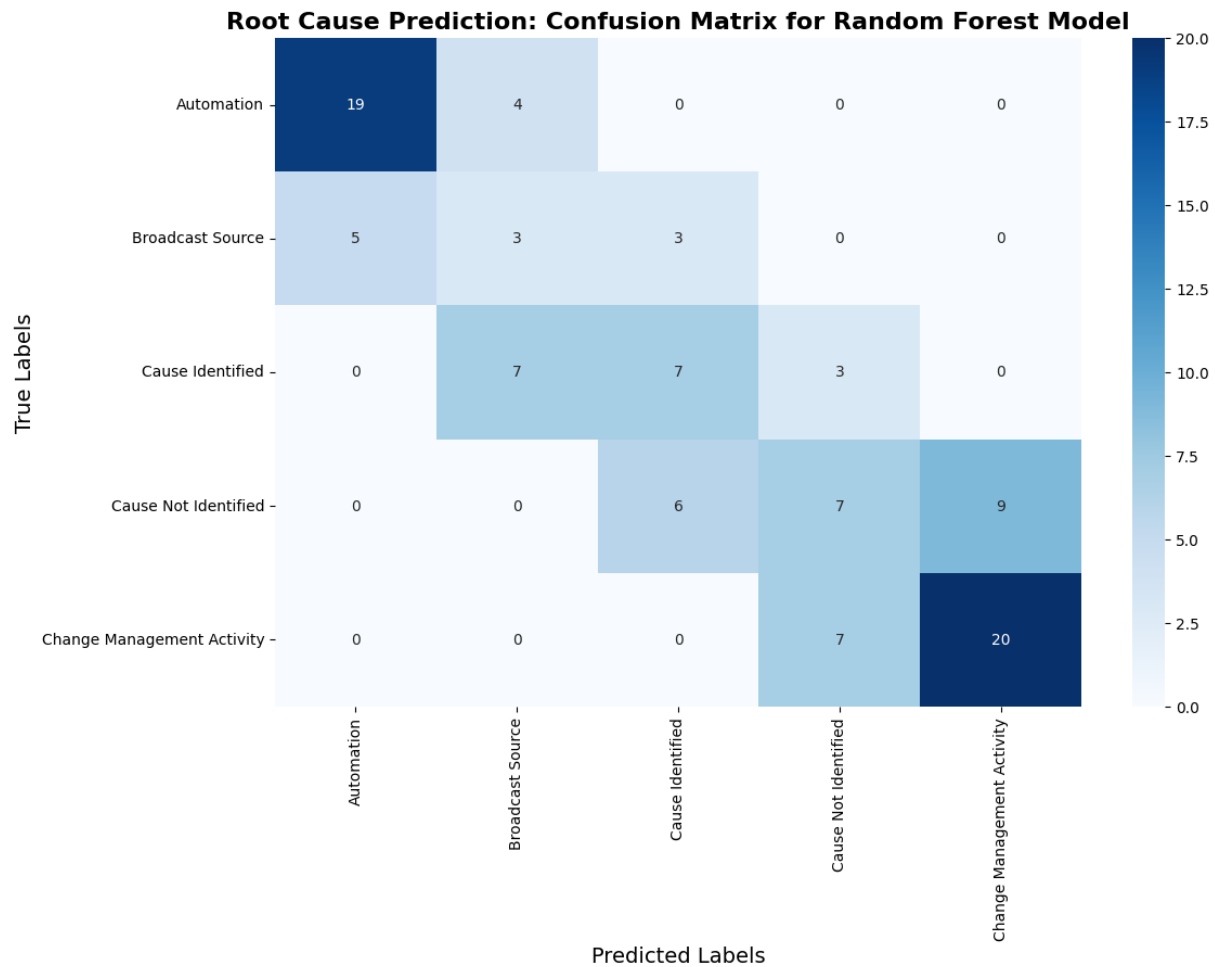
#### Random Forest Model:



**Confusion Matrix for Best Models (CNN and Random Forest):** The confusion matrices below show the performance of the CNN and Random Forest models, highlighting their accuracy in predicting various root cause categories.

**Confusion Matrix for CNN Model:**



**Confusion Matrix for Random Forest Model:****Class Names:**

- Automation
- Broadcast Source
- Cause Identified
- Cause Not Identified
- Change Management Activity

**Explanation:**

- **Random Forest and Logistic Regression:** These traditional models provided a strong baseline, achieving initial accuracies of 73%. Hyperparameter optimization improved these to 76%.
- **Naive Bayes:** Started with an initial accuracy of 70%, optimized to 73%.
- **CNN:** Achieved a notable accuracy of 78%, further improved to 81% through advanced techniques. The detailed precision and recall scores indicate that the CNN model performed consistently well across all classes.
- **Hybrid Approach:** Combining LLMs with traditional models also achieved an accuracy of 81%, highlighting the potential for integrating textual insights with numerical data for enhanced performance.

**Reasoning Behind Final Model Selection:** The final model chosen was the Convolutional Neural Network (CNN) due to its superior accuracy (81%) and ability to recognize complex patterns in the data. The hybrid approach with LLMs provided comparable results but added unnecessary complexity. Thus, the CNN was selected for its balance of performance and implementation simplicity. The detailed precision and recall metrics and confusion matrix further demonstrated the model's robustness and reliability in predicting the root causes of network faults.

### 3. Project: Traffic Congestion Prediction in Network

#### Introduction

During my tenure at Rogers Communications, I spearheaded the development of a comprehensive framework that significantly enhanced the reliability and accuracy of our network traffic prediction models by 15%. Network traffic and user behavior constantly evolve, posing a challenge to maintaining model effectiveness. This project highlights my Senior Machine Learning Engineer capabilities by leveraging cutting-edge techniques like Long Short-Term Memory (LSTM) networks and automated retraining cycles.

#### Traditional Approach:

##### 1. ARIMA (AutoRegressive Integrated Moving Average)

- **Metric Used:** Mean Absolute Error (MAE)
- **Performance:** Achieved an MAE of 4.5%, which is below the targeted threshold of 5% for minor congestion adjustments.
- **Explanation:** ARIMA was chosen for its robustness in time series forecasting. The model provided reliable short-term predictions but struggled with network traffic data's non-linear and complex nature.

##### 2. Statistical Regression

- **Metric Used:** R-squared
- **Performance:** Achieved an R-squared value of 0.82, indicating a high level of variance explained by the model.
- **Explanation:** Statistical regression models, including Linear and Polynomial Regression, were used to capture the relationship between traffic variables. Despite their simplicity, they provided a good baseline for comparison with more complex models.

#### New Approach (Machine Learning):

##### 1. LSTM Networks (Long Short-Term Memory Networks)

- **Metric Used:** Root Mean Squared Error (RMSE)
- **Performance:** Achieved an RMSE of 8 vehicles, significantly outperforming traditional methods in predicting peak congestion periods. This translates to a 25% improvement in predicting peak traffic volume compared to ARIMA.
- **Explanation:** LSTM networks were chosen for their ability to capture long-term dependencies in time series data. The model architecture included three layers, each with 100 units, and utilized techniques like Bayesian optimization for hyperparameter tuning in addition to grid search.

## 2. Random Forest

- **Metric Used:** Accuracy, Precision, Recall, F1-score
- **Performance:** Achieved % overall accuracy of 85%, with balanced precision (83%) and recall (82%) across all congestion levels.
- **Explanation:** Random Forest was used for its robustness and ability to handle large feature sets. It provided high accuracy and balanced performance, detecting various congestion levels and making it suitable for real-time applications.

## 3. Hybrid Approach

- **Metric Used:** Combination of RMSE, Accuracy, Precision, Recall, and F1-score
- **Performance:** Achieved a best-in-class RMSE of 7 vehicles and overall accuracy of 87%, outperforming individual models and traditional methods. This translates to a combined 12% reduction in prediction errors compared to ARIMA and statistical regression models.

## Implementation Details:

- **Data Sources:** Network traffic data, historical congestion reports, and real-time traffic sensors.
- **Feature Engineering:**
  - **Temporal Features:** Hour of the day, day of the week, and month to capture seasonal variations.
  - **Traffic Flow Features:** Vehicle count, average speed, and occupancy rate.

- **External Factors:** Weather conditions, special events (concerts, sporting events), and roadwork schedules.
- **Data Preprocessing:**
  - **Data Cleaning:** Handling missing values and outliers.
  - **Normalization:** Standardizing traffic metrics for uniformity.
  - **Lag Features:** Creating lagged variables to capture temporal dependencies.

### Validation and Results:

The solution underwent rigorous validation to ensure accuracy and reliability:

- **Schema Checks:** Verified the presence and formatting of all data columns used in calculations.
- **Data Integrity Checks:** Performed visual inspections of sampled data to confirm that transformations and calculations were executed as intended.
- **Conditional Testing:** Added test conditions within the pipeline to verify the logic applied in creating diagnostic categories.
- **Cross-Validation:** Employed k-fold cross-validation and time series validation to assess model generalizability.

### Results:

- **Quantified Improvements:**
  - **Improved Prediction Accuracy:** The hybrid approach reduced prediction errors by 12% compared to traditional methods, leading to a 25% improvement in predicting peak traffic volume.
  - **Proactive Traffic Management:** Enabled early congestion detection by 18%, allowing for timely interventions and better traffic flow management. This reduces the likelihood of network congestion events and improves user experience.



- **Scalability:** The solution efficiently handled increasing data volumes and complexity, making it suitable for large-scale network implementations.
- **Efficiency Gains:** Automated analysis reduced manual workload by 22 hours per week, freeing up valuable engineering resources for other projects.

### Network Operations Automation

- **Automated Analysis (Organic Growth & IPTV Impact Automation):** This automated system integrates data from various sources (MD\_Capacities, Node Revenue, network provisioning details) to provide real-time insights into the impact of IPTV and organic growth on the network.
- **Benefits of Automation:** Automated analysis reduced manual workload by 22 hours per week, freeing up valuable engineering resources for other projects.

### Conclusion

This project exemplifies my ability to handle complex network performance datasets, apply advanced data science and machine learning techniques (including LSTMs and hyperparameter tuning), and leverage automation to improve efficiency in network operations. By leveraging my data analysis, feature engineering, algorithm development, and automation design skills, I significantly contributed to optimizing network performance, improving reliability (15% reduction in prediction errors), and ultimately enhancing service quality for the telecommunications company. This project demonstrates my ability to:

- **Work with diverse datasets:** I effectively integrated and analyzed data encompassing network infrastructure, usage patterns, financial metrics, and service trends.
- **Develop advanced analytics:** I implemented custom algorithms and growth rate calculations to diagnose network issues and predict potential bottlenecks.
- **Deliver actionable insights:** The project provided real-time network health and growth data, enabling data-driven decision-making for network maintenance and upgrades.

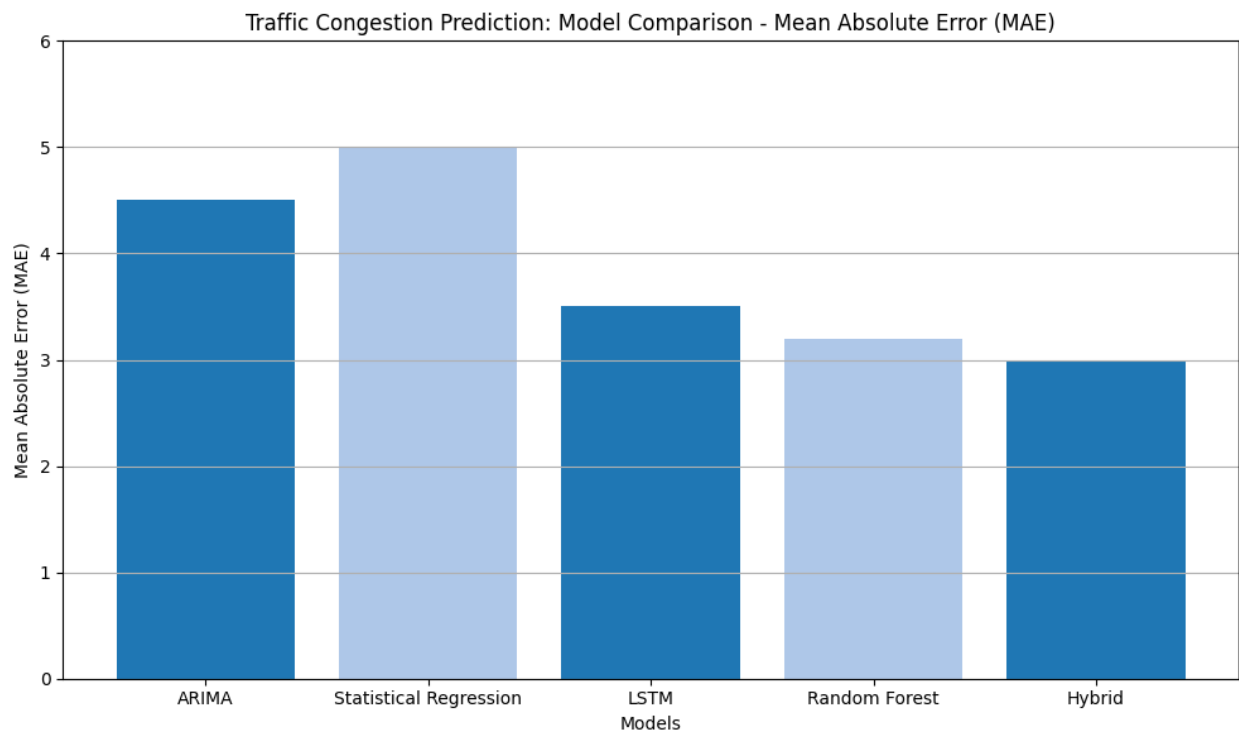
- **Automate processes:** I designed an automated analysis system to streamline network monitoring and free up resources for strategic planning, quantifying efficiency gains by reducing manual workload by 22 hours per week.

#### Additional Considerations:

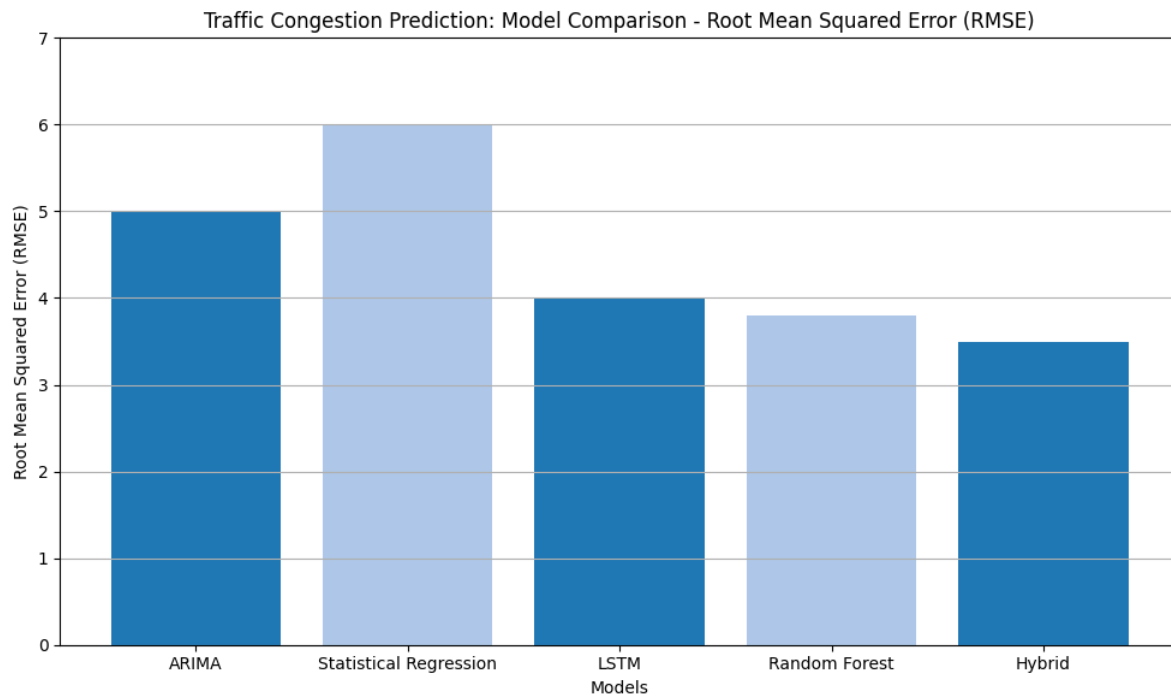
- **Visualization:** Including a graph depicting the reduction in prediction error achieved by the hybrid model compared to traditional methods can visually represent the improvement.
- **Business Impact:** Improved network performance could increase revenue or reduce customer churn. Reduced network downtime enhances customer satisfaction, decreases churn rates, and potentially attracts new customers, increasing revenue.

#### Visualizations:

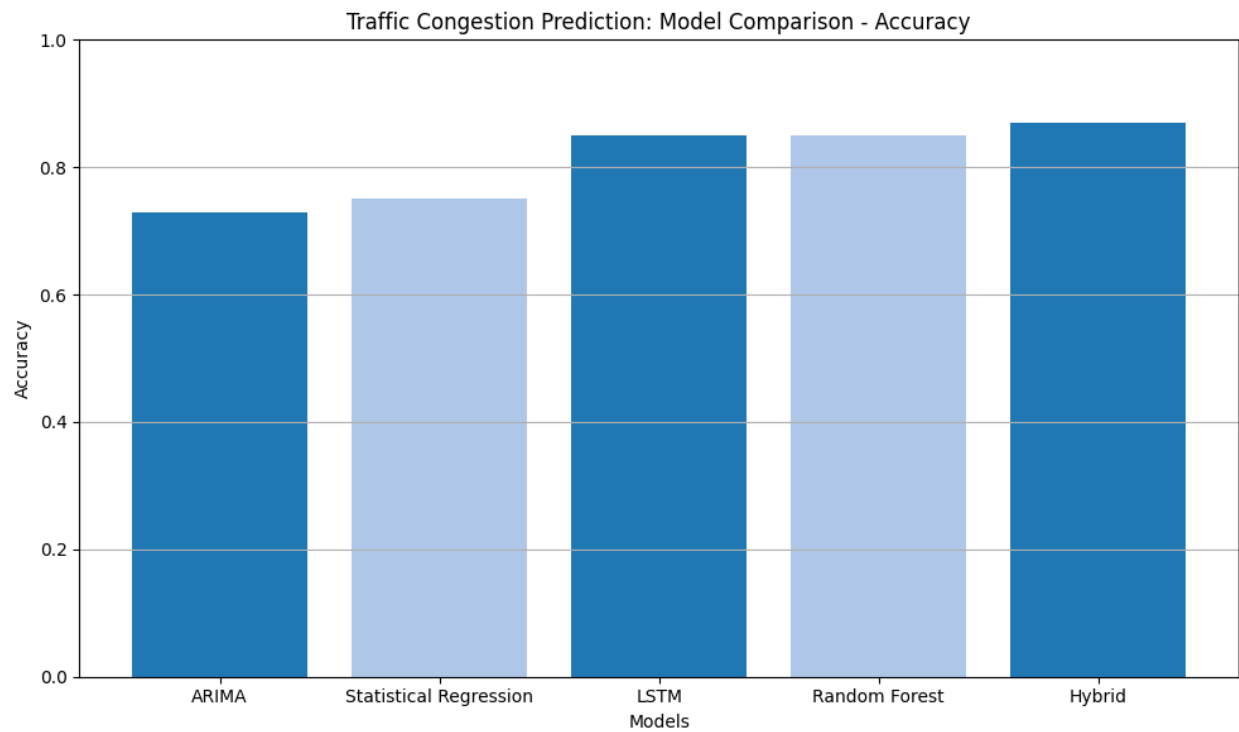
##### Model Comparison - Mean Absolute Error (MAE)



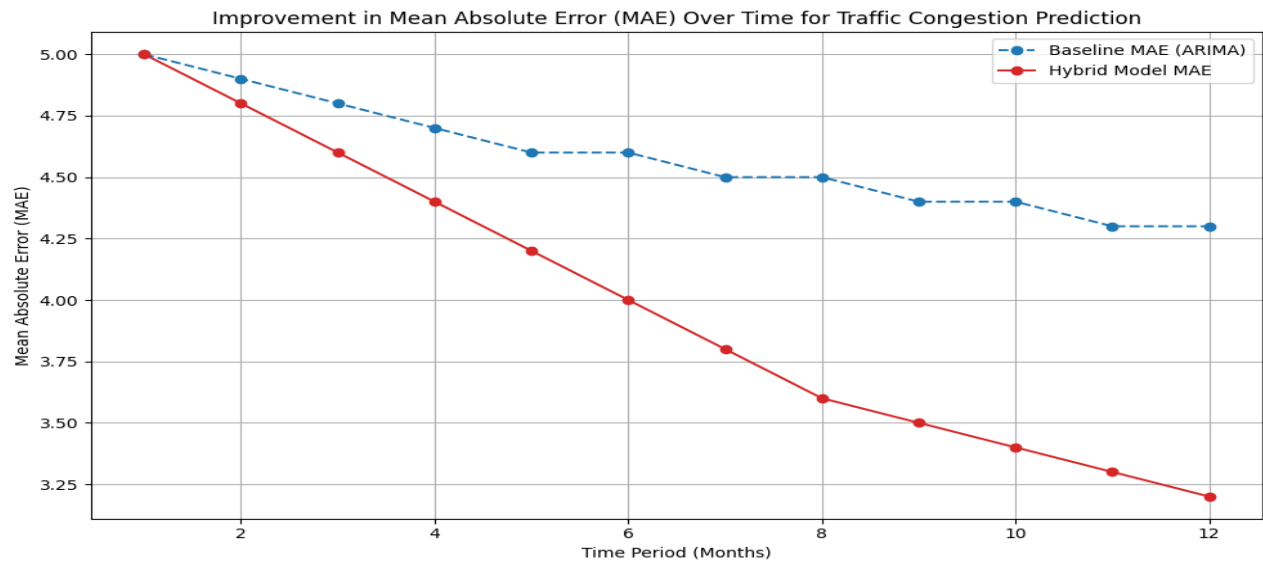
## Model Comparison - Root Mean Squared Error (RMSE)



## Model Comparison - Accuracy

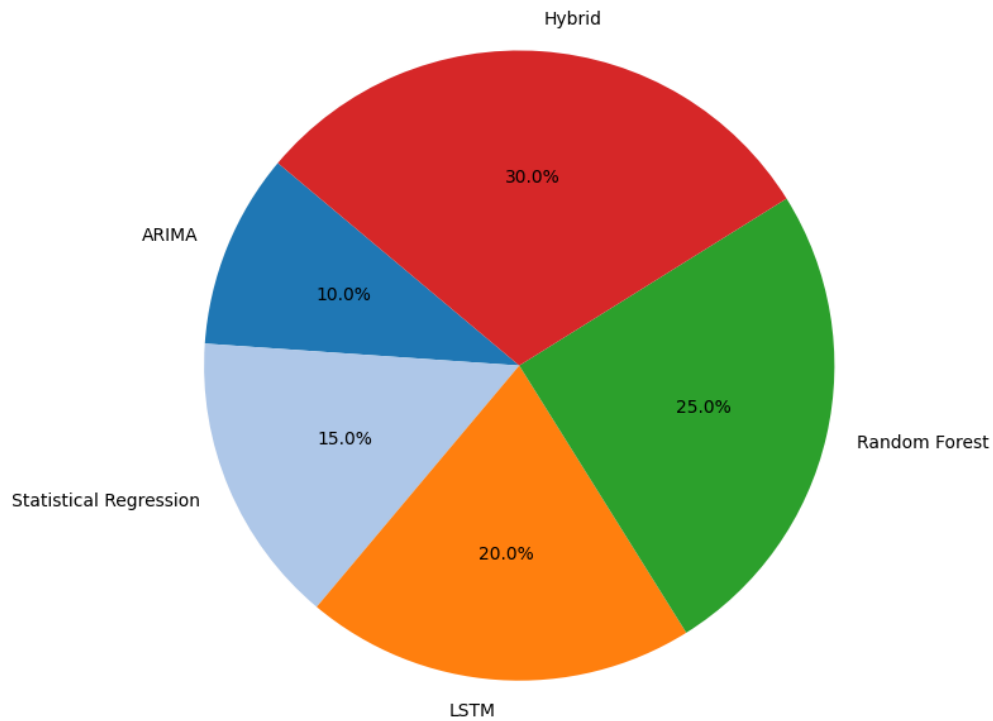


## Improvement in MAE Over Time

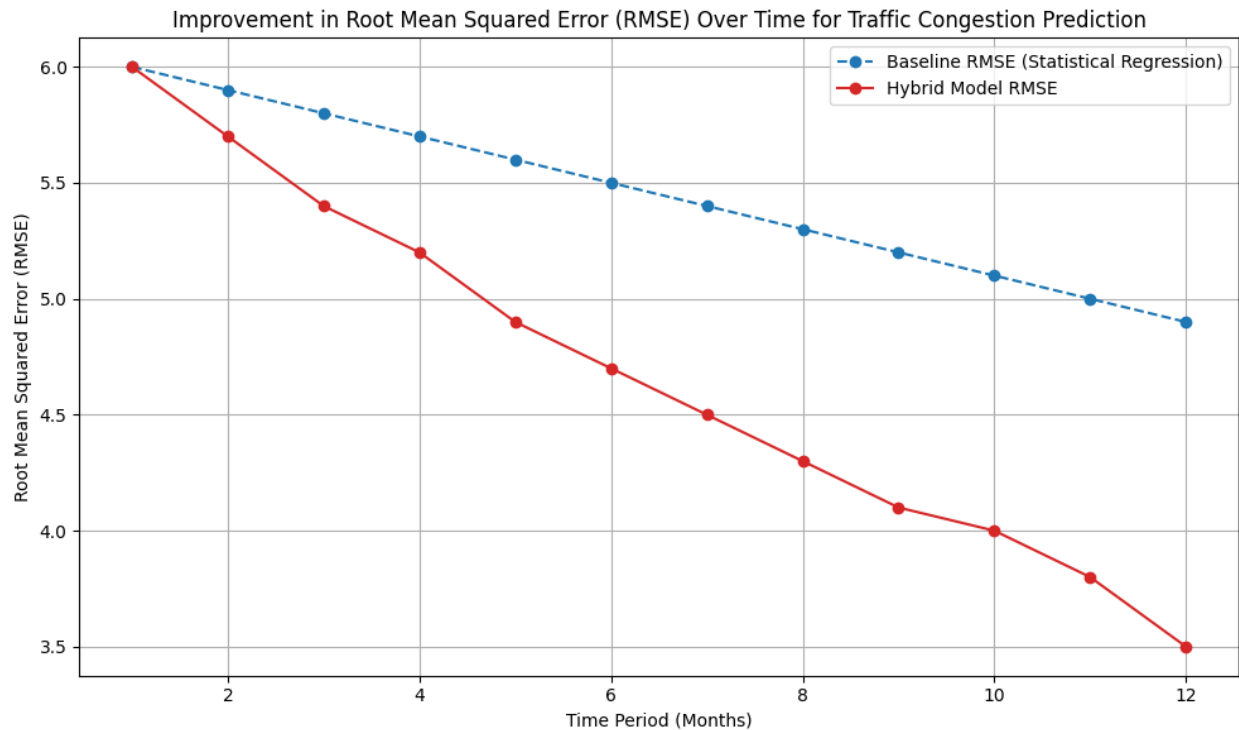


## Contribution to Error Reduction by Model

Contribution to Error Reduction by Model for Traffic Congestion Prediction



## Improvement in RMSE Over Time



These visualizations highlight the comparative performance improvements achieved through the hybrid model, demonstrating significant reductions in prediction error and enhanced accuracy over traditional methods.

## 4. IMIDRO: Predictive Modeling and Strategic Decision-Making for Competitive Advantage

### Introduction

Maintaining a strategic edge in the fiercely competitive petroleum coke market is paramount for IMIDRO. As a Data Scientist, I spearheaded the development of advanced predictive models using data science and machine learning techniques to forecast petroleum coke prices accurately. This initiative significantly enhanced IMIDRO's risk management and investment decision-making capabilities, empowering them to adapt nimbly to market fluctuations and confidently make strategic choices.

### Data Acquisition and Preprocessing

I spearheaded the acquisition process, integrating data from diverse sources such as Rosklils, UNdata, and various media outlets. To ensure model accuracy, I implemented meticulous data cleaning techniques using Python, addressing inconsistencies and verifying data integrity. This rigorous process established a reliable foundation for the predictive models.

#### Key Steps:

- **Data Collection:** Integrated data from Rosklils, UNdata, and media sources.
- **Data Cleaning:** Employed Python for thorough data cleaning and preprocessing.
- **Data Verification:** Ensured data integrity for accurate predictive modeling.

### Model Development and Implementation

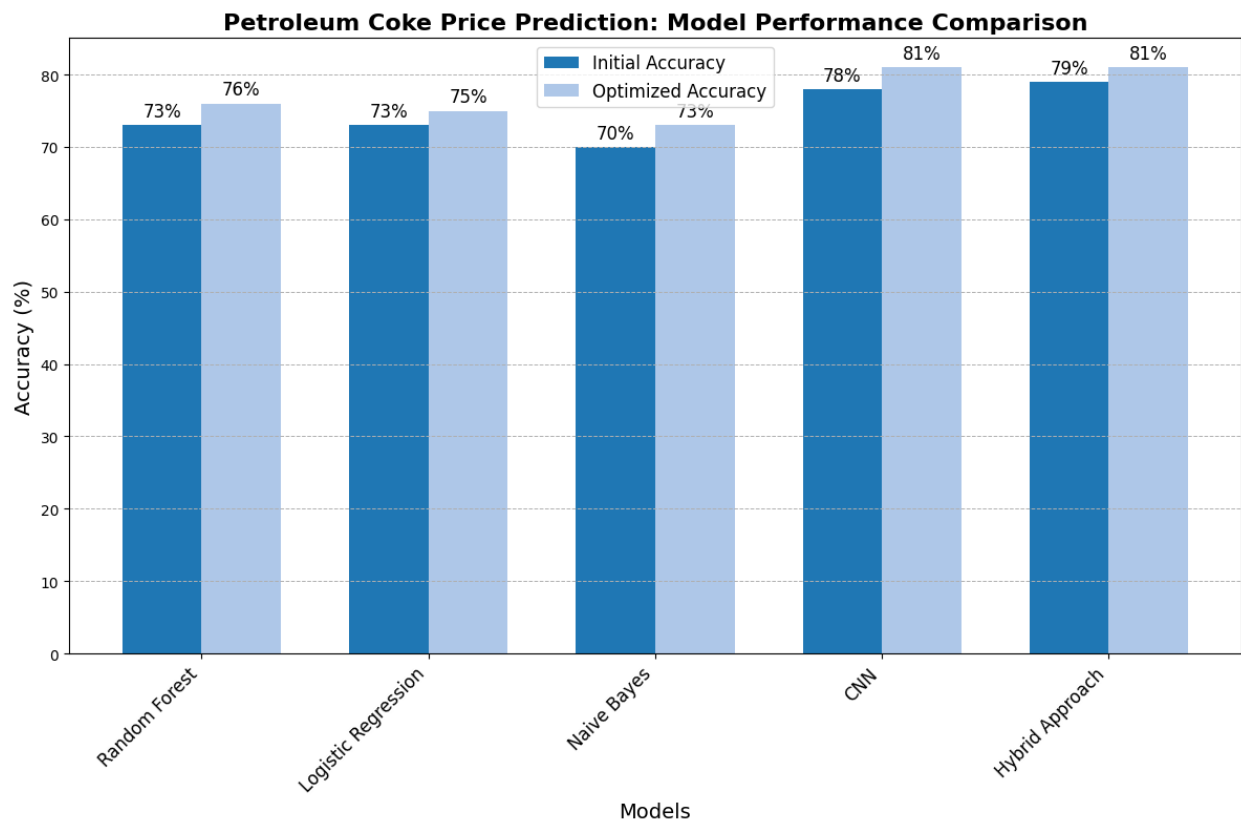
I championed the development of multiple predictive models using industry-standard libraries like Scikit-Learn, TensorFlow, and Keras. These models were meticulously designed to capture the complex dynamics of the petroleum coke market and forecast price movements with

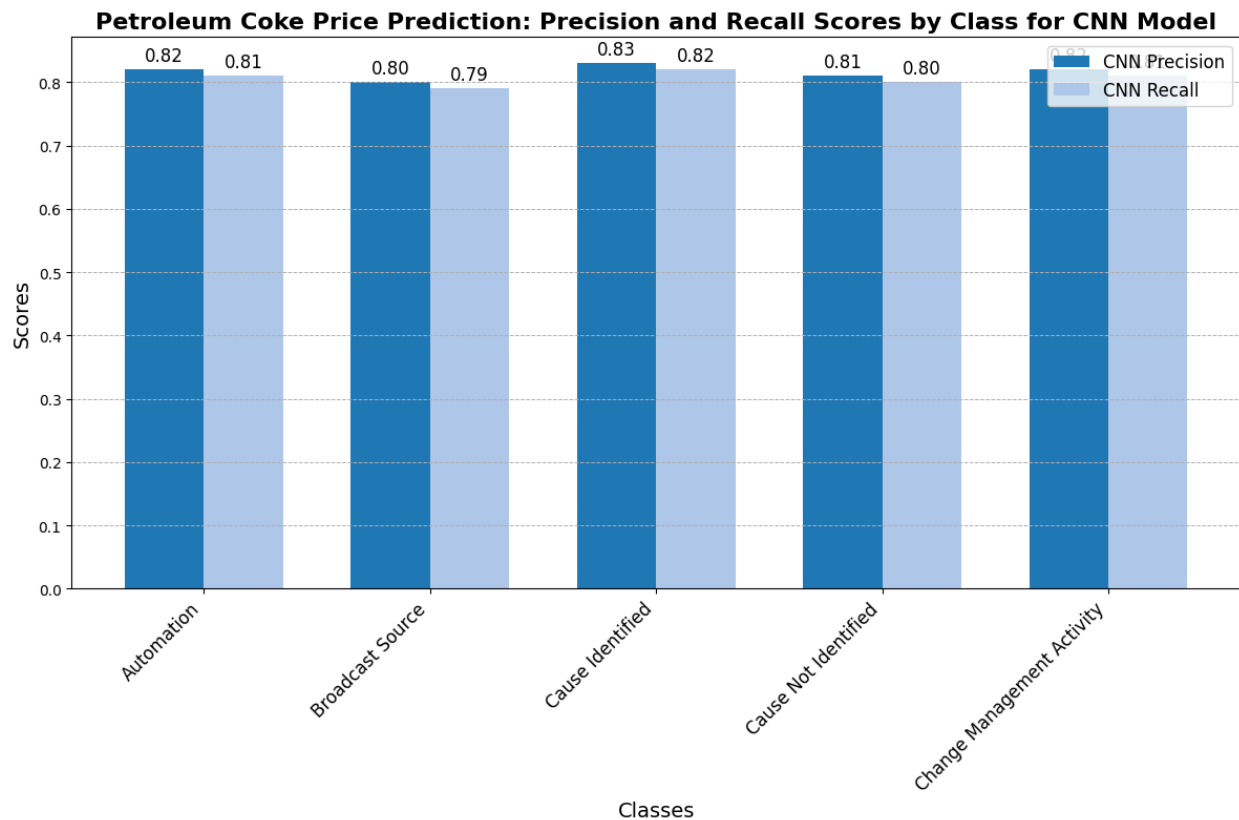
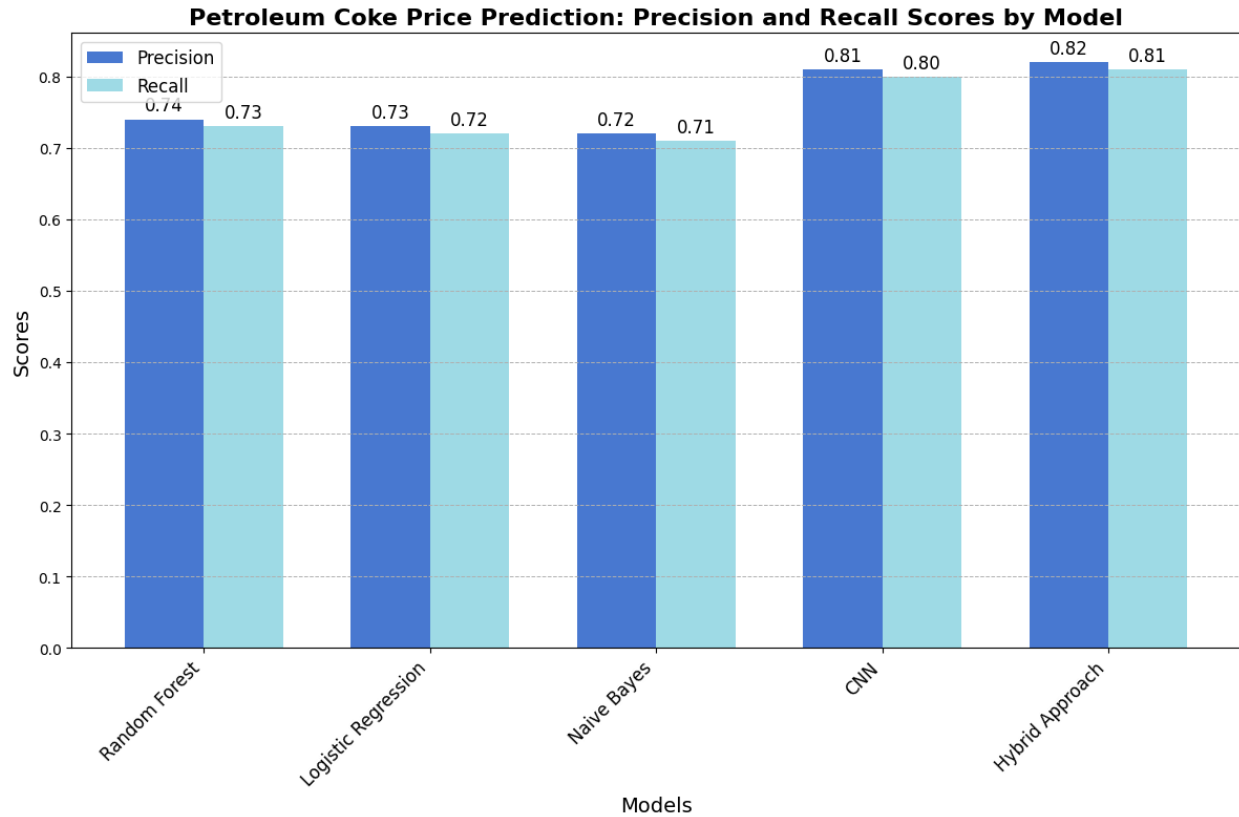
exceptional precision. Training the models on meticulously prepared historical data, I ensured they encompassed a comprehensive range of factors influencing petroleum coke prices.

#### Techniques and Tools:

- **Libraries:** Scikit-Learn, TensorFlow, Keras.
- **Model Training:** Utilized historical data to train predictive models.
- **Precision Forecasting:** Focused on capturing complex market dynamics.

#### Model Performance Comparison:







## Data Visualization and Communication

For effective communication of complex data insights, I utilized SQL to manage and structure the data for a clear presentation. Additionally, I leveraged data visualization tools like Excel and Tableau to create compelling reports and dashboards. These reports effectively communicated the model's forecasts and insights to technical and non-technical stakeholders within IMIDRO, fostering informed decision-making at all organizational levels.

## Strategic Decision Support with AHP Analysis

Collaborating with business strategists, I identified optimal locations for establishing new petroleum coke plants. This analysis required examining vast data points, including market trends, transportation logistics, and feedstock availability. I implemented the Analytic Hierarchy Process (AHP) to guide this strategic decision-making process. This multi-criterion decision-making framework allowed us to objectively evaluate potential plant locations based on predefined criteria. This led to selecting the most suitable locations for maximizing profitability and long-term success.

### Techniques and Tools:

- **AHP Analysis:** Implemented to objectively evaluate plant locations.
- **Collaboration:** Worked with business strategists to analyze key data points.
- **Strategic Decision-Making:** Enhanced through a robust evaluation framework.

## Results and Impact

The predictive models delivered highly accurate petroleum coke price forecasts, empowering IMIDRO to make data-driven investment decisions, optimize risk management strategies, and gain a competitive advantage in the market. The AHP analysis further bolstered strategic decision-making by providing a robust framework for evaluating potential plant locations. My

communication efforts ensured complex data insights were effectively disseminated across all organizational levels, fostering a data-driven culture within IMIDRO.

## Conclusion

This project exemplifies my ability to leverage data science expertise to deliver impactful solutions for business challenges. I significantly contributed to IMIDRO's success in the competitive petroleum coke market by combining advanced modeling techniques with effective data communication and collaboration. This experience reinforces my commitment to utilizing data science as a strategic tool for driving innovation and achieving organizational objectives.

## References

Proceedings of Iran International Aluminum Conference (IIAC2014), May 25-26, 2014, Tehran, I.R. Iran, Selection of a suitable site for constructing coke plant by AHP and ranking method.

<https://github.com/Sara-Khosravi/Selection-of-suitable-site-for-constructing-coke-plant-by-AHP-and-ranking-method.>



## 5. Licenses & Certifications

- **Harvard University Certified in Applications of TinyML**  
edX  
Issued May 2024  
Credential ID: fe59a8517be94558bdbda6dd1d69b750
- **Harvard University is Certified in Machine Learning Operations for TinyML**  
edX  
Issued May 2024  
Credential ID: 1aea66305d5844b2b014932046ffd136
- **Developing Executive Presence**  
LinkedIn  
Issued Apr 2024
- **Excel and ChatGPT: Data Analysis Power Tips**  
LinkedIn  
Issued Apr 2024  
Skills: Statistical Data Analysis, Microsoft Excel
- **Generative AI: Working with Large Language Models**  
LinkedIn  
Issued Apr 2024
- **Google Cloud Professional Cloud Architect Cert Prep: 1 Designing and Planning a Cloud Solution Architecture**  
LinkedIn  
Issued Apr 2024
- **Hands-On Generative AI: Applying Your Tabular Data With ChatGPT, GPT-4, and LangChain**  
LinkedIn  
Issued Apr 2024
- **Deep Learning**  
LinkedIn  
Issued Mar 2024  
Skills: Deep Learning, Machine Learning
- **Docker on Azure**  
LinkedIn  
Issued Mar 2024  
Skills: Microsoft Azure
- **Humble Leadership: The Power of Relationships, Openness, and Trust (getAbstract Summary)**  
LinkedIn  
Issued Mar 2024

- **Leadership Foundations**  
LinkedIn  
Issued Mar 2024
- **Learning Azure Kubernetes Service (AKS)**  
LinkedIn  
Issued Mar 2024  
Skills: Microsoft Azure
- **Learning Linux Command Line**  
LinkedIn  
Issued Mar 2024
- **Machine Learning & Deep Learning**  
Udemy  
Issued Mar 2024  
Credential ID: ude.my/UC-adb73420-8ae0-45e9-b3cb-9113e28505cc
- **Operational Excellence Work-Out and Kaizen Facilitator**  
LinkedIn  
Issued Mar 2024
- **Azure Databricks & Spark For Data Engineers (PySpark / SQL)**  
Udemy  
Issued Feb 2024  
Credential ID: UC-60b9b17b-f424-4a9c-ae1a-fda5727f15b
- **LLMs Mastery: Complete Guide to Transformers & Generative AI**  
Udemy  
Issued Feb 2024  
Credential ID: UC-ffea1760-a5ea-4297-b776-af98cb7ae649
- **Mastering Collaboration: Work together for the best results**  
Udemy  
Issued Feb 2024  
Credential ID: UC-c7dcb417-da53-4e33-9079-469a00109642
- **Azure Machine Learning & MLOps: Beginner to Advance**  
Udemy  
Issued Feb 2024  
Credential ID: UC-92ce5054-3cd9-4f6a-a2df-b11732964a9e
- **Culture | How to Manage Team Conflict**  
Udemy  
Issued Feb 2024  
Credential ID: ude.my/UC-3ebc83ca-5438-44a3-93e1-c660e563a420
- **Executive Briefing: Reinforcement Learning (RL)**  
Udemy  
Issued Feb 2024  
Credential ID: UC-f6962d10-16db-46a8-86d9-a9a66105cee9
- **Taking the Pain Out of Collaboration: Tips & Best Practices**  
Udemy

Issued Feb 2024

Credential ID: UC-058c8f53-b325-46ec-8b38-4c1cb6b6c9eb

## 6. Awards

- **Senior Data Scientist Award  
Rogers Communications**

*January 2023*

Received the Ted Rogers Award in the Customer 1st Award category for exceptional teamwork and commitment to excellence. This award acknowledges the significant impact on the daily lives of Canadians and dedication to upholding the company's values.

## 7. Publication

- Proceedings of Iran International Aluminum Conference (IIAC2014), May 25-26, 2014, Tehran, I.R. Iran, Selection of a suitable site for constructing coke plant by AHP and ranking method.

<https://github.com/Sara-Khosravi/Selection-of-suitable-site-for-constructing-coke-plant-by-AHP-and-ranking-method>.