

## **Names of Project Members:**

**NIHAR MUNIRAJU**

Programming for Machine learning Applications

Project Proposal: Exploring Relationships in Stock Data through Machine Learning

### **1. Introduction:**

In the constantly shifting landscape of financial markets, gaining insights into the behaviour of specific sectors is vital for strategic investment planning. This project proposal is dedicated to exploring the dynamics within the energy and power sector by examining the stock prices of 20 selected companies in this industry. By employing machine learning techniques, the aim is to uncover patterns and relationships that can inform investment decisions and provide a deeper understanding of how these stocks interact with broader market trends. This research could offer valuable perspectives for investors looking to navigate the complexities of the energy and power markets.

### **2. Project Objectives:**

**Examine Historical Data:** Delve into the historical stock data of the top 10 and bottom 10 energy and power companies to understand their market behavior over time.

**Analyze Correlations:** Investigate the correlation between the stock prices of these energy and power companies and the QQQ ETF to identify any potential predictive relationships.

**Linear Regression Modeling:** Implement a Linear Regression model to quantify the linear relationships between the stock prices and the ETF, providing a clear understanding of their interactions.

**Feature Selection with Lasso Regression:** Apply Lasso Regression to select relevant features and perform regularization, enhancing the model's predictive capability and reducing overfitting.

**Data Transformation with PCA:** Utilize Principal Component Analysis (PCA) to transform the data, reducing its dimensionality, and then build another Linear Regression model based on the transformed data.

**Classification Model Exploration:** Explore classification models such as K-Nearest Neighbors (KNN) and Adaboost to predict stock movements based on historical trends and patterns.

**Performance Evaluation:** Assess the performance of the developed models using appropriate metrics to ensure their accuracy and reliability in predicting stock movements.

**Insights Generation:** Extract actionable insights from the models to inform investment strategies and decision-making in the energy and power sector.

By achieving these objectives, the project aims to provide a comprehensive understanding of the relationships between energy and power stocks and the QQQ ETF, offering valuable insights for investors and traders in this sector.

### **3. Data Collection:**

The project will employ Alpaca, a popular financial data platform, to gather historical stock price data for the top 10 and bottom 10 energy and power companies listed on the Nasdaq 100 index. This selection aims to provide a focused analysis on this specific sector, capturing its unique dynamics and trends. Additionally, the corresponding prices of the QQQ ETF will be collected to analyze their relationship with these energy and power stocks. The data will span a considerable time frame to ensure a comprehensive understanding of various market conditions and their impact on the sector. Utilizing Alpaca's API, the project will access accurate and up-to-date information, facilitating a robust analysis of the energy and power market's behavior over time.

### **4. Methodology:**

#### **Data Preprocessing:**

The project will begin with a thorough cleaning and preprocessing of the historical stock data collected from Alpaca. This involves handling missing values, either by filling them with appropriate statistical measures or by removing the affected rows, depending on the context.

Outliers will be identified and treated to ensure they do not skew the analysis. This may involve using techniques such as Z-score or IQR (Interquartile Range) to detect and manage anomalous values.

The data will be normalized or scaled to bring all features to a similar scale, facilitating more accurate model comparisons and interpretations.

Standardization or Min-Max scaling techniques may be employed based on the data distribution.

### **Feature Engineering:**

Relevant features for analysis will be extracted from the historical stock data. This may include technical indicators such as moving averages, volatility measures, and momentum indicators, which are commonly used in stock market analysis.

Additional financial indicators, such as P/E ratios, dividend yields, or sector-specific metrics, may be explored to enrich the feature set and provide deeper insights into the stock behavior.

### **Model Building:**

A Linear Regression model will be implemented initially to establish a baseline understanding of the relationships between the stock prices and the QQQ ETF. This model will serve as a reference point for evaluating the performance of more complex models.

Lasso Regression will be applied to identify key features that have a significant impact on the stock prices. By introducing regularization, Lasso Regression helps in mitigating overfitting and improving model generalization.

PCA (Principal Component Analysis) will be utilized for dimensionality reduction, transforming the original features into a smaller set of uncorrelated components. A new Linear Regression model will be built based on these components to analyze the transformed data.

KNN (K-Nearest Neighbors) classification models will be developed to predict stock movements (up, down, or stable) based on historical data. The choice of 'K' will be optimized through cross-validation to achieve the best predictive performance.

Adaboost, an ensemble classification technique, will be implemented to boost the performance of weak learners. By combining multiple weak models, Adaboost aims to create a stronger classifier with improved accuracy in predicting stock movements.

Throughout the methodology, various performance metrics such as R-squared, RMSE (Root Mean Squared Error), precision, recall, and F1-score will be used to evaluate and compare the models. The analysis will be iterative, with

continuous refinement of models and features to achieve the most accurate and insightful results.

## 5. Evaluation:

To ensure the effectiveness and reliability of the machine learning models developed in this project, a comprehensive evaluation process will be employed:

### **Regression Model Evaluation:**

For regression models like Linear Regression and Lasso Regression, the primary metric used will be R-squared. This metric indicates the proportion of the variance in the dependent variable (stock prices or ETF prices) that is predictable from the independent variables. A higher R-squared value signifies a better fit of the model to the data.

Additionally, the Root Mean Squared Error (RMSE) will be utilized to measure the average magnitude of the errors between the predicted and actual values. A lower RMSE indicates better model performance.

### **Classification Model Evaluation:**

For classification models like K-Nearest Neighbors (KNN) and Adaboost, accuracy will be a key metric. It measures the proportion of correctly predicted instances out of the total instances. While accuracy provides an overall sense of model performance, it may not be sufficient for imbalanced datasets.

Precision and recall will be used to assess the model's ability to correctly identify positive instances (e.g., predicting stock price movements accurately). Precision measures the proportion of true positive predictions among all positive predictions, while recall measures the proportion of true positive predictions among all actual positive instances.

The F1-score, which is the harmonic mean of precision and recall, will be used to provide a balance between the two metrics, especially in cases of class imbalance.

### **Cross-Validation:**

To ensure the robustness of the models, cross-validation will be conducted during training. This technique involves dividing the dataset into multiple subsets (folds) and iteratively training and validating the model on different combinations of these folds. This helps in assessing the model's performance across different subsets of data and reduces the risk of overfitting.

K-fold cross-validation, where the dataset is split into 'k' folds, is commonly used. The model's performance metrics will be averaged across all folds to provide a more reliable estimate of its generalization ability.

By employing these evaluation techniques, the project aims to rigorously assess the performance of each model and ensure their reliability in predicting stock movements and understanding the relationships between energy and power stocks and the QQQ ETF.

## **6. Deliverables:**

### **Analysis Report:**

A detailed report will be prepared, encompassing the entire analysis process, from data collection and preprocessing to model building and evaluation. This report will present the findings and insights derived from the investigation, highlighting the relationships between energy and power stocks and the QQQ ETF.

The report will also discuss the implications of the findings for investors and traders, offering recommendations based on the observed patterns and trends.

### **Visualizations:**

To effectively communicate the relationships between stock prices and the QQQ ETF, a variety of visualizations will be created. These may include line charts showing stock price trends over time, scatter plots illustrating the correlation between stock prices and the ETF, and bar charts comparing model performance metrics.

Visualizations will be designed to be intuitive and informative, aiding in the interpretation of the analysis results.

### **Model Performance Metrics and Comparisons:**

A comprehensive evaluation of the models' performance will be provided, including metrics such as R-squared, RMSE, accuracy, precision, recall, and F1-score. These metrics will be presented in a clear and concise manner, allowing for easy comparison between models.

The report will discuss the strengths and limitations of each model, providing insights into their suitability for different types of analysis and prediction tasks.

### **Codebase with Documentation:**

The complete codebase for the project will be made available, ensuring transparency and reproducibility of the analysis. The code will be well-documented, with clear explanations of each step in the process, from data preprocessing to model implementation and evaluation.

The documentation will also include instructions for running the code and reproducing the results, making it accessible to other researchers and practitioners interested in exploring the energy and power stock market.

## **8. Conclusion:**

In conclusion, this project seeks to unravel the intricate dynamics between the stock prices of energy and power companies and the QQQ ETF, leveraging the power of machine learning techniques. By meticulously analyzing historical data, employing a diverse set of regression and classification models, and rigorously evaluating their performance, this research aims to shed light on the predictive relationships and patterns within the financial markets. The findings of this study are expected to provide valuable insights for investors and financial analysts, enhancing their decision-making process and contributing to a deeper understanding of market behavior. Ultimately, this project aspires to bridge the gap between data-driven analysis and practical investment strategies, empowering stakeholders with a more informed and analytical approach to navigating the complexities of the stock market.