

WESTMINSTER BUSINESS SCHOOL

UNIVERSITY OF
VISION
STRATEGY
OPPORTUNITY
WESTMINSTER 

Machine Learning and Random Trading. A lesson from Intesa SanPaolo S.p.A.

Module Title: **Artificial Intelligence and Machine Learning in Finance
Services**

Module Code: **7FNCE043W**

Module Leader: **Yao, Yumei**

Total Words Count without Tables, Appendices and References: **1954 words.**

Semester 2, 2024/2025

Riccardo Negrisoni
03/03/2025

Table of Contents

1	Company Overview	1
2	Data Collection and Processing	1
2.1	Feature Engineering and Explanatory Data Analysis	1
2.2	Target Variable	3
2.3	A final Overview of all the features	3
3	Machine Learning Analysis	4
4	Confusion Matrices and Classification Reports	4
4.1	Logistic Regression	4
4.2	Extra Trees Classifier	5
5	Cross-Validation and ROC Curves	5
6	Feature Importance and Class-Conditional Means	6
7	Market and Strategy Cumulative Returns Analysis	7
8	Final Results and Literature Review	7
9	Conclusion	7
	References	8
	Appendix A: A Naive Bayes approach	9
	Appendix B: ML to predict High-Variance	10
	Appendix C: Theoretical Foundations	11
	Appendix D: Code	12

1 Company Overview

Intesa Sanpaolo S.p.A., which was formed in 1998, is headquartered in Turin (Italy). It is the backbone of the Italian banking industry since it offers a wide range of financial products and services through six business segments: Banca dei Territori, IMI Corporate and Investment Banking, International Subsidiary Banks, Asset Management, Private Banking, and Insurance. The bank's wide portfolio unequivocally includes lending and deposit operations, structured finance, corporate and investment banking, as well as full-service asset and wealth management products. Pioneering in its approach, Intesa Sanpaolo went the extra mile with digital transformation in order to deliver the best customer experience. Its digital channels, such as the Intesa Sanpaolo Mobile application and the innovative digital bank, Isybank, have secure, easy-to-use products for individuals, families, and companies.

The official [Sanpaolo \(2024\)](#) history page reports that the history of the Group lies in a long tradition of banking innovation and resilience dating back decades. The integration of Banca Intesa with Sanpaolo IMI in 2007 was a turning point, becoming the largest banking group in Italy while upholding its strong regional presence. This solid tradition, with continuous investments in sustainable finance and digital technology, underlies Intesa Sanpaolo's leading position in the international markets as well.

The Bank is also the primary investor in the Research Institute Collegio Carlo Alberto in Turin, Italy. Its involvement in economic and financial research positions it at the forefront of innovation in the 21st century.

2 Data Collection and Processing

The download of Intesa Sanpaolo's historical stock data, from Yahoo Finance, covers the period from January 1st 2004, to December 31st 2024 in order to have 21 full years of daily observations.

I removed missing values to ensure consistency¹. The original columns (**Open**, **High**, **Low**, **Close**, **Volume**) are retained.

2.1 Feature Engineering and Explanatory Data Analysis

Next, several feature-engineering steps are applied, as well as Data analysis of these features.

(1) H-L (High minus Low) and O-C (Close minus Open) are calculated to capture daily price ranges and intraday shifts; [Karpoff \(1987\)](#) discussed the inverse relationship between stock prices

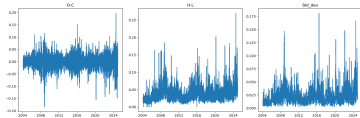


Figure 1: O-C, H-L, and Std.dev Trends Over Time.

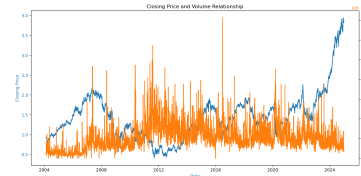


Figure 2: Closing Price (blue) and Volume (orange) Over Time.

and the volume exchanged under specific market conditions. In my analysis, this result seems to be proved (Fig 4).

¹This step is included in the code for completeness. In practice it was not necessary due to the lack of any missing value.

(2) Three moving averages of closing price (3day MA, 10day MA, 30day MA) are created to reveal short- and medium-term trends without introducing look-ahead bias.

As we can see from figure 3, all the MA series are highly correlated. Because of this, including all of them in a machine learning model can introduce redundancy and multicollinearity. In other words, each moving average provides overlapping information about the same underlying price data. This redundancy can distort feature-importance measures and/or may deteriorate predictive performance. In addition, Intesa Sanpaolo has demonstrated strong resilience in the 21st century, showing a solid ability to increase its value after the COVID-19 pandemic.

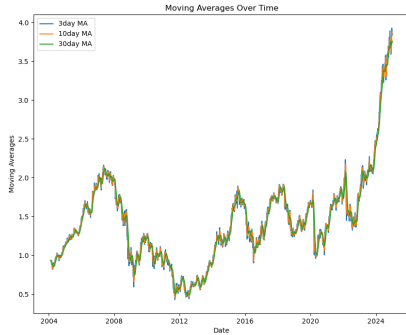


Figure 3: Moving Averages (3-day, 10-day, 30-day) Over Time.

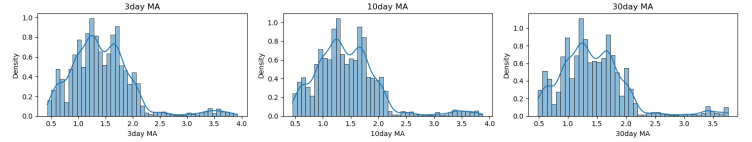


Figure 4: Histograms and Density Estimates of 3-day, 10-day, and 30-day MAs.

Figure 4 shows histograms and kernel density estimates for three moving averages applied to the same underlying time series. Because the 3-day MA reacts quickly to price changes, it yields a broader variable distribution, while the 30-day MA smooths out short-term fluctuations and produces a more peaked shape. The 10-day MA sits between these extremes, smoothing short term dynamics and capturing medium term ones. As expected, all three panels capture the same fundamental behavior.

(3) A rolling standard deviation (`Std_dev`) is derived over a five-day window to reflect short-term volatility. The code also computes daily log-returns, defined as the natural logarithm of the ratio of consecutive closing prices, to stabilize the variance and reduce non-stationarity. Before the analysis of the volatility, it is necessary to make sure computed returns are a stationary series. In figure 5 is clear that returns move around zero.

A Generalized AutoRegressive Conditional Heteroskedasticity (GARCH) model is used to model time-varying volatility and the conditional variance. It assumes that the return series is at least weakly stationary so that the conditional variance process is well-defined and stable over time. If the returns are non-stationary, the GARCH parameters and the estimated variance dynamics can become unreliable or even meaningless. As we have seen the first moment (i.e. the mean) of the returns is constant over time. (`GARCH.t_variance`) is appended to the dataset.

Fitting a Garch (1,1) model with a t-distribution allows us to capture the volatility of the stocks removing excessive sensitivity to the outliers ². Figure 6 shows us the clusters of volatility in our data, reflecting the volatility clustering phenomenon commonly observed in financial time series (Engle, 1982).

²This is clear since the t-distributions allows return to have fatter tails.

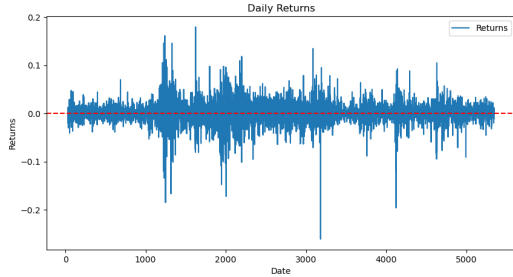


Figure 5: Daily Returns (with zero baseline in red).

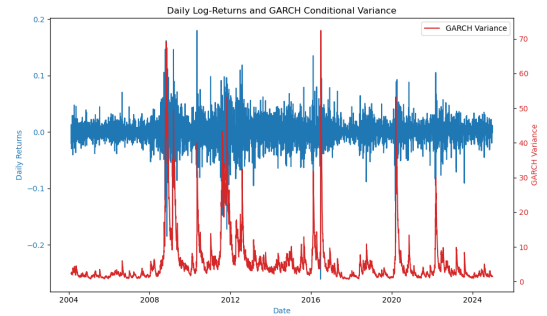


Figure 6: Daily Log-Returns (blue) and GARCH(1,1) Conditional Variance (red).

2.2 Target Variable

A binary target variable (`Price_Rise`) indicates whether the following day’s closing price is higher than the current day’s closing price. All rows with newly introduced NaN values from these rolling operations are dropped.

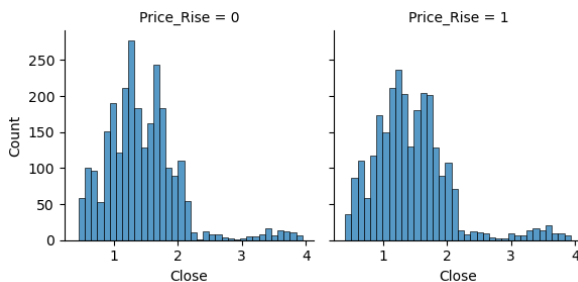


Figure 7: Distribution of Closing Prices for Rising vs. Non-Rising Days.

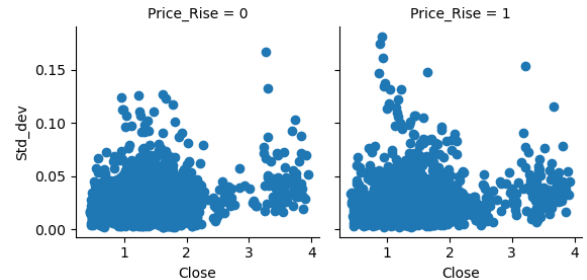


Figure 8: Scatter of Close vs. Std.dev, Separated by Price Rise Categories.

In the first pair of histograms (Fig 7), both distributions peak near the same “Close” value, but when the price rises (`Price_Rise` = 1) there is a slightly greater concentration of higher “Close” values. This suggests that larger closing prices are somewhat more likely to be associated with a price increase on the following day.

Figure 8 shows the second pair of scatter plots, where the “Std_dev” is plotted against “Close” for the two `Price_Rise` categories. Both clouds mostly overlap, even though the points for `Price_Rise` = 1 tend to have lower standard deviations at higher price levels, indicating that when prices are already high and continue to rise, the volatility may be comparatively lower (Nelson, 1991).

2.3 A final Overview of all the features

From the correlation matrix, we observe that `Open`, `High`, `Low`, and `Close` are highly correlated, as expected for price-related variables. Similarly, the moving averages (3day MA, 10day MA, and 30day MA) show strong mutual correlations and closely follow the main price indicators. In contrast, `Volume` and `Returns` show weaker correlations with price variables, suggesting that trading activity and daily return fluctuations do not always move in sync with raw prices. The binary indicator `Price_Rise` has only a modest correlation with other features, indicating that price increases on a given day depend on multiple factors rather than a single variable.

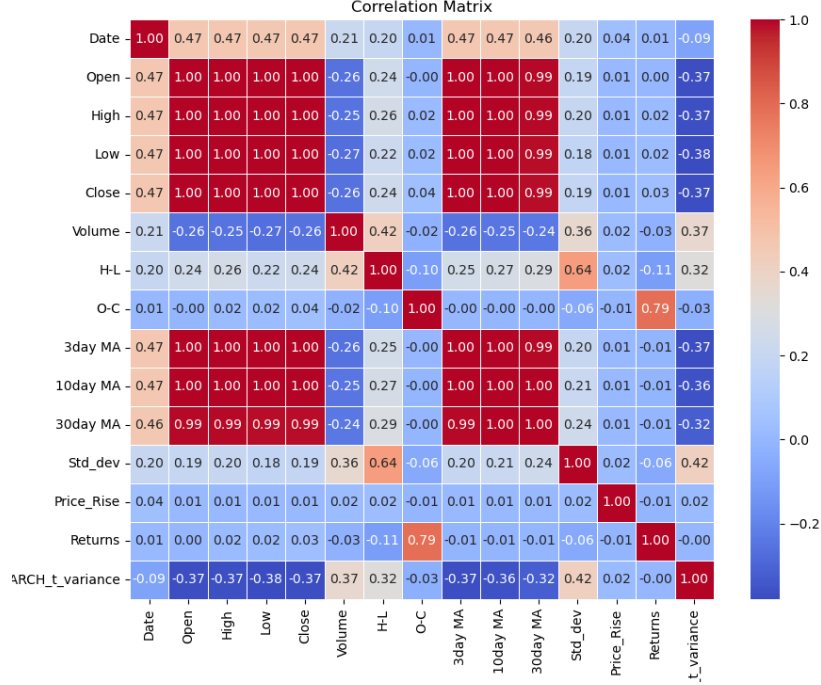


Figure 9: Correlation Matrix of All Features.

3 Machine Learning Analysis

Several classification models were trained to predict whether Intesa Sanpaolo’s stock price would rise on the next trading day. Specifically, I examined Logistic Regression and Extra Trees (Gaussian Naive Bayes was also performed. See Appendix 9). All models used the same feature set: **Volume**, **H-L** (High minus Low), **O-C** (Close minus Open), **Std_dev** (5-day rolling standard deviation of closing price), and **3day MA**. After splitting the dataset into training (85%) and testing (15%)³, **X_train** and **X_test** have been scaled using standard transformation.

The theoretical foundations of the ML models and Cross-Validation can be found in Appendix 9

4 Confusion Matrices and Classification Reports

4.1 Logistic Regression

The confusion matrix for Logistic Regression (Fig 18) shows relatively balanced true negatives (correctly predicting “Down”) and true positives (correctly predicting “Up”), but also the misclassifications (false positives and false negatives) are substantial.

The classification report (Fig 19) shows that the model’s overall accuracy is around 0.51, only slightly above random guessing. Class 0 is identified more effectively in terms of recall, meaning the model catches most of the actual class 0 instances but also misclassifies a fair number of them. Conversely, class 1 has higher precision but lower recall, indicating that while the model is more selective when predicting class 1, it fails to capture many true positives in that category. The macro-averaged F_1 score (0.50) confirms that the model’s performance is relatively balanced between the

³This proportion has been checked and it appears to be the more robust

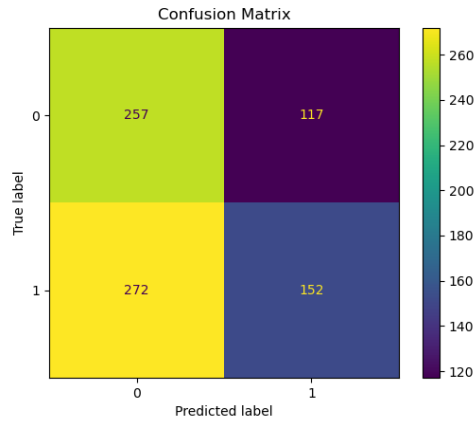


Figure 10: Confusion Matrix for Logistic Regression

Logistic Regression Classification Report:					
	precision	recall	f1-score	support	
0	0.49	0.69	0.57	374	
1	0.57	0.36	0.44	424	
accuracy			0.51	798	
macro avg	0.53	0.52	0.50	798	
weighted avg	0.53	0.51	0.50	798	

Figure 11: Classification Report for Logistic Regression

two classes.

4.2 Extra Trees Classifier

The Extra Trees Classifier exhibits a similar pattern. The distribution of misclassifications is again large, resulting in overall accuracy close to 50%. Despite Extra Trees being a method capable of capturing nonlinearities, the cross-validation scores confirm that it does not outperform the baseline random prediction in this particular application (see Fig 14).

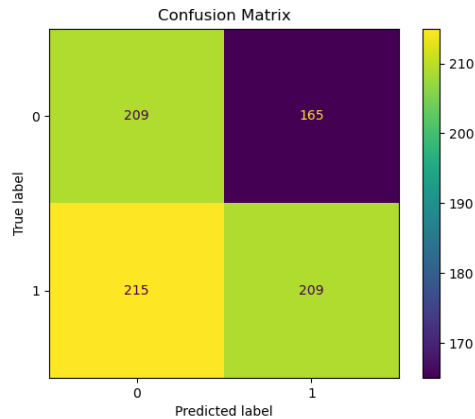


Figure 12: Confusion Matrix for Extra Trees

Extra Trees Classification Report:					
	precision	recall	f1-score	support	
0	0.49	0.57	0.53	374	
1	0.55	0.47	0.51	424	
accuracy			0.52	798	
macro avg	0.52	0.52	0.52	798	
weighted avg	0.52	0.52	0.52	798	

Figure 13: Classification Report for Extra Trees

5 Cross-Validation and ROC Curves

Cross-validation using five folds ($cv=5$) for both Logistic Regression and Extra Trees confirms that mean accuracies remain around 0.50, with standard deviations indicating that performance fluctuates but does not clearly surpass random guessing.

The ROC curve illustrates a model's performance across different classification thresholds. Specifically, it plots the True Positive Rate (how many positive cases are correctly identified) against

```

Logistic Regression - Mean Accuracy: 0.50
Logistic Regression - Standard Deviation: 0.00
Extra Trees - Mean Accuracy: 0.48
Extra Trees - Standard Deviation: 0.01

```

Figure 14: Cross-Validation for Logistic Regression and Extra Trees

the False Positive Rate (how many negative cases are incorrectly classified as positive). In my analysis, the ROC curves for both models exhibit Area Under the Curve (AUC) values near 0.52–0.54, highlighting minimal discriminative power. The models are more or less as good as random guessing.

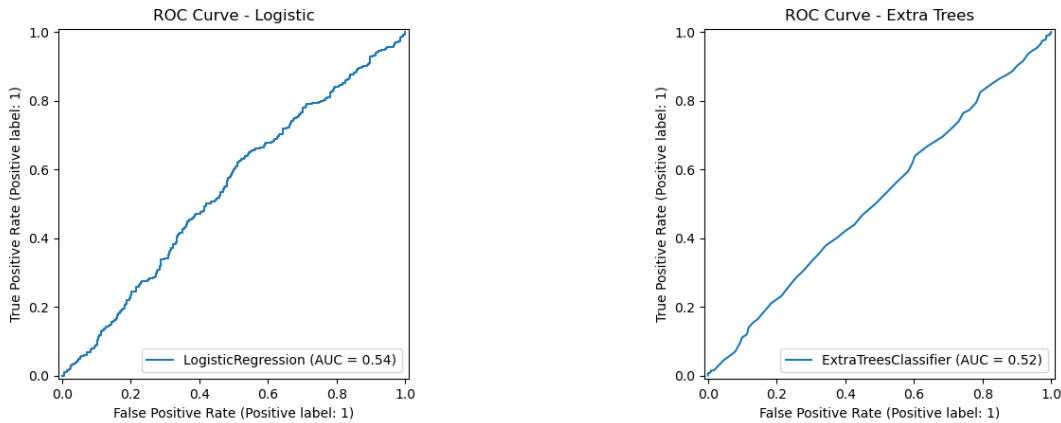


Figure 15: ROC Curves for Logistic Regression (left) and Extra Trees (right).

6 Feature Importance and Class-Conditional Means

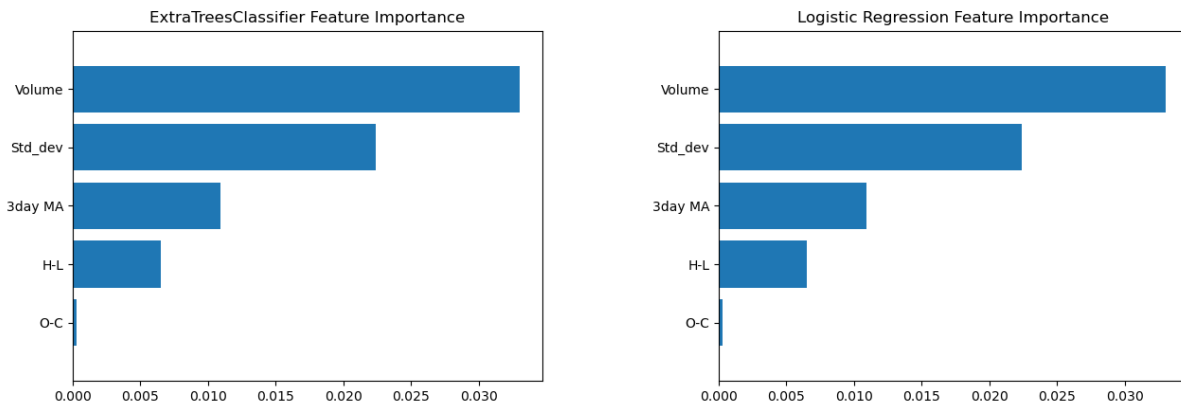


Figure 16: Feature Importance for Extra Trees (left) and Logistic Regression Coefficients (right).

Both Extra Trees and Logistic Regression highlight **Volume** as the most influential predictor for deciding whether the price will rise the next day. **Std_dev** (the short-term volatility measure) also emerges as a key factor. **3day MA** is moderately important, while **H-L** and **O-C** appear to have smaller impacts. Specifically, **O-C** was removed from the regressors after it was found to be unimportant. This step improved the overall accuracy to 53% for the Logistic Regression.

7 Market and Strategy Cumulative Returns Analysis

In figure 17, the red line represents the cumulative market returns calculated by summing the daily log returns of the Intesa Sanpaolo stock. The blue line represents the cumulative returns of a simple long-short strategy driven by model predictions. Specifically, if the model predicts an “Up” day, we go long and capture the market’s next-day return; if it predicts a “Down” day, we effectively short the market and profit when the price falls.



Figure 17: Cumulative Market Returns vs. Strategy Returns (Logistic Regression).

While the strategy can sometimes outperform the market, it remains quite unstable overall. This highlights how challenging it is to build a consistently profitable trading approach based solely on daily price predictions, especially when using a limited set of features and basic models. Without deeper market knowledge or more advanced techniques, short-term forecasts may not be reliable enough to ensure steady profits, particularly when factoring in real-world risks and trading costs.

8 Final Results and Literature Review

Recent studies question whether machine learning (ML) is always the best choice for predicting bank returns. For example, [Makridakis et al. \(2018\)](#) found that many ML models don’t consistently perform better than traditional statistical methods and often struggle with problems like overfitting and limited data. Similarly, [Bluwstein et al. \(2021\)](#) showed that while non-linear ML models can slightly improve banking crisis predictions, their complexity often cancels out those benefits. Lastly, [Antonopoulou et al. \(2023\)](#) pointed out that financial markets are so unpredictable that even advanced ML models can sometimes perform no better than random guessing. These findings suggest that ML’s advantages in financial forecasting are not guaranteed and should be carefully compared to simpler models.

9 Conclusion

Machine learning models struggle to reliably predict the next-day price direction of Intesa Sanpaolo’s stock. The classification reports and cross-validation results consistently show accuracy staying close to 50%, with AUC values ranging from 0.52 to 0.54. While some features like `Volume`, `Std.dev`, and `3day MA` seem to have more influence, the unpredictable nature of short-term price changes and possible shifts in financial data patterns limit the models’ effectiveness. This aligns with a well-known challenge in quantitative finance: predicting daily price movements is extremely difficult, and even small improvements beyond random guessing often require more complex features, larger datasets or advanced modeling techniques. These results emphasize the need to carefully test models on unseen data and be cautious about overfitting or reading too much into slight accuracy gains in noisy financial markets.

References

- Antonopoulou, H. et al. (2023), ‘Utilizing machine learning to reassess the predictability of bank stocks’, *Emerging Science Journal* **7**(3), 724–736.
- Bluwstein, K. et al. (2021), Credit growth, the yield curve and financial crisis prediction: Evidence from a machine learning approach, Technical Report 2614, European Central Bank.
- Boschetti, A. and Massaron, L. (2018), *Python data science essentials: a practitioner’s guide covering essential data science principles, tools, and techniques*, 3rd edn, Packt Publishing, Birmingham, UK.
- Cox, D. R. (1958), ‘The regression analysis of binary sequences’, *Journal of the Royal Statistical Society. Series B (Methodological)* **20**(2), 215–242.
- Dixon, M., Halperin, I. and Bilokon, P. (2020), *Machine learning in finance: from theory to practice*, Springer, Switzerland.
- Engle, R. F. (1982), ‘Autoregressive conditional heteroscedasticity with estimates of the variance of united kingdom inflation’, *Econometrica* **50**(4), 987–1007.
- Geurts, P., Ernst, D. and Wehenkel, L. (2006), ‘Extremely randomized trees’, *Machine Learning* **63**(1), 3–42.
- Karpoff, S. (1987), ‘The relation between price changes and trading volume: A survey’, *Journal of Financial Economics* **17**(1), 153–206.
- Makridakis, S., Spiliotis, E. and Assimakopoulos, V. (2018), ‘Statistical and machine learning forecasting methods: Concerns and ways forward’, *PLOS ONE* **13**(3), e0194889.
- Nelson, D. B. (1991), ‘Conditional heteroskedasticity in asset returns: A new approach’, *Econometrica* **59**(2), 347–370.
- Orlando, G., Chironna, G. and Penikas, H. (2022), A default prediction (pd) model for italian banks: An empirical analysis with econometric and machine learning approaches, in ‘Proceedings of the 3rd International Conference on Modern Management based on Big Data (MMBD)’.
- Sanpaolo, G. I. (2024), ‘Storia’, <https://group.intesasanpaolo.com/it/chi-siamo/storia>. Accessed: 3 March 2025.

Appendix A

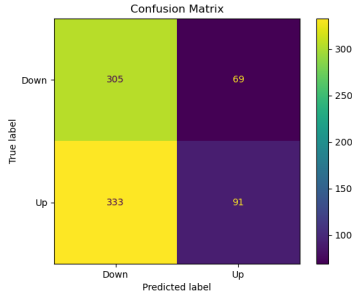


Figure 18: Confusion Matrix for Naive Bayes

	precision	recall	f1-score	support
0	0.48	0.82	0.60	374
1	0.57	0.21	0.31	424
accuracy			0.50	798
macro avg	0.52	0.52	0.46	798
weighted avg	0.53	0.50	0.45	798

Figure 19: Classification Report for Naive Bayes

Gaussian Naive Bayes also struggles to exceed 50% accuracy. These results further validate the difficulty of predicting daily direction of price movement.

In this report, class 0 has a high recall (0.82) but a low precision (0.48), indicating the model correctly identifies most class-0 cases but also misclassifies many class-1 instances as class 0. For class 1, the precision is somewhat better (0.57) but recall is very low (0.21), meaning many true positives for class 1 are missed. Overall accuracy sits at 0.50, close to random guessing, and the F_1 scores reflect that neither class is predicted reliably.

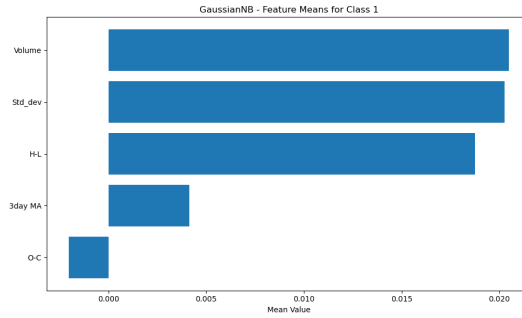


Figure 20: GaussianNB Class-Conditional Means for “Up” Class (Class 1).

Since Gaussian Naive Bayes does not directly provide feature importance via coefficients, the class-conditional means (θ_{1j}) for the “Up” class (Class 1) are used as a proxy. A larger mean in a specific feature suggests that if a day is classified as “Up,” that feature typically has a higher average value. Here, **Volume** and **Std_dev** again rank highest, indicating that “Up” days tend to be associated with higher trading volume and greater recent volatility. **H-L** (daily range) also shows a noticeable mean, implying that upward movements often coincide with a wider trading range.



Figure 21: Cumulative Market Returns vs. Strategy Returns (GaussianNB).

Also this model is not able to beat cumulative market returns.

Appendix B

In Figure 22, the confusion matrix shows that the model very well predicts the positive label, missing most negative cases. Despite this, the ROC curve in Figure 23 yields an AUC of 0.79, suggesting that adjusting the classification threshold could improve separation between the two volatility regimes. Figure 24 indicates a nearly balanced count of correct vs. incorrect predictions, consistent with around 50% accuracy. Finally, Figure 25 reveals that *Volume* dominates the feature importance, overshadowing other inputs in driving the model's volatility-regime predictions. This feature is a pure novelty of this proposed research.

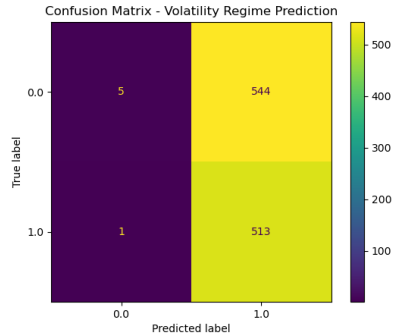


Figure 22: Confusion Matrix — Volatility Regime Prediction

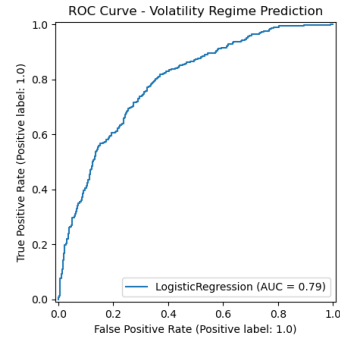


Figure 23: ROC Curve — Volatility Regime Prediction

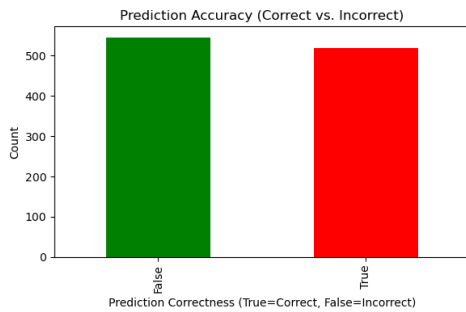


Figure 24: Prediction Accuracy (Correct vs. Incorrect)

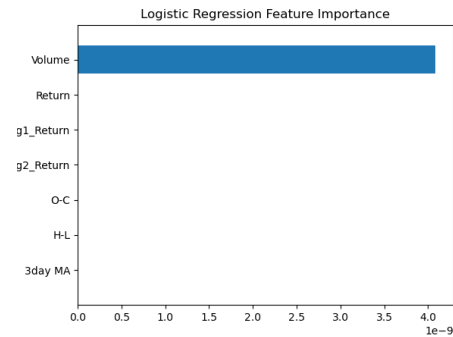


Figure 25: Logistic Regression Feature Importance

Appendix C

Logistic Regression

The logistic function (Cox, 1958) transforms a linear combination of input features into a probability between 0 and 1. Specifically, if x is a vector of predictors and β is the vector of coefficients, then the probability p of a positive class is given by

$$p(y = 1|x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n)}},$$

where β_0 is the intercept and β_1, \dots, β_n are the model parameters. This formulation naturally handles binary outcomes by mapping any real-valued input to a value in $(0, 1)$. Training proceeds by maximizing the likelihood of the observed data under this model.

Extra Trees

Extra Trees method (Geurts et al., 2006) creates multiple decision trees to make predictions. Each tree is built using the original data, but instead of choosing the best split points, it picks them randomly. This randomness, along with selecting features randomly, helps prevent overfitting and makes the model more stable. In the end, all the trees work together by voting on the final result, either by averaging their predictions for numerical values (regression) or by choosing the most common answer (classification).

Cross Validation Analysis

Cross-validation is a technique used to test the performance of a model on new data. In k-fold cross-validation, the data is divided into k smaller sets. The model is trained on k-1 sets and tested on the remaining one. This is done k times, with a different test fold used every time. The performance is found by taking an average across all the folds, based on metrics like accuracy, precision, or recall. Cross-validation, by testing the model on multiple pieces of the data, avoids overfitting and gives a better indication of how the model will perform on new data.

Appendix D

The entire project is available at my [GitHub repository](#)

For completeness and clarity, I report the code in this document:

```
1  # Riccardo Negrisoni
2  # AI and ML
3
4
5
6  # To correctly run the code with no errors PLEASE
7  # use the command : pip install arch
8
9
10
11
12
13  from IPython import get_ipython
14  get_ipython().magic('reset -sf')
15  get_ipython().magic('clear')
16  %%
17  #pip install yfinance pandas openpyxl
18  import warnings
19  warnings.simplefilter(action='ignore', category=FutureWarning)
20  import yfinance as yf
21  import pandas as pd
22  from arch import arch_model
23  import numpy as np
24  import matplotlib.pyplot as plt
25  import seaborn as sns
26  from sklearn import datasets
27  from sklearn.decomposition import PCA
28  from sklearn.preprocessing import StandardScaler
29  from sklearn.model_selection import train_test_split
30  from sklearn.model_selection import cross_val_score
31  from sklearn.linear_model import LogisticRegression
32  from sklearn.linear_model import LinearRegression
33  from sklearn.naive_bayes import GaussianNB
34  from sklearn.ensemble import ExtraTreesClassifier
35  from sklearn.metrics import mean_absolute_error
36  from sklearn.metrics import accuracy_score, roc_curve, classification_report
37  from sklearn.metrics import confusion_matrix, make_scorer, roc_auc_score
38  from sklearn.metrics import RocCurveDisplay, ConfusionMatrixDisplay
39
40
41  from yahooquery import Ticker
42  import os
43
44  %% Section 0
45  ## Relevant Information and preliminary analysis of Intesa
46  # FETCH FUNDAMENTAL DATA FOR INTESA SAN PAOLO
47
48  # Create a Ticker object for Intesa San Paolo (ticker "ISP.MI")
49  ISP = Ticker("ISP.MI")
50
51  # Retrieve and print the ESG scores
```

```

52 print("ESG_Scores:")
53 print(ISP.esg_scores)
54
55 # Retrieve and print key statistics
56 print("\nKey_Statistics:")
57 print(ISP.key_stats)
58
59 # Retrieve and print the summary profile
60 print("\nSummary_Profile:")
61 print(ISP.summary_profile)
62
63 # Retrieve and print institutional ownership details
64 print("\nInstitutional_Ownership:")
65 print(ISP.institution_ownership)
66
67 # Retrieve and print fund ownership details
68 print("\nFund_Ownership:")
69 print(ISP.fund_ownership)
70
71 # Retrieve and print the quarterly balance sheet by specifying frequency="q"
72 print("\nQuarterly_Balance_Sheet:")
73 print(ISP.balance_sheet(frequency="q"))
74
75 # Retrieve and print the cash flow statement
76 print("\nCash_Flow:")
77 print(ISP.cash_flow())
78
79 # Retrieve and print the income statement
80 print("\nIncome_Statement:")
81 print(ISP.income_statement())
82
83 ### IMPORTING DATA
84
85 # Setting my working directory
86 os.chdir(r"C:\Users\Utente\Desktop\WESTMINSTER\AI_and_ML\Intesa_Project")
87 folder_path = r"C:\Users\Utente\Desktop\WESTMINSTER\AI_and_ML\Intesa_Project\
    Figures"
88
89 # Download Intesa Sanpaolo data
90 ticker = "ISP.MI" ### Intesa San Paolo S.p.A. is traded on the Milan Stock
    Exchange with the ticker "ISP.MI"
91 start_date = "2004-01-01"
92 end_date = "2024-12-31"
93
94 ISPdata = yf.download(ticker, start=start_date, end=end_date)
95 ISPdata.dropna(inplace=True)
96
97 # If columns is a MultiIndex (e.g., top level is "ISP.MI"), drop that level
98 if isinstance(ISPdata.columns, pd.MultiIndex):
99     ISPdata.columns = ISPdata.columns.droplevel(1)
100
101 # Move the date index into a normal column
102 ISPdata.reset_index(inplace=True)
103
104 # Keep only the columns you want
105 ISPdata = ISPdata[['Date', 'Open', 'High', 'Low', 'Close', 'Volume']]
106

```

```

107 # Export to Excel without the index
108 FileName = "Intesa_stock_data.xlsx"
109 ISPdata.to_excel(FileName, index=False)
110
111 print(f"Stock_data_saved_to_{FileName}")
112
113
114 %% Section 2 : Feature Engeneering
115
116 # Create new features based on the stock prices
117 ISPdata['H-L'] = ISPdata['High'] - ISPdata['Low']
118 ISPdata['O-C'] = ISPdata['Close'] - ISPdata['Open']
119
120 # Calculate moving averages for the Close price using a 1-day lag to avoid
    lookahead bias
121 ISPdata['3day_MA'] = ISPdata['Close'].shift(1).rolling(window=3).mean()
122 ISPdata['10day_MA'] = ISPdata['Close'].shift(1).rolling(window=10).mean()
123 ISPdata['30day_MA'] = ISPdata['Close'].shift(1).rolling(window=30).mean()
124
125 # Calculate the rolling standard deviation of the Close price over a 5-day
    window
126 ISPdata['Std_dev'] = ISPdata['Close'].rolling(window=5).std()
127
128 # Create a binary target variable: 1 if the next day's Close is higher than
    today's Close, else 0
129 ISPdata['Price_Rise'] = np.where(ISPdata['Close'].shift(-1) > ISPdata['Close'],
    1, 0)
130
131 # Drop any rows with missing values from shifting/rolling operations
132 ISPdata = ISPdata.dropna()
133
134 # Display the first few rows of the enhanced dataset
135 print(ISPdata.head())
136 ISPdata.index.name = "Date"
137 ISPdata.to_excel(FileName, index=False)
138 print(f"Stock_data_saved_to_{FileName}")
139
140
141 %% Extra section 2.1: ANALYSIS OF THE VARIANCE
142
143 # Calculate returns (ensure the series is stationary)
144 returns = np.log(ISPdata['Close'] / ISPdata['Close'].shift(1)).dropna()
145 #Checking for stationarity
146 plt.figure(figsize=(10, 5))
147 plt.plot(returns.index, returns, label='Returns')
148 plt.axhline(0, color='red', linestyle='--')
149 plt.title('Daily_Returns')
150 plt.xlabel('Date')
151 plt.ylabel('Returns')
152 plt.legend()
153 plt.savefig(os.path.join(folder_path, "fig0.ClosingVol.Relationship.png"))
154 plt.show()
155
156 ISPdata['Returns'] = returns
157
158 # Fit a GARCH(1,1) model with t-distribution
159 model = arch_model(returns, vol='Garch', p=1, q=1, dist='t', rescale=True)

```



```

160 fit = model.fit(displ='off')
161
162 # Compute conditional variance (volatility squared)
163 ISPdata = ISPdata.iloc[1:] # Align with returns dropna
164 ISPdata['GARCH_t_variance'] = fit.conditional_volatility ** 2
165
166 print(fit.summary())
167 ISPdata.to_excel(FileName, index=False)
168 print(f"Stock_data_saved_to_{FileName}")
169 fig, ax1 = plt.subplots(figsize=(10, 6))
170
171 # Create a figure and an axis (ax1) for the first plot
172 fig, ax1 = plt.subplots(figsize=(10, 6))
173
174 # Plot Daily Returns on the left y-axis
175 color1 = 'tab:blue'
176 ax1.set_xlabel('Date', color=color1)
177 ax1.set_ylabel('Daily_Returns', color=color1)
178 ax1.plot(ISPdata['Date'], ISPdata['Returns'], color=color1, label='Daily_Returns')
179 ax1.tick_params(axis='y', labelcolor=color1)
180
181 # Create a second axis (ax2) sharing the same x-axis
182 ax2 = ax1.twinx()
183
184 # Plot GARCH Variance on the right y-axis
185 color2 = 'tab:red'
186 ax2.set_ylabel('GARCH_Variance', color=color2)
187 ax2.plot(ISPdata['Date'], ISPdata['GARCH_t_variance'], color=color2, label='GARCH_Variance')
188 ax2.tick_params(axis='y', labelcolor=color2)
189
190 plt.title('Daily_Log-Returns_and_GARCH_Conditional_Variance')
191 fig.tight_layout()
192 plt.legend()
193 plt.savefig(os.path.join(folder_path, "fig20.Garch&Returns.png"))
194 plt.show()
195
196
197
198
199 ### PLOTS
200
201 # Figure 1: Plot for Closing Price and Volume
202 fig, ax1 = plt.subplots(figsize=(12, 6))
203 color = 'tab:blue'
204 ax1.set_xlabel('Date', color=color)
205 ax1.set_ylabel('Closing_Price', color=color)
206 ax1.plot(ISPdata['Date'], ISPdata['Close'], label='Closing_Price', color=color)
207 ax1.tick_params(axis='y', labelcolor=color)
208
209 ax2 = ax1.twinx()
210 color = 'tab:orange'
211 ax2.set_ylabel('Volume', color=color)
212 ax2.plot(ISPdata['Date'], ISPdata['Volume'], label='Volume', color=color)
213 ax2.tick_params(axis='y', labelcolor=color)
214

```

```

215 plt.title('Closing_Price_and_Volume_Relationship')
216 fig.tight_layout()
217 plt.savefig(os.path.join(folder_path, "fig1.ClosingVol.Relationship.png"))
218 plt.show()
219
220 # Figure 2: Plot for 3day, 10day, 30day MAs
221 plt.figure(figsize=(10, 8))
222 plt.plot(ISPdata['Date'], ISPdata['3day_MA'], label='3day_MA')
223 plt.plot(ISPdata['Date'], ISPdata['10day_MA'], label='10day_MA')
224 plt.plot(ISPdata['Date'], ISPdata['30day_MA'], label='30day_MA')
225 plt.title('Moving_Averages_Over_Time')
226 plt.xlabel('Date')
227 plt.ylabel('Moving_Averages')
228 plt.legend()
229 plt.savefig(os.path.join(folder_path, "fig2.MA.png"))
230 plt.show()
231
232 # Figure 3: Plots for O-C, H-L, and Std_dev
233 fig, axes = plt.subplots(1, 3, figsize=(15, 5))
234 axes[0].plot(ISPdata['Date'], ISPdata['O-C'])
235 axes[0].set_title('O-C')
236 axes[1].plot(ISPdata['Date'], ISPdata['H-L'])
237 axes[1].set_title('H-L')
238 axes[2].plot(ISPdata['Date'], ISPdata['Std_dev'])
239 axes[2].set_title('Std_dev')
240
241 plt.tight_layout()
242 plt.savefig(os.path.join(folder_path, "fig3.Comparisons.png"))
243 plt.show()
244
245 # Figure 4: 3day, 10day, 30day MA histogram & density
246 fig, axes = plt.subplots(1, 3, figsize=(15, 3))
247 sns.histplot(data=ISPdata, x='3day_MA', kde=True, stat='density', ax=axes[0])
248 axes[0].set_title('3day_MA')
249
250 sns.histplot(data=ISPdata, x='10day_MA', kde=True, stat='density', ax=axes[1])
251 axes[1].set_title('10day_MA')
252
253 sns.histplot(data=ISPdata, x='30day_MA', kde=True, stat='density', ax=axes[2])
254 axes[2].set_title('30day_MA')
255
256 plt.tight_layout()
257 plt.savefig(os.path.join(folder_path, "fig4.Hist&Density.png"))
258 plt.show()
259
260 # Figure 5: Price Rise & Fall Count Plot + Scatter
261 chart = sns.FacetGrid(ISPdata, col='Price_Rise')
262 chart.map(sns.histplot, 'Close')
263 plt.savefig(os.path.join(folder_path, "figure5_hist.png"))
264 plt.show()
265
266 chart = sns.FacetGrid(ISPdata, col='Price_Rise')
267 chart.map(plt.scatter, 'Close', 'Std_dev')
268 plt.savefig(os.path.join(folder_path, "figure5_scatter.png"))
269 plt.show()
270
271 #Fig 6: Corr Matrix

```

```

272 corr_matrix = ISPdata.corr()
273 plt.figure(figsize=(10, 8))
274 sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt=".2f", linewidths=.5)
275 plt.title("Correlation_Matrix")
276 plt.savefig(os.path.join(folder_path, "fig6.CorrMatrix.png"))
277 plt.show()
278
279
280
281 ###
282 ISPdata.info()
283 ISPdata.describe()
284 ISP.summary_profile
285
286 ISP.institution_ownership
287
288 ISP.fund_ownership
289 ISP.history()
290
291 ### MACHINE LEARNING
292
293 feature_cols = ['Volume', 'H-L', 'O-C', 'Std_dev', '3day_MA']
294 target_col = 'Price_Rise'
295
296 X = ISPdata[feature_cols]
297
298 Y = ISPdata[target_col]
299
300 Y
301
302
303
304 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.15,
305     random_state=240)
306 X_train
307
308 # Standardize features by removing the mean and scaling to unit variance.
309 # StandardScaler transfers data to standard normally distributed data: Gaussian
310 with zero mean and unit variance
311 scaler = StandardScaler()
312 X_train = scaler.fit_transform(X_train)
313 X_test = scaler.transform(X_test)
314
315 regr = LinearRegression()
316 regr.fit(X_train, Y_train)
317 Y_pred = regr.predict(X_test)
318 print ("MAE", mean_absolute_error(Y_test, Y_pred))
319
320 ###
321 # Logistic Regression
322
323 model_logistic = LogisticRegression(random_state=101)
324 model_logistic.fit(X_train, Y_train)
325
326 # Predict on the test set
327 Y_pred_logistic = model_logistic.predict(X_test)
328 print("Logistic_Regression_Classification_Report:")

```

```

327 print(classification_report(Y_test, Y_pred_logistic))
328
329
330 # Extra Trees
331 model_extra_trees = ExtraTreesClassifier(random_state=101)
332 model_extra_trees.fit(X_train, Y_train)
333
334 # Predict on the test set
335 Y_pred_extra = model_extra_trees.predict(X_test)
336 print("Extra_Trees_Classification_Report:")
337 print(classification_report(Y_test, Y_pred_extra))
338
339 ## LOGISTIC REGRESSION PERFORMS BETTER!!!
340
341
342 ### N-fold Cross Validation
343
344 #Cross Validation for Logistic Regression
345 accuracy_scores_log = cross_val_score(model_logistic, X, Y, cv=5, scoring='
    accuracy')
346 print(f"Logistic_Regression_Mean_Accuracy:{accuracy_scores_log.mean():.2f}")
347 print(f"Logistic_Regression_Standard_Deviation:{accuracy_scores_log.std():.2f
    }")
348
349 #Cross Validation for Extra Trees
350 accuracy_scores_et = cross_val_score(model_extra_trees, X, Y, cv=5, scoring='
    accuracy')
351 print(f"Extra_Trees_Mean_Accuracy:{accuracy_scores_et.mean():.2f}")
352 print(f"Extra_Trees_Standard_Deviation:{accuracy_scores_et.std():.2f}")
353
354 ## WE ARE PERFORMING NO BETTER THAN RANDOM GUESSING
355
356
357
358 # -----
359 ### Prediction of Price Rise Using Logistic Regression on X_test Data
360
361 # Re-train or re-use model_extra_trees
362 model_logistic = LogisticRegression(random_state=101)
363 model_logistic.fit(X_train, Y_train)
364
365 Y_pred_logistic_test = model_logistic.predict(X_test)
366 classification_rep = classification_report(Y_test, Y_pred_logistic_test)
367 print("Final_Logistic_Regression_Classification_Report_on_X_test:")
368 print(classification_rep)
369
370
371 ### Confusion Matrix for logistic
372 matrix = ConfusionMatrixDisplay.from_estimator(model_logistic, X_test, Y_test)
373 plt.title('Confusion_Matrix')
374 plt.savefig(os.path.join(folder_path, "fig6.ConfMatrix.png"))
375 plt.show()
376
377
378 ### Confusion Matrix for Extra Trees
379 matrix = ConfusionMatrixDisplay.from_estimator(model_extra_trees, X_test, Y_test
    )

```

```

380 plt.title('Confusion_Matrix')
381 plt.savefig(os.path.join(folder_path, "fig7.ConfMatrix2.png"))
382 plt.show()
383
384 ### ROC Curve for Logistic Regression
385 log_disp = RocCurveDisplay.from_estimator(model_logistic, X_test, Y_test)
386 plt.title("ROC_Curve_Logistic")
387 plt.savefig(os.path.join(folder_path, "roc_curve.Logistic.png"))
388 plt.show()
389
390
391
392 ### BAYESIAN NAIVE
393
394
395 Gauss_Model = GaussianNB()
396 Gauss_Model.fit(X_train, Y_train)
397 Y_pred = Gauss_Model.predict(X_test)
398 target_names=["Down", "Up"]
399
400 print (classification_report(Y_test, Y_pred))
401
402 # Evaluate the model by means of a Confusion Matrix
403
404 matrix = ConfusionMatrixDisplay.from_estimator(Gauss_Model, X_test, Y_test,
405         display_labels=target_names)
406 plt.title('Confusion_Matrix')
407 plt.savefig(os.path.join(folder_path, "fig8.ConfMatrix3.png"))
408 plt.show()
409
410 ### FEATURE IMPORTANCE LOGISTIC
411 classifier = LogisticRegression(random_state=101)
412 classifier.fit(X_train, Y_train)
413
414 # Now retrieve the coefficients from the fitted classifier
415 importance = classifier.coef_[0]
416
417
418 feature_names=X.columns
419
420
421 indices = np.argsort(importance)
422 range1 = range(len(importance[indices]))
423 plt.figure()
424 plt.title("Logistic_Regression_Feature_Importance")
425 plt.barh(range1, importance[indices])
426 plt.yticks(range1, feature_names[indices])
427 plt.ylim([-1, len(range1)])
428 plt.savefig(os.path.join(folder_path, "fig9.Logistic.FeatureImportance.png"))
429 plt.show()
430
431 ### MARKET AND RETURN STRATEGIES WITH LOGISTIC
432
433 # Data Preprocessing
434 ISPdata['Y_pred_logistic_test'] = np.nan
435 ISPdata.iloc[len(ISPdata) - len(Y_pred_logistic_test):, -1] =

```

```

Y_pred_logistic_test # Fill the last rows with predictions
436 trade_ISPdata = ISPdata.dropna() # Drop rows without a prediction
437
438 # Computation of Market Returns
439 trade_ISPdata['Tomorrows_Returns'] = 0.
440 trade_ISPdata['Tomorrows_Returns'] = np.log(trade_ISPdata['Close'] /
441       trade_ISPdata['Close'].shift(1))
442 trade_ISPdata['Tomorrows_Returns'] = trade_ISPdata['Tomorrows_Returns'].shift
443       (-1)
444
445 # Strategy Returns based on Y_pred
446 trade_ISPdata['Strategy_Returns'] = 0.
447 trade_ISPdata['Strategy_Returns'] = np.where(
448     trade_ISPdata['Y_pred_logistic_test'] == True,
449     trade_ISPdata['Tomorrows_Returns'],
450     -trade_ISPdata['Tomorrows_Returns']
451 )
452
453 # ##### Cumulative Market and Strategy Returns
454
455 # Computation of cumulative market and strategy returns
456 trade_ISPdata['Cumulative_Market_Returns'] = np.cumsum(trade_ISPdata['Tomorrows_Returns'])
457 trade_ISPdata['Cumulative_Strategy_Returns'] = np.cumsum(trade_ISPdata['Strategy_Returns'])
458
459 # Plot of cumulative market and strategy returns based on Y_pred
460 plt.figure(figsize=(10, 3))
461 plt.plot(trade_ISPdata['Cumulative_Market_Returns'], color='red', label='Market_Returns')
462 plt.plot(trade_ISPdata['Cumulative_Strategy_Returns'], color='blue', label='Strategy_Returns')
463 plt.title("Cumulative_Market_vs._Strategy_Returns")
464 plt.legend()
465 plt.savefig(os.path.join(folder_path, "fig10.MktStrategy.png"))
466 plt.show()
467
468
469 %% FEATURES IMPORTANCE AND MARKET STRATEGY RETURNS USING NAIVE BAYES
470
471 # FEATURE "IMPORTANCE" with GaussianNB
472 # =====
473 nb_classifier = GaussianNB()
474 nb_classifier.fit(X_train, Y_train)
475
476 # Naive Bayes does not provide coefficients; as a proxy we use the
477     c l a s s conditional means.
478 # (This only makes sense if your task is binary. Here we assume class 1 is the "
479     positive" class.)
480 importance = nb_classifier.theta_[1]
481 feature_names = X.columns
482
483 # Sort the "importance" values for plotting
484 indices = np.argsort(importance)
485 range1 = range(len(importance))

```

```

484 plt.figure(figsize=(10, 6))
485 plt.title("GaussianNB-FeatureMeansforClass1")
486 plt.barh(range1, importance[indices])
487 plt.yticks(range1, feature_names[indices])
488 plt.xlabel("Mean Value")
489 plt.tight_layout()
490 plt.savefig(os.path.join(folder_path, "fig9.GaussianNB_FeatureMeans.png"))
491 plt.show()
492
493
494 # =====
495 # MARKET AND RETURN STRATEGIES WITH GAUSSIAN NAIVE BAYES
496 # =====
497
498 # Get predictions on the test set from the NB model
499 Y_pred_nb = nb_classifier.predict(X_test)
500 print("GaussianNB-ClassificationReportonX_test:")
501 print(classification_report(Y_test, Y_pred_nb))
502
503 # Data Preprocessing: Add NB predictions to ISPdata.
504 # Here we assume that ISPdata originally has all the data in order and that the
505 # predictions
506 # correspond to the last len(Y_pred_nb) rows.
507 ISPdata['Y_pred_nb'] = np.nan
508 ISPdata.iloc[len(ISPdata) - len(Y_pred_nb):, -1] = Y_pred_nb # Fill last rows
509 # with predictions
510 trade_ISPdata = ISPdata.dropna() # Drop rows without a prediction
511
512 # Compute Market Returns as log returns
513 trade_ISPdata['Tomorrows_Returns'] = np.log(trade_ISPdata['Close'] /
514 trade_ISPdata['Close'].shift(1))
515 trade_ISPdata['Tomorrows_Returns'] = trade_ISPdata['Tomorrows_Returns'].shift
516 (-1)
517
518 # Compute Strategy Returns based on the NB predictions:
519 # If Y_pred_nb is True (or 1), take the market return; else, take the negative.
520 trade_ISPdata['Strategy_Returns'] = np.where(
521     trade_ISPdata['Y_pred_nb'] == True, # adjust if your positive class is
522     # represented differently
523     trade_ISPdata['Tomorrows_Returns'],
524     -trade_ISPdata['Tomorrows_Returns']
525 )
526
527 # Compute cumulative returns
528 trade_ISPdata['Cumulative_Market_Returns'] = np.cumsum(trade_ISPdata['Tomorrows_Returns'])
529 trade_ISPdata['Cumulative_Strategy_Returns'] = np.cumsum(trade_ISPdata['Strategy_Returns'])
530
531 # Plot the cumulative returns
532 plt.figure(figsize=(10, 3))
533 plt.plot(trade_ISPdata['Cumulative_Market_Returns'], color='red', label='Market_Returns')
534 plt.plot(trade_ISPdata['Cumulative_Strategy_Returns'], color='blue', label='Strategy_Returns')
535 plt.title("Cumulative_Market_vs.Strategy_Returns_(GaussianNB)")
536 plt.legend()

```

```

532 plt.tight_layout()
533 plt.savefig(os.path.join(folder_path, "fig10.MktStrategy_GaussianNB.png"))
534 plt.show()
535
536 ### ANALYSIS OF VARIANCE WITH ML
537
538 # Use the median of GARCH_t_variance as the threshold
539 vol_threshold = ISPdata['GARCH_t_variance'].median()
540 ISPdata['HighVolatility'] = (ISPdata['GARCH_t_variance'] > vol_threshold).astype
    (int)
541
542 ISPdata['Lag1_Return'] = ISPdata['Return'].shift(1)
543 ISPdata['Lag2_Return'] = ISPdata['Return'].shift(2)
544 ISPdata['HighVolatility_Tomorrow'] = ISPdata['HighVolatility'].shift(-1)
545
546 # Drop rows with missing values caused by pct_change and shifting
547 data_ml = ISPdata.dropna(subset=['Return', 'Lag1_Return', 'Lag2_Return', 'Volume
    ', 'HighVolatility_Tomorrow'])
548
549 # Define Feature Matrix X and Target Vector Y
550 featuresNew = ['Return', 'Lag1_Return', 'Lag2_Return', 'Volume', 'H-L', 'O-C', '3
    day_MA']
551 targetNew = 'HighVolatility_Tomorrow'
552 X = data_ml[featuresNew]
553 Y = data_ml[targetNew]
554
555 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2,
    random_state=101)
556
557 # Train a Logistic Regression Model for Volatility Regime Prediction
558 modelGarch = LogisticRegression(random_state=101)
559 modelGarch.fit(X_train, Y_train)
560
561 # Predict on test set
562 Y_pred = modelGarch.predict(X_test)
563 correct_predictions = (Y_pred == Y_test)
564 print("Classification_Report:")
565 print(classification_report(Y_test, Y_pred))
566 # Plot and Save the Confusion Matrix
567 plt.figure(figsize=(6, 6))
568 cm_disp = ConfusionMatrixDisplay.from_estimator(modelGarch, X_test, Y_test)
569 plt.title("Confusion_Matrix_Volatility_Regime_Prediction")
570 plt.savefig(os.path.join(folder_path, "Volatility_Regime_ConfusionMatrix.png"))
571 plt.show()
572
573 # Plot the ROC Curve
574
575 plt.figure(figsize=(6, 6))
576 roc_disp = RocCurveDisplay.from_estimator(modelGarch, X_test, Y_test)
577 plt.title("ROC_Curve_Volatility_Regime_Prediction")
578 plt.savefig(os.path.join(folder_path, "Volatility_Regime_ROC_Curve.png"))
579 plt.show()
580
581
582 modelGarch.fit(X_train, Y_train)
583 importance = modelGarch.coef_[0]
584 feature_names=X.columns

```



```

585
586
587 indices = np.argsort(importance)
588 range1 = range(len(importance[indices]))
589 plt.figure()
590 plt.title("Logistic Regression Feature Importance")
591 plt.barh(range1, importance[indices])
592 plt.yticks(range1, feature_names[indices])
593 plt.ylim([-1, len(range1)])
594 plt.savefig(os.path.join(folder_path, "fig30.Variance.FeatureImportance.png"))
595 plt.show()
596
597 # Count the number of correct and incorrect predictions.
598 accuracy_counts = correct_predictions.value_counts()
599 print("Prediction Correctness Counts:\n", accuracy_counts)
600
601
602 plt.figure(figsize=(6,4))
603 accuracy_counts.plot(kind='bar', color=['green', 'red'])
604 plt.title("Prediction Accuracy (Correct vs. Incorrect)")
605 plt.xlabel("Prediction Correctness (True=Correct, False=Incorrect)")
606 plt.ylabel("Count")
607 plt.tight_layout()
608 plt.savefig(os.path.join(folder_path, "Volatility_Tomorrow_Prediction_Accuracy.
    png"))
609 plt.show()

```

Listing 1: Complete Python Code