# Data Science Assignment: Maximizing Gain from Market Predictions

### 1) Data Exploration and Pre-processing:

- **Data Understanding**

To ensure a comprehensive understanding of the dataset, the initial phase involves a detailed exploration. This includes examining the key variables and their overall structure, which helps in identifying critical patterns and trends. Specifically, I would focus on detecting time dependencies and seasonality in the market data, as these factors could significantly influence model selection and performance.

Visualization techniques such as time series plots, correlation matrices, and histograms are employed to better understand the distribution of features and the target variable—market changes. These visualizations will reveal underlying patterns, trends, and relationships within the data, which are crucial for subsequent modelling steps.

- **Handling Missing Values**

Missing data can adversely affect model performance, so appropriate imputation methods are necessary. For time series data, I would use forward or backward filling techniques to handle missing values, ensuring that the temporal sequence remains intact. For datasets where missing values are more random, imputation methods like mean or median replacement would be applied. Depending on the complexity and nature of the data, more advanced techniques such as KNN imputation may be considered to enhance data completeness and accuracy.

- **Outlier Detection**

Outliers can skew model performance and lead to inaccurate predictions. Therefore, I would use statistical methods such as Z-score or Interquartile Range (IQR) to identify outliers. If outliers are deemed erroneous, they will be removed to prevent distortion of the model. Conversely, if outliers represent significant market events or volatility spikes, they will be carefully analysed and transformed, as they may hold valuable information for prediction.

- **Feature Engineering**

Feature engineering is crucial for improving model accuracy and interpretability. I would create lag features to capture past values of the target variable (e.g., market change at time t-1, t-2, …..), which can help the model learn temporal dependencies. Additionally, I would generate technical indicators such as moving averages, exponential moving averages (EMA), Relative Strength Index (RSI), and Bollinger Bands, which are commonly used in financial forecasting. These indicators provide insights into market trends and conditions. Furthermore, I would extract temporal features such as the day of the week, month, or holidays, as these can impact market movements and add contextual information to the model.

- **Train-Test Split**

Maintaining temporal order in time series data is essential for realistic model evaluation. Therefore, I would split the dataset by using the earlier portion for training and the later portion for testing. This approach respects the chronological sequence and simulates real-world conditions. To ensure robustness in model validation, I would also implement rolling-window cross-validation. This technique involves training the model on a moving window of past data and testing it on future data, providing a more accurate assessment of model performance over time.

- **Scaling**

Feature scaling is an important step to ensure that the model performs optimally. I would apply <u>StandardScaler</u> or <u>MinMaxScaler</u> to standardize features, particularly when using distance-based models like SVM or neural networks. However, for time series models such as ARIMA, scaling may not be necessary for time-based input variables. Despite this, standardization of technical indicators is still recommended to ensure consistency and improve model accuracy.

By following these steps, I aim to prepare the data effectively, facilitating the development of robust and accurate predictive models.

---

## 2) Model Choice:

- **Model Choice for Market Prediction**

For predicting numerical market changes and the time to reach those changes, a combination of models or a hybrid approach can be employed.

To predict numerical market changes, classical time series models like ARIMA (Auto Regressive Integrated Moving Average) and SARIMA (Seasonal ARIMA) can effectively capture linear trends and seasonal patterns in historical market data. ARIMA models are useful for capturing underlying trends and non-stationarity, while SARIMA extends this capability to include seasonal effects. Alternatively, Long Short-Term Memory (LSTM) networks, a type of recurrent neural network, can capture more complex, non-linear patterns and long-term dependencies, making them suitable for financial forecasting. On the regression front, Gradient Boosting methods such as XGBoost and LightGBM are effective for handling both linear and non-linear relationships and incorporating various engineered features like technical indicators. Random Forest, another ensemble learning method, provides robustness and interpretability by combining multiple decision trees to improve prediction accuracy.

For predicting the time required to achieve a market change, linear regression offers a straightforward baseline model by assuming a linear relationship between predictors and the target variable. This approach is useful for understanding basic trends but may not capture more complex dynamics. Gradient Boosting and Random Forest can be employed here as well, providing the ability to model non-linear relationships and interactions effectively, which might be crucial for more accurate predictions.

A hybrid approach could involve using a multi-task neural network that simultaneously predicts both the numerical market change and the time to reach that change. This model would feature two output layers, one dedicated to the numerical prediction and the other to estimating the time, leveraging shared patterns and features to enhance overall performance. While this approach integrates the

strengths of various models, it also demands a more sophisticated architecture and substantial data for effective training.

In summary, selecting the appropriate model involves balancing complexity, interpretability, and performance. ARIMA and SARIMA are strong candidates for time series analysis, LSTM networks for capturing complex patterns, and Gradient Boosting and Random Forest for their robustness in handling diverse data features. For a more integrated solution, a multi-task neural network could be considered to jointly predict both outcomes.

-------------------------------------------------------------------------------------------------

### 3) Custom Evaluation Metric:

- **Custom Evaluation Metric and it's Influence**

The custom evaluation metric is essential for aligning the model's predictions with the company's specific profit and loss structure. Here's a detailed look at how this metric is constructed and its influence on your modelling approach.

## Custom Metric Construction

### a) Understanding the Profit and Loss Structure:

- **Market Increase > 50 Points:** The company incurs a loss of 20% of the predicted market increase.

- **0 ≤ Market Increase ≤ 50 Points:** The company gains the predicted market increase.

- **Market Increase < 0 Points:** The company incurs a loss of 50 plus the absolute value of the predicted market increase.

### b) Custom Metric Formula:

To capture these conditions, you define a custom function that calculates the overall gain or loss based on predictions:

```python
def custom_profit_loss_metric(y_true, y_pred):
    profit_loss = 0
    for true_val, pred_val in zip(y_true, y_pred):
        if pred_val > 50:
            profit_loss -= 0.2 * pred_val  # Loss of 20% of predicted value
        elif 0 <= pred_val <= 50:
            profit_loss += pred_val  # Gain of predicted value
        else:
            profit_loss -= (50 + abs(pred_val))  # Loss of 50 + absolute value
    return profit_loss
```

This function calculates the total profit or loss based on the predictions for each data point, and then aggregates the results.

**Influence on the Modelling Approach**

**a. Model Selection:**

   - **Focus on Profit Maximization:** Traditional metrics like MSE or MAE are not aligned with the company's profit/loss criteria. By using the custom metric, we focus on maximizing profit rather than minimizing errors.

   - **Choosing Models:** Models that can accurately predict the magnitude of market changes become crucial. For example, regression models (like Gradient Boosting) and time series models (like ARIMA) can be evaluated based on their ability to generate predictions that fall into the profit range or minimize losses.

**b. Training and Evaluation:**

   - **Custom Objective Function:** During model training, use the custom metric as the objective function or evaluation criterion. This means that the model will be optimized to maximize the company's profit rather than just fitting the data.

   - **Cross-Validation:** Implement cross-validation strategies that incorporate the custom metric to ensure that the model generalizes well to unseen data and adheres to the company's profit structure.

**c. Model Tuning:**

   - **Hyperparameter Optimization:** When tuning hyperparameters, use the custom metric to assess the performance of different parameter configurations. This ensures that the chosen hyperparameters lead to the highest potential profit.

   - **Feature Engineering:** Focus on features that impact the profit and loss calculation. For instance, features that influence the likelihood of market changes being in the profitable range (0 to 50 points) or avoid high losses (>50 points) should be prioritized.

**d. Performance Monitoring:**

- **Continuous Evaluation:** Continuously evaluate the model's performance using the custom metric during training and testing phases. Adjust the model or feature set if the profit/loss outcomes are not aligning with expectations.

- **Sensitivity Analysis:** Perform sensitivity analysis to understand how different predictions affect the overall profit or loss. This can help in refining the model and ensuring that it behaves optimally across various market conditions.

## e. Reporting and Insights:

- **Detailed Analysis:** In your report, include a detailed analysis of how the custom metric influenced the model's performance. Compare results with traditional metrics to highlight the differences.

- **Visualization:** Use visualizations to show the impact of different predictions on profit and loss. For example, plot the predicted vs. actual market changes and highlight regions where the model performs well or poorly based on the custom metric.

## Summary

The custom evaluation metric fundamentally shapes our approach by aligning the model's objectives with the company's financial goals. It drives model selection, tuning, and performance evaluation, ensuring that the final model maximizes the company's gain according to the specific profit and loss rules. This focus on profit rather than traditional error metrics helps in building a model that not only fits the data well but also delivers tangible financial benefits.

**4) Model Evaluation:**

To evaluate final model performance, we should use a combination of techniques and metrics:

**a. Hold-out Test Set:** Reserve a portion of your data (typically 20-30%) as a final test set that is not used during model development or tuning. Evaluate your final model on this unseen data to get an unbiased estimate of performance [1][3].

**b. Multiple Performance Metrics:** Use several relevant metrics to get a comprehensive view of model performance [1][2]:

   - **For classification:** Accuracy, precision, recall, F1-score, ROC AUC
   - **For regression**: RMSE, MAE, R-squared

**c. Learning Curves:** Plot training and validation performance across epochs/iterations to assess overfitting/underfitting [2].

**d. Cross-Validation:** Use k-fold cross-validation on your training data to get a robust estimate of model performance across different data splits [3].

**e. Confusion Matrix:** For classification, analyze the confusion matrix to understand error patterns [1].

**f. Business Metrics:** Translate ML metrics to business KPIs that stakeholders care about [2].

**g. Qualitative Evaluation**: Have domain experts review model predictions on sample cases [2].

**h. Baseline Comparison**: Compare your model against simple baselines to ensure it's adding value [3].

**i. Ensemble Methods:** Consider averaging predictions from multiple models trained during cross-validation for improved robustness [3].

**j. Error Analysis:** Examine misclassified examples to gain insights into model weaknesses [1].

The key is to use a combination of quantitative metrics and qualitative analysis to thoroughly assess your model's strengths and limitations before deployment. Regularly re-evaluate performance after deployment to detect any degradation over time.

## Citations:

[1] https://towardsdatascience.com/various-ways-to-evaluate-a-machine-learning-models-performance-230449055f15?gi=18384328f258

[2] https://neptune.ai/blog/ml-model-evaluation-and-selection

[3] https://mindfulmodeler.substack.com/p/how-to-get-from-evaluation-to-final

[4] https://www.fiddler.ai/model-evaluation-in-model-monitoring/what-is-model-performance-evaluation

[5] https://graphite-note.com/a-comprehensive-guide-to-model-evaluation-in-machine-learning/

[6] https://www.jeremyjordan.me/evaluating-a-machine-learning-model/

[7] https://sebastianraschka.com/faq/docs/evaluate-a-model.html

[8] https://machinelearning101.readthedocs.io/en/latest/_pages/07_model_performance.html

----------------------------------------------------------------------------------------------------

## 5) Suggestions for Further Improvement:

- **Hyperparameter Tuning:** Use *GridSearchCV* or *RandomizedSearchCV* for hyperparameter optimization. For models like Gradient Boosting or Random Forest, tuning parameters like learning rate, maximum depth, or the number of trees can improve performance.

- **Feature Selection:** Apply techniques like recursive feature elimination (RFE) to identify the most important features influencing the market prediction. This could simplify your model and improve interpretability.

- **Ensemble Learning:** Combine predictions from multiple models (e.g., using stacking or voting classifiers) to reduce variance and bias, potentially increasing overall performance.

- **Incorporate External Data:** Financial markets are influenced by global events, economic indicators, and sentiment analysis from news and social media. Adding features like *sentiment analysis scores*, economic indicators (GDP growth, inflation rates), or *market volatility indices* can enhance model performance.

- **Adaptive Models:** Financial markets evolve, so models need to adapt. Consider using online learning models or regularly retraining the model with recent data to capture changing patterns in the market.

-------------------------------------------------------------------------------------------------