# Data Science Canvas

| Project: | |
|---|---|
| Team: | |

## Problem Statement | Execution & Evaluation | Data Collection & Preparation

| **Business Case & Value Added** | **Model Selection** | **Model Requirements** | **Skills** | **Model Evaluation** | **Data Storytelling** | **Data Selection & Cleansing** | **Data Collection** |
|---|---|---|---|---|---|---|---|
| Which business case should be analyzed and what added value does it generate?<br><br>The case is predicting each trading day's high and low for NIFTY 500 stocks and generating early buy/sell signals. Added value: improves intraday decision-making, reduces trader bias, and enables realistic profit estimation through systematic backtesting including costs and slippage. | Which analysis methods can be considered on the basis of the specific data landscape and the business case?<br><br>Minute-level OHLCV data with highly non-linear and noisy patterns requires models that can capture complex dependencies. Thus, tree-based gradient boosting models XGBoost, LightGBM, CatBoost , neural network, Stacking Ensemble ensembled together with probabilities[0.3, 0.2, 0.2, 0.2, 0.1] were used for BUY/SELL classification, whiletwo separated sequential3 layered (128 → 64 → 32 ) LSTM models were used for predicting daily high and low prices. These models balance accuracy, speed, and the ability to generalize across noisy intraday data. | Which model requirements must be complied with in order to obtain a valid model?<br><br>Model must train using only past data (strict time-based split) and avoid look-ahead when generating indicators. Predictions must be early enough for trading and robust under different market regimes. | What skills are needed to provide the data and model development?<br><br>Time-series preprocessing, indicator engineering, ML modeling with tree ensembles and simple neural nets, backtesting logic, and performance evaluation including transaction-cost modeling. | Which indicators require quality control and validation and how should they be interpreted? Is real-time monitoring necessary?<br><br>Classification signals must be evaluated using **Accuracy, F2 Score, and AUC**, ensuring the model correctly detects directional movement patterns. For high/low regression, **R² and RMSE** must be monitored to verify the quality and stability of numerical predictions. Cross-validation is needed to ensure generalization across stocks and days. | What requirements does the target group have for the presentation of the results and how do I effectively communicate this data?<br><br>Stakeholders need clear, visual summaries of intraday trends, cluster behavior, model accuracy, and prediction errors. Charts comparing actual vs predicted highs/lows, F2-score bars, ROC curves, and cluster-wise EDA help them quickly evaluate reliability. | Which of the available data is relevant? Do the data have to be cleaned up?<br><br>Use regular-hours minute data only. Clean duplicates, erroneous ticks, missing minutes, and negative or zero-price anomalies. Align timestamps to uniform intervals and normalise inconsistent formatting. | How and with which methods should additionally required data be collected? What properties has this data to fulfil?<br><br>Data is downloaded from the Kaggle Algo Trading dataset (2015–2025) containing minute OHLCV per stock. All data must have synchronized timestamps, consistent trading-session boundaries, and maintain minute-level granularity without gaps. |
| **Data Landscape**<br>Which data is required for this and which is already available? Which additional data has to be collected?<br><br>Available: minute-level OHLCV data for all NIFTY 500 stocks. Required: engineered indicators (returns, ATR, rolling highs/lows, VWAP). | | **Software & Libraries**<br>Which software should be used? Is there already a standard solution? Which libraries are used?<br><br>Python-based workflow using pandas, NumPy, scikit-learn, XGBoost/LightGBM, and statsmodels. Matplotlib/seaborn/plotly for visualization and Dash for reporting. | | | | **Data Integration**<br>In which system should the data from different sources be migrated?<br><br>Combine all stock files into a unified columnar dataset (parquet). Integrate indicators and additional market features in a feature store ensuring strict time alignment. | **Explorative Data Analysis**<br>Are there outliers or structures to be considered? Creation of descriptive key figures for the first assessment of the data.<br><br>Check for missing minutes, extreme price spikes, zero-volume entries, and duplicated ticks. |