# Forecasting Unit Trust Sales and Redemption

Leveraging on Data Science

25th July 2024

AIML Knowledge Sharing

# About Me

Aswadi Abdul Rahman

✉ aswadi.abdulrahman@gmail.com

in https://www.linkedin.com/in/aswadiabdulrahman/

💼 Head of Analytic Department, ASNB

**Other Working Experiences:**

**slido**

Please download and install the Slido app on all computers you use

# How would you rate your proficiency with AI and machine learning tools/method?

ⓘ Start presenting to display the poll results on this slide.

# Why we need to have unit trust forecast



**Enhanced Strategic Planning:** Prepare for potential sales or redemption trends for next year.



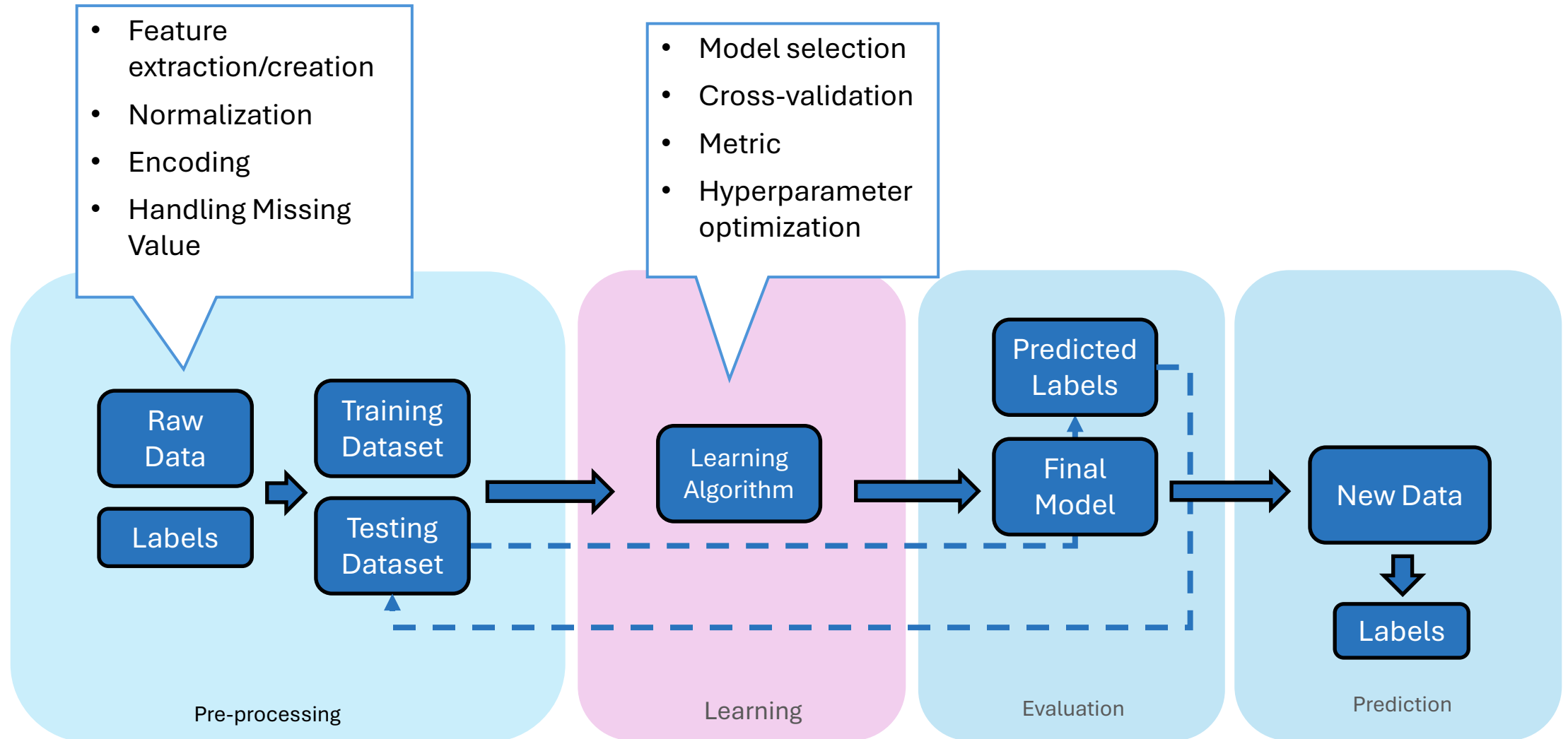**Investment Decision:** Enable fund managers to manage effectively.



**Risk Management:** Identify market trends and risks, allowing proactive measures.



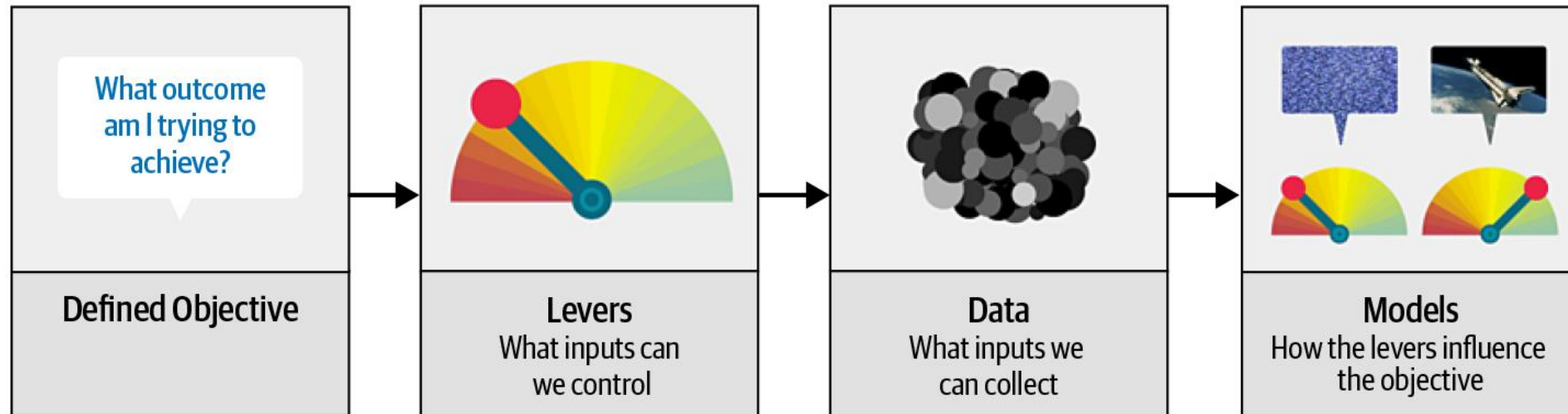**Product Management:** Ensure proper allocation of investment products.

# Flow Chart of Predictive Modeling

- Feature extraction/creation
- Normalization
- Encoding
- Handling Missing Value

- Model selection
- Cross-validation
- Metric
- Hyperparameter optimization

**Raw Data**

**Labels**

**Training Dataset**

**Testing Dataset**

**Learning Algorithm**

**Predicted Labels**

**Final Model**

**New Data**

**Labels**

Pre-processing

Learning

Evaluation

Prediction

# Details on the Pre-processing

- **Feature Selection:**
    - Training: Include all relevant features you can think of.
    - Forecast: Consider features you might have access to in the future.



"The Drivetrain Approach" - Designing Great Data Products (Jeremy, Margit, Mike)

# Details on the Pre-processing

- **Feature Extraction:**
  - Break down dates into components like year, month, and day.
  - Understanding features that having the same correlation (multi-collinearity) using rank correlation
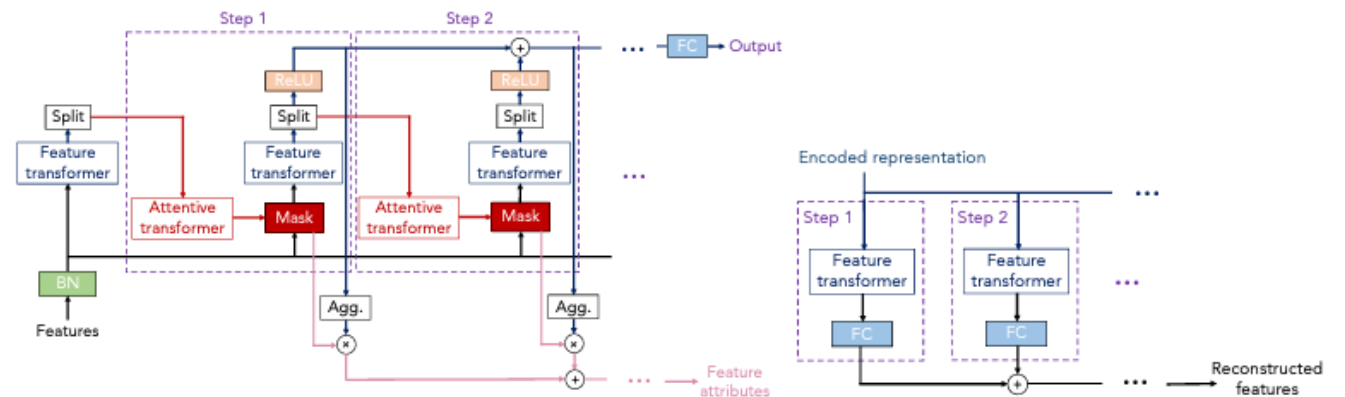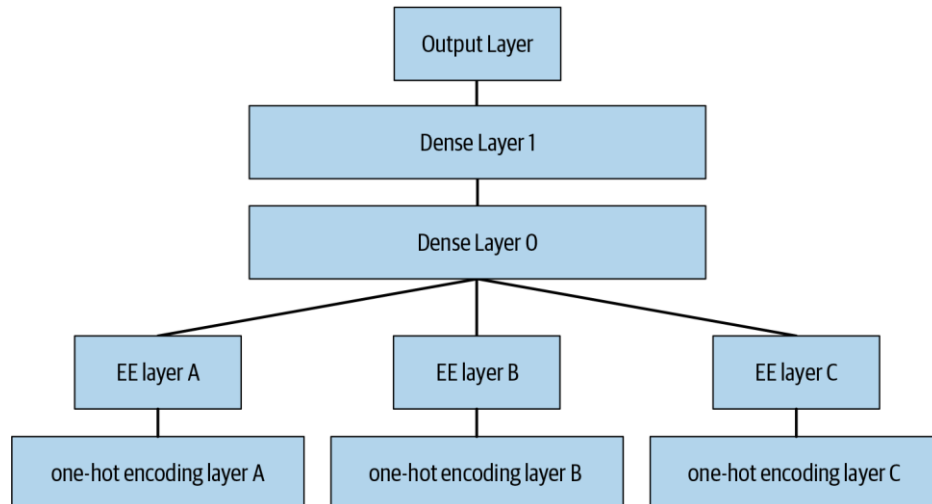
# Details on the Pre-processing

- **Encoding:**
  - Primarily use label encoding.
  - Occasionally use embedding techniques like entity embedding or unsupervised like TabNet.

# Details on the Pre-processing
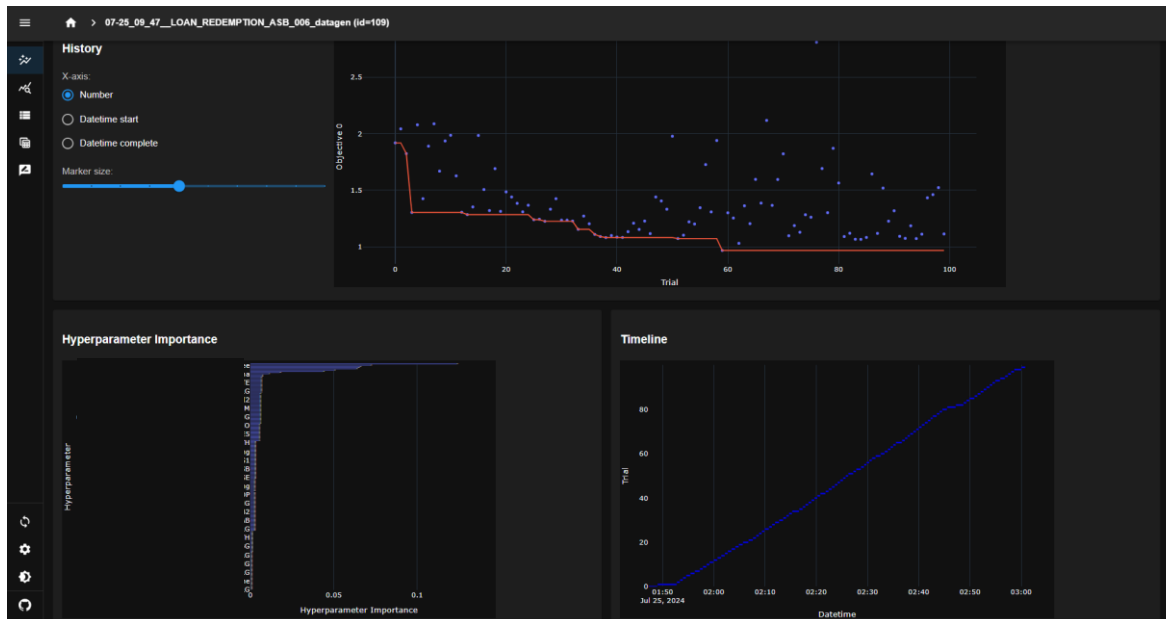
- **Normalization:**
  - Typically, use the log1p transformation for normalization.

- **Handling Missing Values:**
  - Create new features to label missing values.
  - Usually, fill missing values with -1.

# Details on the Learning

- **Hyperparameter optimization:**
  - In the beginning we use grid search for random forest.
  - Now we use Optuna using Bayesian optimization on XGBoost.
  - Then store all the experiment in ML Flow

# Details on the Learning

- **Evaluation Metric:**



### Mean Squared Error (MSE):

**Definition**: Measures the average squared difference between predicted and actual values.

Use for Subscription forecast

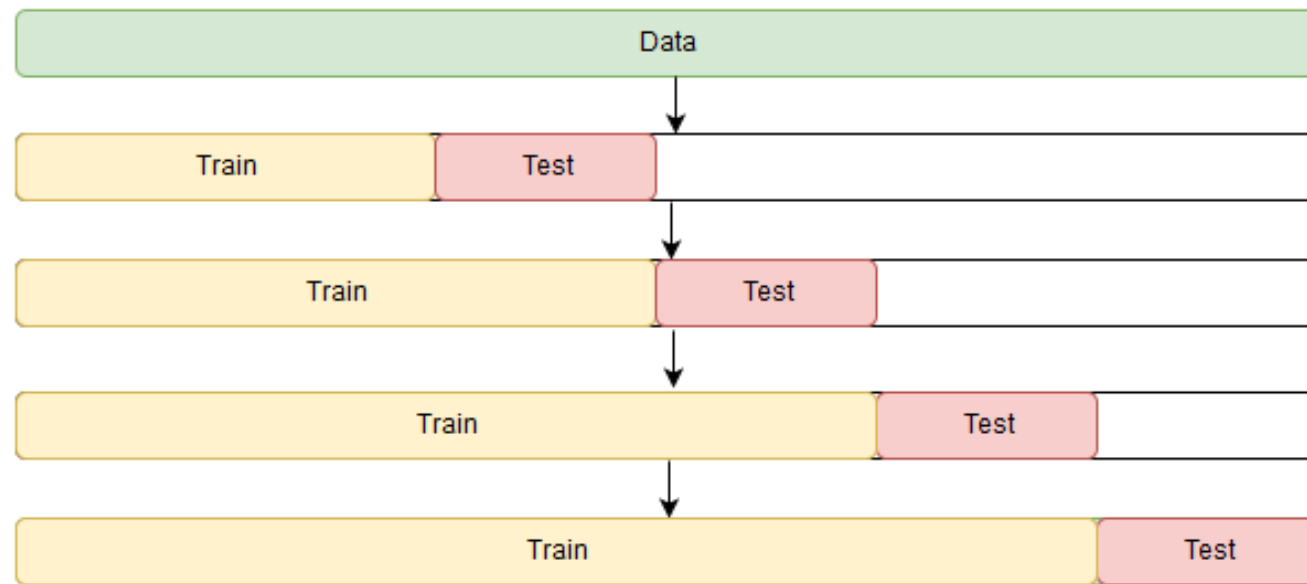### Mean Squared Error (MSE) penalized under Predict:

**Definition**: Measures the average squared difference between predicted and actual with penalized errors under zero.

Use for Redemption forecast

# Details on the Learning

- **Model selection:** Our baseline is random forest, then move up to gradient boost tree.

- **Cross-validation:** K-fold rolling windowing

# Details on the Learning

By forecasting fund sales and redemptions, we can identify potential risks and take steps to mitigate them.

**Historical Monthly data:**

- Each fund Redemption
- Lockdown
- Fund Dividend
- Marketing Campaign
- Fund Net Asset Value
- Macroeconomics

**Forecast Monthly data:**

- Fund Net Asset Value
- Macroeconomics
- Fund Dividend
- Marketing Campaign

ML Model Development

Forecast using Trained model

# What matters for others?

- Each prediction value was determined by its trained features and weight which using interpreter (tree interpreter)
- What happen on each forecasted value? <u>Discuss with subject matter expert</u>

# What matters for others?
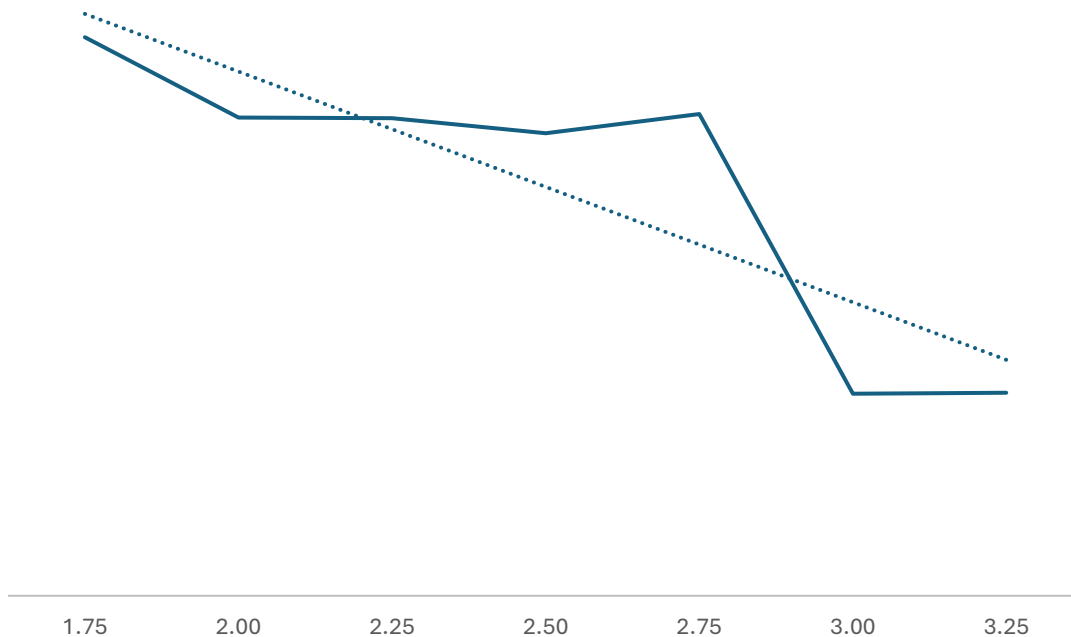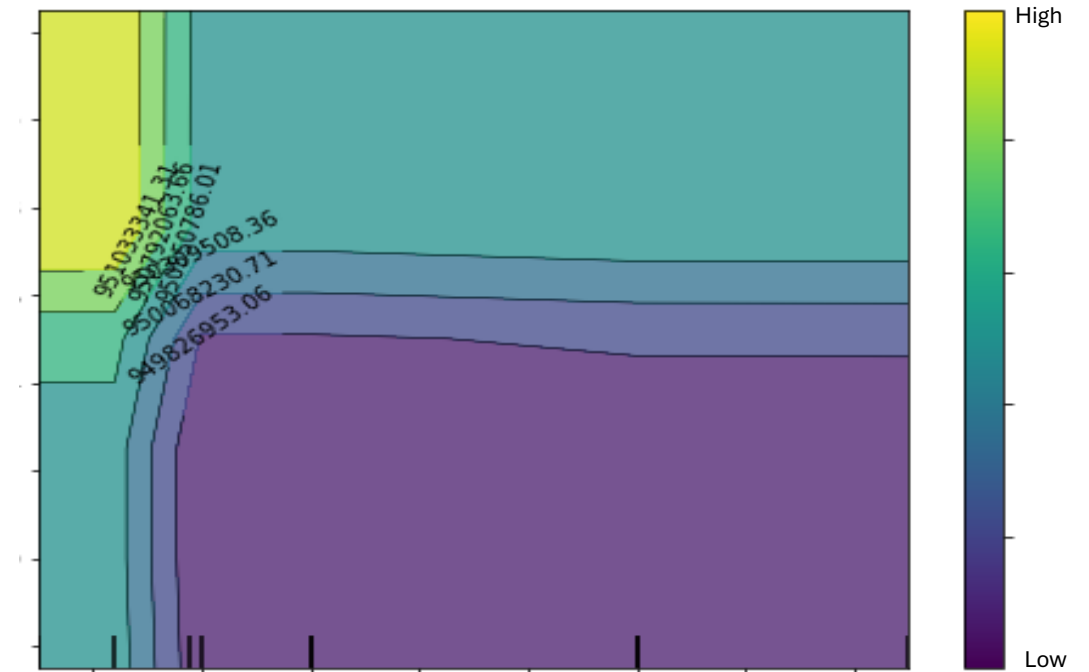
By training the data using the ML models, other than forecast we will also extract the sub output which are features what-if analysis using partial dependent plot (PDP).
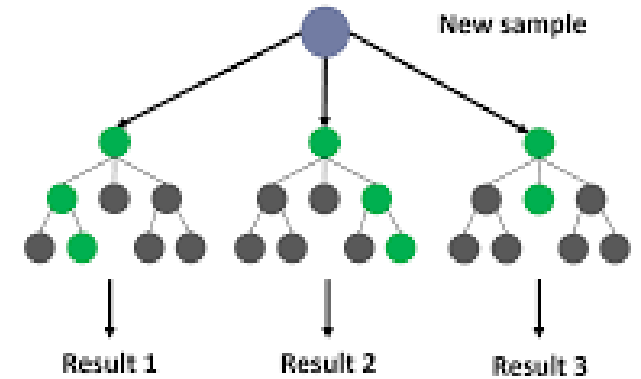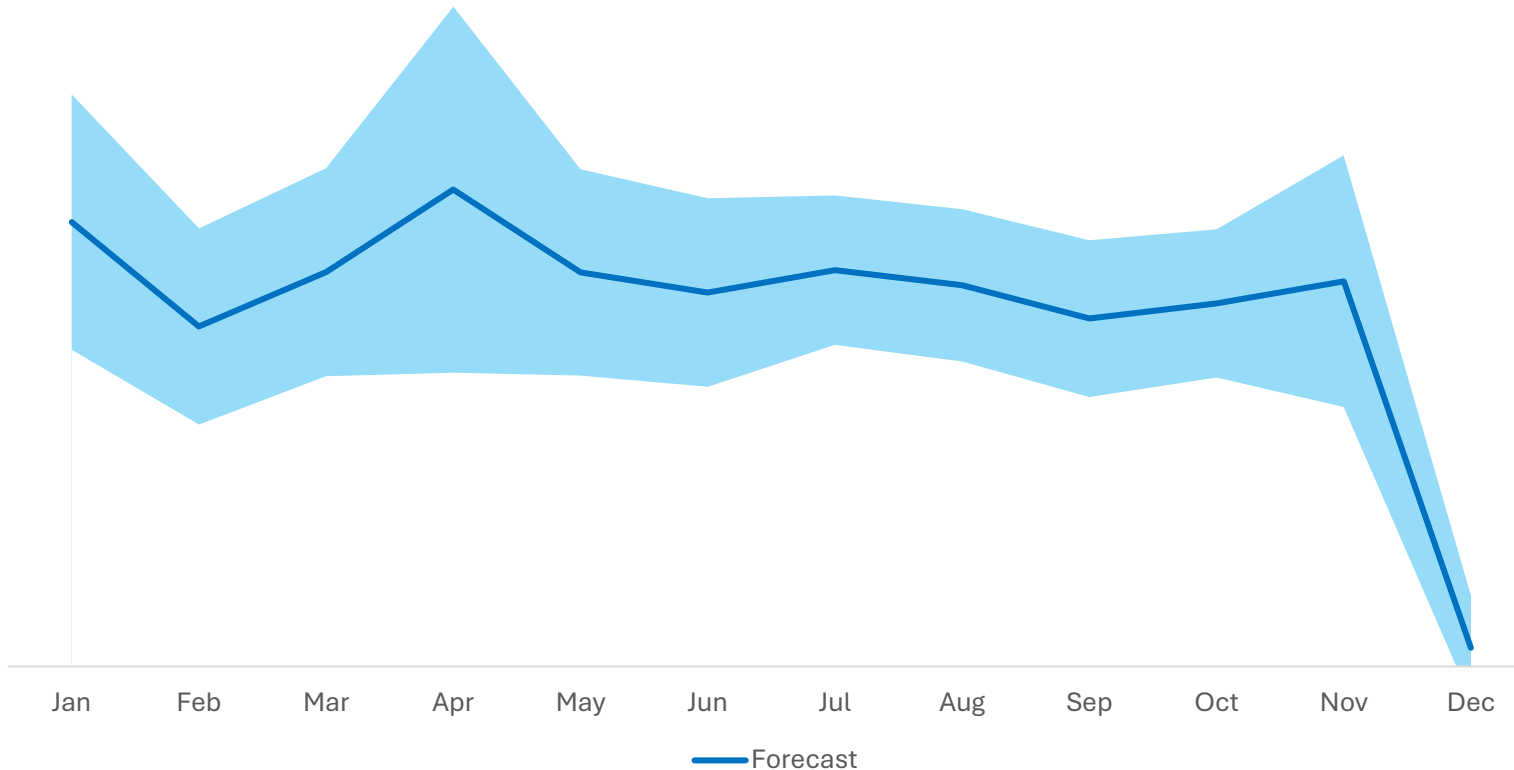


PDP for Macroeconomic 1



PDP for Macroeconomic 1 vs features 1

# What matters for others?

Forecast with confidence interval will have affect on their decision. We are using standard deviation of the prediction tree.



Conformal Prediction?

# Forecasting Challenges

- **Retrain with new data:**
  - To avoid data drift
  - NEVER train the whole dataset for forecast.
  - Need to consider future data availability especially external data.
- **Strategic consideration:**
  - Some fund that does not have enough data
  - And new initiatives should be done outside the model.
- **Future accuracy expectation:** Find stability instead of accuracy
- **Retrain with new method:** Only when your model is not stable & outside of the confidence interval.

# Other Use Cases

- Unsupervised: Topic Modelling Clustering

- Predicting Customer Unitholding

- AI Chatbot – Leverage on OpenAI LLM

- AI Avatar Chatbot

- Lookalike model for Customer Segment

- Instant Analysis using LLM

- AMLA Transaction Classification

# We are hiring!!

**AMLA Data Analytic Manager**

Experience Needed

- \>5 years of total working experiences
- \>3 years  on Compliance and AMLA fields
- \>3 years hands on working on ML models in python and ETL

What will you do

- Supervised 2 data analytics/ data science
- Working together with compliances department and security commission for relevant guidelines

Email to aswadi@pnb.com.my

Thank You