# perplexity

# Explanation of Your Python Code

Your Python code is a complete pipeline for:

- Reading customer churn data,
- Training multiple Random Forest models to find the best,
- Predicting customer churn on new data,
- Saving the predicted churn customers to an Excel file,
- And visualizing feature importance.

Let me walk you through each part and explain the libraries you used.

## Predefined Libraries Used

### 1. pandas

- Used for data manipulation and analysis.
- Reads Excel files into DataFrames, handles tabular data easily.

### 2. numpy

- Provides support for numerical operations and arrays.
- Used mainly here for sorting feature importance.

### 3. matplotlib.pyplot & seaborn

- matplotlib.pyplot is a plotting library for creating graphs.
- seaborn builds on matplotlib and provides nicer, more statistical graphics.
- You use them to plot feature importance.

### 4. sklearn.model_selection.train_test_split

- Utility to split your total dataset into training and testing parts.
- Helps in evaluating the model fairly.

### 5. sklearn.ensemble.RandomForestClassifier

- Implements the Random Forest algorithm, an ensemble of decision trees.
- Used for classification tasks like predicting customer churn.

### 6. sklearn.metrics (classification_report, confusion_matrix, precision_score)

- Tools to evaluate model performance:
  - Confusion matrix: counts true/false positives/negatives.
  - Classification report: precision, recall, f1-score per class.
  - Precision score: ratio of true positives over predicted positives.

### 7. sklearn.preprocessing.LabelEncoder

- Converts categorical (text) variables into numeric codes for the model.

### 8. joblib

- Usually used to save/load trained models, but not explicitly used here.

### Code Breakdown

### Reading and Preparing Training Data

- Reads `'vw_ChurnData'` sheet from Excel file into a DataFrame `data`.
- Removes columns not relevant for prediction.
- Encodes categorical columns into numbers because ML models require numeric input.
- Converts churn status (`Stayed`/`Churned`) into 0/1.

### Splitting Data

- Splits the dataset into 80% for training, 20% for testing.

### Training Multiple Models

- Defines function `train_and_evaluate` to:
  - Train a Random Forest with a given random seed,
  - Predict churn on test data,
  - Print performance metrics, and return the model with precision score.
- Loops over different random seeds to train multiple models,
- Selects the model with the highest precision.

### Visualizing Feature Importance

- Uses the best model's feature importance scores to plot how much each feature contributes to predictions.

### Predicting on New Data

- Reads `'vw_JoinData'` sheet for new customer data to predict churn.

- Copies the original data to keep unencoded columns.

- Prepares new data by dropping non-predictive columns.

- Encodes categorical columns using the encoders fit on training data.

- Uses the best trained model to predict churn.

- Adds predictions back to the original new data.

### Extracting and Saving Churned Customers

- Filters customers predicted to churn.

- Saves this filtered data to a new Excel file.

- Prints confirmation.

### Summary

Your script automates the machine learning workflow for churn prediction:

- Reads data,

- Preprocesses and encodes it,

- Trains multiple models to maximize precision,

- Applies the best model to new data,

- And prepares results in Excel for business use or dashboards.

If you want, I can help explain any specific part in more detail or on how to extend this pipeline with automation or integration with Power BI.