# Project Synopsis: AutoFeatureEngineering

## Overview

The AutoFeatureEngineering package automates essential feature engineering tasks for machine learning,

from handling missing values to encoding categorical variables, generating polynomial features, and ranking feature

importance. This tool aims to simplify data preparation, making it faster and more efficient to build robust machine

learning models.

## Key Features

1. Data Type Detection: Automatically identifies and classifies columns in the dataset as numerical

or categorical, enabling customized processing based on data type.

2. Missing Value Handling: Provides several strategies

to manage missing values:

 - Mean Imputation: Fills missing values with the column mean.

 - Median Imputation: Fills

missing values with the column median.

 - Row Dropping: Removes rows with any missing values.

 - Target Column Imputation:

Handles missing values in the target column separately, filling with mean for numerical targets or mode for categorical

targets.

3. Encoding Categorical Variables: Offers flexible options for encoding categorical features:

 - One-Hot Encoding: Converts

categorical variables into binary columns for each category.

 - Ordinal Encoding: Encodes categories as ordered numbers.

- Target Encoding: Encodes categories based on the mean value of the target variable for each category.

4. Polynomial Feature

   Engineering: Generates polynomial and interaction terms for numerical columns to capture non-linear relationships within the data.

5. Feature Ranking: Uses Random Forest models to rank features based on importance, aiding in informed feature selection.

6. Feature Selection: Retains only relevant features by removing low-importance features, always ensuring that the target

   column is included in the final dataset.

## Technologies Used

Python Libraries:

 - pandas: For data manipulation and handling.

 - scikit-learn: Provides

   preprocessing, encoding, and feature selection tools.

 - category_encoders: Specifically used for target encoding.

 - Random Forest Models: For ranking and selecting important features based on feature importance scores.

## Use Cases

1. Automated Data Preparation for Machine Learning: The package is ideal for streamlining data preprocessing,

making datasets ready for model training by addressing missing values, encoding, and feature selection.

2. Prototyping Machine

  Learning Models: It's highly suitable for creating quick prototypes by automating multiple feature engineering tasks, allowing

  data scientists to focus on model development and evaluation.

3. Handling Mixed Data Types: The package can process datasets

  with a combination of numerical and categorical data, making it useful across diverse fields, including finance,

healthcare, and

  retail.

4. Feature Selection for High-Dimensional Datasets: Particularly valuable for high-dimensional datasets, this tool ranks

  and retains only relevant features, enhancing model interpretability and reducing training time.

## Usage Instructions

Installation: The package can be installed via PyPI and imported for immediate use in any Python

  environment.

1. Initialization: Initialize the AutoFeatureEngineering class by specifying:

  - target_col: The column to be

  predicted.

  - null_strategy: Method for handling missing values, such as mean or median.

  - encoding_method: The method to encode

  categorical variables, such as one-hot, ordinal, or target encoding.

2. Feature Engineering Pipeline: After initialization, the

fit_transform method applies the complete feature engineering pipeline:

- Detects and classifies column types in the dataset.

- Manages missing values based on the specified strategy.

- Encodes categorical features using the selected encoding method.

- Generates polynomial features for numerical columns.

- Ranks and retains relevant features based on importance, ensuring the

  target column is included.

3. Results: The output is a processed dataset, optimized with engineered features, ready for direct use

  in machine learning model training.

## Conclusion

AutoFeatureEngineering is a comprehensive tool for automated feature engineering, designed to streamline data

  preparation and improve model-building efficiency. With its flexible options for handling missing values, encoding, and

feature

  selection, this package is an invaluable asset for data scientists working on complex machine learning projects.