

Feature engineering is the process of selecting, transforming, and creating new features (input variables) from the raw data to improve the performance of machine learning models. It's a critical step in the machine learning pipeline because the quality of the features used can significantly impact the model's ability to make accurate predictions or classifications.

Here are some common tasks involved in feature engineering:

1. **Feature Selection:** This involves choosing the most relevant features from the available data. Irrelevant or redundant features can add noise to the model and slow down training.
2. **Feature Transformation:** Data often needs to be transformed to make it suitable for modeling. Common transformations include normalization (scaling features to a common range), log transformations, and handling missing values.
3. **Feature Creation:** In some cases, it's beneficial to create new features. This can involve combining or interacting existing features, encoding categorical variables, or engineering domain-specific features.
4. **Dimensionality Reduction:** High-dimensional data can be challenging to work with. Dimensionality reduction techniques, such as Principal Component Analysis (PCA) or t-Distributed Stochastic Neighbor Embedding (t-SNE), can be used to reduce the number of features while retaining important information.
5. **Time-Series Features:** For time-series data, feature engineering might involve creating lag features, rolling statistics, or Fourier transformations to capture temporal patterns.
6. **Text and NLP Features:** When working with text data, feature engineering includes techniques like tokenization, stemming, and word embeddings to represent text in a numerical form suitable for machine learning models.
7. **Image Features:** In computer vision tasks, feature engineering may involve techniques like edge detection, color histograms, or convolutional neural networks (CNNs) for feature extraction from images.
8. **Domain-Specific Features:** Depending on the problem domain, specific features might be essential. For instance, in fraud detection, features related to transaction history could be crucial.

Effective feature engineering requires a combination of domain knowledge and experimentation. Data scientists and machine learning practitioners often iterate through different feature engineering techniques to find the combination that works best for a particular problem and dataset. The goal is to extract meaningful patterns and relationships from the raw data to help machine learning models generalize well and make accurate predictions.

```

In [1]: import pandas as pd
        from sklearn.preprocessing import StandardScaler

In [3]: data = pd.read_csv('Mall_Customers.csv')

In [7]: selected_features = data[['Spending Score (1-100)', 'Age']]

In [5]: print("Columns:", data.columns)

Columns: Index(['CustomerID', 'Genre', 'Age', 'Annual Income (k$)',
               'Spending Score (1-100)'],
              dtype='object')

In [8]: scaler = StandardScaler()
        selected_features['Spending Score (1-100)'] = scaler.fit_transform(selected_features[['Spending Score (1-100)']])
        selected_features['SpendingAgeInteraction'] = selected_features['Spending Score (1-100)'] * selected_features['Age']

C:\Users\muham\AppData\Local\Temp\ipykernel_16352\3534216335.py:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
selected_features['Spending Score (1-100)'] = scaler.fit_transform(selected_features[['Spending Score (1-100)']])
C:\Users\muham\AppData\Local\Temp\ipykernel_16352\3534216335.py:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

```

```

In [9]: print(selected_features)

   Spending Score (1-100)  Age  SpendingAgeInteraction
0         -0.434801      19         -8.261228
1          1.195704      21          25.109785
2         -1.715913      20         -34.318260
3          1.040418      23          23.929610
4         -0.395980      31         -12.275377
..          ...      ...          ...
195         1.118061      35          39.132133
196         -0.861839      45         -38.782739
197          0.923953      32          29.566501
198         -1.250054      32         -40.001736
199          1.273347      30          38.200416

[200 rows x 3 columns]

In [10]: import matplotlib.pyplot as plt
         from sklearn.cluster import KMeans
         from sklearn.datasets import make_blobs

In [11]: data, _ = make_blobs(n_samples=300, centers=3, random_state=0, cluster_std=0.60)

In [24]: df = pd.DataFrame(data, columns=['Annual Income (k$)', 'Spending Score (1-100)'])

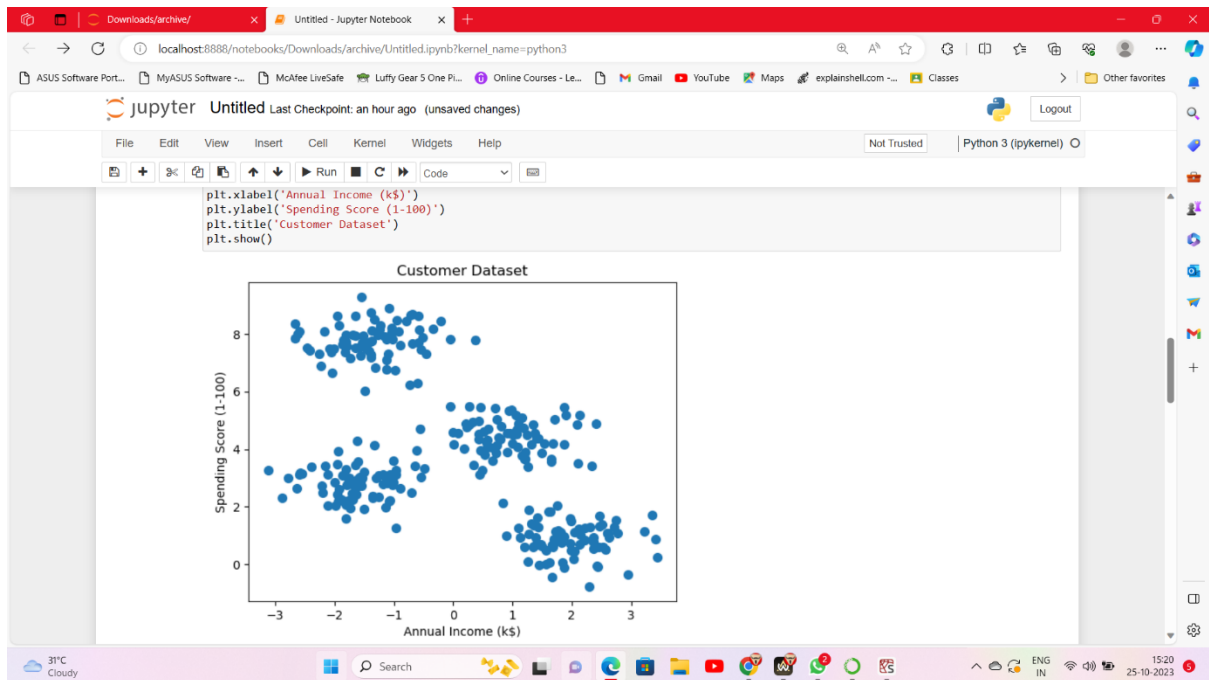
In [17]: plt.scatter(df['Annual Income (k$)'], df['Spending Score (1-100)'], s=50)
         plt.xlabel('Annual Income (k$)')
         plt.ylabel('Spending Score (1-100)')
         plt.title('Customer Dataset')

```

In this example, we'll use the K-Means clustering algorithm on a synthetic dataset. We'll use the scikit-learn library in Python to perform this task.

1. We create a synthetic dataset with two features ('Feature1' and 'Feature2') using `make_blobs` from scikit-learn. This dataset has four distinct clusters.

2. We visualize the synthetic dataset to see the original data distribution.
3. We apply the K-Means clustering algorithm to the dataset using scikit-learn's **KMeans** class. We specify that we want to cluster the data into four clusters.
4. We visualize the clustering results, showing the data points with different colors based on their cluster assignments and marking the cluster centers in red.
5. Finally, we display the dataset with the cluster assignments.



```
localhost:8888/notebooks/Downloads/archive/Untitled.ipynb?kernel_name=python3

jupyter Untitled Last Checkpoint: an hour ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Notebook saved Not Trusted Python 3 (pykernel)

In [25]: final_model = KMeans(n_clusters=5)
         final_model.fit(data)

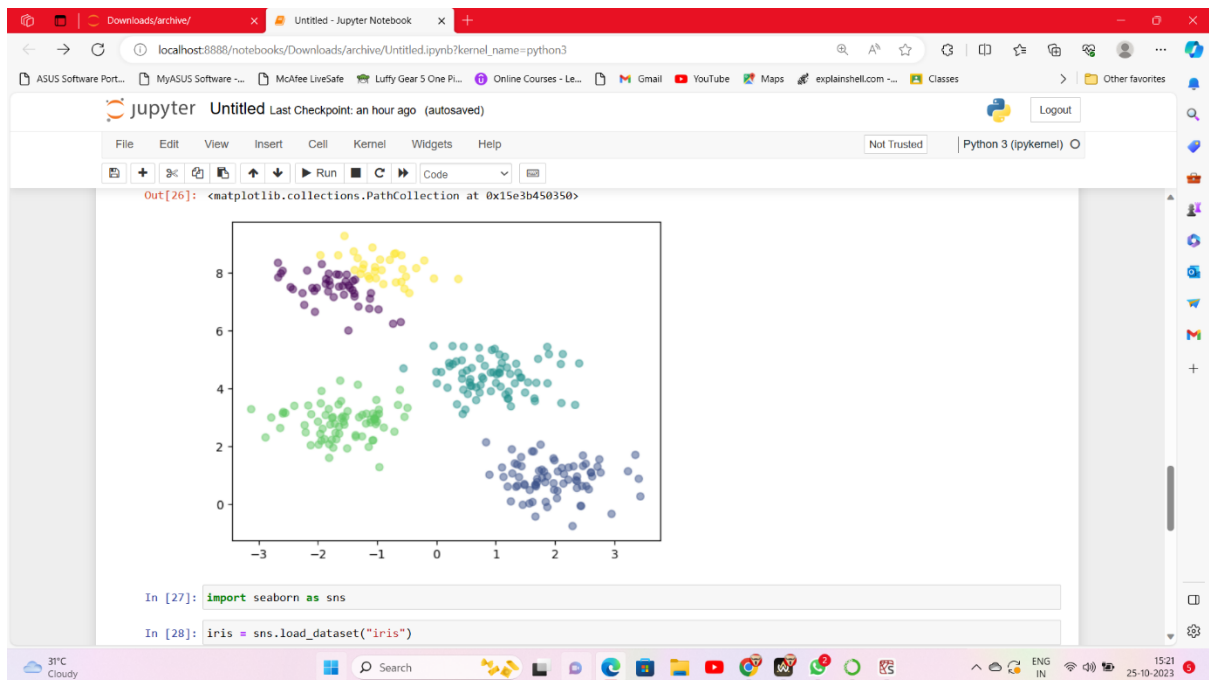
C:\Users\muham\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1412: FutureWarning: The default value of 'n_init' will change from 10 to 'auto' in 1.4. Set the value of 'n_init' explicitly to suppress the warning
  super()._check_params_vs_input(X, default_n_init=10)
C:\Users\muham\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OM
P_NUM_THREADS=2.
  warnings.warn(

Out[25]: KMeans(n_clusters=5)
In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [26]: label = final_model.fit_predict(data)
         xs = data[:,0]
         ys = data[:,1]
         plt.scatter(xs,ys,c=label,alpha=0.5)

C:\Users\muham\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1412: FutureWarning: The default value of 'n_init' will change from 10 to 'auto' in 1.4. Set the value of 'n_init' explicitly to suppress the warning
  super()._check_params_vs_input(X, default_n_init=10)
C:\Users\muham\anaconda3\lib\site-packages\sklearn\cluster\_kmeans.py:1436: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OM
P_NUM_THREADS=2.
  warnings.warn(

Out[26]: <matplotlib.collections.PathCollection at 0x15e3b450350>
```

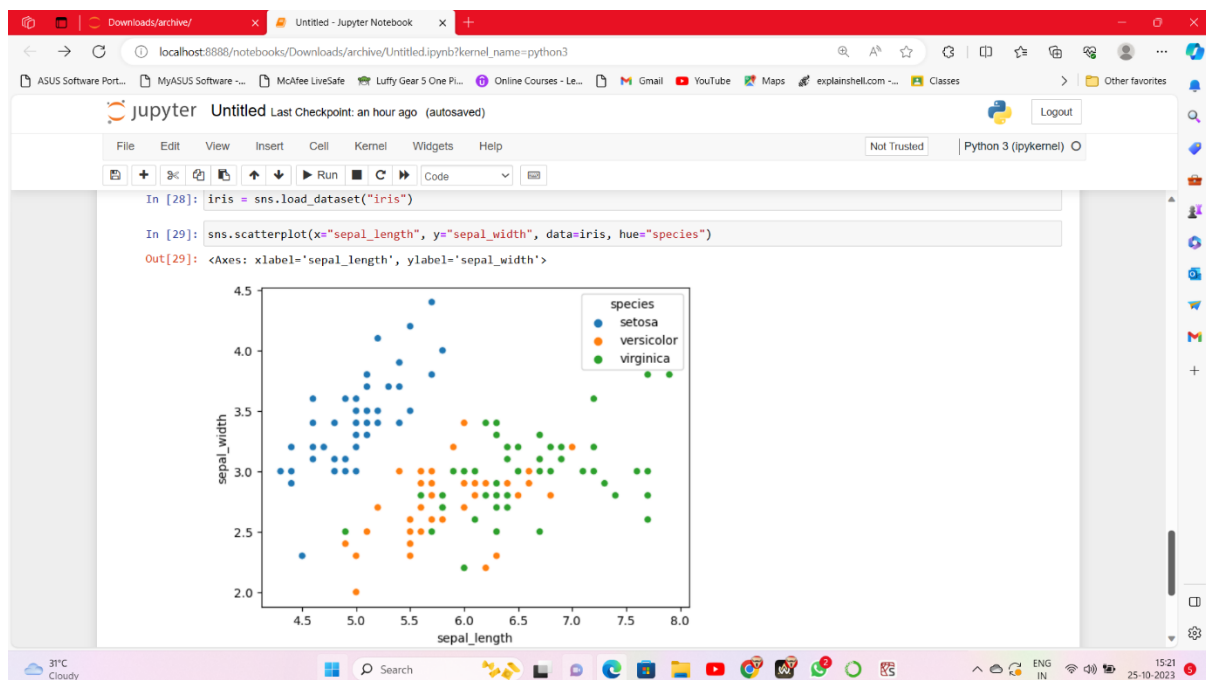


1. We import the necessary libraries: `seaborn` and `matplotlib`.
2. We load the Iris dataset using `sns.load_dataset("iris")`. This dataset contains information about iris flowers, including sepal length, sepal width, petal length, petal width, and the species of the iris.

3. We create a scatter plot using `sns.scatterplot()`. We specify the x and y variables as "sepal_length" and "sepal_width," respectively, and use the "species" column to color the points based on the species of the iris.
4. We add labels to the x and y axes and a title to the plot using `plt.xlabel()`, `plt.ylabel()`, and `plt.title()`.
5. Finally, we display the plot with `plt.show()`.

This code will generate a scatter plot showing the relationship between sepal length and sepal width for the different species of iris flowers, with each species represented by a different color.

You can customize the plot further, add legends, change colors, or explore different types of plots depending on your dataset and the information you want to convey. Matplotlib and Seaborn provide extensive customization options for creating informative and visually appealing data visualizations.



Name: Athul Krishnan V R

College: JCT College of Engineering and Technology

NM id: au720921244009