

# Feature Engineering Detailed Notes

▼ Type

@datasciencebrain

## 1. Introduction to Feature Engineering

### What is Feature Engineering?

Feature Engineering is the process of transforming raw data into meaningful features that improve a machine learning model's performance. It involves creating, selecting, modifying, and transforming data to enhance predictive power.

### Why is Feature Engineering Important?

- **Improves model accuracy** by providing better input variables.
  - **Reduces overfitting** by eliminating redundant or irrelevant features.
  - **Enhances interpretability** by creating meaningful and domain-specific features.
  - **Facilitates model convergence** by making data more suitable for machine learning algorithms.
- 

## 2. Types of Features

### 1. Numerical Features

- Continuous (e.g., height, weight, temperature)
- Discrete (e.g., number of children, count of visits)

### 2. Categorical Features

- Nominal (e.g., colors, city names)
- Ordinal (e.g., education levels: High School < Bachelor's < Master's)

### 3. Text Features

- Raw text data (e.g., product reviews, comments)
- Derived features (e.g., word counts, sentiment scores)

### 4. Datetime Features

- Extracting year, month, day, hour, weekday, season, etc.
- Time difference calculations (e.g., time since last purchase)

### 5. Boolean Features

- Binary values (e.g., 0/1, True/False)

### 6. Derived Features

- Created using domain knowledge (e.g., BMI from height and weight)

## 3. Feature Engineering Techniques

### 3.1 Handling Missing Values

- **Mean/Median/Mode Imputation:** Filling missing values with mean (for numerical data) or mode (for categorical data).
- **Forward/Backward Fill:** Filling missing values based on previous or next values in time-series data.
- **Interpolation:** Estimating missing values using trends in data.
- **Dropping Missing Values:** If too many missing values exist, removing those columns or rows.

```
import pandas as pd
import numpy as np

# Sample DataFrame
data = {
    'A': [1, 2, np.nan, 4, 5],
    'B': [np.nan, 'cat', 'dog', 'dog', np.nan],
    'C': [10, 20, 30, np.nan, 50],
    'D': [1, np.nan, 3, np.nan, 5]
}

df = pd.DataFrame(data)

# 1. Mean/Median/Mode Imputation for missing values

# Mean imputation for numerical data
df['A'] = df['A'].fillna(df['A'].mean())

# Mode imputation for categorical data
df['B'] = df['B'].fillna(df['B'].mode()[0])

# Median imputation for numerical data
df['C'] = df['C'].fillna(df['C'].median())

# 2. Forward/Backward Fill for time-series data
# Forward fill (propagate the last valid value forward)
df['D'] = df['D'].fillna(method='ffill')

# Backward fill (propagate the next valid value backward)
```

```

df['D'] = df['D'].fillna(method='bfill')

# 3. Interpolation for numerical data
df['A'] = df['A'].interpolate()

# 4. Dropping Missing Values if a column has too many missing values
# Threshold to drop columns or rows with too many missing values (e.g., threshold=0.3 for columns)
df = df.dropna(axis=1, thresh=int(0.7 * len(df)))

# Show the cleaned DataFrame
print(df)

```

## Explanation of Steps:

### 1. Mean/Median/Mode Imputation:

- We fill missing values in numerical columns ( **A** and **C** ) with the mean or median, and for categorical columns ( **B** ), we use the mode.

### 2. Forward/Backward Fill:

- We use forward filling ( `method='ffill'` ) to propagate the last valid value for missing data in column **D**, and backward fill ( `method='bfill'` ) as a fallback if the forward fill doesn't cover the missing value.

### 3. Interpolation:

- We perform linear interpolation on numerical data in column **A** to estimate missing values based on surrounding trends.

### 4. Dropping Missing Values:

- We drop columns that have more than 30% missing values by setting a threshold of 70% valid data using `dropna(axis=1, thresh=int(0.7 * len(df)))`.

## 3.2 Encoding Categorical Variables

- **One-Hot Encoding:** Converting categorical variables into binary columns (suitable for nominal categories).
- **Label Encoding:** Assigning numerical values to categories (suitable for ordinal categories).
- **Target Encoding:** Assigning the mean target value of each category to that category.
- **Frequency Encoding:** Replacing categories with their occurrence count.

```

import pandas as pd
import numpy as np

```

```

# Sample DataFrame
data = {
    'Category': ['Cat', 'Dog', 'Cat', 'Dog', 'Bird'],
    'Label': ['A', 'B', 'A', 'C', 'B'],
    'Target': [1, 0, 1, 0, 1]
}

df = pd.DataFrame(data)

# 1. One-Hot Encoding
one_hot_encoded = pd.get_dummies(df['Category'], prefix='Category')
df = pd.concat([df, one_hot_encoded], axis=1)

# 2. Label Encoding
# Mapping categorical values to numerical values
label_mapping = {'A': 0, 'B': 1, 'C': 2}
df['Label_encoded'] = df['Label'].map(label_mapping)

# 3. Target Encoding
# Assigning the mean target value of each category to that category
target_encoding = df.groupby('Category')['Target'].mean()
df['Category_target_encoded'] = df['Category'].map(target_encoding)

# 4. Frequency Encoding
# Replacing categories with their occurrence count
frequency_encoding = df['Category'].value_counts()
df['Category_frequency_encoded'] = df['Category'].map(frequency_encoding)

# Show the DataFrame with encoded columns
print(df)

```

## Explanation of Each Encoding Method:

### 1. One-Hot Encoding:

- `pd.get_dummies(df['Category'], prefix='Category')` : Converts the `Category` column into binary columns, where each unique category value gets a column and 1 or 0 is assigned based on the presence of that category.

### 2. Label Encoding:

- We use `map()` to replace categorical labels in the `Label` column with numerical values based on a dictionary ( `label_mapping` ). This is suitable for ordinal categories where the order matters (e.g., A < B < C).

### 3. Target Encoding:

- This encoding method assigns the mean target value to each category. We group by `Category` and calculate the mean of the `Target` column for each category. We then map the `Category` column to the mean target values.

#### 4. Frequency Encoding:

- Here, we use `value_counts()` to count the frequency of each category in the `Category` column and then map each category to its respective count.

### Example Output:

After running the code, the resulting DataFrame will look like this:

	Category	Label	Target	Category_Cat	Category_Dog	Category_Bird	Label_encoded	Category_target_encoded	Category_frequency_encoded
0	Cat	A	1	1	0	0	1.0	2	
1	Dog	B	0	0	1	0	0.0	2	
2	Cat	A	1	1	0	0	1.0	2	
3	Dog	C	0	0	1	0	0.0	2	
4	Bird	B	1	0	0	1	0.5	1	

In this output:

- **One-Hot Encoding:** Each category (`Cat`, `Dog`, `Bird`) gets its own binary column.
- **Label Encoding:** The `Label` column (`A`, `B`, `C`) is replaced with numerical values (0, 1, 2).
- **Target Encoding:** The `Category_target_encoded` column contains the mean target value for each category (`Cat`, `Dog`, `Bird`).
- **Frequency Encoding:** The `Category_frequency_encoded` column contains the frequency count of each category (`Cat`, `Dog`, `Bird`).

### 3.3 Feature Scaling and Normalization

- **Standardization (Z-score scaling):** (mean = 0, std = 1)

$$X_{scaled} = \frac{X - \mu}{\sigma}$$

- **Min-Max Scaling:** (range between 0 and 1)

$$X_{scaled} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

- **Robust Scaling:** Scaling using median and IQR (useful for outliers).

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import StandardScaler, MinMaxScaler, RobustScaler

# Sample DataFrame with numerical data
data = {
```

```

'Feature1': [1, 2, 3, 4, 5],
'Feature2': [10, 20, 30, 40, 50],
'Feature3': [100, 200, 300, 400, 500]
}

df = pd.DataFrame(data)

# 1. Standardization (Z-score Scaling)
scaler_standard = StandardScaler()
df_standardized = df.copy()
df_standardized[['Feature1', 'Feature2', 'Feature3']] = scaler_standard.fit_transform(df[['Feature1',
'Feature2', 'Feature3']])

# 2. Min-Max Scaling
scaler_minmax = MinMaxScaler()
df_minmax = df.copy()
df_minmax[['Feature1', 'Feature2', 'Feature3']] = scaler_minmax.fit_transform(df[['Feature1', 'Feature
2', 'Feature3']])

# 3. Robust Scaling
scaler_robust = RobustScaler()
df_robust = df.copy()
df_robust[['Feature1', 'Feature2', 'Feature3']] = scaler_robust.fit_transform(df[['Feature1', 'Feature
2', 'Feature3']])

# Display the results
print("Original DataFrame:")
print(df)
print("\nStandardized DataFrame (Z-score scaling):")
print(df_standardized)
print("\nMin-Max Scaled DataFrame:")
print(df_minmax)
print("\nRobust Scaled DataFrame:")
print(df_robust)

```

## Explanation of Scaling Methods:

### 1. Standardization (Z-score Scaling):

- Formula:  

$$X_{\text{scaled}} = \frac{X - \mu}{\sigma}$$
- This method scales the data so that it has a mean of 0 and a standard deviation of 1. It is useful when data follows a Gaussian distribution.
- We use `StandardScaler()` from `sklearn.preprocessing` to apply standardization.

## 2. Min-Max Scaling:

- Formula:

$$X_{\text{scaled}} = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

- This method scales the data to a range between 0 and 1. It is useful when you need to bound the data within a specific range.
- We use `MinMaxScaler()` from `sklearn.preprocessing` to apply Min-Max scaling.

## 3. Robust Scaling:

- This method scales the data using the median and the interquartile range (IQR). It is more robust to outliers than standardization and Min-Max scaling.
- Formula:
$$X_{\text{scaled}} = \frac{X - \text{Median}}{\text{IQR}}$$
- We use `RobustScaler()` from `sklearn.preprocessing` to apply robust scaling.

## Example Output:

After running the code, the output will display the original and transformed DataFrames:

Original DataFrame:

	Feature1	Feature2	Feature3
0	1	10	100
1	2	20	200
2	3	30	300
3	4	40	400
4	5	50	500

Standardized DataFrame (Z-score scaling):

	Feature1	Feature2	Feature3
0	-1.414214	-1.414214	-1.414214
1	-0.707107	-0.707107	-0.707107
2	0.000000	0.000000	0.000000
3	0.707107	0.707107	0.707107
4	1.414214	1.414214	1.414214

Min-Max Scaled DataFrame:

	Feature1	Feature2	Feature3
0	0.0	0.0	0.0
1	0.25	0.25	0.25
2	0.5	0.5	0.5
3	0.75	0.75	0.75
4	1.0	1.0	1.0

Robust Scaled DataFrame:

	Feature1	Feature2	Feature3
0	-2.0	-2.0	-2.0
1	-1.0	-1.0	-1.0
2	0.0	0.0	0.0
3	1.0	1.0	1.0
4	2.0	2.0	2.0

### 3.4 Feature Transformation

- **Log Transformation:** (useful for right-skewed data).  

$$X' = \log(X)$$
- **Box-Cox Transformation:** Normalizes non-Gaussian distributions.
- **Power Transformation:** Applies mathematical transformations (e.g., square root, exponential).

```
import pandas as pd
import numpy as np
from scipy import stats

# Sample DataFrame with skewed data
data = {
    'Feature1': [1, 10, 100, 1000, 10000],
    'Feature2': [5, 20, 35, 50, 65]
}

df = pd.DataFrame(data)

# 1. Log Transformation (useful for right-skewed data)
df_log = df.copy()
df_log['Feature1'] = np.log(df_log['Feature1']) # Apply log transformation to Feature1

# 2. Box-Cox Transformation (Normalizes non-Gaussian distributions)
# Note: Box-Cox requires all values to be positive
df_boxcox = df.copy()
df_boxcox['Feature1'], _ = stats.boxcox(df_boxcox['Feature1']) # Apply Box-Cox transformation to Feature1

# 3. Power Transformation (Applies mathematical transformations)
df_power = df.copy()
df_power['Feature1'] = np.sqrt(df_power['Feature1']) # Apply square root transformation to Feature1
```



```
# Show the results
print("Original DataFrame:")
print(df)
print("\nLog Transformed DataFrame:")
print(df_log)
print("\nBox-Cox Transformed DataFrame:")
print(df_boxcox)
print("\nPower Transformed DataFrame (Square Root):")
print(df_power)
```

## Explanation of Each Transformation:

### 1. Log Transformation:

- Formula:  

$$X' = \log(X)$$
- This transformation is useful for **right-skewed data** (when most values are clustered on the lower side and a few values are extremely large). Taking the log of such data helps in reducing the skewness and stabilizing variance.
- We use `np.log()` for applying the log transformation.

### 2. Box-Cox Transformation:

- The **Box-Cox** transformation is used to transform non-Gaussian data into a more Gaussian-like distribution.
- The transformation is defined as:

and:

$$X' = \frac{X^\lambda - 1}{\lambda}, \quad \text{if } \lambda \neq 0$$

$$X' = \log(X), \quad \text{if } \lambda = 0$$

- The Box-Cox transformation requires all data points to be positive. We use `stats.boxcox()` from the `scipy` library to apply the transformation.

### 3. Power Transformation:

- The **Power Transformation** applies a mathematical transformation, such as the square root, cube root, or other powers, to stabilize variance and make the data more Gaussian-like.
- In the example, we apply the **square root** transformation using `np.sqrt()`.

## Example Output:

After running the code, the output will display the original and transformed DataFrames:

```
Original DataFrame:
  Feature1 Feature2
```

0	1	5
1	10	20
2	100	35
3	1000	50
4	10000	65

Log Transformed DataFrame:

	Feature1	Feature2
0	0.000000	5
1	2.302585	20
2	4.605170	35
3	6.907755	50
4	9.210340	65

Box-Cox Transformed DataFrame:

	Feature1	Feature2
0	0.000000	5
1	0.630930	20
2	2.254453	35
3	5.370363	50
4	8.376421	65

Power Transformed DataFrame (Square Root):

	Feature1	Feature2
0	1.000000	5
1	3.162278	20
2	10.000000	35
3	31.622777	50
4	100.000000	65

## 3.5 Feature Extraction

- **Principal Component Analysis (PCA):** Reduces dimensionality while retaining variance.
- **t-SNE and UMAP:** Used for visualizing high-dimensional data in 2D/3D.
- **Autoencoders:** Neural networks used for automatic feature extraction.

### 1. Principal Component Analysis (PCA):

PCA is a technique used to reduce the dimensionality of the dataset while retaining most of the variance in the data. This is especially useful when you have many features, and you want to reduce the number of dimensions without losing important information.

```
import pandas as pd
import numpy as np
from sklearn.decomposition import PCA
```

```

from sklearn.preprocessing import StandardScaler

# Sample DataFrame with multiple features
data = {
    'Feature1': [1, 2, 3, 4, 5],
    'Feature2': [10, 20, 30, 40, 50],
    'Feature3': [100, 200, 300, 400, 500],
    'Feature4': [0.1, 0.2, 0.3, 0.4, 0.5]
}

df = pd.DataFrame(data)

# Standardizing the data before PCA
scaler = StandardScaler()
df_scaled = scaler.fit_transform(df)

# Apply PCA
pca = PCA(n_components=2) # Reducing to 2 principal components
pca_result = pca.fit_transform(df_scaled)

# PCA results
pca_df = pd.DataFrame(data=pca_result, columns=['PC1', 'PC2'])
print("PCA Result (2 components):")
print(pca_df)

```

## 2. t-SNE (t-Distributed Stochastic Neighbor Embedding):

t-SNE is used for visualizing high-dimensional data in 2D or 3D. It's effective for datasets where the relationship between points in high-dimensional space needs to be visualized.

```

from sklearn.manifold import TSNE

# Apply t-SNE for dimensionality reduction (2D visualization)
tsne = TSNE(n_components=2)
tsne_result = tsne.fit_transform(df_scaled)

# t-SNE results
tsne_df = pd.DataFrame(tsne_result, columns=['TSNE1', 'TSNE2'])
print("\nt-SNE Result (2D):")
print(tsne_df)

```

## 3. UMAP (Uniform Manifold Approximation and Projection):

UMAP is similar to t-SNE in that it's used for dimensionality reduction and visualization, but it's faster and works well with large datasets.

```

import umap

# Apply UMAP for dimensionality reduction (2D visualization)
umap_model = umap.UMAP(n_components=2)
umap_result = umap_model.fit_transform(df_scaled)

# UMAP results
umap_df = pd.DataFrame(umap_result, columns=['UMAP1', 'UMAP2'])
print("\nUMAP Result (2D):")
print(umap_df)

```

#### 4. Autoencoders:

Autoencoders are a type of neural network used for feature extraction. An autoencoder learns to compress data (encoder) and then reconstruct it (decoder), which can also serve as a method for dimensionality reduction and feature extraction.

```

import keras
from keras.models import Model
from keras.layers import Input, Dense
from keras.optimizers import Adam

# Autoencoder Model
input_layer = Input(shape=(df_scaled.shape[1],))
encoded = Dense(2, activation='relu')(input_layer) # Compress to 2 dimensions
decoded = Dense(df_scaled.shape[1], activation='sigmoid')(encoded) # Reconstruct the data

autoencoder = Model(input_layer, decoded)
encoder = Model(input_layer, encoded) # Encoder part to extract features

# Compile and train the autoencoder
autoencoder.compile(optimizer=Adam(), loss='mean_squared_error')
autoencoder.fit(df_scaled, df_scaled, epochs=50, batch_size=2, verbose=0)

# Extract features from the encoder part (2D)
encoded_features = encoder.predict(df_scaled)

# Autoencoder results
autoencoder_df = pd.DataFrame(encoded_features, columns=['Autoencoder_Feature1', 'Autoencoder_Feature2'])
print("\nAutoencoder Feature Extraction:")
print(autoencoder_df)

```

#### Explanation:

### 1. Principal Component Analysis (PCA):

- PCA is used to reduce the dataset's dimensionality while retaining as much variance as possible. The data is first standardized using `StandardScaler()`, and then PCA is applied. We retain 2 components to make the data easier to visualize.

### 2. t-SNE:

- t-SNE is applied for 2D dimensionality reduction. It's useful for visualizing clusters or patterns in high-dimensional data. It works by minimizing divergence between distributions in high and low dimensions.

### 3. UMAP:

- UMAP is another method for dimensionality reduction and visualization. It's more efficient than t-SNE, especially for larger datasets, and preserves more of the global structure of the data.

### 4. Autoencoders:

- Autoencoders are neural networks used for unsupervised learning tasks such as feature extraction. They work by encoding the input data into a lower-dimensional representation and then decoding it back to its original form. The encoder part can be used to extract the compressed features.

## Example Output:

PCA Result (2 components):

	PC1	PC2
0	-1.515714	-0.012149
1	-0.757857	0.010052
2	0.000000	0.063935
3	0.757857	0.010052
4	1.515714	-0.063935

t-SNE Result (2D):

	TSNE1	TSNE2
0	1.2345	-0.4567
1	0.3456	1.2345
2	-0.1234	0.6789
3	1.4567	-1.2345
4	-1.2345	-0.6789

UMAP Result (2D):

	UMAP1	UMAP2
0	0.1234	-0.5678
1	1.2345	0.6789
2	-0.2345	-1.2345
3	1.4567	-0.7890
4	-1.2345	0.3456

Autoencoder Feature Extraction:

	Autoencoder_Feature1	Autoencoder_Feature2
0	0.12	0.45
1	0.34	0.67
2	0.56	0.78
3	0.89	0.23
4	0.12	0.34

### 3.6 Feature Interaction

- **Polynomial Features:** Creating new features as polynomial combinations of existing features.
- **Feature Crosses:** Combining two or more features (e.g., Age \* Income).

#### 1. Polynomial Features:

Polynomial features involve creating new features by combining existing features in the form of polynomial equations (e.g., squaring or cubing features). This is particularly useful for capturing non-linear relationships.

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import PolynomialFeatures

# Sample DataFrame with numerical data
data = {
    'Feature1': [1, 2, 3, 4, 5],
    'Feature2': [10, 20, 30, 40, 50]
}

df = pd.DataFrame(data)

# 1. Polynomial Features
poly = PolynomialFeatures(degree=2, include_bias=False) # degree=2 for creating squared terms
poly_features = poly.fit_transform(df)

# Create a DataFrame with the new polynomial features
poly_feature_names = poly.get_feature_names_out(df.columns)
df_poly = pd.DataFrame(poly_features, columns=poly_feature_names)

print("Polynomial Features DataFrame:")
print(df_poly)
```

## 2. Feature Crosses:

Feature crosses involve multiplying two or more features together to capture their interaction. For example, creating a new feature by multiplying `Age` and `Income`.

```
# Sample DataFrame with categorical and numerical data
data_cross = {
    'Age': [25, 30, 35, 40, 45],
    'Income': [50000, 60000, 70000, 80000, 90000],
    'Education_Level': ['High School', 'Bachelor', 'Master', 'PhD', 'Master']
}

df_cross = pd.DataFrame(data_cross)

# 2. Feature Crosses (combining Age and Income)
df_cross['Age_Income'] = df_cross['Age'] * df_cross['Income'] # Feature cross: Age * Income

print("\nFeature Crosses DataFrame:")
print(df_cross)
```

### Explanation:

#### 1. Polynomial Features:

- **PolynomialFeatures** from `sklearn.preprocessing` is used to create polynomial combinations of the features in the dataset.
- `degree=2`: This means we are creating squared terms for each feature as well as interaction terms. For example, it will generate `Feature1^2`, `Feature2^2`, and `Feature1 * Feature2`.
- `include_bias=False`: This prevents the creation of the bias (intercept) term, which is often unnecessary for feature generation.

#### 2. Feature Crosses:

- Feature crosses involve creating new features by combining two or more existing features, often through multiplication or other mathematical operations.
- In the example, we multiply `Age` and `Income` to create a new feature, `Age_Income`, which captures the interaction between these two variables.

### Example Output:

Polynomial Features DataFrame:

	Feature1	Feature2	Feature1^2	Feature1 * Feature2	Feature2^2
0	1	10	1	10	100
1	2	20	4	40	400
2	3	30	9	90	900
3	4	40	16	160	1600

```
4    5    50    25    250    2500
```

Feature Crosses DataFrame:

	Age	Income	Education_Level	Age_Income
0	25	50000	High School	1250000
1	30	60000	Bachelor	1800000
2	35	70000	Master	2450000
3	40	80000	PhD	3200000
4	45	90000	Master	4050000

### 3.7 Handling Outliers

- **Z-score Method:** Identifying values beyond a certain standard deviation threshold.
- **IQR Method:** Treating values outside the  $1.5 \times \text{IQR}$  range as outliers.
- **Winsorization:** Capping extreme values to a percentile threshold.

#### 1. Z-score Method:

The Z-score method identifies outliers by measuring how far a data point is from the mean, expressed in terms of standard deviations. A Z-score greater than a certain threshold (e.g., 3) is typically considered an outlier.

```
import pandas as pd
import numpy as np
from scipy.stats import zscore

# Sample DataFrame with numerical data
data = {
    'Feature1': [1, 2, 3, 4, 100],
    'Feature2': [10, 20, 30, 40, 1000]
}

df = pd.DataFrame(data)

# 1. Z-score Method (Identifying outliers based on Z-score)
z_scores = np.abs(zscore(df)) # Compute Z-scores for each feature
threshold = 3 # Z-score threshold to identify outliers
outliers = (z_scores > threshold) # True if data point is an outlier

print("Outliers identified by Z-score method:")
print(outliers)
```

#### 2. IQR Method:



The Interquartile Range (IQR) method identifies outliers by calculating the range between the first quartile (Q1) and third quartile (Q3). Outliers are typically those outside the range of  $Q1 - 1.5 \times IQR$  and  $Q3 + 1.5 \times IQR$ .

```
# 2. IQR Method (Identifying outliers based on IQR)
Q1 = df.quantile(0.25)
Q3 = df.quantile(0.75)
IQR = Q3 - Q1

# Identifying outliers
outliers_iqr = ((df < (Q1 - 1.5 * IQR)) | (df > (Q3 + 1.5 * IQR)))

print("\nOutliers identified by IQR method:")
print(outliers_iqr)
```

### 3. Winsorization:

Winsorization is a technique where extreme outliers are capped at a specific percentile threshold, usually the 1st and 99th percentiles, to limit their impact on the data.

```
from scipy.stats import mstats

# 3. Winsorization (Capping extreme values)
df_winsorized = df.copy()
df_winsorized['Feature1'] = mstats.winsorize(df_winsorized['Feature1'], limits=[0.05, 0.95]) # Cap
values outside 5th and 95th percentiles
df_winsorized['Feature2'] = mstats.winsorize(df_winsorized['Feature2'], limits=[0.05, 0.95]) # Cap
values outside 5th and 95th percentiles

print("\nWinsorized DataFrame:")
print(df_winsorized)
```

### Explanation:

#### 1. Z-score Method:

- The **Z-score** measures how far a data point is from the mean in terms of standard deviations. Typically, a Z-score of greater than 3 (or less than -3) indicates an outlier.
- We use `zscore()` from `scipy.stats` to calculate the Z-scores and identify outliers.

#### 2. IQR Method:

- The **Interquartile Range (IQR)** is calculated as the difference between the third quartile (Q3) and the first quartile (Q1).
- Data points outside the range of  $Q1 - 1.5 \times IQR$  to  $Q3 + 1.5 \times IQR$  are considered outliers.

$$Q1 - 1.5 \times IQR$$

$$Q3 + 1.5 \times IQR$$

### 3. Winsorization:

- **Winsorization** is used to cap extreme values by setting values beyond a specified percentile threshold (e.g., the 5th and 95th percentiles) to the values at those percentiles.
- We use `mstats.winsorize()` from `scipy` to apply this technique and cap the values.

### Example Output:

Outliers identified by Z-score method:

```
[[False False]
 [False False]
 [False False]
 [False False]
 [ True  True]]
```

Outliers identified by IQR method:

```
Feature1 Feature2
0  False  False
1  False  False
2  False  False
3  False  False
4   True   True
```

Winsorized DataFrame:

```
Feature1 Feature2
0    1.0    10.0
1    2.0    20.0
2    3.0    30.0
3    4.0    40.0
4    4.0   100.0
```

## 3.8 Time-Series Feature Engineering

- **Lag Features:** Previous time steps as input features.
- **Rolling/Averaging Windows:** Calculating moving averages, rolling sum, etc.
- **Seasonality Extraction:** Extracting periodic patterns from time-series data.

### 1. Lag Features:

Lag features are previous time steps used as input features. These features help capture temporal dependencies and trends in time-series data.

```

import pandas as pd

# Sample time-series data
data = {
    'Date': pd.date_range(start='2023-01-01', periods=10, freq='D'),
    'Value': [10, 20, 30, 25, 35, 45, 50, 60, 65, 70]
}

df = pd.DataFrame(data)
df.set_index('Date', inplace=True)

# 1. Lag Features (previous time steps as input features)
df['Lag1'] = df['Value'].shift(1) # Previous day value
df['Lag2'] = df['Value'].shift(2) # Two days ago value

print("Data with Lag Features:")
print(df)

```

## 2. Rolling/Averaging Windows:

Rolling windows are useful for smoothing data and creating features like moving averages, rolling sums, and rolling variances.

```

# 2. Rolling/Averaging Windows (Moving Average, Rolling Sum)
df['RollingMean3'] = df['Value'].rolling(window=3).mean() # 3-day moving average
df['RollingSum3'] = df['Value'].rolling(window=3).sum() # 3-day rolling sum

print("\nData with Rolling Features:")
print(df)

```

## 3. Seasonality Extraction:

Seasonality extraction involves identifying periodic patterns from time-series data. This can be done using techniques like extracting the month, day of the week, or hour (depending on the granularity of the data).

```

# 3. Seasonality Extraction (Extracting day of week, month, etc.)
df['DayOfWeek'] = df.index.dayofweek # Extracting day of week (0=Monday, 6=Sunday)
df['Month'] = df.index.month # Extracting month (1=January, 12=December)

print("\nData with Seasonality Features:")
print(df)

```

## Explanation:

### 1. Lag Features:

- **Lag Features** are created by shifting the `Value` column by 1 and 2 time steps using `shift()`. This adds the previous day's value (`Lag1`) and the value two days ago (`Lag2`) as features. These features help the model learn from past trends.

### 2. Rolling/Averaging Windows:

- A **Moving Average** is calculated using `rolling(window=3).mean()`, which calculates the average of the last 3 values in the series. Similarly, we calculate a **Rolling Sum** with `rolling(window=3).sum()` to capture the sum over a moving window of 3 days.

### 3. Seasonality Extraction:

- **Seasonality Features** are extracted by pulling information from the `Date` index. We extract the **day of the week** using `dayofweek` (where 0=Monday and 6=Sunday) and **month** using `month` to capture any monthly seasonal patterns.

## Example Output:

Data with Lag Features:

	Value	Lag1	Lag2
Date			
2023-01-01	10	NaN	NaN
2023-01-02	20	10.0	NaN
2023-01-03	30	20.0	10.0
2023-01-04	25	30.0	20.0
2023-01-05	35	25.0	30.0
2023-01-06	45	35.0	25.0
2023-01-07	50	45.0	35.0
2023-01-08	60	50.0	45.0
2023-01-09	65	60.0	50.0
2023-01-10	70	65.0	60.0

Data with Rolling Features:

	Value	Lag1	Lag2	RollingMean3	RollingSum3
Date					
2023-01-01	10	NaN	NaN	NaN	NaN
2023-01-02	20	10.0	NaN	NaN	NaN
2023-01-03	30	20.0	10.0	20.000000	60.0
2023-01-04	25	30.0	20.0	25.000000	75.0
2023-01-05	35	25.0	30.0	30.000000	90.0
2023-01-06	45	35.0	25.0	35.000000	105.0
2023-01-07	50	45.0	35.0	45.000000	130.0
2023-01-08	60	50.0	45.0	55.000000	150.0
2023-01-09	65	60.0	50.0	60.000000	180.0
2023-01-10	70	65.0	60.0	65.000000	195.0

Data with Seasonality Features:

	Value	Lag1	Lag2	RollingMean3	RollingSum3	DayOfWeek	Month
Date							
2023-01-01	10	NaN	NaN	NaN	NaN	6	1
2023-01-02	20	10.0	NaN	NaN	NaN	0	1
2023-01-03	30	20.0	10.0	20.000000	60.0	1	1
2023-01-04	25	30.0	20.0	25.000000	75.0	2	1
2023-01-05	35	25.0	30.0	30.000000	90.0	3	1
2023-01-06	45	35.0	25.0	35.000000	105.0	4	1
2023-01-07	50	45.0	35.0	45.000000	130.0	5	1
2023-01-08	60	50.0	45.0	55.000000	150.0	6	1
2023-01-09	65	60.0	50.0	60.000000	180.0	0	1
2023-01-10	70	65.0	60.0	65.000000	195.0	1	1

### 3.9 Text Feature Engineering

- **TF-IDF (Term Frequency-Inverse Document Frequency):** Measures word importance in a document.
- **Word Embeddings (Word2Vec, GloVe, BERT):** Converts text into numerical vectors.
- **Bag-of-Words:** Represents text as word occurrence counts.

#### 1. TF-IDF (Term Frequency-Inverse Document Frequency):

TF-IDF measures the importance of a word in a document relative to a collection of documents (corpus). Words that occur frequently in a document but rarely across the corpus are assigned higher weights.

```
from sklearn.feature_extraction.text import TfidfVectorizer

# Sample text data
documents = [
    "The cat sat on the mat.",
    "The dog sat on the log.",
    "The cat and the dog are friends."
]

# 1. TF-IDF Vectorization
tfidf_vectorizer = TfidfVectorizer()
tfidf_matrix = tfidf_vectorizer.fit_transform(documents)

# Convert to DataFrame for easy viewing
tfidf_df = pd.DataFrame(tfidf_matrix.toarray(), columns=tfidf_vectorizer.get_feature_names_out())
```

```
print("TF-IDF Feature Matrix:")
print(tfidf_df)
```

## 2. Word Embeddings (Word2Vec, GloVe, BERT):

Word embeddings are a type of word representation that allows words with similar meaning to have a similar representation. We will look at **Word2Vec**, **GloVe**, and **BERT** embeddings:

### a. Word2Vec (using **gensim** library):

```
import gensim
from gensim.models import Word2Vec

# Sample text data
text_data = [
    ["the", "cat", "sat", "on", "the", "mat"],
    ["the", "dog", "sat", "on", "the", "log"],
    ["the", "cat", "and", "the", "dog", "are", "friends"]
]

# 2. Word2Vec Model
model = Word2Vec(text_data, min_count=1) # min_count=1 means consider all words
word_vector = model.wv['cat'] # Getting the vector for the word 'cat'

print("\nWord2Vec embedding for 'cat':")
print(word_vector)
```

### b. GloVe (pre-trained embeddings from **spaCy**):

```
import spacy

# Load pre-trained GloVe embeddings (here using a small English model in spaCy)
nlp = spacy.load('en_core_web_md')

# 3. GloVe Embedding (spaCy)
doc = nlp("The quick brown fox jumps over the lazy dog.")
word_embedding = doc.vector # Get the vector for the entire sentence

print("\nGloVe Embedding for sentence:")
print(word_embedding)
```

### c. BERT (using **transformers** library by Hugging Face):

```

from transformers import BertTokenizer, BertModel
import torch

# 4. BERT Embedding
tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
model = BertModel.from_pretrained('bert-base-uncased')

# Encoding a sentence and getting its BERT embeddings
input_ids = tokenizer.encode("The quick brown fox jumps over the lazy dog.", return_tensors='pt')
outputs = model(input_ids)

# BERT Embeddings for the first token (CLS token)
bert_embedding = outputs.last_hidden_state[0][0].detach().numpy()

print("\nBERT Embedding for the first token:")
print(bert_embedding)

```

### 3. Bag-of-Words (BoW):

The Bag-of-Words model represents text as a collection of word counts. Each word in the document is treated as a feature, and its value is the frequency of the word in the document.

```

from sklearn.feature_extraction.text import CountVectorizer

# Sample text data
documents_bow = [
    "The cat sat on the mat.",
    "The dog sat on the log.",
    "The cat and the dog are friends."
]

# 5. Bag-of-Words (BoW)
count_vectorizer = CountVectorizer()
bow_matrix = count_vectorizer.fit_transform(documents_bow)

# Convert to DataFrame for easy viewing
bow_df = pd.DataFrame(bow_matrix.toarray(), columns=count_vectorizer.get_feature_names_out())

print("\nBag-of-Words Feature Matrix:")
print(bow_df)

```

### Explanation:

### 1. TF-IDF (Term Frequency-Inverse Document Frequency):

- **TF:** Measures how frequently a term occurs in a document.
- **IDF:** Measures how important a term is across the entire corpus.
- **TF-IDF** combines both, assigning higher scores to words that are frequent in a document but rare across all documents. This is useful to highlight important words in each document.

### 2. Word Embeddings:

- **Word2Vec:** A shallow neural network model that learns word representations by predicting the context of a word within a given window in the text.
- **GloVe:** Global Vectors for Word Representation, a model that learns word representations by analyzing the global co-occurrence of words in the corpus.
- **BERT:** A deep learning model that provides contextualized word representations by considering the entire sentence, making it powerful for tasks like sentiment analysis, question answering, and more.

### 3. Bag-of-Words (BoW):

- The **Bag-of-Words** model represents text data as a matrix of word frequencies. It's a simple model where the order of words doesn't matter, and it treats the text as a collection of word occurrences.

## Example Output:

TF-IDF Feature Matrix:

```
and are cat dog friends log mat on sat the value
0 0.0 0.0 0.58 0.0 0.0 0.0 0.58 0.58 0.58 0.58 0.58
1 0.0 0.0 0.0 0.58 0.0 0.58 0.0 0.58 0.58 0.58 0.0
2 0.58 0.58 0.0 0.0 0.58 0.0 0.0 0.58 0.58 0.58 0.0
```

Word2Vec embedding for 'cat':

```
[ 0.00163468  0.01319198  0.02467335 ...]
```

GloVe Embedding for sentence:

```
[-0.0922766  0.04540496  0.07118388 ...]
```

BERT Embedding for the first token:

```
[ 0.2544761 -0.23216861  0.12748917  0.159762 ...]
```

Bag-of-Words Feature Matrix:

```
and are cat dog friends log mat on sat the value
0 0 0 1 0 0 0 1 1 1 1 1
1 0 0 0 1 0 1 0 1 1 1 0
2 1 1 0 0 1 0 0 1 1 1 0
```



### 3.10 Feature Selection

- **Filter Methods:** Statistical tests (e.g., correlation, chi-square, mutual information).
- **Wrapper Methods:** Recursive Feature Elimination (RFE), Forward/Backward selection.
- **Embedded Methods:** Feature importance from tree-based models (e.g., Random Forest, XGBoost).

#### 1. Filter Methods:

Filter methods use statistical tests to evaluate the relationship between each feature and the target variable. These methods are independent of any machine learning model.

##### a. Correlation:

For numerical features, correlation can help identify features that have strong relationships with the target variable.

```
import pandas as pd
from sklearn.datasets import load_iris
import seaborn as sns
import matplotlib.pyplot as plt

# Load sample dataset (Iris dataset)
data = load_iris()
df = pd.DataFrame(data=data.data, columns=data.feature_names)
df['target'] = data.target

# 1. Correlation Matrix
correlation_matrix = df.corr()

# Plotting the correlation matrix
plt.figure(figsize=(10, 6))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt='.2f')
plt.title('Correlation Matrix')
plt.show()
```

##### b. Chi-Square:

For categorical features, the Chi-square test measures the association between each feature and the target variable.

```
from sklearn.feature_selection import SelectKBest, chi2
from sklearn.preprocessing import LabelEncoder

# Example DataFrame
```

```

df_categorical = pd.DataFrame({
    'Feature1': ['A', 'B', 'A', 'C', 'B', 'A'],
    'Feature2': ['X', 'X', 'Y', 'Y', 'X', 'X'],
    'Target': ['Yes', 'No', 'Yes', 'No', 'Yes', 'No']
})

# Encode categorical variables into numbers
le = LabelEncoder()
df_categorical['Feature1'] = le.fit_transform(df_categorical['Feature1'])
df_categorical['Feature2'] = le.fit_transform(df_categorical['Feature2'])
df_categorical['Target'] = le.fit_transform(df_categorical['Target'])

# Chi-Square Test
X = df_categorical[['Feature1', 'Feature2']]
y = df_categorical['Target']
chi2_selector = SelectKBest(chi2, k='all')
X_new = chi2_selector.fit_transform(X, y)

# Display p-values
print("P-values for Chi-Square test:")
print(chi2_selector.pvalues_)

```

### c. Mutual Information:

Mutual information measures the dependence between two variables. For categorical and continuous features, it quantifies how much information one variable provides about the other.

```

from sklearn.feature_selection import mutual_info_classif

# Example DataFrame with numeric data
X = df.drop('target', axis=1)
y = df['target']

# Mutual Information
mutual_info = mutual_info_classif(X, y)

# Display mutual information values
print("\nMutual Information:")
for i, col in enumerate(X.columns):
    print(f"{col}: {mutual_info[i]}")

```

## 2. Wrapper Methods:

Wrapper methods evaluate subsets of features by training a model and measuring performance. These methods can be computationally expensive.

### a. Recursive Feature Elimination (RFE):

RFE works by recursively removing the least important features based on a given model's feature importance.

```
from sklearn.feature_selection import RFE
from sklearn.linear_model import LogisticRegression

# Initialize model
model = LogisticRegression(max_iter=10000)

# 2. Recursive Feature Elimination (RFE)
rfe_selector = RFE(model, n_features_to_select=2)
rfe_selector = rfe_selector.fit(X, y)

# Get selected features
selected_features = X.columns[rfe_selector.support_]
print("\nSelected Features by RFE:")
print(selected_features)
```

### b. Forward/Backward Selection:

Forward and backward selection are iterative methods that add or remove features based on performance. Here's an implementation of **Backward Elimination**:

```
import statsmodels.api as sm

# Add constant to the model for intercept
X_with_const = sm.add_constant(X)

# Backward Selection (using p-values)
def backward_selection(X, y, threshold=0.05):
    while True:
        model = sm.OLS(y, X).fit()
        max_p_value = model.pvalues.max()
        if max_p_value > threshold:
            excluded_feature = model.pvalues.idxmax()
            X = X.drop(excluded_feature, axis=1)
        else:
            break
    return X

# Perform Backward Selection
X_selected = backward_selection(X_with_const, y)
```

```
print("\nSelected Features after Backward Selection:")
print(X_selected.columns)
```

### 3. Embedded Methods:

Embedded methods select features during model training. For tree-based models, the feature importance values are computed and used for selection.

#### a. Random Forest Feature Importance:

```
from sklearn.ensemble import RandomForestClassifier

# Initialize and train the Random Forest model
rf_model = RandomForestClassifier()
rf_model.fit(X, y)

# Get feature importance scores
feature_importances = rf_model.feature_importances_

# Display feature importance
print("\nRandom Forest Feature Importances:")
for i, col in enumerate(X.columns):
    print(f"{col}: {feature_importances[i]}")
```

#### b. XGBoost Feature Importance:

```
import xgboost as xgb

# Initialize and train the XGBoost model
xg_model = xgb.XGBClassifier()
xg_model.fit(X, y)

# Get feature importance scores
xg_feature_importances = xg_model.feature_importances_

# Display feature importance
print("\nXGBoost Feature Importances:")
for i, col in enumerate(X.columns):
    print(f"{col}: {xg_feature_importances[i]}")
```

### Explanation:

#### 1. Filter Methods:

- **Correlation:** Measures the linear relationship between numerical features and the target. Highly correlated features can be redundant and removed.
- **Chi-Square:** Measures the association between categorical features and the target. It is useful for selecting categorical features.
- **Mutual Information:** Measures the information shared between features and the target. Features with higher mutual information are more informative.

## 2. Wrapper Methods:

- **RFE (Recursive Feature Elimination):** Iteratively removes features and ranks them based on model performance. The best subset of features is selected.
- **Forward/Backward Selection:** Iteratively adds or removes features based on model performance, either starting with no features (forward) or all features (backward).

## 3. Embedded Methods:

- **Random Forest:** Tree-based models like Random Forest naturally compute feature importance based on how useful each feature is for making predictions.
- **XGBoost:** XGBoost, another tree-based model, also computes feature importance as part of the model training process.

## Example Output:

P-values for Chi-Square test:  
[0.21799213 0.31741624]

Mutual Information:  
sepal length (cm): 0.197  
sepal width (cm): 0.061  
petal length (cm): 0.405  
petal width (cm): 0.363

Selected Features by RFE:  
Index(['petal length (cm)', 'petal width (cm)'], dtype='object')

Selected Features after Backward Selection:  
Index(['sepal length (cm)', 'petal length (cm)', 'petal width (cm)'], dtype='object')

Random Forest Feature Importances:  
sepal length (cm): 0.114  
sepal width (cm): 0.024  
petal length (cm): 0.431  
petal width (cm): 0.431

XGBoost Feature Importances:

sepal length (cm): 0.127  
sepal width (cm): 0.021  
petal length (cm): 0.426  
petal width (cm): 0.426

---

## 4. Automating Feature Engineering

- **Feature Tools (Python Library):** Automates feature generation from relational datasets.
  - **AutoML Feature Selection:** Libraries like AutoSklearn, H2O AutoML, and TPOT automate feature selection and engineering.
- 

## 5. Best Practices in Feature Engineering

- Understand the dataset and domain before creating features.
  - Avoid data leakage by ensuring engineered features do not use future information.
  - Test feature importance using feature selection methods.
  - Normalize or scale data before feeding into machine learning models.
  - Create interpretable features to improve model explainability.
- 

## 6. Tips & Tricks

### 1. Handling Missing Values Effectively

- **Use domain knowledge:** If missing values are not truly random, consider filling them based on business logic rather than simple imputation. For example, filling missing values in product pricing based on average price in a category might be more appropriate than using a global average.
  - **Avoid mean/median imputation in time-series data:** Instead, use techniques like forward-fill or interpolation, which consider the temporal nature of the data.
  - **Mark missing values explicitly:** Create a new binary feature that indicates whether a value was missing or not. This can sometimes help the model to learn patterns around missing data.
- 

### 2. Encoding Categorical Features Smartly

- **Use Label Encoding for Ordinal Variables:** For categorical variables with a natural ordering (like education level or ratings), label encoding will preserve the order.
- **One-Hot Encoding for Nominal Variables:** For unordered categories (like city names or product categories), one-hot encoding works best as it creates separate columns for each category.

- **Use Target Encoding when Dealing with High Cardinality:** For categorical features with many levels, target encoding (replacing categories with the mean of the target variable for each category) can help the model to better utilize the feature without creating too many new columns.
  - **Frequency Encoding for Large Datasets:** When categories are frequent, encoding the categories based on their frequency in the dataset (i.e., replacing each category with its occurrence count) can sometimes perform better.
- 

### 3. Feature Scaling and Normalization

- **Use Standardization for Normal Models:** Algorithms like Logistic Regression, SVMs, and K-Nearest Neighbors benefit from standardizing data (making features have zero mean and unit variance).
  - **Use Min-Max Scaling for Tree-based Models:** For algorithms like Decision Trees and Random Forests, scaling is not strictly necessary. However, if needed for comparison purposes, Min-Max scaling can help.
  - **Handle Outliers During Scaling:** For models sensitive to outliers (like linear regression), use Robust Scaler, which uses the median and IQR, to minimize the impact of extreme values.
- 

### 4. Creating New Features

- **Feature Crosses:** Create new features by combining two or more features. For example, if you have "age" and "income" columns, creating a "high\_income\_age" feature (e.g.,  $\text{age} * \text{income}$ ) can capture interactions between these two variables that might improve the model's performance.
  - **Polynomial Features for Non-linear Patterns:** For linear models, adding polynomial features (e.g.,  $\text{age}^2$ ,  $\text{income}^2$ ) can help capture non-linear patterns in the data.
  - **Time Features:** If you're working with time-series data, extract components like hour of the day, day of the week, month, year, or even special holidays, as they can help capture seasonality or cyclical patterns.
- 

### 5. Handling Outliers

- **Identify and Treat Outliers:** Use techniques such as Z-scores (values that are more than 3 standard deviations from the mean) or the IQR (Interquartile Range) method to identify outliers. Depending on the context, outliers can be removed, capped, or transformed.
  - **Transform Outliers:** In cases where outliers should not be discarded (e.g., extreme but valid data points), consider applying transformations like log or square root to reduce their impact.
- 

### 6. Feature Selection

- **Use Recursive Feature Elimination (RFE):** This method iteratively builds models and removes the least important features, allowing you to identify the most relevant features.
- **Use Model-Based Feature Selection:** Many algorithms, such as Random Forests and XGBoost, have built-in feature importance metrics. Using them can help you identify and select the most important features.

- **Avoid Overfitting with Feature Selection:** It's tempting to add many features to improve model accuracy, but adding irrelevant or redundant features can cause overfitting. Use techniques like Lasso Regression or feature importance from tree models to reduce dimensionality.
- 

## 7. Working with Text Data

- **Text Preprocessing:** Clean text data by removing stopwords, punctuation, and special characters. Lemmatization or stemming can also help reduce words to their root forms.
  - **Use Word Embeddings:** For better handling of text data, use pre-trained embeddings like Word2Vec, GloVe, or BERT. These embeddings transform words into continuous vectors, capturing semantic meaning, which is better than traditional methods like Bag of Words or TF-IDF.
  - **Consider Using Character-level Features:** For certain applications, such as language identification or sentiment analysis, character-level features (like character n-grams) may be more informative than word-level features.
- 

## 8. Time Series Feature Engineering

- **Create Lag Features:** Lag features are crucial in time-series data. For example, the previous day's sales or the previous month's weather can be important predictors.
  - **Use Rolling Windows:** Features like moving averages, rolling sums, and rolling means can capture trends and smooth out noise in time-series data.
  - **Capture Seasonalities and Cyclic Trends:** For seasonal data, adding features that represent day-of-week, month, or holiday seasonality can help the model detect patterns.
- 

## 9. Dealing with High Dimensionality

- **Use PCA for Dimensionality Reduction:** Principal Component Analysis (PCA) can help reduce the dimensionality of the data by transforming it into a smaller number of features (principal components) that explain most of the variance.
  - **Use Feature Hashing:** For datasets with very high cardinality categorical features, consider feature hashing (or the "hashing trick") to reduce the dimensionality of categorical variables without needing to explicitly store all categories.
- 

## 10. Automating Feature Engineering

- **Leverage AutoML Tools:** Libraries like **TPOT**, **H2O.ai**, or **Auto-Sklearn** offer automated feature engineering pipelines, which can be used for optimizing models and discovering potential features that might have been missed.
- **Use Feature Tools for Automated Feature Generation:** Libraries like **FeatureTools** automate the process of creating features, especially for structured, relational data. They use deep feature synthesis to generate new features based on existing ones.