# Feature Engineering Overview

Feature engineering

Feature engineering is the process of selecting, transforming, and creating input features from raw data to improve the performance of machine learning models. It involves techniques to extract meaningful insights, eliminate irrelevant information, and enhance the predictive power of models. Here's a breakdown of key aspects:

## Key Steps in Feature Engineering

1. **Understanding the Data**:

   - Explore the dataset to identify data types, missing values, and distributions.
   - Analyze correlations and relationships between features and target variables.

2. **Handling Missing Data**:

   - Impute missing values using strategies like mean, median, or mode.
   - Drop rows or columns with excessive missing data if appropriate.

3. **Encoding Categorical Variables**:

   - Use techniques like one-hot encoding, label encoding, or target encoding to handle non-numeric categories.

4. **Scaling and Normalization**:

   - Apply standardization (z-score) or normalization (min-max) to bring features to the same scale.
   - Essential for algorithms sensitive to feature magnitudes, like SVM or K-means.

5. **Feature Transformation**:

   - Apply mathematical transformations (log, square root, or polynomial) to address skewness or add non-linear relationships.
   - Use binning to group continuous data into categories.

6. **Feature Creation**:

   - Create new features by combining existing ones (e.g., ratios, differences, or aggregations).
   - Engineer domain-specific features based on problem understanding.

7. **Dimensionality Reduction**:

   - Use PCA, t-SNE, or feature selection techniques to reduce the number of features while retaining essential information.

8. **Outlier Treatment**:

   - Identify and handle outliers using z-scores, IQR, or domain knowledge.
   - Treating outliers can improve model robustness.

9. **Feature Interaction**:

   - Explore interactions between features using multiplication, division, or conditional logic.

- Useful in datasets where combined effects influence outcomes.

## Tools for Feature Engineering

- **Python Libraries**:
  - **Pandas**: Data manipulation and preprocessing.
  - **Scikit-learn**: Feature scaling, encoding, and selection.
  - **Feature-engine**: Specialized feature engineering tools.
- **Automated Feature Engineering**:
  - Libraries like **Featuretools** or **AutoFeat** generate features automatically.

## Example in Python

```python
import pandas as pd
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.impute import SimpleImputer
from sklearn.decomposition import PCA

# Sample dataset
data = pd.DataFrame({
    'Age': [25, 30, None, 35],
    'Salary': [50000, 60000, 55000, None],
    'Gender': ['Male', 'Female', 'Female', 'Male']
})

# Handle missing values
imputer = SimpleImputer(strategy='mean')
data['Age'] = imputer.fit_transform(data[['Age']])
data['Salary'] = imputer.fit_transform(data[['Salary']])

# Encode categorical variables
encoder = OneHotEncoder(sparse=False)
gender_encoded = encoder.fit_transform(data[['Gender']])
gender_df = pd.DataFrame(gender_encoded, columns=encoder.get_feature_names_out(['Gender']))

# Combine processed data
data = pd.concat([data, gender_df], axis=1).drop(columns=['Gender'])

# Scale numerical features
scaler = StandardScaler()
data[['Age', 'Salary']] = scaler.fit_transform(data[['Age', 'Salary']])

# Apply dimensionality reduction (optional)
pca = PCA(n_components=1)
data['PCA_Feature'] = pca.fit_transform(data[['Age', 'Salary']])

print(data)
```

## Benefits of Feature Engineering

- Improves model accuracy and generalization.
- Reduces overfitting by eliminating redundant features.
- Simplifies models with fewer, more meaningful variables.

Properly engineered features often outperform raw data, making this step crucial for machine learning success.

Let's say i have 10 independent columns and target variable .

It non linear regression problem.

Please advice me  some best exploratory data analysis

For a non-linear regression problem with 10 independent columns and a target variable, a well-structured exploratory data analysis (EDA) is crucial to understand the data's patterns and relationships. Here's a step-by-step guide tailored to your scenario:

## 1. Understand the Data Structure

- **Data Overview**:
  - Check for the number of rows and columns, data types, and missing values.

```python
df.info()
df.describe()
df.isnull().sum()
```

- **Visualize Missing Values**:
  - Use a heatmap or bar plot to identify missing data patterns.

```python
import seaborn as sns
import matplotlib.pyplot as plt
sns.heatmap(df.isnull(), cbar=False, cmap="viridis")
```

## 2. Univariate Analysis

- **For Continuous Variables**:
  - Plot histograms, box plots, and kernel density estimates (KDE) to observe distributions, outliers, and skewness.

```python
for col in df.select_dtypes(include='number').columns:
    sns.histplot(df[col], kde=True)
    plt.title(f"Distribution of {col}")
    plt.show()
```

- **For Categorical Variables**:
  - Use bar plots or pie charts to understand the frequency distribution.

```python
for col in df.select_dtypes(include='object').columns:
    df[col].value_counts().plot(kind='bar')
    plt.title(f"Frequency of {col}")
    plt.show()
```

---

## 3. Bivariate Analysis

- **Target vs. Independent Variables**:
  - Use scatter plots to visualize relationships between numerical features and the target variable.

```python
for col in df.drop(columns='target').columns:
    sns.scatterplot(x=df[col], y=df['target'])
    plt.title(f"{col} vs Target")
    plt.show()
```

- **For Categorical Variables**:
  - Use box plots or violin plots to analyze the target distribution across categories.

```python
for col in df.select_dtypes(include='object').columns:
    sns.boxplot(x=df[col], y=df['target'])
    plt.title(f"{col} vs Target")
    plt.show()
```

- **Pairwise Relationships**:
  - Use pair plots to identify potential interactions and clusters.

```python
sns.pairplot(df, vars=df.drop(columns='target').columns, hue='target')
```

---

## 4. Multivariate Analysis

- **Correlation Matrix**:
  - Identify correlations between features and the target variable. This is useful for detecting multicollinearity and relationships.

```python
corr = df.corr()
sns.heatmap(corr, annot=True, cmap='coolwarm')
```

- **Feature Interaction Analysis**:
  - Use scatter plots or 3D plots to explore feature interactions with the target variable.

```python
from mpl_toolkits.mplot3d import Axes3D
fig = plt.figure()
ax = fig.add_subplot(111, projection='3d')
ax.scatter(df['feature1'], df['feature2'], df['target'])
plt.title("Feature1 vs Feature2 vs Target")
plt.show()
```

- **PCA for Visualization**:
  - Reduce dimensions to visualize feature interactions in 2D or 3D.

```python
from sklearn.decomposition import PCA
pca = PCA(n_components=2)
pca_result = pca.fit_transform(df.drop(columns='target'))
sns.scatterplot(x=pca_result[:, 0], y=pca_result[:, 1], hue=df['target'])
```

---

## 5. Outlier Detection

- **Box Plot for Outliers**:

```python
for col in df.select_dtypes(include='number').columns:
    sns.boxplot(df[col])
    plt.title(f"Boxplot of {col}")
    plt.show()
```

- **Z-Score Method**:
  - Identify features with extreme values.

```python
from scipy.stats import zscore
df['zscore'] = zscore(df['feature'])
outliers = df[(df['zscore'] > 3) | (df['zscore'] < -3)]
```

---

## 6. Feature Importance (Initial Insights)

- **Feature Correlation with Target**:

```python
target_corr = df.corr()['target'].sort_values(ascending=False)
print(target_corr)
```

- **Tree-Based Feature Importance**:
  - Train a simple tree-based model like Random Forest to observe feature importance.

```python
from sklearn.ensemble import RandomForestRegressor
model = RandomForestRegressor()
X = df.drop(columns='target')
y = df['target']
model.fit(X, y)
feature_importances = pd.Series(model.feature_importances_, index=X.columns)
feature_importances.sort_values(ascending=False).plot(kind='bar')
```

## 7. Non-linear Patterns

- **Polynomial Relationships**:
  - Plot quadratic or cubic fits to check for non-linear trends.

```python
for col in df.drop(columns='target').columns:
    sns.regplot(x=df[col], y=df['target'], order=2)
    plt.title(f"Polynomial Fit: {col} vs Target")
    plt.show()
```

- **Interaction Effects**:
  - Use interaction plots or partial dependency plots to explore combined effects.

```python
from sklearn.inspection import PartialDependenceDisplay
from sklearn.ensemble import GradientBoostingRegressor
model = GradientBoostingRegressor()
model.fit(X, y)
PartialDependenceDisplay.from_estimator(model, X, features=[0, 1])  # Feature indices
```

## 8. Handling Multicollinearity

- Use **Variance Inflation Factor (VIF)** to detect and remove multicollinear features.

```python
from statsmodels.stats.outliers_influence import variance_inflation_factor
X = df.drop(columns='target')
vif_data = pd.DataFrame()
vif_data['Feature'] = X.columns
vif_data['VIF'] = [variance_inflation_factor(X.values, i) for i in range(X.shape[1])]
print(vif_data)
```

## 9. Identify Clusters (Optional)

If clustering is suspected, use K-Means or DBSCAN to find patterns.

```python
from sklearn.cluster import KMeans
kmeans = KMeans(n_clusters=3)
clusters = kmeans.fit_predict(X)
sns.scatterplot(x=X['feature1'], y=X['feature2'], hue=clusters)
```

This approach ensures you explore the data comprehensively, laying the foundation for robust feature engineering and model building. Let me know if you'd like any part of this tailored further!