# Data Preprocessing Techniques

**StandardScaler:**

$$\frac{x_i - mean(x)}{stdev(x)}$$

- The StandardScaler ensures that all numerical attributes are scaled to have a mean of 0 and a standard deviation of 1 before they are fed to the machine learning model.

- It is used to normalize the data, ensuring that all features contribute equally to the analysis.

- This is particularly important for algorithms sensitive to the scale of input data, like SVM, k-NN, or logistic regression.

- **Example**: Consider a dataset with two features: age (ranging from 20 to 60) and income (ranging from 30,000 to 100,000). The Standard Scaler will transform these features to have a mean of 0 and a standard deviation of 1, making them comparable despite their different units.

# MinMaxScaler:

$$\frac{x_i - min(x)}{max(x) - min(x)}$$

- used to rescale the features of a dataset so that they lie within a specific range, usually between 0 and 1.

- It transforms each feature individually by scaling the data based on its minimum and maximum values.

- used to ensure that all features contribute equally to a model's performance, particularly in algorithms that rely on the magnitude of data, such as gradient descent optimization in neural networks or k-nearest neighbors (KNN).

- **Example**: Suppose you have a dataset with the feature values [10, 20, 30, 40, 50]. After applying MinMaxScaler (scaled to the range 0 to 1), the transformed feature values would be [0, 0.25, 0.5, 0.75, 1], where 10 maps to 0 and 50 maps to 1.

# RobustScaler:

- scales features using statistics that are robust to outliers, specifically the median and the interquartile range (IQR).

- used to mitigate the impact of outliers by centering the data around the median and scaling it based on the IQR, making it effective for datasets where outliers significantly affect the mean and standard deviation.

- **Example**: Given a dataset with feature values [1,2,2,3,100], applying RobustScaler would transform it by subtracting the median (2) and dividing by the IQR (3-2=1), resulting in [−1,0,0,1,98], which reduces the influence of the outlier (100).

$$\frac{x_i - Q_1(x)}{Q_3(x) - Q_1(x)}$$

# Normalization:

- used to scale the features of a dataset so that they fall within a specific range, typically [0, 1] or [-1, 1]. This ensures that no single feature dominates due to its scale when applying machine learning algorithms.

- used to improve the performance and training stability of ML models, especially those that rely on gradient-based optimization, such as neural networks, as it prevents features with larger ranges from disproportionately influencing the model.

- **Example**: If you have a dataset with features like age (ranging from 20 to 70) and income (ranging from 30,000 to 120,000), without normalization, the model might prioritize income over age due to its larger scale. After normalization, both features will be on a similar scale, ensuring a balanced influence on the model's predictions.

$$\frac{x_i}{\sqrt{x_i^2 + y_i^2 + z_i^2}}$$

# Binarization:

- process of converting numerical or categorical data into a binary format (0s and 1s).

- used to transform variables that have two possible states or to simplify data for certain algorithms.

- used to prepare data for ML models that require or perform better with binary inputs, such as some types of classification algorithms, or to make the data more interpretable by reducing it to a binary form.

- **Example**: Suppose you have a dataset with a feature **"Temperature"** where values above 20°C are considered "High" and values below or equal to 20°C are "Low." Binarization would convert this feature into a binary variable: "1" for "High" and "0" for "Low."

# Encoding Categorical (Ordinal & Nominal) Features:

- involves converting categorical variables into a numerical format.

- **Ordinal encoding** is used for categorical data with a meaningful order (e.g., "Low," "Medium," "High"), while **Nominal encoding** is for categories without a specific order (e.g., "Red," "Blue," "Green").

- essential because most algorithms require numerical input, and categorical data in its original form cannot be directly used by these algorithms.

- **Example**: Consider a dataset with a categorical feature "Education Level" with values "High School," "Bachelor's," and "Master's." For ordinal encoding, these could be converted to 1, 2, and 3, respectively, reflecting their order. For a nominal feature like "Color," one-hot encoding would create separate binary columns for "Red," "Blue," and "Green," with a value of 1 in the column corresponding to the color and 0 in the others.

# Imputation:

- process of replacing missing data with substitute values to ensure that a dataset is complete and can be used for analysis.

- Without imputation, many ML models and statistical analyses cannot handle missing values effectively.

- It is used to maintain data integrity and avoid biases that might occur if missing data points are simply ignored or removed, which could lead to inaccurate conclusions or predictions.

- **For example**, in a dataset of customer ages, if some ages are missing, imputation might involve replacing those missing values with the average age of the remaining customers, so the dataset remains usable for further analysis.

**Polynomial Features:**

- generates new features by taking the original features and raising them to different powers (degrees).

- It expands the feature space to include interactions between variables, allowing for more complex relationships in the data to be captured.

- used to enhance the model's ability to capture non-linear relationships between features and the target variable.

- By including polynomial terms, the model can fit a wider variety of curves, improving its predictive power for non-linear data.

- **Example**: If you have a dataset with a single feature x, generating polynomial features up to the 2nd degree would create new features like $x^2$. For instance, if $x = 2$, the polynomial features might include x and $x^2$ as 2 and 4, respectively. This can help a linear model fit a curve rather than just a straight line.

# Custom Transformer:

- used to apply specific, user-defined transformations to data, allowing for tailored data preparation that isn't covered by standard transformers.

- provides flexibility for custom processing steps.

- often employed when existing transformers in libraries  don't fit the specific needs of your data. For instance, if you need to create new features from existing ones in a particular way or handle data in a unique format, a Custom Transformer can be designed to address these requirements.

- **Example**: Suppose you have a dataset with a column of timestamps, and you want to extract features like the hour of the day and the day of the week from these timestamps. A Custom Transformer can be created to extract these features and add them to your dataset, making it easier for a model to leverage this temporal information.

# Text Processing:

- **Cleaning and Normalizing**: Removes punctuation, converts text to lowercase, and eliminates stop words to clean and normalize the text for consistency and reduces noise in the data.

- **Tokenization**: splits text into smaller units, such as words or phrases (tokens). Helps in analyzing and processing text at a granular level, making it easier to extract meaningful patterns.

- **Stemming and Lemmatization**: These techniques reduce words to their base or root forms (e.g., "running" to "run" or "better" to "good"). This helps in grouping similar words together, improving the accuracy of text analysis.

**Example**: For the sentence "The quick brown foxes are jumping over the lazy dog":
  - **Cleaning**: "the quick brown foxes are jumping over the lazy dog" (converted to lowercase, punctuation removed).
  - **Tokenization**: ["the", "quick", "brown", "foxes", "are", "jumping", "over", "the", "lazy", "dog"].
  - **Stemming/Lemmatization**: ["the", "quick", "brown", "fox", "be", "jump", "over", "the", "lazy", "dog"].

# CountVectorizer:

- **Text to Numerical Conversion**:
  - used to convert text data into a matrix of token counts.
  - creates a vocabulary of all unique words in the dataset and then counts the frequency of each word in each document.

- **Feature Extraction**:
  - helps in transforming textual data into a format that can be used for ML models, by turning text into a numerical representation that captures word occurrences.

- **Example**: Consider two sentences: **"I love programming"** and **"Programming is fun"**. CountVectorizer would create a vocabulary like {'I': 0, 'love': 1, 'programming': 2, 'is': 3, 'fun': 4} and generate a matrix like:
  - Sentence 1: [1, 1, 1, 0, 0]
  - Sentence 2: [0, 0, 1, 1, 1]

  This matrix can then be used as input for machine learning algorithms.

# Term Frequency Inverse Document Frequency (TfIdf):

- a statistical measure used to evaluate the importance of a word in a document relative to a collection of documents (corpus).

- It combines two metrics:
  - **Term Frequency (TF)**: How often a term appears in a document.
  - **Inverse Document Frequency (IDF)**: How common or rare a term is across all documents.

- helps to highlight words that are important to a particular document while reducing the weight of commonly occurring words across many documents.

- helps in distinguishing relevant terms in text analysis and information retrieval.

# Term Frequency Inverse Document Frequency (TfIdf):

**Example**: Suppose you have two documents:
- **Doc 1: "Cats are great pets"**
- **Doc 2: "Dogs are great pets"**

- **TF**: In Doc 1, the term "cats" has a TF of 1 (it appears once out of 4 words). In Doc 2, the term "dogs" has a TF of 1 (also appears once out of 4 words).

- **IDF**: Both "cats" and "dogs" appear in only one of the documents, so their IDF will be higher compared to more common words like "are" or "great".

- **TfIdf Calculation**: "Cats" and "dogs" will have a higher TfIdf score in their respective documents than "are" or "great," highlighting their significance in distinguishing between the documents.

- useful in text processing tasks like search engines and text classification, where identifying the most relevant words is crucial.

# HashingVectorizer:

- **Feature Extraction**:
    - HashingVectorizer is used in text data preprocessing to convert text documents into numerical feature vectors.
    - It does this by applying a hashing function to the text, which maps each word to a fixed-size feature space.

- **Efficiency**:
    - does not require storing a vocabulary of all possible words. This makes it memory-efficient and scalable, especially for large datasets.

- **Hash Collisions**: Since it uses hashing, there can be collisions where different words map to the same feature. However, this trade-off is acceptable in many applications due to the efficiency gains.

# HashingVectorizer:

**Example**: Assume we have the following two text documents:
- **"Machine learning is amazing"**
- **"I love learning and AI"**

We use HashingVectorizer to convert these documents into numerical vectors with a fixed number of features (let's say 5 for simplicity). Here's how it works:

**Tokenization and Hashing**:
**Document 1: Tokens:** ["Machine", "learning", "is", "amazing"]
**Document 2: Tokens:** ["I", "love", "learning", "and", "AI"]

HashingVectorizer applies a hash function to each token to map them into one of the 5 feature indices. This means each token is hashed into one of 5 possible slots.

# HashingVectorizer:

**Example**: **Feature Vector Creation**:

**Hashing Function**: The hash function converts each token into an integer value that determines which feature index (0 through 4) the token will be placed into.

**Example Mapping**: Let's assume the hash function maps tokens to indices as follows: "Machine" -> index 0, "learning" -> index 1, "is" -> index 2, "amazing" -> index 3, "I" -> index 1, "love" -> index 4, "and" -> index 2, "AI" -> index 3

**Document 1 Vector**: "Machine" maps to index 0, "learning" maps to index 1, "is" maps to index 2, "amazing" maps to index 3. The resulting vector might look like [1, 1, 1, 1, 0], where each index shows the count of tokens hashed to that position.

**Document 2 Vector:** "I" maps to index 1, "love" maps to index 4, "learning" maps to index 1, "and" maps to index 2, "AI" maps to index 3
The resulting vector might look like [0, 2, 1, 1, 1].

# HashingVectorizer:

**Example**:

Final Vectors:

    **Document 1 Vector: [1, 1, 1, 1, 0]**
    **Document 2 Vector: [0, 2, 1, 1, 1]**

Here's what each number in the vector represents:

Index 0: Count of "Machine" in Document 1 (1), and 0 in Document 2.

Index 1: Count of "learning" in Document 1 (1) and "I" + "learning" in Document 2 (2).

Index 2: Count of "is" in Document 1 (1) and "and" in Document 2 (1).

Index 3: Count of "amazing" in Document 1 (1) and "AI" in Document 2 (1).

Index 4: Count of "love" in Document 2 (1), and 0 in Document 1.

# Image using skimage:

- In skimage, Image refers to the data structure used to represent and manipulate images.

- It is essentially a multi-dimensional array (e.g., 2D for grayscale or 3D for color images) that stores pixel values.

- **Usage in Data Preprocessing**:
  - **Normalization**: Images can be normalized to a specific range, improving the performance of machine learning models.
  - **Transformation**: Operations such as **resizing, cropping, or rotating** images can be performed to standardize input data.
  - **Filtering**: Techniques like **denoising or edge detection** help in enhancing the image quality for better analysis.

- **Example**: Suppose you have a grayscale image of size 256x256 pixels. You might use skimage to resize it to 128x128 pixels for consistency across your dataset. Additionally, you could apply a Gaussian filter to smooth out noise, preparing the image for further analysis.

# Plotting

various types of plots and their common uses:

- **Bar Plot:**
  - Use: To compare categorical data or show counts of items within different categories.
  - Example: Comparing sales numbers for different products.

- **Histogram:**
  - Use: To display the distribution of a single continuous variable and show the frequency of data within certain intervals (bins).
  - Example: Visualizing the distribution of ages in a population.

- **Line Plot:**
  - Use: To show trends over time or continuous data points.
  - Example: Tracking stock prices over a year.

various types of plots and their common uses:

- **Scatter Plot:**
  - Use: To explore the relationship or correlation between two continuous variables.
  - Example: Analyzing the relationship between study hours and exam scores.

- **Box Plot (Box-and-Whisker Plot):**
  - Use: To summarize the distribution of a dataset, showing median, quartiles, and potential outliers.
  - Example: Comparing the test scores of students from different classes.

- **Pie Chart:**
  - Use: To represent proportions or percentages of a whole in categorical data.
  - Example: Displaying the market share of different companies.

- **Heatmap:**
  - Use: To show the intensity of data values across two variables with color coding.
  - Example: Visualizing the correlation matrix between multiple financial indicators.

various types of plots and their common uses:

- **Violin Plot:**
    - Use: To show the distribution of data across different categories, similar to a box plot but with a density estimation.
    - Example: Comparing the distribution of exam scores between different classes.

- **Pair Plot:**
    - Use: To show pairwise relationships between several variables in a dataset.
    - Example: Exploring relationships between different features in a dataset of house prices.

- **Contour Plot:**
    - Use: To represent three-dimensional data in two dimensions, showing levels of a third variable.
    - Example: Plotting temperature variations across a geographical area.

# KDE Plotting:

- Kernel Density Estimation (KDE): is a non-parametric way to estimate the probability density function of a continuous random variable.

- It smooths out the data to create a continuous curve that represents the distribution.

- used to visualize the distribution of data and identify patterns such as peaks, valleys, and overall shape of the data distribution. Unlike histograms, KDE provides a smooth estimate that can help in understanding the underlying distribution.

- **Example**: Imagine you have a dataset of heights for a group of people. Using KDE, you can create a smooth curve to show the distribution of heights, revealing if the data is normally distributed, bimodal, or skewed, which can provide insights into the characteristics of the population.