

✓ Dimensionality Reduction Techniques 1

Agenda

1. Introduction to Dimensionality Reduction

- Importance of Dimensionality Reduction
- Common Challenges in Dimensionality Reduction
- Types of Dimensionality Reduction

2. Feature Selection

- Filter Methods
- Wrapper Methods
- Embedded Methods

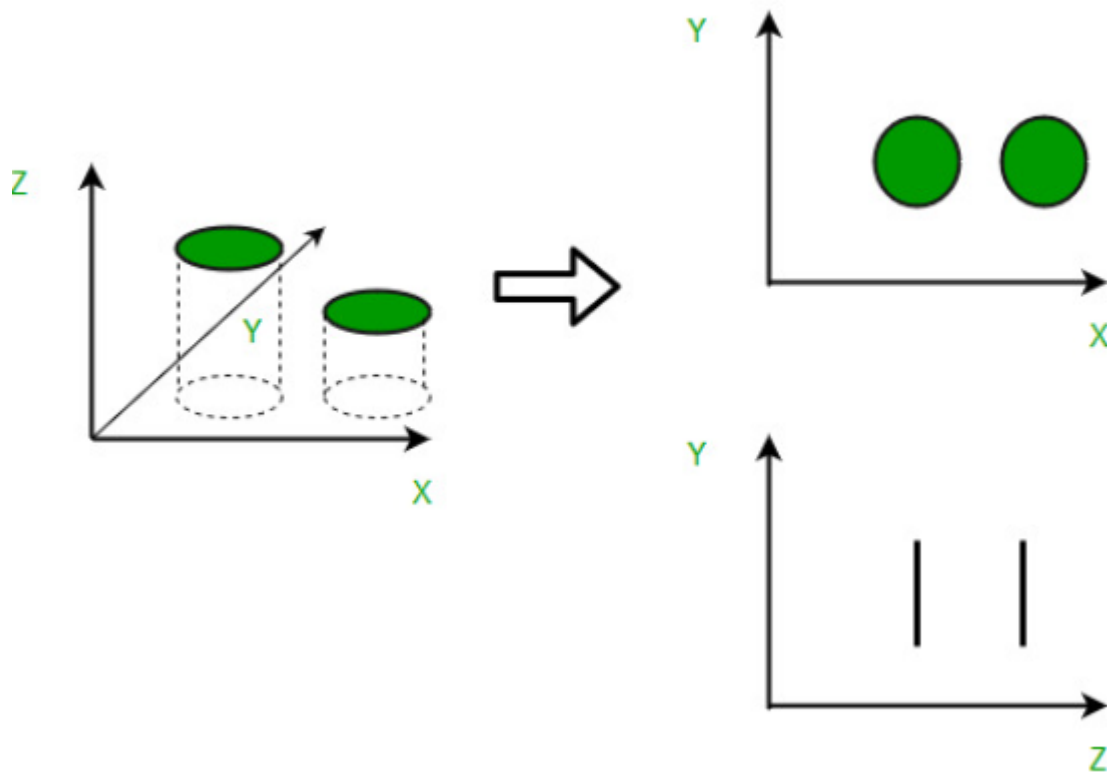
3. Feature Extraction

- Principal Component Analysis(PCA)
- Kernel Principal Component Analysis (Kernel PCA)
- Linear Discriminant Analysis

✓ Introduction to Dimensionality Reduction

- Dimensionality reduction is a critical preprocessing technique in machine learning and data science, where we reduce the number of features (variables) in a dataset while preserving the essential structure, patterns, and variability.
- This technique helps overcome the challenges posed by high-dimensional data, including computational inefficiency, overfitting, and the "curse of dimensionality".
- Dimensionality reduction simplifies the data without losing significant information, improving model performance and interpretability.

Dimensionality Reduction



✓ Importance of Dimensionality Reduction

- **Curse of Dimensionality:**

- As the number of dimensions (features) increases, the volume of the feature space grows exponentially, making data sparse and difficult to analyze.

- **Overfitting:**

- High-dimensional data may lead to overfitting, where models learn the noise in the data instead of general patterns.

- **Visualization:**

- Reducing data to two or three dimensions can make it easier to visualize and understand.

- **Noise Reduction:**

- Many features in high-dimensional data may be irrelevant or redundant.
- Dimensionality reduction can filter out noise and focus on the most important features.

- **Faster Computation:**

- Reducing the number of features reduces the computational cost for training models, improving both speed and efficiency.

- **Improved Performance:**

- By eliminating irrelevant features, models can focus on the most meaningful ones, often improving accuracy and generalization.

✓ Common Challenges in Dimensionality Reduction

- **Loss of Information:**

- When reducing dimensions, there is always a trade-off between retaining variance and losing information.

- **Interpretability:**

- Some techniques, like PCA, may result in components that are linear combinations of the original features, making it difficult to interpret them in a meaningful way.

- **Non-linearity:**

- Linear techniques like PCA may not capture complex, non-linear relationships in the data. In such cases, non-linear techniques like t-SNE or autoencoders are more appropriate.

✓ Types of Dimensionality Reduction Techniques

Dimensionality reduction techniques can be broadly categorized into feature selection and feature extraction.

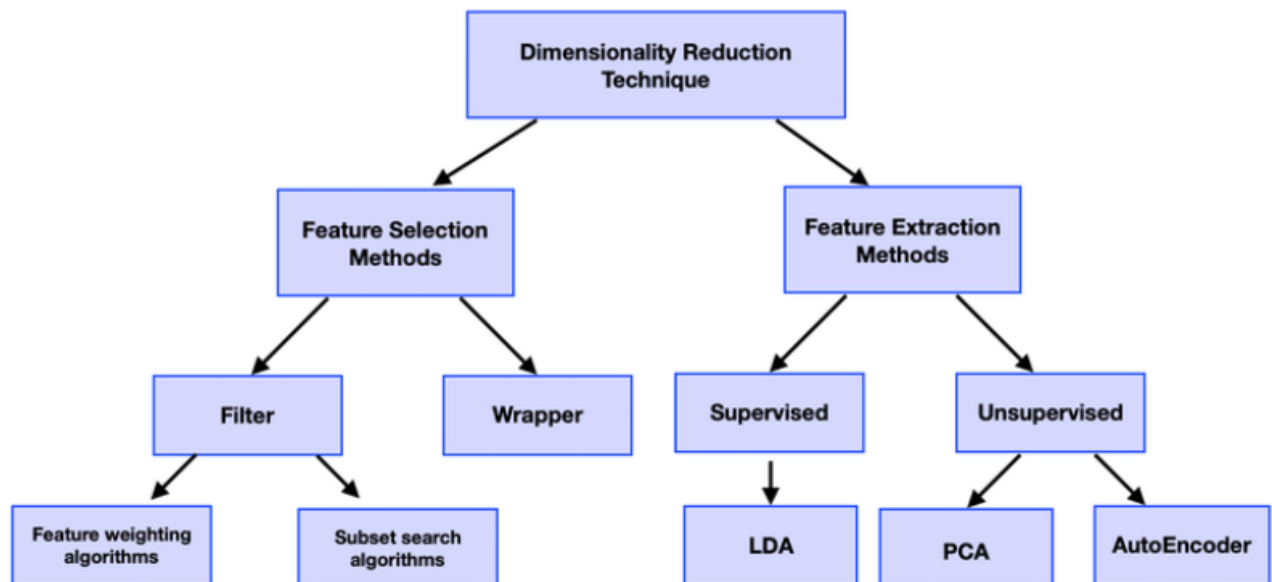
1. Feature Selection

- Feature selection involves selecting a subset of the original features based on certain criteria, without transforming the feature space. The goal is to retain the most important variables and discard irrelevant or redundant ones.
 - **Filter Methods:** Select features based on statistical techniques such as correlation, mutual information, or statistical tests like Chi-square.
 - **Wrapper Methods:** These methods evaluate the performance of subsets of features using a specific model and iteratively add or remove features (e.g., recursive feature elimination).
 - **Embedded Methods:** Feature selection is performed within the training process of the model (e.g., Lasso regression, decision trees).

2. Feature Extraction

- Feature extraction transforms the original high-dimensional features into a new, lower-dimensional set of features. This transformation captures the essence of the original data while reducing redundancy.

- **Principal Component Analysis (PCA):** PCA transforms data by projecting it onto the principal components, which are the directions of maximum variance in the data. It's a linear method.
- **Singular Value Decomposition (SVD):** SVD is a matrix factorization technique that decomposes a matrix into three other matrices, capturing the latent structure of the data.
- **Linear Discriminant Analysis (LDA):** LDA focuses on maximizing the separation between different classes by projecting the data onto a lower-dimensional space that emphasizes class separability.
- **t-SNE (t-distributed Stochastic Neighbor Embedding):** t-SNE is a non-linear technique that is useful for visualizing high-dimensional data by reducing it to two or three dimensions while preserving the local structure.
- **Autoencoders:** These are neural networks that learn a compressed representation (encoding) of the input data, and then reconstruct the data from this reduced form.



✓ Feature Selection

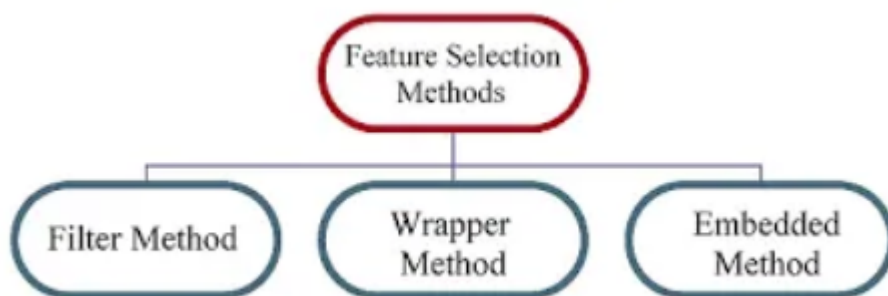
- Feature selection is a critical technique in machine learning where we select a subset of relevant features (or variables) from the dataset to train the model.
- The goal is to retain only the most important features, reducing dimensionality and improving the model's performance by eliminating redundant, irrelevant, or noisy data.
- Unlike dimensionality reduction, which transforms the features into a new feature space, feature selection retains the original features and discards the less useful ones.

Importance of Feature Selection

- **Model Performance:** By eliminating irrelevant or redundant features, feature selection helps improve model accuracy and generalization.
- **Reduction of Overfitting:** Fewer features reduce the risk of overfitting by simplifying the model and preventing it from learning noise in the data.
- **Faster Computation:** Reducing the number of features decreases the computational complexity of the algorithm, making it faster and more efficient.
- **Improved Interpretability:** Models with fewer features are easier to interpret and understand, which is critical in domains like healthcare, finance, and scientific research.
- **Noise Reduction:** By filtering out irrelevant features, feature selection helps reduce noise in the data, improving model robustness.
- **Avoiding the Curse of Dimensionality:** Feature selection helps mitigate the curse of dimensionality, which refers to the problems associated with high-dimensional data.

✓ Types of Feature Selection Methods

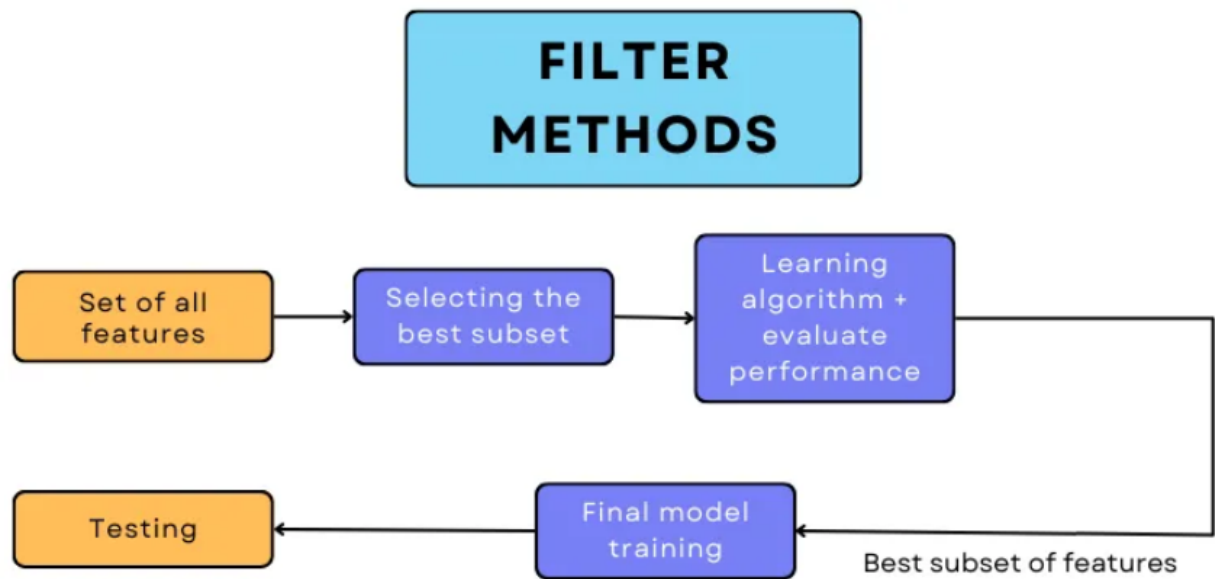
Feature selection techniques can be divided into three main categories: filter methods, wrapper methods, and embedded methods.



✓ 1. Filter Methods

- Filter methods assess the relevance of features by analyzing the intrinsic properties of the data independently of any machine learning model.
- These techniques utilize statistical measures, such as correlation coefficients, chi-squared tests, or mutual information, to rank features based on their ability to predict the target variable.
- By focusing solely on the data characteristics, filter methods provide a straightforward and computationally efficient approach to feature selection. Once the features are ranked, the top-ranked ones are selected for further modeling.
- This approach helps in reducing dimensionality and improving model performance while maintaining interpretability.

- However, filter methods may overlook interactions between features that could be captured by more complex modeling techniques.



Common Filter Methods:

- Correlation Coefficient: This method examines the correlation between each feature and the target variable. Features with a higher correlation to the target are selected.
- Chi-Square Test: This test is used to determine the independence between two categorical variables. It measures the statistical relationship between features and the target variable.
- ANOVA (Analysis of Variance): ANOVA is used to measure the statistical difference between the means of different groups, often used for continuous target variables.
- Mutual Information: Mutual information measures the amount of information that one variable provides about another. It is useful for determining non-linear relationships between features and the target variable.
- Variance Threshold: This method removes features that have low variance across the dataset. If a feature has low variance, it provides little information and can be discarded.

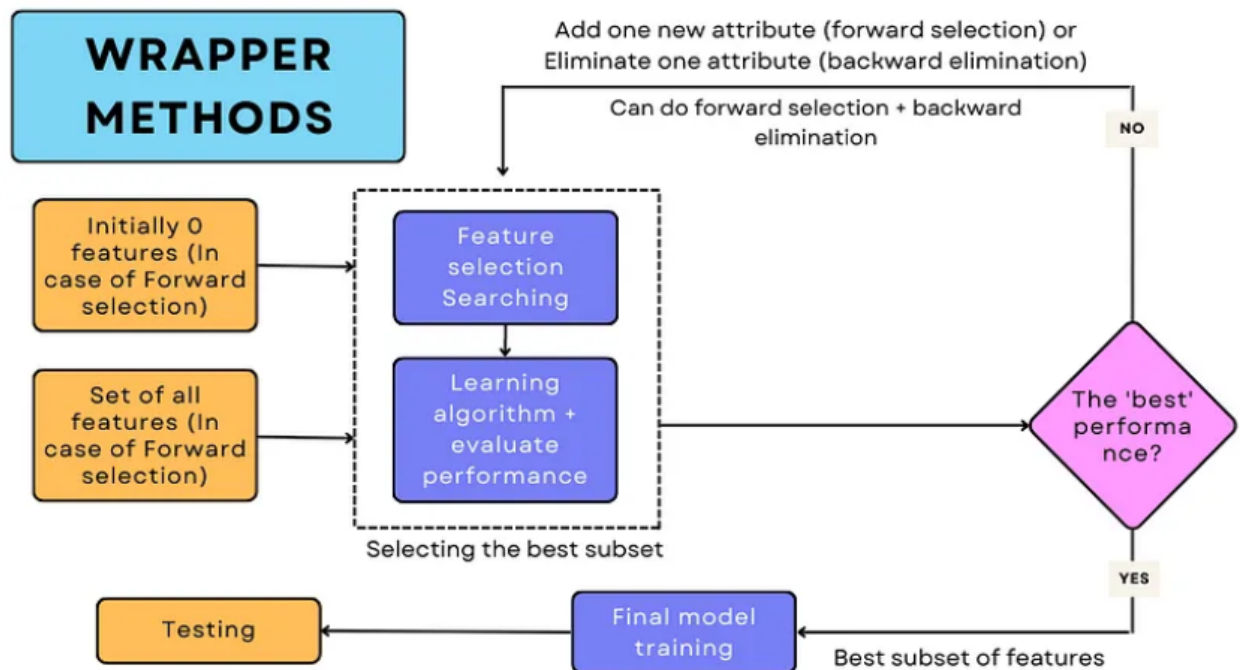
Example of Filter Method:

- Correlation-based Selection: For a regression task, we might calculate the Pearson correlation coefficient between each feature and the target variable. Features with a high correlation (close to +1 or -1) are considered important, while those with low correlation (close to 0) are discarded.

✓ 2. Wrapper Methods

- Wrapper methods function like a collaborative team of analysts working closely with a machine learning model.

- These techniques assess subsets of features based on their performance with the selected model, systematically evaluating different combinations.
- By iteratively adding or removing features, wrapper methods aim to identify the optimal feature set that enhances model performance. This hands-on approach provides a tailored selection process, allowing the model to focus on the most informative features.
- As a result, wrapper methods can lead to improved accuracy and efficiency in machine learning tasks, although they may require more computational resources compared to other feature selection techniques.



Common Wrapper Methods:

- **Forward Selection:** This method starts with an empty set of features and adds one feature at a time based on the improvement in model performance. The process continues until adding more features no longer improves the performance.
- **Backward Elimination:** This method starts with all features and removes one feature at a time based on the decrease in model performance. It continues until removing more features significantly affects performance.
- **Recursive Feature Elimination (RFE):** RFE is a popular wrapper method where the model is trained iteratively, and the least important features (based on the model's weights or importance scores) are removed. This process repeats until the desired number of features is reached.

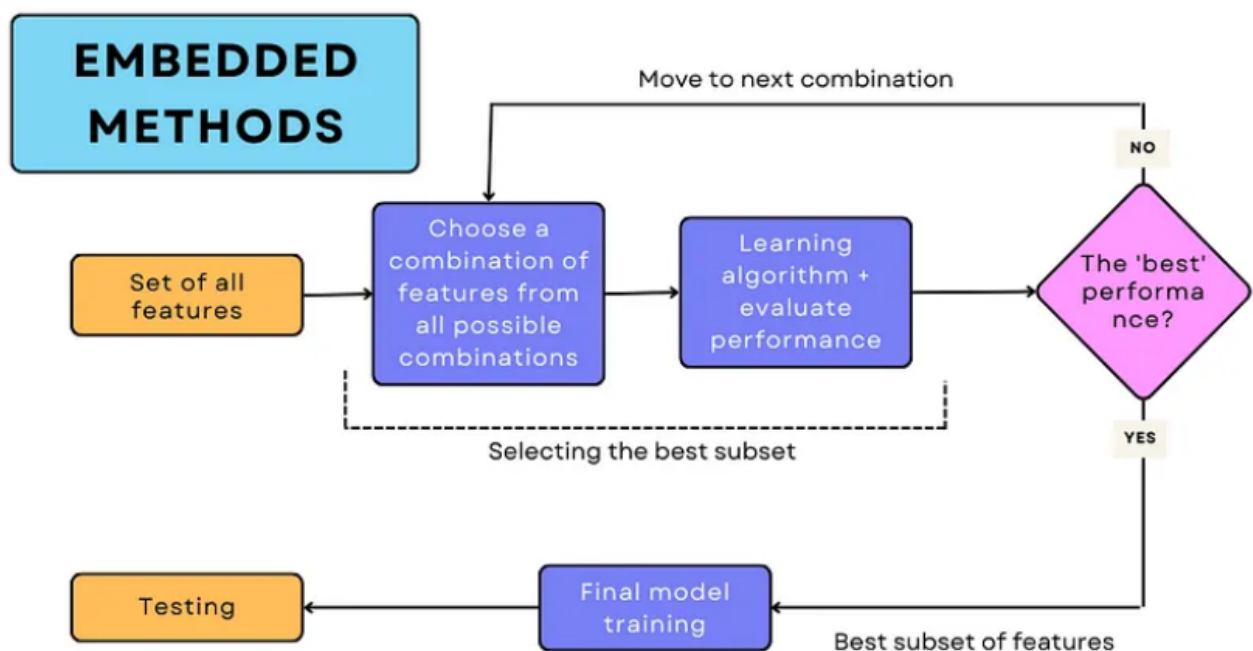
Example of Wrapper Method:

- **Recursive Feature Elimination (RFE):** In RFE, a classifier like Random Forest is trained with all features. The model identifies the least important feature based on its feature

importance score and removes it. This process is repeated until the desired number of features is reached, and the model performance is evaluated at each iteration.

✓ 3. Embedded Methods

- Embedded methods combine feature selection with the model training process, adopting a "why not both?" approach.
- By integrating feature selection directly into the modeling framework, these methods enable the model to simultaneously learn the relationships between features and the target variable while identifying the most relevant features.
- This dual functionality enhances the model's efficiency and performance, as it focuses on the most significant variables without the need for separate preprocessing steps.
- As a result, embedded methods offer a streamlined approach to feature selection, making them particularly useful in complex machine learning tasks where interpretability and accuracy are crucial.



Common Embedded Methods:

- Lasso (L1 Regularization): Lasso adds a penalty to the loss function proportional to the absolute value of the coefficients of the features. It drives the coefficients of irrelevant features to zero, effectively selecting a subset of features.
- Ridge (L2 Regularization): Ridge regularization adds a penalty proportional to the square of the coefficients. While it doesn't lead to exact zero coefficients, it reduces the impact of irrelevant features.
- Elastic Net: Elastic Net is a combination of L1 and L2 regularization, which balances the advantages of both Lasso and Ridge.

- **Tree-based Methods:** Decision trees, Random Forest, and Gradient Boosting automatically rank features based on their importance during the model training process.

Example of Embedded Method:

- **Lasso Regression:** Lasso performs feature selection by adding an L1 penalty to the model's cost function, shrinking the coefficients of irrelevant features to zero. This effectively selects a sparse subset of features that are most important for predicting the target variable.

Feature selection can be easily implemented in Python using libraries like Scikit-learn. Here's an example of using a filter method (Variance Threshold) and an embedded method (Lasso) in Python:

```
import numpy as np
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.linear_model import Lasso
from sklearn.feature_selection import VarianceThreshold, SelectFromModel, RFE
from sklearn.tree import DecisionTreeClassifier
import pandas as pd

# Load Iris dataset
iris = load_iris()
X = iris.data
y = iris.target
feature_names = iris.feature_names

# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_s

# --- Filter Method (Variance Threshold) ---
# Apply variance threshold
filter_selector = VarianceThreshold(threshold=0.1)
X_train_filter = filter_selector.fit_transform(X_train)
X_test_filter = filter_selector.transform(X_test)

# Get selected feature indices
filter_indices = filter_selector.get_support(indices=True)
selected_filter_features = [feature_names[i] for i in filter_indices]

print("Selected Features (Filter Method):", selected_filter_features)

# --- Embedded Method (Lasso Regression) ---
# Fit Lasso regression model
lasso = Lasso(alpha=0.1)
lasso.fit(X_train_filter, y_train)

# Select important features based on Lasso coefficients
embedded_selector = SelectFromModel(lasso, prefit=True)
X_train_embedded = embedded_selector.transform(X_train_filter)

# Get selected feature indices
embedded_indices = embedded_selector.get_support(indices=True)
selected_embedded_features = [selected_filter_features[i] for i in embedded_indic

print("Selected Features (Embedded Method):", selected_embedded_features)

# --- Wrapper Method (Recursive Feature Elimination) ---
# Create a Decision Tree Classifier
model = DecisionTreeClassifier()

# Recursive Feature Elimination
rfe = RFE(estimator=model, n_features_to_select=2) # Select the top 2 features
rfe.fit(X_train, y_train)

# Get the selected features
X_train_wrapper = rfe.transform(X_train)
```

```
# Get selected feature indices
wrapper_indices = rfe.get_support(indices=True)
selected_wrapper_features = [feature_names[i] for i in wrapper_indices]

print("Selected Features (Wrapper Method):", selected_wrapper_features)
```

➡ Selected Features (Filter Method): ['sepal length (cm)', 'sepal width (cm)',
Selected Features (Embedded Method): ['petal length (cm)']
Selected Features (Wrapper Method): ['petal length (cm)', 'petal width (cm)']

✓ Feature Extraction

- Feature extraction is a process in machine learning and data preprocessing where raw data is transformed into a set of features that can be used by an algorithm to improve model performance.
- Unlike feature selection, which aims to select a subset of the original features, feature extraction involves creating new features from the existing data by applying certain transformations.
- The main objective of feature extraction is to capture the most informative aspects of the data while reducing its dimensionality.

Importance of Feature Extraction

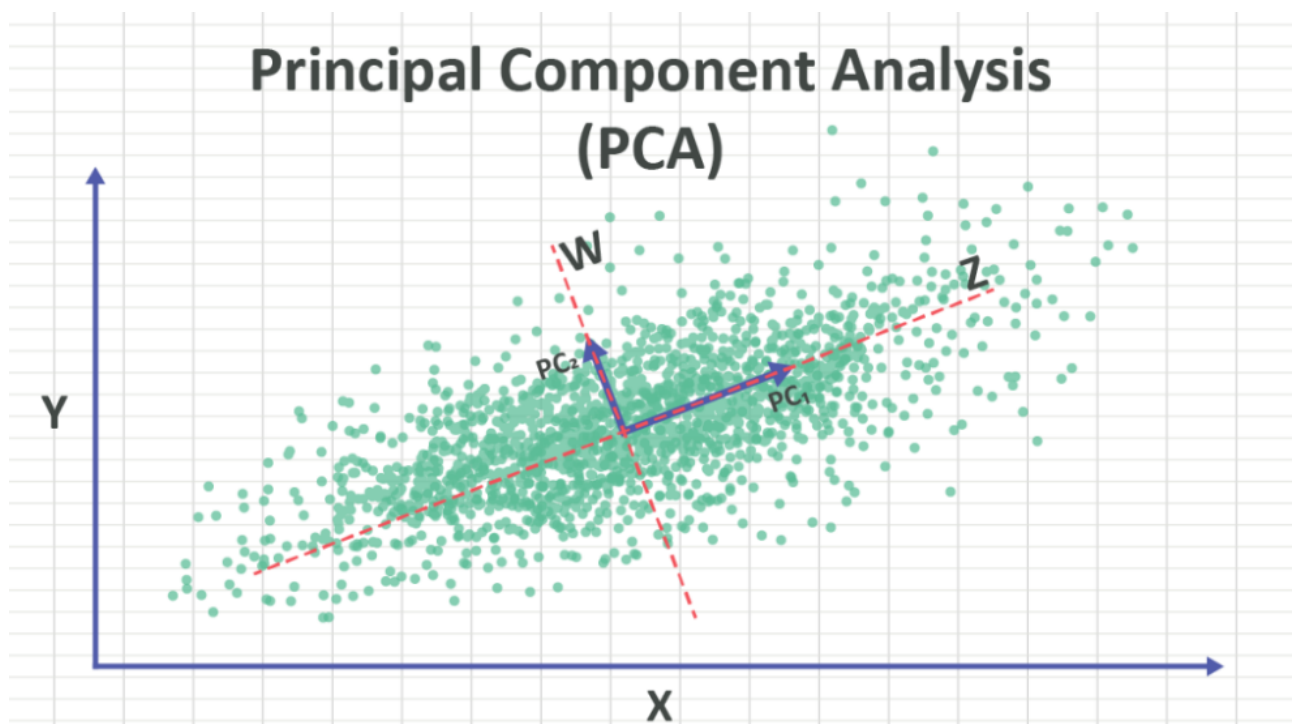
- Dimensionality Reduction: By transforming raw data into a lower-dimensional feature set, feature extraction helps in reducing the dimensionality of the dataset, making computations faster and models simpler.
- Improving Model Performance: Extracted features often provide a more informative representation of the data, which can lead to better model accuracy and generalization.
- Data Representation: High-dimensional data can be difficult for algorithms to handle. Feature extraction helps in creating compact representations of the data while retaining essential information.
- Handling Complex Data: For unstructured or semi-structured data (e.g., text, images), feature extraction is essential to convert data into a form suitable for machine learning algorithms.

✓ Key Feature Extraction Techniques

Feature extraction methods can be broadly categorized into manual methods (designed by human experts) and automated methods (learned by algorithms such as neural networks). Let's explore some common feature extraction techniques.

✓ Principal Component Analysis (PCA)

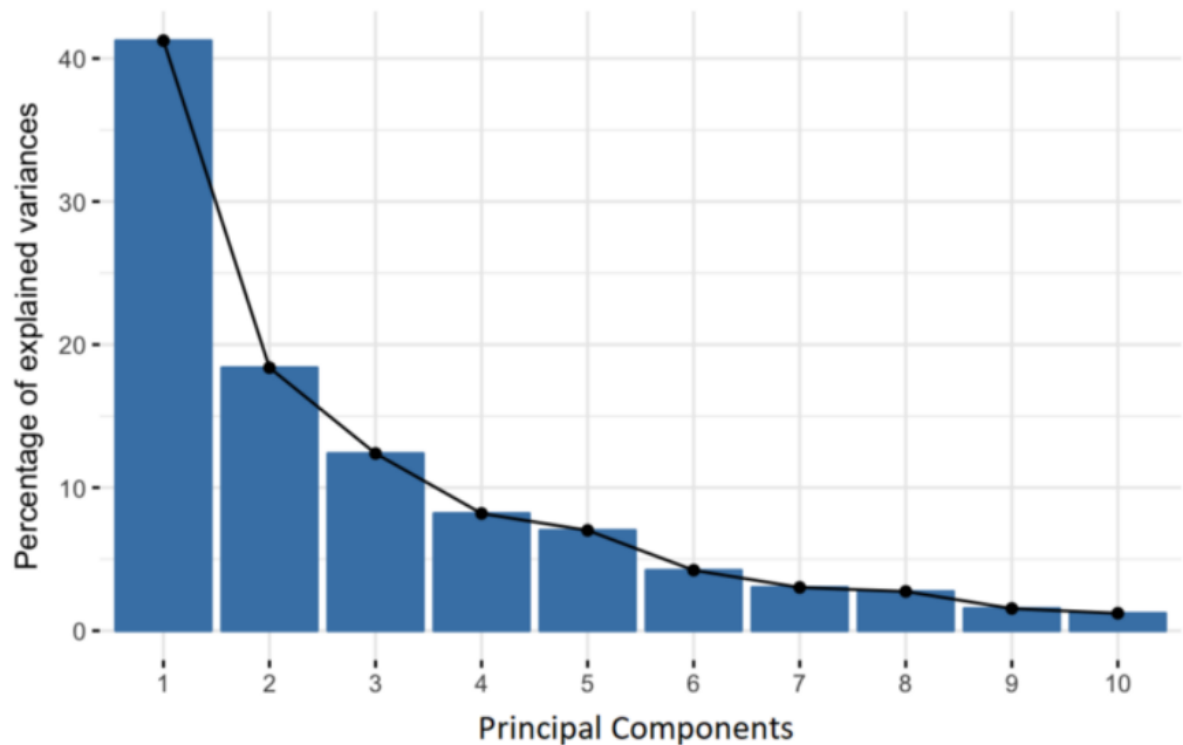
- Principal Component Analysis (PCA) is a widely utilized dimensionality reduction technique that converts high-dimensional data into a lower-dimensional representation while preserving as much variance as possible.
- By identifying the directions of maximum variance in the data, PCA generates new variables known as principal components, which are linear combinations of the original features. This transformation allows for the simplification of complex datasets, making them easier to analyze and visualize.
- PCA is particularly beneficial in various applications, including machine learning, statistics, and data visualization, where it aids in enhancing model performance and reducing computational costs.
- By eliminating less informative features, PCA helps to mitigate issues such as overfitting and improves the interpretability of the results.
- Furthermore, it is often used as a preprocessing step before applying other machine learning algorithms, enabling better insights into the underlying structure of the data.



✓ What Are Principal Components?

- Principal components are new variables created by combining the original variables in a way that ensures these new components are uncorrelated.
- The goal of this transformation is to capture as much information as possible from the original variables into the first few components.

- For example, when dealing with 10-dimensional data, PCA generates 10 principal components.
- However, PCA prioritizes placing the maximum amount of information into the first principal component, then distributing the next largest portion of information into the second component, and so forth.
- This process continues until all components are derived, as illustrated in the accompanying scree plot.



✓ Steps in PCA

1. Standardization

- PCA is sensitive to the relative scales of features, which can skew the results if some features have larger ranges than others.
- To address this, it is crucial to standardize the data so that each feature has a mean of 0 and a standard deviation of 1.
- This normalization prevents features with larger scales from disproportionately influencing the principal components, allowing for a more accurate representation of the underlying data structure. - Standardization ensures that all features contribute equally to the PCA analysis.
- Formula for standardization:

$$z = (x - \mu) / \sigma$$

where x is the data point, μ is the mean, and σ is the standard deviation.

2. Covariance Matrix Computation

- The covariance matrix is a crucial component in understanding the relationships between features in a dataset.
- It quantifies how much each pair of features varies together, revealing the degree of correlation between them.
- A positive covariance indicates that as one feature increases, the other tends to increase as well, while a negative covariance suggests an inverse relationship.
- By analyzing the covariance matrix, one can identify patterns and dependencies among features, which is essential for techniques like Principal Component Analysis (PCA).
- This matrix serves as the foundation for determining the principal components, helping to highlight the directions of maximum variance in the data.
- The covariance matrix is defined as:

$$\text{Cov}(X) = \frac{1}{n-1} X^T X$$

where X is the data matrix, and X^T is the transpose of X .

3. Eigenvalue and Eigenvector Computation

- Eigenvalues and eigenvectors are derived from the covariance matrix and play a key role in dimensionality reduction techniques like Principal Component Analysis (PCA).
- The eigenvectors indicate the directions (principal components) in which the data exhibits the highest variability, effectively capturing the most informative features of the dataset.
- Correspondingly, the eigenvalues measure the magnitude of variance explained by each eigenvector.
- A higher eigenvalue signifies that its associated eigenvector accounts for more variance in the data, guiding the selection of principal components.
- This relationship helps prioritize the components that contribute most significantly to the data's structure and variability.
- Eigenvalue equation:

$$\text{Cov}(X) v = \lambda v$$

where v is the eigenvector and λ is the eigenvalue.

4. Selecting Principal Components

- In Principal Component Analysis (PCA), the eigenvectors, which represent the principal components, are sorted in descending order based on their corresponding eigenvalues.
- This sorting helps identify the most significant directions of variance within the data.
- The top k eigenvectors are then selected, where k denotes the desired number of dimensions to which the data will be reduced.
- These selected principal components capture the maximum variance, ensuring that the most informative features are retained while reducing the overall dimensionality of the dataset.
- This process facilitates a more efficient representation of the data, enhancing interpretability and model performance.

5. Projection onto Principal Components

- Once the top k principal components have been selected, the original data is projected onto these components to create a reduced-dimensional dataset.
- This projection involves calculating the dot product between the original data and the selected eigenvectors, transforming the data into the new feature space defined by the principal components.
- The resulting lower-dimensional representation retains the most significant variance from the original dataset while discarding less informative features.
- This reduced-dimensional dataset is easier to analyze, visualize, and utilize in machine learning models.
- Projection formula:

$$Z = XW$$

where Z is the projected data, X is the standardized data, and W is the matrix of selected eigenvectors.

✓ Applications of PCA

- Data Visualization: PCA is used to reduce high-dimensional data to 2D or 3D for visualization purposes.
- Noise Filtering: By removing components that explain very little variance (noise), PCA can be used to denoise data.
- Feature Extraction: PCA can identify the most important features of a dataset, which can be useful for feature engineering and improving model performance.
- Compression: PCA can reduce the size of datasets for storage and computational efficiency.

Example:

```
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler
import numpy as np

# Sample data
data = np.array([[2.5, 2.4], [0.5, 0.7], [2.2, 2.9], [1.9, 2.2], [3.1, 3.0]])

# Standardizing the data
scaler = StandardScaler()
data_scaled = scaler.fit_transform(data)

# Applying PCA
pca = PCA(n_components=2)
pca_transformed = pca.fit_transform(data_scaled)

# Print the transformed data
print("PCA Transformed Data:\n", pca_transformed)
print("Explained Variance Ratio:", pca.explained_variance_ratio_)

⇒ PCA Transformed Data:
[[ 0.5124457  0.23853092]
 [-2.57528445  0.06114533]
 [ 0.69555387 -0.43434461]
 [-0.1485184  -0.0800397 ]
 [ 1.51580328  0.21470806]]
Explained Variance Ratio: [0.96982031 0.03017969]
```

The explained variance ratio is high (close to 1), PCA has successfully captured most of the important variance in the data.

✓ Kernel Principal Component Analysis (Kernel PCA)

- Kernel PCA is an extension of the classical Principal Component Analysis (PCA) that uses the kernel trick to handle non-linear relationships in data.
- In contrast to standard PCA, which only captures linear correlations between features, Kernel PCA enables PCA to work on data that is not linearly separable by implicitly mapping the data into a higher-dimensional space where linear techniques (like PCA) can be applied.
- It is a form of dimensionality reduction that allows capturing complex, non-linear patterns in the data.

Why Kernel PCA?

- Classical PCA works by identifying the directions (principal components) that maximize the variance in the data, but it assumes that the data lies in a linear subspace.

- This assumption is restrictive for many real-world datasets where the underlying relationships between features are non-linear.
- For example:
 - Data such as concentric circles or moon-shaped clusters cannot be separated by a linear transformation, and PCA fails in these scenarios.
 - Kernel PCA solves this issue by using a kernel function to implicitly map the data to a higher-dimensional space where the data becomes linearly separable.

✓ How Kernel PCA Works

- Kernel PCA uses the kernel trick to transform the original feature space into a higher-dimensional space, where standard PCA can be applied.
- The key idea is that instead of directly computing the principal components from the data in the original feature space, Kernel PCA computes them from the kernel matrix, which represents pairwise similarities between data points in the higher-dimensional space.

Steps Involved:

- **Mapping the data to a higher-dimensional space:**
 - Using a kernel function ϕ , the original data X is mapped to a higher-dimensional feature space $\Phi(X)$.
 - The mapping ϕ is implicit, meaning that you don't compute the transformation explicitly, but you work with the kernel matrix that contains dot products in the higher-dimensional space.
- **Kernel Function:**
 - A kernel function $K(x_i, x_j) = \phi(x_i)^\top \phi(x_j)$ is used to compute the similarity between pairs of data points in the high-dimensional space.
 - Common kernel functions include:
 - Polynomial Kernel: Captures polynomial relationships.
 - Gaussian (RBF) Kernel: Captures non-linear, radial relationships.
 - Sigmoid Kernel: Models data with a structure similar to neural networks.
- **Construct the Kernel Matrix:**
 - Compute the kernel matrix K , which contains pairwise similarities between all data points: $K_{ij} = \phi(x_i)^\top \phi(x_j)$
 - This kernel matrix is symmetric and captures the relationships in the higher-dimensional space without explicitly computing the transformation.
- **Center the Kernel Matrix:**

- Center the kernel matrix K to ensure that the data has zero mean in the feature space. This step is crucial for PCA to work correctly.

- **Eigen Decomposition:**

- Perform eigen decomposition on the centered kernel matrix:

$$Kv = \lambda v$$

Here, v are the eigenvectors, and λ are the corresponding eigenvalues.

- The eigenvectors corresponding to the largest eigenvalues represent the principal components in the higher-dimensional space.

- **Project the Data:**

- The original data points are projected onto the principal components derived from the kernel matrix to obtain the new lower-dimensional representation of the data.

✓ Mathematical Formulation of Kernel PCA

1. Kernel Function:

- The kernel function $K(x_i, x_j)$ computes the similarity between two data points in the higher-dimensional feature space without explicitly computing the coordinates.

Examples:

- Polynomial Kernel:

$$K(x_i, x_j) = (x_i^T x_j + c)^d$$

where c is a constant, and d is the degree of the polynomial.

- RBF (Gaussian) Kernel:

$$K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$$

where σ is a parameter controlling the width of the Gaussian function.

2. Compute Kernel Matrix:

- Given a dataset $X = [x_1, x_2, \dots, x_n]$, the kernel matrix K is computed as:

$$K_{ij} = \phi(x_i)^T \phi(x_j)$$

- The matrix K contains pairwise kernel evaluations between all data points.

3. Center the Kernel Matrix:

- To center the kernel matrix K , use the following formula:

$$K' = K - \frac{1}{n}K - \frac{1}{n}K + \frac{1}{n}K\frac{1}{n}$$

- where $\frac{1}{n}$ is a matrix of ones, scaled appropriately.

4. Eigen Decomposition:

- Perform eigen decomposition on the centered kernel matrix

$$K'v = \lambda v$$

- Here, v are the eigenvectors, and λ are the eigenvalues. The eigenvectors corresponding to the largest eigenvalues define the principal components.

5. Dimensionality Reduction:

- To reduce the dimensionality of the data, project the original data onto the first k principal components:

$$Z = Kv$$

- where Z is the new, lower-dimensional representation of the data.

Example

```

import numpy as np
from sklearn.datasets import make_moons
from sklearn.decomposition import KernelPCA
import matplotlib.pyplot as plt

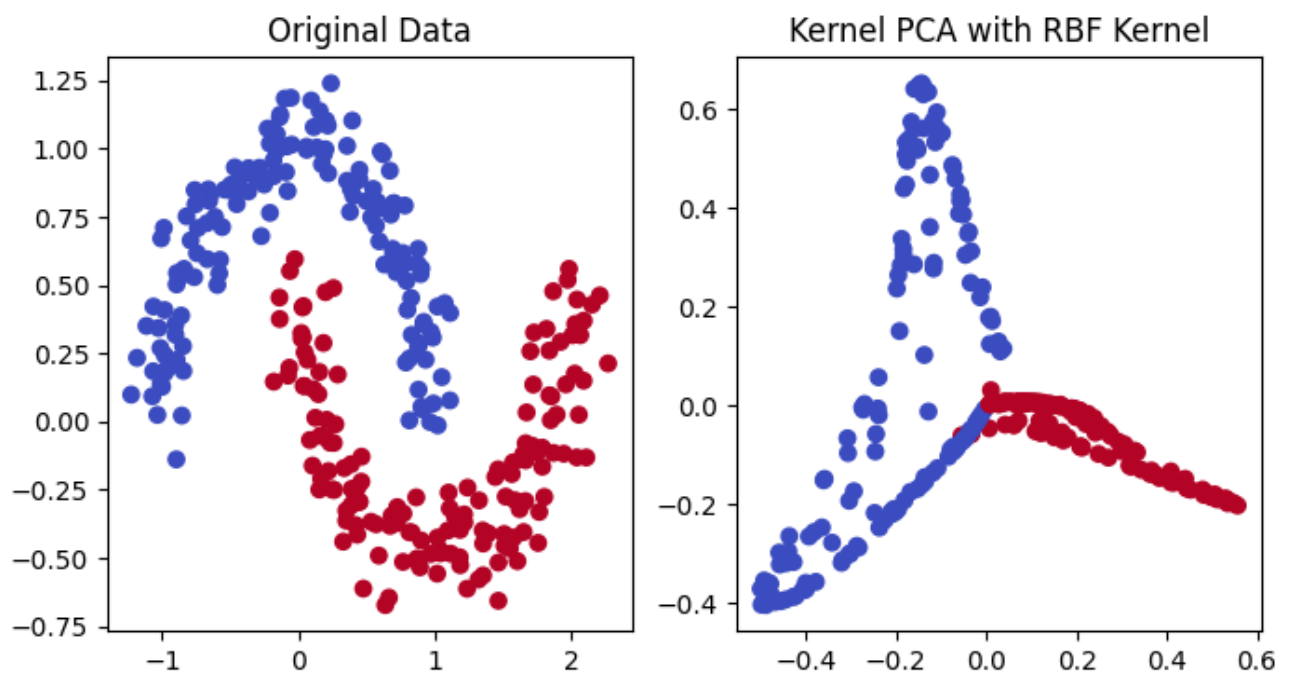
# Create a non-linear dataset (moon-shaped data)
X, y = make_moons(n_samples=300, noise=0.1, random_state=42)

# Apply Kernel PCA with RBF kernel
kpca = KernelPCA(n_components=2, kernel='rbf', gamma=15)
X_kpca = kpca.fit_transform(X)

# Plotting the original data
plt.figure(figsize=(8,4))
plt.subplot(1, 2, 1)
plt.scatter(X[:, 0], X[:, 1], c=y, cmap='coolwarm')
plt.title('Original Data')

# Plotting the transformed data
plt.subplot(1, 2, 2)
plt.scatter(X_kpca[:, 0], X_kpca[:, 1], c=y, cmap='coolwarm')
plt.title('Kernel PCA with RBF Kernel')
plt.show()

```



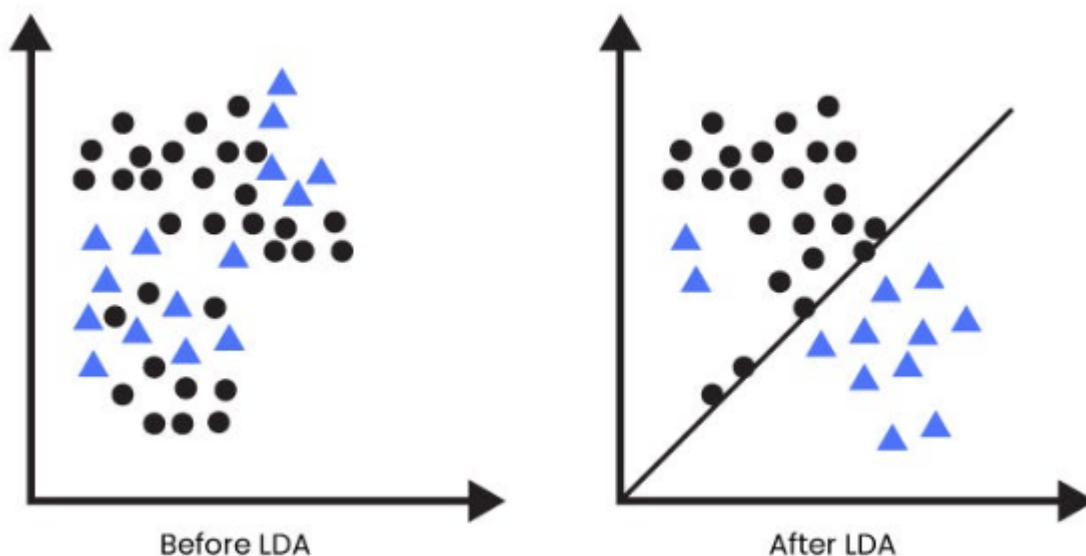
- `make_moons` generates a non-linear dataset.
- `KernelPCA` applies the RBF (Gaussian) kernel to perform non-linear dimensionality reduction.
- The output visualizes the transformed data in the reduced feature space, where the moon-shaped clusters become linearly separable.

✓ Linear Discriminant Analysis (LDA)

- Linear Discriminant Analysis (LDA) is a powerful technique used for both classification and dimensionality reduction, focusing on finding linear combinations of features that optimally separate different classes in the data.
- Unlike PCA, which emphasizes variance, LDA aims to maximize the distance between the means of different classes while minimizing the variance within each class.
- This makes LDA particularly effective for supervised classification tasks, as it enhances the discriminative power of the features.
- By projecting the data onto a lower-dimensional space defined by the discriminant axes, LDA not only reduces dimensionality but also improves the performance of subsequent models.
- It is widely applied in various domains, including image recognition, bioinformatics, and marketing analytics, making it a valuable tool in the machine learning toolkit.

Concept of LDA

- LDA aims to project the data into a lower-dimensional space while maximizing the class separability. It seeks to find a linear transformation of the feature space such that the ratio of between-class variance to within-class variance is maximized, ensuring that classes are well separated.
- In simpler terms, LDA tries to maximize the distance between the means of different classes (inter-class variance) while minimizing the spread (variance) within each class (intra-class variance).



✓ Steps in Linear Discriminant Analysis

1. Compute the mean vectors for each class

- The mean of each class in the dataset is calculated. If there are k classes, there will be k mean vectors, one for each class.

2. Compute the scatter matrices

- Within-class scatter matrix S_W Measures how much each class's samples deviate from their respective mean.
- Between-class scatter matrix S_B Measures how much the class means deviate from the overall mean of the dataset.

3. Solve the eigenvalue problem

- Find the eigenvalues and eigenvectors for the matrix $S_W^{-1} S_B$. The eigenvectors corresponding to the largest eigenvalues provide the directions along which the classes are best separated.

4. Select the top eigenvectors

- Select the top $k-1$ eigenvectors (where k is the number of classes). This step reduces the dimensionality of the data and helps with class separation