

SUPPORT VECTOR MACHINE (SVM)

A Support Vector Machine (SVM) is a powerful machine learning algorithm. It helps us solve problems like classifying data into groups, predicting outcomes, and even spotting unusual data points. SVMs can handle both simple (linear) and complex (nonlinear) data patterns, making them useful in many areas. It helps us solve problems like classifying data into groups, predicting outcomes, and even spotting unusual data points. SVMs can handle both simple (linear) and complex (nonlinear) data patterns, making them useful in many areas. The key strength of SVMs lies in their ability to identify the optimal decision boundary, known as the hyperplane, that separates different classes in the data. This ensures the separation is as clear as possible by maximising the margin between the hyperplane and the closest data points from each class, called support vectors.

By transforming data into higher dimensions using a technique called the kernel trick, SVMs can also classify complex, nonlinear data patterns. This makes them suitable for both binary (two-class) and multi-class classification problems, as well as regression tasks.

The dimension of the hyperplane depends on the number of features. For instance, if there are two input features, the hyperplane is simply a line, and if there are three input features, the hyperplane becomes a 2-D plane. As the number of features increases beyond three, the complexity of visualising the hyperplane also increases. The hyperplane that maximizes this margin is called the maximum-margin hyperplane or hard margin. In real-world scenarios, data is not always perfectly separable. SVM introduces soft margins, which allow some data points to fall on the wrong side of the margin. For such cases, SVM adds a penalty for each misclassified or margin-violating point to balance between maximizing the margin and minimizing errors.

The model minimizes a combination of margin size and the penalty, commonly using hinge loss.

Hyperplane: The hyperplane is the decision boundary used to separate data points of different classes in a feature space. For linear classification, this is a linear equation represented as $w \cdot x + b = 0$.

Support Vectors: Support vectors are the closest data points to the hyperplane. These points are critical in determining the hyperplane and the margin in Support Vector Machine (SVM).

Margin: The margin refers to the distance between the support vector and the hyperplane. The primary goal of the SVM algorithm is to maximize this margin, as a wider margin typically results in better classification performance.

Kernel: The kernel is a mathematical function used in SVM to map input data into a higher-dimensional feature space. This allows the SVM to find a hyperplane in cases where data points are not linearly separable in the original space.

Hard Margin: A hard margin refers to the maximum-margin hyperplane that perfectly separates the data points of different classes without any misclassifications.

Soft Margin: When data contains outliers or is not perfectly separable, SVM introduces a slack variable for each data point to allow some misclassifications while balancing between maximizing the margin and minimizing violations.

PCA (Principle Component Analysis)

Principal Component Analysis (PCA) is a statistical technique used for dimensionality reduction. It simplifies complex datasets by transforming them into a smaller number of uncorrelated variables called principal components, while retaining as much of the original dataset's variability as possible.

Principal Component Analysis (PCA) is a statistical procedure that uses an orthogonal transformation that converts a set of correlated variables to a set of uncorrelated variables. PCA is the most widely used tool in exploratory data analysis and in machine learning for predictive models. Moreover, Principal Component Analysis (PCA) is an unsupervised learning algorithm technique used to examine the interrelations among a set of variables. It is also known as a general factor analysis where regression determines a line of best fit. The main goal of Principal Component Analysis (PCA) is to reduce the dimensionality of a dataset while preserving the most important patterns or relationships between the variables without any prior knowledge of the target variables.

1. **Dimensionality Reduction:** PCA reduces the dataset's features while keeping the essential information, making it easier to analyze high-dimensional data.
2. **Principal Components:** These are new variables formed as linear combinations of the original ones. They are uncorrelated, with the first component capturing the most variance, followed by the second, and so on.
3. **Variance:** PCA prioritizes features with higher variance, as they carry more significant information.
4. **Eigenvalues and Eigenvectors:** PCA uses these mathematical tools to identify principal components. Eigenvalues show how much variance a component captures, while eigenvectors indicate the direction.
5. **Covariance Matrix:** This matrix shows how features vary together. PCA uses it to identify directions of maximum variance in the data.

Steps in PCA

1. **Standardization:** Scale the data so all features have a mean of zero and standard deviation of one.
2. **Compute Covariance Matrix:** Understand relationships between features.
3. **Calculate Eigenvalues and Eigenvectors:** Determine the importance and direction of each component.
4. **Select Principal Components:** Choose the components that capture the most variance.
5. **Transform Data:** Project the dataset onto these components for dimensionality reduction.

Clustering (K-Means & DB Scan)

Clustering is a type of unsupervised learning used to group data points based on their similarities

K-Means Clustering

K-Means organises data into clusters by assigning each point to the nearest center, called a centroid. The algorithm begins by initializing centroids and then iteratively adjusts them to minimize the variation within each cluster. The number of clusters must be specified beforehand, making this method straightforward but dependent on prior knowledge of the data.

K-Means works well when the clusters are evenly distributed and shaped like spheres. It's computationally efficient, making it suitable for larger datasets. However, it is sensitive to outliers, as these points can significantly affect the centroid's position. Additionally, it may struggle with data that doesn't conform to regular shapes or densities.

DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

DBSCAN identifies clusters based on the density of data points. It groups points that are densely packed together, while points in low-density regions are classified as noise. DBSCAN requires two parameters: ϵ (neighbourhood radius) and minimum points required to form a dense region.

DBSCAN is particularly effective when clusters have irregular shapes or when noise and outliers are present. However, it can be sensitive to the choice of parameters and may struggle with clusters of varying densities.

Key Differences

1. Approach:

- K-Means minimizes distance to centroids, making it more geometric.
- DBSCAN focuses on density, which makes it effective for irregular shapes.

2. Outlier Treatment:

- K-Means integrates outliers into clusters.
- DBSCAN identifies and labels outliers as noise.
-

Both K-Means and DBSCAN are versatile tools for exploring data. In marketing, they can segment customers based on purchasing behavior. In biology, they help group genes or species based on shared traits. In geospatial analysis, they can identify regions of interest or patterns in spatial data. While their approaches differ, both algorithms aim to make sense of complex datasets by organizing them into meaningful groups.