

Linear Regression

Definition: Linear Regression is a statistical method used to model the relationship between a dependent variable and one or more independent variables using a linear equation.

When to Use:

- When there is a linear relationship between the independent and dependent variables.
- When interpretability and simplicity are important.
- For problems with continuous output.

When to Avoid:

- When the relationship between variables is non-linear.
- When the dataset contains a lot of outliers.
- When there is high multicollinearity among independent variables.

Advantages:

- Simple and easy to implement.
- Computationally efficient.
- Coefficients are easily interpretable.
- Works well with linearly separable data.

Disadvantages:

- Assumes a linear relationship between features and target.
- Sensitive to outliers.
- May not perform well with complex datasets.
- Assumes independence of errors.

Typical Use Cases:

- Predicting housing prices based on features like size, location, and number of rooms.
- Estimating sales and revenue based on marketing spend.
- Analysing the relationship between variables in scientific studies.

Scenarios Where It May Not Perform Well:

- Non-linear relationships, such as predicting stock prices based on historical data.
- Datasets with many outliers, like detecting fraudulent transactions.
- When there is multicollinearity, such as predicting a target variable with highly correlated features.

Ridge Regression

Definition: Ridge Regression is a type of linear regression that includes an L2 regularization term to shrink the coefficients and reduce model complexity.

When to Use:

- When there is multicollinearity.
- To prevent overfitting in a linear model.

When to Avoid:

- When feature selection is crucial.
- When there is no multicollinearity.

Advantages:

- Reduces model complexity and prevents overfitting.
- Handles multicollinearity well.

Disadvantages:

- Requires tuning of the regularization parameter (λ).
- Less interpretable due to shrinkage of coefficients.

Typical Use Cases:

- Financial modelling.
- Predictive modelling in datasets with many features.

Scenarios Where It May Not Perform Well:

- Sparse datasets where feature selection is needed.

Lasso Regression

Definition: Lasso Regression is a type of linear regression that includes an L1 regularization term to enforce sparsity in the coefficients.

When to Use:

- When feature selection is important.
- When there is high dimensionality.

When to Avoid:

- When all features are important.
- When there is no multicollinearity.

Advantages:

- Can reduce the number of features.
- Helps in feature selection and interpretation.

Disadvantages:

- May underperform if no true sparsity.
- Requires tuning of the regularization parameter (λ).

Typical Use Cases:

- High-dimensional datasets like genetics.
- Situations requiring feature selection.

Scenarios Where It May Not Perform Well:

- Non-sparse datasets where all features are important.

Elastic Net Regression

Definition: Elastic Net Regression combines L1 and L2 regularization to improve both Ridge and Lasso Regression.

When to Use:

- When there is multicollinearity and feature selection is important.
- When both Ridge and Lasso Regression are needed.

When to Avoid:

- When the dataset is not high-dimensional.
- When there is no multicollinearity.

Advantages:

- Combines benefits of Ridge and Lasso.
- Flexible regularization.

Disadvantages:

- Requires tuning of two regularization parameters (λ_1 and λ_2).

Typical Use Cases:

- Genomics and other high-dimensional datasets.
- Predictive modelling with feature selection.

Scenarios Where It May Not Perform Well:

- Small or low-dimensional datasets.

Polynomial Regression

Definition: Polynomial Regression extends linear regression by considering polynomial relationships between the independent and dependent variables.

When to Use:

- When the relationship between variables is non-linear.

When to Avoid:

- When there is a risk of overfitting.
- When the dataset is too small.

Advantages:

- Can model non-linear relationships.

Disadvantages:

- High risk of overfitting with high-degree polynomials.
- Requires careful tuning.

Typical Use Cases:

- Curve fitting in engineering.
- Predicting growth trends.

Scenarios Where It May Not Perform Well:

- High-dimensional datasets with many features.

Support Vector Regression (SVR)

Definition: SVR uses Support Vector Machines to perform regression by finding a hyperplane that best fits the data points within a specified margin of tolerance.

When to Use:

- When the data is high-dimensional.
- When robust and flexible models are needed.

When to Avoid:

- When the dataset is very large.
- When interpretability is crucial.

Advantages:

- Effective in high-dimensional spaces.
- Robust to overfitting.

Disadvantages:

- Computationally intensive.
- Requires careful tuning of parameters.

Typical Use Cases:

- Financial time series forecasting.
- Predictive modelling with non-linear relationships.

Scenarios Where It May Not Perform Well:

- Very large datasets due to computational complexity.

Decision Tree Regression

Definition: Decision Tree Regression uses a tree-like model of decisions to predict the value of a target variable based on several input variables.

Formula: Decision trees partition the data into subsets based on the value of input features, but there is no single formula.

When to Use:

- When interpretability is important.
- For datasets with non-linear relationships.

When to Avoid:

- When the data is noisy.
- When a very precise prediction is required.

Advantages:

- Easy to understand and interpret.
- Handles both numerical and categorical data.

Disadvantages:

- Prone to overfitting.
- Unstable with small variations in data.

Typical Use Cases:

- Customer segmentation.
- Predicting outcomes in healthcare.

Scenarios Where It May Not Perform Well:

- Datasets with a lot of noise or small variations.

Random Forest Regression

Definition: Random Forest Regression is an ensemble method that uses multiple decision trees to improve the predictive performance and control overfitting.

Formula: Combines the predictions of several decision trees, but there is no single formula.

When to Use:

- When high accuracy is required.
- When the data has a lot of features.

When to Avoid:

- When interpretability is crucial.
- When computational resources are limited.

Advantages:

- Reduces overfitting compared to a single decision tree.
- Handles large datasets well.

Disadvantages:

- Less interpretable than a single decision tree.
- Requires more computational resources.

Typical Use Cases:

- Credit scoring.
- Predictive modeling in marketing.

Scenarios Where It May Not Perform Well:

- Situations requiring model interpretability.

Gradient Boosting Regression

Definition: Gradient Boosting Regression builds an ensemble of trees sequentially, where each tree corrects the errors of the previous one.

Formula: Sequentially adds models to minimize the loss function, but there is no single formula.

When to Use:

- When high predictive performance is needed.
- For structured/tabular data.

When to Avoid:

- When the dataset is too large.
- When computational speed is a concern.

Advantages:

- High predictive accuracy.
- Can handle complex datasets well.

Disadvantages:

- Computationally expensive.
- Prone to overfitting if not properly tuned.

Typical Use Cases:

- Predictive maintenance.
- Fraud detection.

Scenarios Where It May Not Perform Well:

- Very large datasets due to computational time.

AdaBoost Regression

Definition: AdaBoost Regression is an ensemble method that combines multiple weak learners to form a strong learner by focusing on the errors of previous learners.

Formula: Weights are adjusted to emphasize misclassified instances, but there is no single formula.

When to Use:

- When boosting is needed for better accuracy.
- When dealing with weak learners.

When to Avoid:

- When the dataset is very noisy.
- When computational resources are limited.

Advantages:

- Improves weak learners.
- Can handle different types of weak learners.

Disadvantages:

- Sensitive to noisy data.
- Requires careful tuning.

Typical Use Cases:

- Spam detection.
- Customer churn prediction.

Scenarios Where It May Not Perform Well:

- Noisy datasets with a lot of errors.

XGBoost Regression

Definition: XGBoost Regression is an efficient and scalable implementation of gradient boosting that improves speed and performance.

Formula: Uses gradient boosting framework, but there is no single formula.

When to Use:

- When high accuracy and speed are required.
- For large-scale datasets.

When to Avoid:

- When interpretability is crucial.
- When the dataset is very small.

Advantages:

- High performance and efficiency.
- Handles large datasets well.

Disadvantages:

- Less interpretable.
- Can be prone to overfitting if not properly tuned.

Typical Use Cases:

- Competition datasets.
- Large-scale machine learning problems.

Scenarios Where It May Not Perform Well:

- Small datasets or those requiring high interpretability.

LightGBM Regression

Definition: LightGBM Regression is a gradient boosting framework that uses tree-based learning algorithms and is designed for high performance and efficiency.

Formula: Uses histogram-based algorithms, but there is no single formula.

When to Use:

- When high performance and speed are needed.
- For large-scale datasets.

When to Avoid:

- When interpretability is crucial.
- When the dataset is very small.

Advantages:

- Very fast and efficient.
- Handles large datasets well.

Disadvantages:

- Less interpretable.
- Requires careful parameter tuning.

Typical Use Cases:

- Real-time prediction systems.
- Large-scale machine learning tasks.

Scenarios Where It May Not Perform Well:

- Small datasets or those requiring high interpretability.

CatBoost Regression

Definition: CatBoost Regression is a gradient boosting framework that handles categorical features automatically.

Formula: Uses gradient boosting framework, but there is no single formula.

When to Use:

- When handling categorical data.
- For high performance and efficiency.

When to Avoid:

- When interpretability is crucial.
- When the dataset is very small.

Advantages:

- Handles categorical features well.
- High performance and efficiency.

Disadvantages:

- Less interpretable.
- Requires parameter tuning.

Typical Use Cases:

- E-commerce recommendations.
- Predictive modelling with categorical data.

Scenarios Where It May Not Perform Well:

- Small datasets or those requiring high interpretability.

Logistic Regression

Definition: Logistic Regression is a statistical model used for binary classification that predicts the probability of a binary outcome based on one or more predictor variables.

When to Use:

- When the target variable is binary.
- For linearly separable data.

When to Avoid:

- When there is a non-linear relationship between variables.
- For multiclass classification without modification.

Advantages:

- Simple and easy to implement.
- Outputs probabilities.
- Interpretable coefficients.

Disadvantages:

- Assumes a linear relationship between features and log odds.
- Can be less effective with complex data.

Typical Use Cases:

- Spam detection.
- Credit scoring.
- Medical diagnosis.

Scenarios Where It May Not Perform Well:

- Non-linear relationships.
- High-dimensional data without regularization.

K-Nearest Neighbors (KNN)

Definition: K-Nearest Neighbors is a non-parametric algorithm that classifies a data point based on the majority class of its k nearest neighbors.

Formula: No explicit formula; based on distance metrics like Euclidean distance.

When to Use:

- When the data is small and well-labeled.
- For non-linear decision boundaries.

When to Avoid:

- With large datasets (computationally expensive).
- When feature scaling is not done.

Advantages:

- Simple and intuitive.
- No training phase.

Disadvantages:

- Computationally expensive during prediction.
- Sensitive to the choice of k and distance metric.

Typical Use Cases:

- Image recognition.
- Recommender systems.
- Anomaly detection.

Scenarios Where It May Not Perform Well:

- Large datasets.
- High-dimensional data.

Support Vector Machine (SVM)

Definition: SVM is a supervised learning algorithm that finds the hyperplane which best separates different classes in the feature space.

When to Use:

- For high-dimensional data.
- When a clear margin of separation is required.

When to Avoid:

- With very large datasets.
- When interpretability is crucial.

Advantages:

- Effective in high-dimensional spaces.
- Robust to overfitting with proper kernel.

Disadvantages:

- Computationally intensive.
- Requires tuning of parameters (C, kernel).

Typical Use Cases:

- Text categorization.
- Image classification.
- Bioinformatics.

Scenarios Where It May Not Perform Well:

- Large datasets.
- Noisy data.

Decision Tree Classification

Definition: Decision Tree Classification uses a tree-like model of decisions to classify data points.

Formula: No explicit formula; based on splitting criteria like Gini impurity or entropy.

When to Use:

- For interpretable models.
- When feature interactions are important.

When to Avoid:

- With noisy data.
- When a precise prediction is needed.

Advantages:

- Easy to understand and interpret.
- Handles both numerical and categorical data.

Disadvantages:

- Prone to overfitting.
- Unstable with small variations in data.

Typical Use Cases:

- Customer segmentation.
- Fraud detection.
- Medical diagnosis.

Scenarios Where It May Not Perform Well:

- Noisy data.
- Small data variations.

Random Forest Classification

Definition: Random Forest Classification is an ensemble method that uses multiple decision trees to improve predictive performance and control overfitting.

Formula: Combines predictions of several decision trees.

When to Use:

- When high accuracy is required.
- When data has a lot of features.

When to Avoid:

- When interpretability is crucial.
- When computational resources are limited.

Advantages:

- Reduces overfitting.
- Handles large datasets well.

Disadvantages:

- Less interpretable.
- Requires more computational resources.

Typical Use Cases:

- Credit scoring.
- Predictive modeling in marketing.
- Healthcare predictions.

Scenarios Where It May Not Perform Well:

- Situations requiring model interpretability.

Gradient Boosting Classification

Definition: Gradient Boosting Classification builds an ensemble of trees sequentially, where each tree corrects the errors of the previous one.

Formula: Sequentially adds models to minimize the loss function.

When to Use:

- When high predictive performance is needed.
- For structured/tabular data.

When to Avoid:

- When the dataset is too large.
- When computational speed is a concern.

Advantages:

- High predictive accuracy.
- Handles complex datasets well.

Disadvantages:

- Computationally expensive.
- Prone to overfitting if not properly tuned.

Typical Use Cases:

- Predictive maintenance.
- Fraud detection.
- Customer churn prediction.

Scenarios Where It May Not Perform Well:

- Very large datasets.

AdaBoost Classification

Definition: AdaBoost Classification is an ensemble method that combines multiple weak learners to form a strong learner by focusing on errors of previous learners.

Formula: Weights are adjusted to emphasize misclassified instances.

When to Use:

- When boosting is needed for better accuracy.
- When dealing with weak learners.

When to Avoid:

- When the dataset is very noisy.
- When computational resources are limited.

Advantages:

- Improves weak learners.
- Can handle different types of weak learners.

Disadvantages:

- Sensitive to noisy data.
- Requires careful tuning.

Typical Use Cases:

- Spam detection.
- Customer churn prediction.
- Medical diagnostics.

Scenarios Where It May Not Perform Well:

- Noisy datasets.

XGBoost Classification

Definition: XGBoost Classification is an efficient and scalable implementation of gradient boosting that improves speed and performance.

Formula: Uses gradient boosting framework.

When to Use:

- When high accuracy and speed are required.
- For large-scale datasets.

When to Avoid:

- When interpretability is crucial.
- When the dataset is very small.

Advantages:

- High performance and efficiency.
- Handles large datasets well.

Disadvantages:

- Less interpretable.
- Can be prone to overfitting if not properly tuned.

Typical Use Cases:

- Competition datasets.
- Large-scale machine learning problems.
- Sales forecasting.

Scenarios Where It May Not Perform Well:

- Small datasets.

LightGBM Classification

Definition: LightGBM Classification is a gradient-boosting framework that uses tree-based learning algorithms and is designed for high performance and efficiency.

Formula: Uses histogram-based algorithms.

When to Use:

- When high performance and speed are needed.
- For large-scale datasets.

When to Avoid:

- When interpretability is crucial.
- When the dataset is very small.

Advantages:

- Very fast and efficient.
- Handles large datasets well.

Disadvantages:

- Less interpretable.
- Requires careful parameter tuning.

Typical Use Cases:

- Real-time prediction systems.
- Large-scale machine learning tasks.
- Web traffic prediction.

Scenarios Where It May Not Perform Well:

- Small datasets.

CatBoost Classification

Definition: CatBoost Classification is a gradient-boosting framework that handles categorical features automatically.

Formula: Uses gradient boosting framework.

When to Use:

- When handling categorical data.
- For high performance and efficiency.

When to Avoid:

- When interpretability is crucial.
- When the dataset is very small.

Advantages:

- Handles categorical features well.
- High performance and efficiency.

Disadvantages:

- Less interpretable.
- Requires parameter tuning.

Typical Use Cases:

- E-commerce recommendations.
- Predictive modeling with categorical data.
- Customer behavior analysis.

Scenarios Where It May Not Perform Well:

- Small datasets.

Naive Bayes (Gaussian, Multinomial, Bernoulli)

Definition: Naive Bayes is a probabilistic classifier based on Bayes' theorem, assuming independence between features.

When to Use:

- When the independence assumption holds.
- For text classification.

When to Avoid:

- When features are highly correlated.

Advantages:

- Simple and fast.
- Works well with high-dimensional data.

Disadvantages:

- Assumes independence of features.
- Less accurate when the independence assumption is violated.

Typical Use Cases:

- Spam detection.
- Text classification.
- Sentiment analysis.

Scenarios Where It May Not Perform Well:

- Datasets with highly correlated features.

Quadratic Discriminant Analysis (QDA)

Definition: QDA is a classification algorithm that models the class-conditional distribution of the data using a quadratic decision boundary.

Formula: Based on the quadratic form of the Mahalanobis distance.

When to Use:

- When classes have distinct covariance structures.
- For datasets with non-linear class boundaries.

When to Avoid:

- When the dataset is small.
- When the covariance matrices are similar.

Advantages:

- Handles non-linear decision boundaries.
- More flexible than LDA.

Disadvantages:

- Requires large datasets.
- Sensitive to outliers.

Typical Use Cases:

- Medical diagnostics.
- Image recognition.

Scenarios Where It May Not Perform Well:

- Small datasets.
- Datasets with similar covariance structures.

Linear Discriminant Analysis (LDA)

Definition: Linear Discriminant Analysis (LDA) is a classification and dimensionality reduction technique that finds a linear combination of features that best separates two or more classes.

When to Use:

- When classes are linearly separable.
- For dimensionality reduction while preserving class separability.

When to Avoid:

- When classes have different covariance structures.
- For non-linear class boundaries.

Advantages:

- Simple and interpretable.
- Effective for dimensionality reduction and visualization.

Disadvantages:

- Assumes that features are normally distributed within each class.
- Can be less effective if classes are not linearly separable.

Typical Use Cases:

- Face recognition.
- Pattern recognition.
- Medical diagnostics where class separability is important.

Scenarios Where It May Not Perform Well:

- When classes have significantly different covariance matrices.
- In cases with non-linear relationships between features.

K-Means Clustering

Definition: K-means clustering is an iterative algorithm that partitions data into k clusters, where each data point belongs to the cluster with the nearest centroid.

When to Use:

- When the number of clusters is known in advance.
- For well-separated, spherical clusters.

When to Avoid:

- When clusters are not spherical or have varying sizes and densities.
- When the number of clusters is not known.

Advantages:

- Simple and efficient.
- Works well for large datasets.

Disadvantages:

- Sensitive to initial centroid placement.
- Assumes clusters are spherical and equally sized.

Typical Use Cases:

- Market segmentation.
- Document clustering.
- Image compression.

Scenarios Where It May Not Perform Well:

- Non-spherical clusters.
- Clusters with varying densities.

K-Medoids Clustering

Definition: K-Medoids Clustering is similar to K-Means but uses actual data points as cluster centers (medoids) instead of centroids.

Formula/Key Concepts:

- **Objective Function:** Minimize the sum of dissimilarities between points and their respective medoids.

When to Use:

- When the data has outliers or is noisy.
- When you need cluster centers to be actual data points.

When to Avoid:

- For very large datasets (computationally expensive).

Advantages:

- More robust to outliers than K-Means.
- Centers are actual data points.

Disadvantages:

- Computationally intensive for large datasets.
- Not suitable for high-dimensional data.

Typical Use Cases:

- Clustering with noisy data.
- When you need interpretable cluster centers.

Scenarios Where It May Not Perform Well:

- Large datasets or high-dimensional data.

Hierarchical Clustering (Agglomerative, Divisive)

Definition: Hierarchical Clustering builds a hierarchy of clusters either by iteratively merging smaller clusters (agglomerative) or by splitting a large cluster (divisive).

Formula/Key Concepts:

- **Agglomerative:** Start with individual points and merge clusters based on distance metrics.
- **Divisive:** Start with all points in one cluster and recursively split.

When to Use:

- When you want to understand the data hierarchy.
- For small to medium-sized datasets.

When to Avoid:

- With very large datasets due to high computational complexity.
- When you need a predefined number of clusters.

Advantages:

- Does not require specifying the number of clusters.
- Produces a dendrogram (tree-like diagram) showing cluster relationships.

Disadvantages:

- Computationally expensive.
- Difficult to handle large datasets.

Typical Use Cases:

- Taxonomy and classification problems.
- Data exploration.

Scenarios Where It May Not Perform Well:

- Large datasets.
- When exact cluster numbers are needed.

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

Definition: DBSCAN is a density-based clustering algorithm that groups together points that are close to each other based on a distance threshold and a minimum number of points.

Formula/Key Concepts:

- **Core Point:** A point with at least a minimum number of neighbors within a given radius.
- **Reachable Points:** Points within the neighborhood of a core point.
- **Noise:** Points that do not belong to any cluster.

When to Use:

- When clusters are of arbitrary shape.
- When dealing with noise and outliers.

When to Avoid:

- When clusters have varying densities.
- When the distance metric is not appropriate.

Advantages:

- Can find arbitrarily shaped clusters.
- Robust to noise and outliers.

Disadvantages:

- Performance depends on the choice of parameters (epsilon and minPts).
- Not suitable for datasets with varying density.

Typical Use Cases:

- Spatial data analysis.
- Anomaly detection.

Scenarios Where It May Not Perform Well:

- Clusters with varying density.
- Large datasets with high dimensionality.

OPTICS (Ordering Points To Identify the Clustering Structure)

Definition: OPTICS is a density-based clustering algorithm that extends DBSCAN by creating an ordering of data points to identify clusters of varying density.

Formula/Key Concepts:

- **Reachability Plot:** A plot that helps visualize the clustering structure based on reachability distance.

When to Use:

- When dealing with clusters of varying density.
- For datasets where the density of clusters varies significantly.

When to Avoid:

- For very large datasets without sufficient computational resources.

Advantages:

- Can handle clusters with varying densities.
- Provides a reachability plot for detailed clustering insights.

Disadvantages:

- More complex and computationally intensive than DBSCAN.
- Requires parameter tuning.

Typical Use Cases:

- Complex spatial clustering.
- Multi-scale clustering problems.

Scenarios Where It May Not Perform Well:

- Very large datasets.
- Situations where parameter tuning is not feasible.

Mean Shift Clustering

Definition: Mean Shift Clustering is a non-parametric clustering algorithm that finds clusters by iteratively shifting data points towards the mode (highest density) of the data.

Formula/Key Concepts:

- **Mean Shift Vector:** Moves each data point towards the mean of points within a given radius.

When to Use:

- When clusters are of arbitrary shape.
- For data with complex, unknown cluster structures.

When to Avoid:

- When the choice of bandwidth (radius) is challenging.
- For very large datasets.

Advantages:

- Does not require the number of clusters to be specified.
- Can find arbitrarily shaped clusters.

Disadvantages:

- Computationally expensive for large datasets.
- Performance heavily depends on bandwidth selection.

Typical Use Cases:

- Image segmentation.
- Data with unknown cluster shapes.

Scenarios Where It May Not Perform Well:

- Large datasets with high dimensionality.
- When choosing the bandwidth is difficult.

Gaussian Mixture Models (GMM)

Definition: Gaussian Mixture Models (GMM) is a probabilistic model that assumes data is generated from a mixture of several Gaussian distributions with unknown parameters.

Formula/Key Concepts:

- **Likelihood Function:** Estimates the probability of data points belonging to each Gaussian component.
- **Expectation-Maximization (EM) Algorithm:** Used to estimate the parameters.

When to Use:

- When you assume the data comes from a mixture of Gaussian distributions.
- For soft clustering where data points can belong to multiple clusters.

When to Avoid:

- When the data does not fit Gaussian assumptions.
- For very high-dimensional data.

Advantages:

- Provides probabilistic cluster assignments.
- Can model complex data distributions.

Disadvantages:

- Assumes data is Gaussian, which may not always be true.
- Can be computationally intensive.

Typical Use Cases:

- Image segmentation.
- Anomaly detection in finance.

Scenarios Where It May Not Perform Well:

- Non-Gaussian data.
- High-dimensional spaces.

BIRCH (Balanced Iterative Reducing and Clustering using Hierarchies)

Definition: BIRCH is a clustering algorithm that incrementally and dynamically clusters incoming data using a tree structure.

Formula/Key Concepts:

- **Clustering Feature Tree (CF Tree):** A tree structure used to summarize and cluster data efficiently.

When to Use:

- For large datasets that do not fit in memory.
- When you need to perform clustering incrementally.

When to Avoid:

- For very small datasets.
- When the data has highly irregular clusters.

Advantages:

- Efficient for large datasets.
- Incremental and dynamic clustering.

Disadvantages:

- Less effective for data with highly irregular clusters.
- Requires parameter tuning for optimal performance.

Typical Use Cases:

- Large-scale data clustering.
- Streaming data analysis.

Scenarios Where It May Not Perform Well:

- Small datasets.
- Data with highly irregular clusters.

Principal Component Analysis (PCA)

Definition: Principal Component Analysis (PCA) is a linear dimensionality reduction technique that transforms data into a set of orthogonal components, capturing the maximum variance in the data.

Formula/Key Concepts:

- **Objective:** Find the eigenvectors and eigenvalues of the covariance matrix to identify principal components.
- **Transformation:** Data is projected onto the principal components.

When to Use:

- When you want to reduce dimensionality while retaining as much variance as possible.
- For linear relationships between features.

When to Avoid:

- When data has non-linear relationships.
- When the interpretability of principal components is not required.

Advantages:

- Reduces dimensionality while preserving variance.
- Helps in noise reduction and visualization.

Disadvantages:

- Assumes linear relationships.
- Principal components may not be easily interpretable.

Typical Use Cases:

- Data preprocessing and visualization.
- Feature reduction for machine learning.

Scenarios Where It May Not Perform Well:

- Non-linear data relationships.
- When features are not linearly correlated.

Linear Discriminant Analysis (LDA)

Definition: Linear Discriminant Analysis (LDA) is a technique for both classification and dimensionality reduction that finds a linear combination of features that best separates multiple classes.

Formula/Key Concepts:

- **Objective:** Maximize the ratio of between-class variance to within-class variance.
- **Projection:** Projects data onto a lower-dimensional space that maximizes class separability.

When to Use:

- When you want to reduce dimensionality while preserving class separability.
- For problems where classes are linearly separable.

When to Avoid:

- When classes have different covariance structures.
- For non-linear class boundaries.

Advantages:

- Reduces dimensionality while preserving class structure.
- Simple and interpretable.

Disadvantages:

- Assumes normally distributed data.
- Less effective with non-linear relationships.

Typical Use Cases:

- Pattern recognition.
- Classification tasks with high-dimensional data.

Scenarios Where It May Not Perform Well:

- Classes with differing covariance matrices.
- Non-linear class boundaries.

Kernel PCA

Definition: Kernel PCA is an extension of PCA that uses kernel methods to perform non-linear dimensionality reduction by mapping data into a higher-dimensional space.

Formula/Key Concepts:

- **Kernel Trick:** Uses a kernel function to compute the principal components in a higher-dimensional feature space.
- **Objective:** Finds principal components in the transformed space.

When to Use:

- When data has non-linear relationships.
- For capturing complex patterns not achievable with linear PCA.

When to Avoid:

- When computational resources are limited (can be computationally expensive).
- For very large datasets.

Advantages:

- Captures non-linear structures.
- Flexible with different kernel functions.

Disadvantages:

- Computationally intensive.
- Choice of kernel and parameters can be challenging.

Typical Use Cases:

- Non-linear feature extraction.
- Complex data transformations.

Scenarios Where It May Not Perform Well:

- Large datasets with high computational demands.
- When the choice of kernel is not optimal.

t-Distributed Stochastic Neighbor Embedding (t-SNE)

Definition: t-Distributed Stochastic Neighbor Embedding (t-SNE) is a non-linear dimensionality reduction technique designed for visualizing high-dimensional data in lower dimensions.

Formula/Key Concepts:

- **Objective:** Minimize the divergence between probability distributions in high-dimensional and low-dimensional spaces.
- **Optimization:** Uses gradient descent to find a lower-dimensional representation.

When to Use:

- For visualizing complex, high-dimensional data.
- When data has non-linear relationships.

When to Avoid:

- For large datasets due to computational constraints.
- When interpretability of the transformed features is needed.

Advantages:

- Effective for visualizing high-dimensional data.
- Captures non-linear relationships and complex structures.

Disadvantages:

- Computationally expensive for large datasets.
- Results can be sensitive to parameter choices (e.g., perplexity).

Typical Use Cases:

- Data visualization.
- Exploratory data analysis.

Scenarios Where It May Not Perform Well:

- Very large datasets.
- When results need to be reproducible or highly interpretable.

Independent Component Analysis (ICA)

Definition: Independent Component Analysis (ICA) is a technique used to separate a multivariate signal into additive, independent components.

Formula/Key Concepts:

- **Objective:** Find components that are statistically independent from each other.
- **Transformation:** Uses statistical methods to estimate independent components.

When to Use:

- When you need to separate mixed signals into independent sources.
- For applications involving blind source separation.

When to Avoid:

- When components are not independent.
- For data with complex non-linear relationships.

Advantages:

- Useful for separating mixed signals.
- Can reveal hidden factors in data.

Disadvantages:

- Assumes components are statistically independent.
- Requires large datasets for accurate results.

Typical Use Cases:

- Signal processing (e.g., separating audio sources).
- Feature extraction in machine learning.

Scenarios Where It May Not Perform Well:

- Non-independent sources.
- When data size is insufficient for reliable separation.

Factor Analysis

Definition: Factor Analysis is a technique used to identify underlying relationships between variables by modeling the data as a linear combination of potential factors.

Formula/Key Concepts:

- **Objective:** Identify factors that explain the variance and correlations among observed variables.
- **Model:** Data is modeled as $X = \Lambda F + \epsilon$, where Λ is the factor loading matrix, F is the factor matrix, and ϵ is the error term.

When to Use:

- When you want to uncover underlying factors driving the observed variables.
- For data with potential latent variables.

When to Avoid:

- For datasets with insufficient sample size.
- When factors cannot be easily interpreted.

Advantages:

- Identifies underlying structure in data.
- Reduces dimensionality by combining correlated variables.

Disadvantages:

- Results can be difficult to interpret.
- Assumes linear relationships among variables.

Typical Use Cases:

- Psychology and social sciences (e.g., personality assessments).
- Market research (e.g., identifying underlying factors in consumer behavior).

Scenarios Where It May Not Perform Well:

- Small datasets or data with high complexity.
- When factors are not easily interpretable.

Bagging (Bootstrap Aggregating)

Definition: Bagging (Bootstrap Aggregating) is an ensemble method that improves model stability and accuracy by combining predictions from multiple models trained on different subsets of the data.

Formula/Key Concepts:

- **Bootstrap Sampling:** Randomly sample subsets of data with replacement.
- **Aggregation:** Combine predictions (e.g., majority voting for classification, averaging for regression).

When to Use:

- When you want to reduce model variance and improve stability.
- For high-variance models like decision trees.

When to Avoid:

- When you need to handle high bias issues.
- For very small datasets.

Advantages:

- Reduces variance and helps prevent overfitting.
- Simple and effective with decision trees.

Disadvantages:

- Can be computationally expensive with large base models.
- Does not address high bias.

Typical Use Cases:

- Improving accuracy of decision trees.
- Ensemble methods in complex machine learning pipelines.

Scenarios Where It May Not Perform Well:

- Low-bias, high-variance scenarios.
- Extremely high-dimensional data without sufficient samples.

Boosting

Definition: Boosting is an ensemble technique that sequentially builds models, each correcting errors made by the previous models, to improve overall performance.

Key Types of Boosting:

1. AdaBoost:

- **Definition:** Adaptive Boosting focuses on training a series of models where each model is trained to correct the mistakes of its predecessors.
- **When to Use:** When you need to reduce both bias and variance.
- **Advantages:** Simple and effective for many types of data.
- **Disadvantages:** Sensitive to noisy data and outliers.

2. Gradient Boosting:

- **Definition:** Builds models sequentially, with each new model correcting the residual errors of the previous models.
- **When to Use:** When you need to improve model performance by reducing bias.
- **Advantages:** Handles a variety of data types and can model complex relationships.
- **Disadvantages:** Computationally intensive and sensitive to hyperparameters.

3. XGBoost:

- **Definition:** Extreme Gradient Boosting is an optimized version of gradient boosting with improved performance and scalability.
- **When to Use:** For large datasets and complex models needing high performance.
- **Advantages:** High accuracy, efficiency, and scalability.
- **Disadvantages:** Can be complex to tune and interpret.

4. LightGBM:

- **Definition:** Light Gradient Boosting Machine is a faster variant of gradient boosting with optimized performance for large datasets.
- **When to Use:** When speed and scalability are critical.
- **Advantages:** Faster training and lower memory usage.
- **Disadvantages:** Requires careful parameter tuning.

5. CatBoost:

- **Definition:** Categorical Boosting is designed to handle categorical features efficiently in boosting.
- **When to Use:** When dealing with categorical data.
- **Advantages:** Handles categorical features natively, reduces the need for preprocessing.
- **Disadvantages:** Less mature compared to XGBoost and LightGBM.

When to Avoid:

- When models are already low-bias and high-variance.
- For very noisy data without proper regularization.

Advantages:

- Reduces both bias and variance.
- Often leads to state-of-the-art performance.

Disadvantages:

- Can be computationally intensive.
- Sensitive to hyperparameters and noisy data.

Typical Use Cases:

- Competitions like Kaggle.
- Complex datasets needing high predictive power.

Scenarios Where It May Not Perform Well:

- Very noisy datasets without proper tuning.
- Situations requiring very fast model training.

Stacking (Stacked Generalization)

Definition: Stacking is an ensemble method that combines multiple models by training a meta-model to learn the best way to combine the predictions from base models.

Formula/Key Concepts:

- **Base Models:** Train several different models on the same dataset.
- **Meta-Model:** Trains on the predictions of base models to make the final prediction.

When to Use:

- When you want to leverage the strengths of multiple models.
- For complex datasets where single models might not perform well.

When to Avoid:

- For very simple problems where a single model might suffice.
- When interpretability is crucial and the stacking model becomes too complex.

Advantages:

- Can significantly improve model performance by combining diverse models.
- Helps capture different aspects of the data.

Disadvantages:

- Computationally expensive.
- Can be complex to implement and tune.

Typical Use Cases:

- Model stacking in competitions and complex predictive tasks.
- Combining different model types to improve accuracy.

Scenarios Where It May Not Perform Well:

- Simple problems where stacking does not provide significant benefits.
- Scenarios with limited computational resources.

Voting Classifier

Definition: Voting Classifier is an ensemble method that combines predictions from multiple models using voting strategies to determine the final classification.

Formula/Key Concepts:

- **Voting Strategies:** Majority voting (hard voting) or average probabilities (soft voting).

When to Use:

- When you want to combine predictions from multiple models to improve accuracy.
- For tasks where individual models have different strengths.

When to Avoid:

- When the base models are very similar and do not add additional value.
- For very small datasets.

Advantages:

- Simple to implement and understand.
- Can improve performance by leveraging diverse models.

Disadvantages:

- Performance depends on the diversity and quality of base models.
- Can be less effective if models are highly correlated.

Typical Use Cases:

- Combining different classifiers to improve accuracy.
- Simple ensemble approach for diverse model types.

Scenarios Where It May Not Perform Well:

- When all base models are very similar or overfitted.
- For very high-dimensional data where base models may be less effective.

Artificial Neural Networks (ANN)

Definition: Artificial Neural Networks (ANNs) are computational models inspired by the human brain, consisting of interconnected nodes (neurons) organized in layers, used for learning patterns from data.

Formula/Key Concepts:

- **Layers:** Input layer, hidden layers, and output layer.
- **Activation Functions:** Sigmoid, ReLU, Tanh, etc.
- **Training:** Uses backpropagation and gradient descent to update weights.

When to Use:

- When modeling complex relationships in data.
- For tasks requiring learning non-linear patterns.

When to Avoid:

- For very small datasets where overfitting is a concern.
- When interpretability is crucial.

Advantages:

- Flexible and can model complex functions.
- Suitable for a variety of tasks including classification, regression, and pattern recognition.

Disadvantages:

- Requires a large amount of data and computational power.
- Can be prone to overfitting.

Typical Use Cases:

- Image and speech recognition.
- Predictive analytics and data classification.

Scenarios Where It May Not Perform Well:

- Limited data availability.
- Tasks requiring high interpretability.

Convolutional Neural Networks (CNN)

Definition: Convolutional Neural Networks (CNNs) are specialized neural networks designed to process data with a grid-like topology, such as images, using convolutional layers to capture spatial hierarchies.

Formula/Key Concepts:

- **Convolutional Layers:** Apply filters to input data to detect features.
- **Pooling Layers:** Reduce dimensionality while retaining important features.
- **Activation Functions:** Often ReLU.

When to Use:

- For tasks involving image or spatial data.
- When capturing hierarchical patterns and features is important.

When to Avoid:

- For non-image data where spatial relationships are not relevant.
- When computational resources are limited.

Advantages:

- Excellent for image recognition and computer vision tasks.
- Can learn spatial hierarchies and features automatically.

Disadvantages:

- Computationally intensive and requires substantial memory.
- May not perform well on non-image data.

Typical Use Cases:

- Image classification and object detection.
- Video analysis and medical image analysis.

Scenarios Where It May Not Perform Well:

- Non-visual data where spatial patterns are not relevant.
- Limited computational resources.

Recurrent Neural Networks (RNN)

Definition: Recurrent Neural Networks (RNNs) are designed to handle sequential data by maintaining a hidden state that captures information from previous time steps.

Formula/Key Concepts:

- **Hidden States:** Maintain memory of previous time steps.
- **Activation Functions:** Often ReLU or Tanh.

When to Use:

- For sequential or time-series data.
- When capturing temporal dependencies is crucial.

When to Avoid:

- For tasks where long-term dependencies are important (due to vanishing gradient issues).
- When computational efficiency is a concern.

Advantages:

- Can process sequences of variable length.
- Captures temporal dependencies in data.

Disadvantages:

- Struggles with long-term dependencies (vanishing gradients).
- Training can be slow and computationally intensive.

Typical Use Cases:

- Natural language processing (NLP).
- Time-series forecasting and speech recognition.

Scenarios Where It May Not Perform Well:

- Long sequences with complex dependencies.
- When efficient training is required.

Long Short-Term Memory Networks (LSTM)

Definition: Long Short-Term Memory Networks (LSTMs) are a type of RNN designed to overcome the vanishing gradient problem by using gating mechanisms to control the flow of information.

Formula/Key Concepts:

- **Gates:** Input gate, forget gate, and output gate to regulate information flow.
- **Cell State:** Maintains long-term dependencies.

When to Use:

- For sequential data where long-term dependencies are important.
- When traditional RNNs struggle with gradient issues.

When to Avoid:

- For simple sequential tasks where LSTMs are not necessary.
- When computational resources are limited.

Advantages:

- Effective at capturing long-term dependencies in sequences.
- Handles complex sequential data well.

Disadvantages:

- Computationally intensive.
- Requires careful tuning of hyperparameters.

Typical Use Cases:

- Machine translation and text generation.
- Complex time-series forecasting.

Scenarios Where It May Not Perform Well:

- Simple sequential tasks.
- Limited computational resources.

Gated Recurrent Units (GRU)

Definition: Gated Recurrent Units (GRUs) are a variant of LSTMs that simplify the architecture by combining the input and forget gates into a single gate, improving efficiency.

Formula/Key Concepts:

- **Gates:** Update gate and reset gate.
- **Hidden State:** Directly updated without a separate cell state.

When to Use:

- When you need a more computationally efficient alternative to LSTMs.
- For tasks requiring capturing temporal dependencies.

When to Avoid:

- For very simple sequential tasks where GRUs may be overkill.
- When the additional complexity of GRUs is not justified.

Advantages:

- Fewer parameters compared to LSTMs, making training faster.
- Efficient at capturing temporal dependencies.

Disadvantages:

- May not capture long-term dependencies as effectively as LSTMs.
- Requires careful tuning of hyperparameters.

Typical Use Cases:

- Time-series forecasting and NLP tasks.
- Situations where computational efficiency is a concern.

Scenarios Where It May Not Perform Well:

- Extremely long sequences with very complex dependencies.
- When training efficiency is less critical.

Autoencoders

Definition: Autoencoders are neural networks designed for unsupervised learning, learning efficient representations (encoding) of data and reconstructing the original input from these representations.

Formula/Key Concepts:

- **Encoder:** Maps input to a lower-dimensional representation.
- **Decoder:** Reconstructs the original input from the representation.
- **Loss Function:** Typically mean squared error (MSE) between input and reconstruction.

When to Use:

- For dimensionality reduction and feature extraction.
- When learning unsupervised representations of data.

When to Avoid:

- For tasks requiring supervised learning.
- When interpretability of the encoding is not possible.

Advantages:

- Effective for feature extraction and data compression.
- Can be used for anomaly detection and denoising.

Disadvantages:

- Limited by the quality of the learned representations.
- Requires careful tuning of network architecture.

Typical Use Cases:

- Data compression and denoising.
- Anomaly detection and feature learning.

Scenarios Where It May Not Perform Well:

- When a high level of interpretability is needed.
- For data with very complex structures that autoencoders cannot capture.

Generative Adversarial Networks (GANs)

Definition: Generative Adversarial Networks (GANs) are a class of neural networks designed to generate new data samples by training two networks, a generator and a discriminator, in a game-theoretic framework.

Formula/Key Concepts:

- **Generator:** Creates fake data samples.
- **Discriminator:** Distinguishes between real and fake samples.
- **Objective:** Generator tries to fool the discriminator, while the discriminator tries to correctly classify samples.

When to Use:

- For generating realistic synthetic data.
- When exploring new data generation and enhancement techniques.

When to Avoid:

- For tasks where data generation is not needed.
- When computational resources are limited.

Advantages:

- Can generate high-quality synthetic data.
- Useful for data augmentation and simulation.

Disadvantages:

- Training can be unstable and require careful tuning.
- Can be computationally intensive.

Typical Use Cases:

- Image and video generation.
- Data augmentation and simulation.

Scenarios Where It May Not Perform Well:

- Limited resources and unstable training environments.
- When generating data is not necessary or useful.

Isolation Forest

Definition: Isolation Forest is an anomaly detection method that isolates observations by randomly selecting a feature and then randomly selecting a split value between the maximum and minimum values of the selected feature.

Formula/Key Concepts:

- **Isolation Trees:** Construct multiple trees where anomalies are isolated faster than normal observations.
- **Anomaly Score:** Based on the average path length in the trees; shorter path lengths indicate anomalies.

When to Use:

- For datasets with high-dimensional features.
- When the goal is to detect anomalies in large datasets.

When to Avoid:

- For very small datasets.
- When interpretability of detected anomalies is crucial.

Advantages:

- Efficient with high-dimensional data.
- Fast training and prediction.

Disadvantages:

- May not be effective with very small or very imbalanced datasets.
- Anomalies may not be easily interpretable.

Typical Use Cases:

- Fraud detection in financial transactions.
- Network intrusion detection.

Scenarios Where It May Not Perform Well:

- Small datasets where tree-based methods may overfit.
- When high interpretability of results is required.

Local Outlier Factor (LOF)

Definition: Local Outlier Factor (LOF) is an anomaly detection method that identifies outliers based on their local density compared to their neighbors.

Formula/Key Concepts:

- **Local Density:** Measures how isolated a data point is relative to its neighbors.
- **Anomaly Score:** Calculated using the ratio of the local density of the point to the local density of its neighbors.

When to Use:

- For datasets with varying densities.
- When local deviations from the density are of interest.

When to Avoid:

- For datasets where global anomalies are more relevant than local anomalies.
- When the dataset is very large and computational efficiency is a concern.

Advantages:

- Effective for detecting anomalies in datasets with varying densities.
- Can capture anomalies that are not globally outlying.

Disadvantages:

- Computationally intensive for large datasets.
- May require tuning of parameters for optimal performance.

Typical Use Cases:

- Anomaly detection in spatial and temporal data.
- Outlier detection in datasets with varying densities.

Scenarios Where It May Not Perform Well:

- Large datasets where computational resources are a constraint.
- Datasets with uniform density where local anomalies are less relevant.

One-Class SVM

Definition: One-Class SVM is an anomaly detection method that learns a decision boundary around the normal data points to identify outliers as those that fall outside this boundary.

Formula/Key Concepts:

- **Objective:** Maximize the margin around the normal data points while minimizing the number of data points outside this margin.
- **Kernel Trick:** Can use different kernel functions to handle non-linearly separable data.

When to Use:

- For datasets where only normal data is available (unsupervised anomaly detection).
- When you need to identify outliers in a high-dimensional space.

When to Avoid:

- For datasets with a large number of outliers, as it may become less effective.
- When computational resources are limited.

Advantages:

- Effective for high-dimensional data.
- Can handle non-linear boundaries with kernel functions.

Disadvantages:

- Can be computationally expensive, especially with large datasets.
- Performance can be sensitive to the choice of kernel and hyperparameters.

Typical Use Cases:

- Fraud detection in scenarios with only normal transaction data.
- Outlier detection in high-dimensional spaces.

Scenarios Where It May Not Perform Well:

- Very large datasets or datasets with numerous outliers.
- Scenarios where computational resources are constrained.

Q-Learning

Definition: Q-learning is a model-free reinforcement learning algorithm that learns the value of action-reward pairs to derive an optimal policy for decision-making.

When to Use:

- For problems with a discrete and relatively small state-action space.
- When you need a model-free approach to learn optimal policies.

When to Avoid:

- In environments with large or continuous state-action spaces.
- When the environment changes frequently or is non-stationary.

Advantages:

- Simple to implement and understand.
- Effective for small to medium-sized state-action spaces.

Disadvantages:

- Limited scalability to large state-action spaces.
- Requires a large number of iterations for convergence.

Typical Use Cases:

- Simple game environments like Grid World.
- Basic control tasks with discrete states and actions.

Scenarios Where It May Not Perform Well:

- High-dimensional or continuous state-action spaces.
- Dynamic environments where policies need frequent updates.

Deep Q-Networks (DQN)

Definition: Deep Q-Networks (DQN) extend Q-Learning by using deep neural networks to approximate the Q-function, allowing for handling larger and more complex state spaces.

Formula/Key Concepts:

- **Neural Network:** Used to approximate the Q-function, $Q(s,a;\theta)$.
- **Experience Replay:** Stores past experiences and samples randomly for training to improve stability.
- **Target Network:** A separate network used to stabilize training by periodically updating its weights from the main network.

When to Use:

- For problems with large or high-dimensional state spaces.
- When traditional Q-Learning struggles due to large state-action spaces.

When to Avoid:

- For very simple environments where a tabular Q-Learning approach is sufficient.
- When computational resources are limited.

Advantages:

- Can handle large and complex state spaces.
- Utilizes deep learning to capture intricate patterns.

Disadvantages:

- Computationally expensive and requires extensive training.
- Training can be unstable and sensitive to hyperparameters.

Typical Use Cases:

- Complex video games (e.g., Atari games).
- Robotics and autonomous systems with high-dimensional sensory inputs.

Scenarios Where It May Not Perform Well:

- Environments with very limited computational resources.
- Tasks where quick, simple solutions are preferable.

Policy Gradient Methods

Definition: Policy Gradient Methods directly optimize the policy by updating the policy parameters based on the gradient of the expected reward, rather than estimating value functions.

When to Use:

- For problems where the state or action space is continuous.
- When you need to learn complex policies that cannot be easily represented by value functions.

When to Avoid:

- For problems where value-based methods are more straightforward or efficient.
- When dealing with very large action spaces without sufficient data.

Advantages:

- Can handle high-dimensional and continuous action spaces.
- Capable of learning complex policies directly.

Disadvantages:

- Training can be slow and requires careful tuning.
- May suffer from high variance in gradient estimates.

Typical Use Cases:

- Continuous control tasks (e.g., robotic arm manipulation).
- Complex policy learning where action spaces are large or continuous.

Scenarios Where It May Not Perform Well:

- Environments with discrete and small action spaces where value-based methods are sufficient.
- Tasks requiring fast training and convergence.

Actor-Critic Methods

Definition: Actor-critic methods combine policy gradient methods (actor) with value-based methods (critic) to leverage both policy learning and value estimation for better performance.

When to Use:

- For problems requiring both policy optimization and value estimation.
- When dealing with large and complex state-action spaces.

When to Avoid:

- For problems where simpler methods provide sufficient performance.
- When computational resources are limited.

Advantages:

- Combines benefits of both policy and value-based approaches.
- Can improve learning efficiency and stability.

Disadvantages:

- More complex to implement and tune.
- Requires careful balance between actor and critic updates.

Typical Use Cases:

- Complex reinforcement learning tasks (e.g., advanced game environments, robotics).
- Scenarios where both policy and value estimation are needed for effective learning.

Scenarios Where It May Not Perform Well:

- Simple tasks where single-method approaches are adequate.
- Environments with limited computational resources or requiring rapid prototyping.