# Machine Learning Revision

## Important Basic Terms

1. **Machine Learning:** Machine learning is a subset of artificial intelligence that focuses on the development of algorithms and statistical models that enable computers to learn from and make predictions or decisions based on data without being explicitly programmed. It involves the use of algorithms to identify patterns or relationships within data and make predictions or decisions based on those patterns.

2. **Supervised Learning:** Supervised learning is a type of machine learning where the model is trained on a labeled dataset, which means that each example in the dataset is associated with a target label or outcome. The goal of supervised learning is to learn a mapping from input features to target labels. Examples include:
   Classification: Predicting whether an email is spam or not spam.
   Regression: Predicting the price of a house based on its features.

3. **Unsupervised Learning:** Unsupervised learning is a type of machine learning where the model is trained on an unlabeled dataset, which means that there are no target labels or outcomes provided. The goal of unsupervised learning is to find hidden patterns or structures within the data. Examples include:
   Clustering: Grouping similar customers based on their purchasing behavior.
   Dimensionality Reduction: Reducing the number of features in a dataset while preserving its underlying structure.

4. **Reinforcement Learning:** Reinforcement learning is a type of machine learning where an agent learns to make decisions by interacting with an environment. The agent receives feedback in the form of rewards or penalties based on its actions, and the goal is to learn a policy that maximizes the cumulative reward over time. Examples include:
   Training an AI agent to play chess or Go.
   Teaching a robot to navigate a maze.

5. **Parameters:** Parameters are the internal variables of a machine learning model that are learned from training data. They define the behavior or configuration of the model and are adjusted during the training process to minimize the error between the model's predictions and the actual target values.

6. **Hyperparameters:** Hyperparameters are configuration settings that are external to the model and are set before the training process begins. They control aspects of the learning process such as the model's complexity, regularization strength, and optimization

algorithm. Examples include learning rate, regularization parameter, and number of hidden layers in a neural network.

7. **Weights:** Weights are the coefficients that multiply the input features in a machine learning model. They determine the importance or contribution of each feature to the model's predictions and are learned during the training process.

8. **Costs:** Costs, also known as loss or error, represent the discrepancy between the predicted output of a machine learning model and the actual target values in the training data. The goal of training a model is to minimize the cost function, which measures the average error over all training examples.

9. **Bias:** Bias is the error introduced by approximating a real-world problem with a simplified model. It represents the difference between the average prediction of the model and the true value of the target variable. High bias can lead to underfitting, where the model is too simple to capture the underlying patterns in the data.

10. **Variance:** Variance is the variability of model predictions for a given input when trained on different datasets. It represents the model's sensitivity to fluctuations in the training data. High variance can lead to overfitting, where the model captures noise or random fluctuations in the training data.

11. **Overfitting:** Overfitting occurs when a machine learning model learns the training data too well, capturing noise or random fluctuations in the data rather than underlying patterns. As a result, the model performs well on the training data but poorly on unseen or test data.

12. **Underfitting:** Underfitting occurs when a machine learning model is too simple to capture the underlying patterns in the data. It typically occurs when the model is not complex enough to represent the true relationship between the input features and the target variable, resulting in poor performance on both training and test data.

13. **Regularization:** Regularization is a technique used to prevent overfitting by adding a penalty term to the cost function that discourages large parameter values. It helps to control the complexity of the model and reduce the risk of overfitting by imposing constraints on the weights or parameters of the model.

14. **Activation Function:** An activation function is a mathematical function that determines the output of a neuron in a neural network based on its input. It introduces non-linearity into the network, allowing it to learn complex patterns in the data. Common activation functions include sigmoid, tanh, ReLU, and softmax.

15. **Cost Function:** A cost function, also known as a loss function or objective function, measures the difference between the predicted output of a machine learning model and the actual target values. The goal of training a model is to minimize the cost function, which represents the error between the predicted and actual values.

16. **Learning Rate:** The learning rate is a hyperparameter that controls the size of the steps taken during the optimization process. It determines how quickly or slowly the model learns from the training data. A higher learning rate can lead to faster convergence but may cause the model to overshoot the optimal solution, while a lower learning rate may result in slower convergence but more stable training.

17. **Iterations, Epochs, and Batches:**
    a. **Iterations:** In machine learning, an iteration refers to a single update of the model's parameters using a batch of training examples. It involves computing the gradients of the cost function with respect to the parameters and updating the parameters using an optimization algorithm such as gradient descent.
    b. **Epochs:** An epoch refers to one complete pass through the entire training dataset during the training process. It consists of multiple iterations, where each iteration updates the model's parameters using a different batch of training examples. Training typically involves running multiple epochs until the model converges or a stopping criterion is met.
    c. **Batches:** A batch is a subset of the training data used to compute the gradients of the cost function during each iteration. Batching allows for more efficient computation and memory usage, especially for large datasets. Common batch sizes include 32, 64, or 128 examples per batch.

## Linear Regression

Linear regression is a fundamental supervised learning algorithm used for modeling the relationship between a dependent variable and one or more independent variables. It is a simple yet powerful technique that assumes a linear relationship between the input variables and the output.

**Process**

**Data Collection:** Gather a dataset containing the independent variables (features) and the dependent variable (target).

**Data Preprocessing:** Clean the data by handling missing values, outliers, and encoding categorical variables. Split the data into training and testing sets.

**Feature Selection/Extraction:** Choose relevant features that have a significant impact on the dependent variable. Feature engineering techniques such as polynomial features, interaction terms, and dimensionality reduction may be applied.

**Model Training:** Fit the linear regression model to the training data by finding the optimal coefficients (weights) that minimize the difference between the predicted and actual values of the target variable. This is often done using the method of least squares or gradient descent.

**Model Evaluation:** Assess the performance of the model using evaluation metrics such as mean squared error (MSE), root mean squared error (RMSE), mean absolute error (MAE), or R-squared. Evaluate the model on the testing set to check for overfitting.

**Prediction:** Use the trained model to make predictions on new or unseen data. Evaluate the model's performance on real-world data and refine as needed.

**Gradient Descent**

Gradient descent is an optimization algorithm used to minimize the cost function (error) of a machine learning model. In linear regression, gradient descent iteratively updates the model's parameters (coefficients) by moving in the direction of the steepest descent of the cost function. The learning rate controls the size of the steps taken during each iteration.

**Lasso and Ridge Regularization**

**Lasso (L1 Regularization):** Lasso regularization adds a penalty term to the cost function that is proportional to the absolute values of the model's coefficients. It encourages sparsity by shrinking some coefficients to zero, effectively performing feature selection and reducing model complexity.

**Ridge (L2 Regularization):** Ridge regularization adds a penalty term to the cost function that is proportional to the squared magnitudes of the model's coefficients. It penalizes large coefficients and tends to shrink them towards zero without eliminating them entirely, thus reducing the model's variance and mitigating overfitting.

**Example Projects**
- **Predicting House Prices:** Use linear regression to predict the prices of houses based on features such as size, number of bedrooms, location, etc.

- **Stock Price Prediction:** Apply linear regression to predict future stock prices based on historical stock data and relevant market indicators.

- **Salary Prediction:** Use linear regression to predict the salaries of employees based on features such as years of experience, education level, job title, etc.

- **Customer Lifetime Value Prediction:** Predict the lifetime value of customers using linear regression based on their historical purchasing behavior, demographics, and other relevant factors.

- **Medical Diagnosis:** Apply linear regression to predict medical outcomes or diagnose diseases based on patient data such as symptoms, test results, and medical history.

## Multiple Linear Regression

Multiple linear regression is an extension of simple linear regression that allows for modeling the relationship between a dependent variable and two or more independent variables. It assumes a linear relationship between the independent variables and the dependent variable, with each independent variable having a linear impact on the dependent variable while holding other variables constant.

- **Multicollinearity:** Multicollinearity occurs when two or more independent variables in a multiple linear regression model are highly correlated with each other. It can lead to unstable coefficient estimates and inflated standard errors. Techniques such as variance inflation factor (VIF) are used to detect and address multicollinearity.

- **Adjusted R-squared:** R-squared (coefficient of determination) measures the proportion of variance in the dependent variable that is explained by the independent variables in the model. Adjusted R-squared is a modified version of R-squared that adjusts for the number of predictors in the model. It penalizes the addition of unnecessary variables and provides a more accurate measure of model fit.

- **Interaction Terms:** Interaction terms are additional variables created by multiplying two or more independent variables together. They allow the model to capture non-additive effects and interactions between variables. For example, in a marketing study, an interaction term between advertising spend and customer income could capture how the effect of advertising varies with income levels.

- **Feature Scaling:** Feature scaling is the process of standardizing or normalizing the independent variables in a multiple linear regression model to ensure that all variables have the same scale. It prevents variables with larger magnitudes from dominating the model and improves convergence during training. Common techniques include min-max scaling and z-score normalization.

- **Model Assumptions:** Multiple linear regression relies on several assumptions, including linearity, independence of errors, homoscedasticity (constant variance of errors), and normality of errors. Violations of these assumptions can lead to biased estimates and inaccurate predictions. Diagnostic tests such as residual analysis and Q-Q plots are used to assess the validity of these assumptions.

# Decision Trees

Decision trees are a popular supervised learning algorithm used for classification and regression tasks. They create a tree-like structure of decisions based on features of the data to predict the target variable.

**Basic Concepts**

**Tree Structure:** Decision trees consist of nodes, branches, and leaves. Each node represents a feature of the dataset, each branch represents a decision based on that feature, and each leaf node represents the predicted outcome.

**Splitting Criteria:** Decision trees split the dataset into subsets based on the values of features to maximize the homogeneity of the target variable within each subset. Common splitting criteria include Gini impurity and information gain (entropy).

**Decision Making:** At each internal node of the tree, a decision is made based on the value of a specific feature. The dataset is split into subsets based on this decision, and the process repeats recursively until a stopping criterion is met.

**Leaf Nodes and Predictions:** When the tree reaches a leaf node, it makes a prediction based on the majority class (for classification) or the mean value (for regression) of the target variable in that subset.

**Pruning:** Pruning is a technique used to reduce the size of the decision tree by removing unnecessary branches. It helps prevent overfitting and improves the generalization ability of the model.

**Some Important Concepts**

**Entropy and Information Gain:** Entropy measures the impurity or randomness of a dataset. Information gain quantifies the reduction in entropy achieved by splitting the dataset based on a particular feature. Decision trees aim to maximize information gain at each split to create more homogeneous subsets.

**Gini Impurity:** Gini impurity is another measure of impurity used in decision trees. It represents the probability of incorrectly classifying a randomly chosen element if it were labeled according to the distribution of classes in the subset. Decision trees aim to minimize Gini impurity at each split.

**Overfitting:** Decision trees are prone to overfitting, especially when they become too deep and complex. Overfitting occurs when the model captures noise and patterns specific to the training data, leading to poor generalization on unseen data. Techniques such as pruning, limiting tree depth, and using ensemble methods (e.g., random forests) help mitigate overfitting.

**Feature Importance:** Decision trees provide a measure of feature importance based on how much each feature contributes to reducing impurity or increasing information gain. Features with higher importance values are considered more influential in predicting the target variable.

**Categorical and Numerical Variables:** Decision trees can handle both categorical and numerical variables. For categorical variables, the tree splits the dataset based on each category. For numerical variables, the tree selects an optimal threshold to split the data into two subsets.

# Logistic Regression

Logistic regression is a statistical method used for binary classification tasks, where the target variable has two possible outcomes (e.g., yes/no, spam/not spam). Despite its name, logistic regression is a classification algorithm, not a regression algorithm.

**Basic Concepts**

**Sigmoid Function:** Logistic regression uses the sigmoid function (also known as the logistic function) to transform the output of a linear equation into a probability value between 0 and 1. The sigmoid function is defined as $\sigma(z) = 1/1 + e^{-z}$, where (z) is the linear combination of the input features and model coefficients.

**Binary Classification:** Logistic regression models the probability that a given input belongs to a particular class. If the probability is greater than a threshold (typically 0.5), the input is classified as belonging to the positive class; otherwise, it is classified as belonging to the negative class.

**Maximum Likelihood Estimation:** Logistic regression estimates the model coefficients (weights) using maximum likelihood estimation (MLE). The goal is to find the set of coefficients that maximize the likelihood of observing the actual target classes given the input features.

**Decision Boundary:** The decision boundary separates the input space into regions corresponding to the positive and negative classes. In logistic regression, the decision boundary is linear, meaning it is a straight line (or hyperplane in higher dimensions) defined by the equation $(\theta^T X = 0 \)$, where $(\theta)$ is the vector of model coefficients and $(X)$ is the vector of input features.

**Important Concepts**

**Regularization:** Regularization techniques such as L1 (Lasso) and L2 (Ridge) regularization can be applied to logistic regression to prevent overfitting and improve generalization performance. Regularization adds a penalty term to the cost function, which penalizes large coefficients and encourages simpler models.

**Confusion Matrix and Evaluation Metrics:** Confusion matrix is a table that summarizes the performance of a classification model by comparing the predicted classes with the actual classes. Evaluation metrics derived from the confusion matrix include accuracy, precision, recall (sensitivity), specificity, F1-score, and ROC-AUC (Receiver Operating Characteristic - Area Under Curve).

**Multiclass Classification:** Logistic regression can be extended to handle multiclass classification tasks using techniques such as one-vs-rest (OvR) or one-vs-one (OvO) classification. In OvR, separate binary logistic regression models are trained for each class, while in OvO, pairwise comparisons are made between all pairs of classes.

**Feature Importance:** Similar to linear regression, logistic regression provides a measure of feature importance based on the magnitude of the coefficients. Features with higher absolute coefficient values are considered more influential in predicting the target class.

**Example Projects**
**Email Spam Detection:** Use logistic regression to classify emails as spam or not spam based on features extracted from the email content, sender information, and email header.

**Customer Churn Prediction:** Apply logistic regression to predict whether customers will churn (cancel their subscription or leave the service) based on their demographics, usage patterns, and interactions with the company.

**Credit Risk Assessment:** Use logistic regression to assess the risk of default on credit loans by classifying applicants as low-risk or high-risk based on their credit history, income, and other financial indicators.

**Disease Diagnosis:** Apply logistic regression to assist in medical diagnosis by predicting the likelihood of a patient having a particular disease based on their symptoms, medical history, and diagnostic test results.

# Naive Bayes

Naive Bayes is a probabilistic classification algorithm based on Bayes' theorem with an assumption of independence between features. Despite its simplicity, Naive Bayes is effective and widely used for text classification, spam filtering, sentiment analysis, and other tasks.

**Basic Concepts**
**Bayes' Theorem:** Naive Bayes is based on Bayes' theorem, which describes the probability of an event given prior knowledge of conditions that might be related to the event. Bayes' theorem is expressed as:

$$P(A|B) = P(B|A) . P(A)/P(B)$$

where $P(A|B)$ is the probability of event A given event B, $P(B|A)$ is the probability of event B given event A, $P(A)$ is the prior probability of event A, and $P(B)$ is the prior probability of event B.

**Independence Assumption:** Naive Bayes assumes that the presence of a particular feature in a class is independent of the presence of other features. This assumption simplifies the calculation of probabilities and makes the algorithm computationally efficient.

**Probability Estimation:** Naive Bayes calculates the probabilities of each class based on the occurrence of features in the training data. It uses the training data to estimate the prior probabilities of each class and the conditional probabilities of each feature given the class.

**Classification Rule:** Naive Bayes uses a simple classification rule known as the maximum a posteriori (MAP) rule to assign the most probable class label to a given instance. According to

the MAP rule, the class with the highest posterior probability given the input features is selected as the predicted class.

**Important Concepts**

**Types of Naive Bayes**

There are different variations of Naive Bayes classifiers, including:
- **Gaussian Naive Bayes:** Assumes that continuous features follow a Gaussian (normal) distribution.
- **Multinomial Naive Bayes:** Suitable for text classification tasks with discrete features (e.g., word counts).
- **Bernoulli Naive Bayes:** Applicable when features are binary (e.g., presence/absence of a word).
-

**Handling Zero Probabilities:** Naive Bayes can encounter zero probabilities for features not present in the training data. Techniques such as Laplace smoothing or add-one smoothing are used to handle zero probabilities and prevent the model from assigning zero probability to unseen features.

**Feature Independence vs. Correlation:** While Naive Bayes assumes feature independence, it may still perform well even when this assumption is violated. It can capture some degree of feature correlation and still provide effective classification results, especially in text classification tasks.

**Scalability and Efficiency:** Naive Bayes is computationally efficient and scalable, making it suitable for large-scale datasets and real-time applications. It requires minimal parameter tuning and can handle high-dimensional feature spaces efficiently.