

1 - Key concepts

jueves, 2 de mayo de 2024 8:35

Properties:

Homoscedasticity:

Homoscedasticity is a term used in statistics to describe a situation where the variability of the data points around the regression line is about the same across all levels of the independent variable(s). In simpler terms, it means that the spread or scatter of the data points is consistent across the range of values of the independent variable(s).

Imagine you're fitting a regression line to a scatter plot of data points. Homoscedasticity means that the vertical distance between each data point and the regression line is roughly the same, regardless of where you are along the x-axis. In other words, the "thickness" of the "cloud" of points around the line stays the same from left to right.

Multicollinearity:

Key ML concepts:

□ **Hypothesis** — The hypothesis is noted h_{θ} and is the model that we choose. For a given input data $x^{(i)}$ the model prediction output is $h_{\theta}(x^{(i)})$.

□ **Likelihood** — The likelihood of a model $L(\theta)$ given parameters θ is used to find the optimal parameters θ through likelihood maximization. We have:

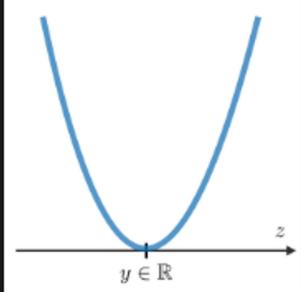
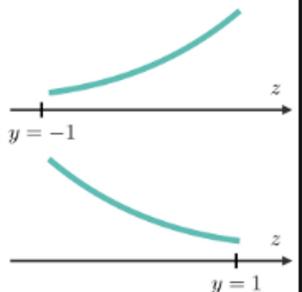
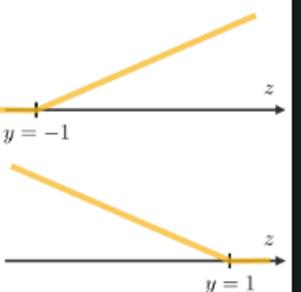
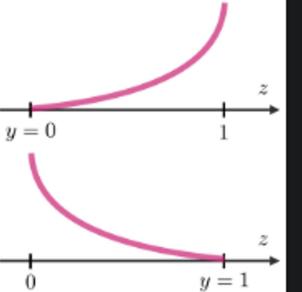
$$\theta^{\text{opt}} = \arg \max_{\theta} L(\theta)$$

Remark: in practice, we use the log-likelihood $\ell(\theta) = \log(L(\theta))$ which is easier to optimize.

It measures the probability of observing the given data under a specific model or hypothesis.

In short and simple terms, likelihood tells us how probable the observed data is, given the assumptions of our model. It's like asking, "How likely is it that we would see this data if our model were true?" Likelihood plays a crucial role in tasks such as parameter estimation, where the goal is to find the model parameters that maximize the likelihood of observing the data.

□ Loss function — A loss function is a function $L : (z, y) \in \mathbb{R} \times Y \mapsto L(z, y) \in \mathbb{R}$ that takes as inputs the predicted value z corresponding to the real data value y and outputs how different they are. The common loss functions are summed up in the table below:

Least squared error	Logistic loss	Hinge loss	Cross-entropy
$\frac{1}{2}(y - z)^2$	$\log(1 + \exp(-yz))$	$\max(0, 1 - yz)$	$-(y \log(z) + (1 - y) \log(1 - z))$
			
Linear regression	Logistic regression	SVM	Neural Network

A loss function is a way to measure how well a machine learning model is performing on a task. It calculates the difference between the predicted output of the model and the actual target output. The goal of the model during training is to minimize this difference, or "loss," by adjusting its parameters. In simple terms, the loss function tells the model how wrong it is, and the model tries to minimize this wrongness during training.

□ Cost function — The cost function J is commonly used to assess the performance of a model, and is defined with the loss function L as follows:

$$J(\theta) = \sum_{i=1}^m L(h_\theta(x^{(i)}), y^{(i)})$$

The cost function is similar to a loss function—it measures how well a model is performing, but it **calculates the overall error across all training examples rather than just a single example**. The model's goal is still to minimize this cost function during training by adjusting its parameters.

A smooth and well-behaved cost function is one that doesn't have sudden changes or discontinuities in its gradients, has a single global minimum, and has bounded gradients.

Supervised Learning: It's like teaching a machine to recognize patterns by showing it labeled examples. For instance, you teach it to recognize cats by showing it pictures of cats labeled "cat" and pictures of other things labeled "not cat."

Unsupervised Learning: Here, the machine has to find patterns on its own without labeled data. It's like exploring a new city without a map or guide - you try to group similar things together or find underlying structures without being told what they are.

Underfitting: Imagine not studying enough for a test and performing poorly as a result. Underfitting occurs when a model is too simple to capture the underlying structure of the data, leading to poor performance on both training and test data.

Overfitting: Imagine memorizing answers to specific questions rather than understanding the topic. Overfitting happens when a model learns the training data too well, including noise or irrelevant details, and performs poorly on new, unseen data.

Learning rate: The learning rate is a hyperparameter that controls the step size or rate at which a machine learning model updates its parameters during training. In optimization algorithms such as gradient descent, the learning rate determines the magnitude of the updates applied to the model parameters in each iteration.

Cross-validation: Cross-validation is a resampling technique used to assess the performance of machine learning models. It involves splitting the dataset into multiple subsets, called folds, to train and evaluate the model multiple times. The main goal of cross-validation is to estimate how well a model will generalize to new, unseen data.

Transfer Learning: It's like leveraging knowledge gained from one task to improve performance on another. In transfer learning, a model pre-trained on a large dataset (source task) is adapted to a new, related task (target task) by transferring learned representations. This can save time and resources by starting with a model that has already learned useful features.

- **Fine-Tuning:** It's like making small adjustments to a well-trained athlete's technique to excel in a specific event. Fine-tuning is a transfer learning technique where a pre-trained model is further trained on a new dataset specific to the target task. The model's parameters are adjusted slightly to adapt to the new data while retaining the knowledge gained from the original training. Fine-tuning allows the model to specialize in the nuances of the new task while benefiting from the general features learned during pre-training.

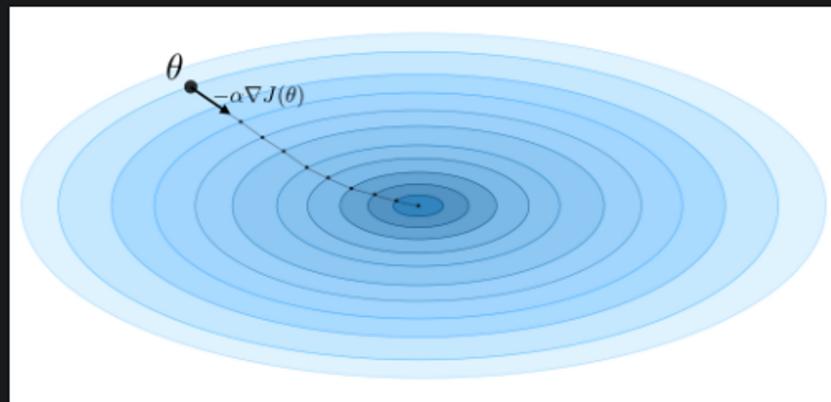
Ensembling: Ensembling is like combining opinions from different people to make a decision. There are two main types: Bagging (Bootstrap Aggregating), where multiple models are trained independently and then combined (e.g., Random Forest), and Boosting, where models are trained sequentially, and each new model learns from the mistakes of the previous ones (e.g., AdaBoost, Gradient Boosting).

Incremental Learning: It's like learning new things gradually without forgetting what you've learned before. Incremental learning refers to continuously updating a model's knowledge with new data over time, allowing it to adapt and improve without retraining from scratch.

Optimization methods:

□ **Gradient descent** — By noting $\alpha \in \mathbb{R}$ the learning rate, the update rule for gradient descent is expressed with the learning rate and the cost function J as follows:

$$\theta \leftarrow \theta - \alpha \nabla J(\theta)$$



Remark: Stochastic gradient descent (SGD) is updating the parameter based on each training example, and batch gradient descent is on a batch of training examples.

Gradient descent is an optimization algorithm used to minimize the cost or loss function of a machine learning model. It works by iteratively adjusting the parameters of the model in the direction that reduces the value of the cost function.

Gradient descent is more suitable in cases where:

- The cost function is not necessarily smooth or well-behaved.
- The optimization problem involves a large number of parameters.
- Computational resources are limited.
- Robustness to noise or ill-conditioned problems is important.

□ Newton's algorithm — Newton's algorithm is a numerical method that finds θ such that $\ell'(\theta) = 0$. Its update rule is as follows:

$$\boxed{\theta \leftarrow \theta - \frac{\ell'(\theta)}{\ell''(\theta)}}$$

Remark: the multidimensional generalization, also known as the Newton-Raphson method, has the following update rule:

$$\theta \leftarrow \theta - (\nabla_{\theta}^2 \ell(\theta))^{-1} \nabla_{\theta} \ell(\theta)$$

Newton's algorithm, also known as Newton's method or Newton-Raphson method, is an optimization algorithm used to find the minimum (or maximum) of a function. It's particularly useful for finding the roots of equations or minimizing the cost or loss function in machine learning.

In short and simple terms, Newton's algorithm is like finding the lowest point in a hilly terrain by taking into account both the steepness of the hill (gradient) and the curvature of the hill (second derivative). It uses this information to iteratively update the current guess of the minimum until it converges to a solution. This makes Newton's algorithm efficient for optimization tasks where the objective function is smooth and well-behaved.

Newton's algorithm is more suitable in cases where:

- The cost function is smooth and well-behaved.
- The optimization problem involves a small or moderate number of parameters.
- Computational resources are not a limiting factor.
- Fast convergence is desired, especially in regions with steep gradients or complex curvature.

Model results improvement methods:

Normalization:

Normalization is a technique used to rescale numeric features to a similar scale. There are several types of normalization commonly used in machine learning:

- **Min-Max Normalization (Feature Scaling):** This method scales the data to a fixed range, typically between 0 and 1. It's calculated using the formula:

$$X_{\text{normalized}} = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$
 where X is the original value, X_{\min} is the minimum value of the feature, and X_{\max} is the maximum value of the feature.
- **Z-Score Normalization (Standardization):** This method standardizes the data to have a mean of 0 and a standard deviation of 1. It's calculated using the formula:

$$X_{\text{normalized}} = \frac{X - \mu}{\sigma}$$
 where X is the original value, μ is the mean of the feature, and σ is the standard deviation of the feature.
- **Robust Scaling:** This method scales the data to be robust to outliers by using the median and interquartile range (IQR) instead of the mean and standard deviation. It's calculated using the formula:

$$X_{\text{normalized}} = \frac{X - \text{median}}{\text{IQR}}$$
 where X is the original value, the median is the median of the feature, and the IQR is the interquartile range of the feature.

- **Unit Vector Normalization:** This method scales each feature so that the magnitude of the feature vector is 1. It's calculated using the formula:

$$X_{\text{normalized}} = \frac{X}{\|X\|}$$

where X is the original feature vector, and $\|X\|$ is the Euclidean norm (magnitude) of the feature vector.

Transformations:

Several techniques can transform a variable distribution to resemble a Gaussian (normal) distribution:

- **Log Transformation:** Taking the logarithm of the variable often helps in reducing right skewness (positive skewness) and making the distribution more symmetrical.
- **Square Root Transformation:** Similar to the log transformation, the square root function can help reduce skewness, especially for variables with positive skewness.
- **Box-Cox Transformation:** The Box-Cox transformation is a parametric method that can be applied to stabilize the variance and make the distribution more Gaussian-like. It includes a parameter, lambda (λ), which determines the transformation applied to the data.
- **Yeo-Johnson Transformation:** Similar to the Box-Cox transformation, the Yeo-Johnson transformation is a modification that allows for both positive and negative values of the variable.
- **Power Transformation:** Power transformations involve raising the variable to a power (other than 1 or 0), such as squaring or cubing, to adjust the distribution.
- **Quantile Transformation:** This technique maps the variable's distribution to a specified probability distribution, often a standard normal distribution. It's useful for transforming the entire distribution rather than just reducing skewness.
- **Rank Transformation:** Rank transformation assigns ranks to the observations, effectively converting the variable into one with a uniform distribution. This can help in situations where other transformations are not effective.

Regularization:

Regularization techniques are used in machine learning to prevent overfitting and improve the generalization performance of models. Some common regularization techniques include:

- **L1 Regularization (Lasso):** This technique adds a penalty term to the loss function proportional to the absolute value of the model's coefficients. It encourages sparsity in the model by shrinking less important features' coefficients to zero, effectively performing feature selection.
- **L2 Regularization (Ridge):** L2 regularization adds a penalty term proportional to the square of the model's coefficients to the loss function. It penalizes large coefficients and encourages them to be small, leading to smoother model weights and reducing the model's sensitivity to individual data points.
- **Elastic Net Regularization:** Elastic Net combines both L1 and L2 regularization by adding penalties for both the absolute value and the square of the model's coefficients to the loss function. It combines the benefits of Lasso (feature selection) and Ridge (reducing multicollinearity) regularization.
- **Dropout:** Dropout is a regularization technique commonly used in neural networks. It randomly drops (sets to zero) a fraction of the input units during training to prevent complex co-adaptations of neurons and reduce overfitting. Dropout effectively creates an ensemble of multiple subnetworks and improves generalization.
- **Early Stopping:** Early stopping is a simple regularization technique that stops training the model when the performance on a validation dataset starts to degrade, thus preventing the model from overfitting to the training data.
- **Data Augmentation:** Data augmentation is a regularization technique commonly used in image classification tasks. It involves creating additional training examples by applying random transformations to the existing data, such as rotation, scaling, cropping, or flipping, to increase

the diversity of the training dataset and improve generalization.

- **Batch Normalization:** Batch normalization is a regularization technique used in deep neural networks. It normalizes the activations of each layer in the network by adjusting and scaling them based on the mean and variance of the activations within each mini-batch during training. Batch normalization helps stabilize the training process and reduces internal covariate shift, leading to faster convergence and better generalization.

Hyperparameter optimization:

Hyperparameter optimization techniques are methods used to systematically search for the optimal hyperparameters of a machine learning model. These techniques help automate the process of tuning hyperparameters, which are parameters that are set before the learning process begins and control aspects of the learning process itself.

Some common hyperparameter optimization techniques include:

- **Grid Search:** Grid search is a basic hyperparameter optimization technique that exhaustively searches through a specified subset of hyperparameter values. It evaluates the model performance for each combination of hyperparameters and selects the one that yields the best performance.
- **Random Search:** Random search randomly samples hyperparameter values from specified distributions and evaluates the model performance for each sampled combination. Random search is more computationally efficient than grid search and often finds good hyperparameter values with fewer evaluations.
- **Bayesian Optimization:** Bayesian optimization is a sequential model-based optimization technique that builds a probabilistic model of the objective function (model performance) and uses it to intelligently select the next set of hyperparameters to evaluate. Bayesian optimization typically requires fewer evaluations compared to grid search or random search, making it suitable for expensive-to-evaluate objective functions.
- **Gradient-Based Optimization:** Gradient-based optimization methods, such as gradient descent, can be used to optimize hyperparameters by treating the objective function (model performance) as a differentiable function of the hyperparameters. However, this approach requires access to the gradients of the objective function with respect to the hyperparameters, which may not always be available.
- **Evolutionary Algorithms:** Evolutionary algorithms, such as genetic algorithms or particle swarm optimization, mimic natural selection processes to iteratively evolve a population of hyperparameter configurations over multiple generations. These algorithms can explore a wide range of hyperparameter values and are suitable for complex optimization problems.
- **Meta-Learning or AutoML:** Meta-learning approaches, also known as AutoML (Automated Machine Learning), use machine learning algorithms to learn the optimal hyperparameter search strategy from previous model evaluations. These techniques can adaptively select hyperparameters based on the characteristics of the dataset and the model architecture.
- **Hyperband:** Hyperband is a bandit-based hyperparameter optimization technique that uses a multi-armed bandit algorithm to allocate computational resources to different hyperparameter configurations. It dynamically allocates more resources to promising configurations and discards underperforming ones, leading to efficient hyperparameter optimization.

2 - End-to-End

jueves, 2 de mayo de 2024 13:34

1. Data Collection and Understanding:

- **Data sources:** Identify and gather data from diverse sources such as databases, APIs, spreadsheets, or files.
- **Quality assessment:** Assess the quality of the data by checking for completeness, consistency, accuracy, and relevance.
- **Data documentation:** Document metadata, data dictionaries, or README files to provide context and understanding of the dataset.

2. Data Cleaning and Preprocessing:

- **Missing values handling:** Implement strategies like mean imputation, mode imputation, or advanced imputation techniques (e.g., KNN imputation) to deal with missing data.
- **Data formatting:** Convert data types, handle datetime objects, and ensure consistency in coding categorical variables.
- **Outlier detection:** Use statistical methods or visualization techniques to identify outliers and decide on appropriate treatment methods.
- **Normalization/Standardization:** Normalize or standardize numerical features to ensure they have similar scales, which can improve the performance of certain algorithms.

3. Exploratory Data Analysis (EDA):

- **Distribution exploration:** Analyze the distribution of variables to understand their characteristics and identify potential data transformations.
- **Correlation analysis:** Examine correlations between variables to identify relationships and multicollinearity issues.
- **Feature importance:** Assess the importance of features using techniques like feature importance plots or correlation-based methods.
- **Visualization techniques:** Utilize various visualization tools and libraries (e.g., matplotlib, seaborn) to create informative and insightful plots.

4. Feature Selection and Dimensionality Reduction:

- **Feature importance methods:** Employ techniques such as recursive feature elimination (RFE), feature importance scores, or domain knowledge to select relevant features.
- **Dimensionality reduction algorithms:** Apply algorithms like PCA, t-SNE, or LDA to reduce the number of features while preserving the most important information.

5. Model Building and Evaluation:

- **Algorithm selection:** Choose appropriate machine learning algorithms based on the problem type, dataset size, and complexity.
- **Cross-validation:** Use techniques like k-fold cross-validation to estimate model performance and assess generalization ability.
- **Hyperparameter tuning:** Optimize model hyperparameters using techniques like grid search, random search, or Bayesian optimization to improve performance.
- **Evaluation metrics:** Select evaluation metrics (e.g., accuracy, precision, recall, F1-score, ROC-AUC) based on the problem domain and requirements.

6. Interpretation and Insights:

- **Model interpretation:** Interpret model coefficients, feature importances, or decision boundaries to understand how the model makes predictions.
- **Insights extraction:** Extract actionable insights and recommendations from the analysis to inform business decisions or guide further investigations.
- **Communication skills:** Effectively communicate findings to diverse audiences, including stakeholders, clients, or non-technical users, using clear and concise language.

7. Iteration and Refinement:

- **Continuous improvement:** Iterate on the analysis process by incorporating feedback, refining methodologies, and exploring alternative approaches.
- **Documentation and reproducibility:** Document analysis steps, code, assumptions, and decisions to ensure transparency and reproducibility.
- **Feedback loop:** Establish a feedback loop with stakeholders to validate findings, address concerns, and iterate on analysis outputs.

ML pipelines:

A machine learning (ML) pipeline is a sequence of steps or processes that automate and streamline the end-to-end machine learning workflow, from data preprocessing and feature engineering to model training, evaluation, and deployment. It encapsulates the entire process of building, training, and deploying machine learning models in a structured and repeatable manner.

Here's a breakdown of the components typically included in an ML pipeline:

- **Data Ingestion:** This step involves collecting, loading, and accessing data from various sources, such as databases, files, APIs, or streaming platforms.
- **Data Preprocessing:** Data preprocessing includes cleaning, transforming, and formatting the raw data to make it suitable for analysis. Tasks may include handling missing values, encoding categorical variables, scaling numerical features, and splitting data into training and testing sets.
- **Feature Engineering:** Feature engineering involves creating new features or transforming existing ones to improve model performance. This step may include techniques like dimensionality reduction, feature selection, or generating new features based on domain knowledge.
- **Model Selection:** Model selection involves choosing the appropriate machine learning algorithm(s) for the task at hand, based on factors such as the problem type (e.g., classification, regression), data characteristics, and performance requirements.
- **Model Training:** In this step, the selected model(s) are trained on the training dataset using the chosen algorithm(s). The training process involves optimizing model parameters or hyperparameters to minimize the chosen loss function and improve predictive accuracy.
- **Model Evaluation:** Once trained, the model(s) are evaluated using the test dataset to assess their performance and generalization ability. Evaluation metrics such as accuracy, precision, recall, F1-score, or ROC-AUC are calculated to measure model performance.
- **Model Deployment:** After successful evaluation, the trained model(s) are deployed into production environments to make predictions on new, unseen data. Deployment may involve creating APIs, web services, or batch processing pipelines to serve model predictions to end-users or downstream applications.
- **Monitoring and Maintenance:** Once deployed, the model(s) need to be monitored regularly to ensure they continue to perform well over time. Monitoring involves tracking key performance metrics, detecting drifts or changes in data distribution, and retraining or updating models as needed to maintain their effectiveness.

3 - Supervised Learning

jueves, 2 de mayo de 2024 8:35

Resource: <https://stanford.edu/~shervine/teaching/cs-229/cheatsheet-supervised-learning>

Supervised learning is a type of machine learning where the model learns patterns from labeled data. It's like having a teacher guiding the learning process. The model is trained on input-output pairs, and it learns to make predictions or decisions based on the input data.

Type of prediction — The different types of predictive models are summed up in the table below:

	Regression	Classification
Outcome	Continuous	Class
Examples	Linear regression	Logistic regression, SVM, Naive Bayes

For classification there is:

Binary Classification: In binary classification, there are only two possible classes or categories. Examples include spam detection (spam or not spam), medical diagnosis (disease or healthy), and sentiment analysis (positive or negative sentiment).

Multiclass Classification: In multiclass classification, there are more than two possible classes or categories. Examples include image classification (classifying images into different types of objects or scenes), document classification (classifying documents into different topics or categories), and handwritten digit recognition (classifying digits from 0 to 9).

Type of model — The different models are summed up in the table below:

	Discriminative model	Generative model
Goal	Directly estimate $P(y x)$	Estimate $P(x y)$ to then deduce $P(y x)$
What's learned	Decision boundary	Probability distributions of the data
Illustration		
Examples	Regressions, SVMs	GDA, Naive Bayes

Discriminative models focus on predicting the output directly from the input, while generative models aim to understand the underlying data distribution and generate new samples based on it.

1. Discriminative models:

Discriminative models learn the decision boundary directly from the data without explicitly modeling the probability distributions of the features and the classes.

1.1 Linear Models:

- Regression -

Linear Regression:

Linear regression is a statistical **method used to model the relationship between a dependent variable (often denoted as y) and one or more independent variables (often denoted as x)**. It assumes that this relationship is linear, meaning that changes in the independent variables are associated with proportional changes in the dependent variable.

In simple linear regression, there is only one independent variable, and the relationship between x and y is modeled using a straight line:

$$y = mx + b$$

Where:

- y is the dependent variable (the variable we want to predict).
- x is the independent variable (the variable used to make predictions).
- m is the slope of the line (representing the relationship between x and y).
- b is the y-intercept (the value of y when $x=0$).

The goal of linear regression is to find the best-fitting line that minimizes the difference between the predicted values (based on the line) and the actual values of the dependent variable. This is often done by minimizing a cost function, such as the mean squared error (MSE), which measures the average squared difference between the predicted and actual values.

Linear regression can be extended to multiple independent variables, known as multiple linear regression, where the relationship between the dependent variable and multiple independent variables is modeled using a linear equation:

$$y = b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$$

Where:

- y is the dependent variable.
- x_1, x_2, \dots, x_n are the independent variables.
- b_0 is the y-intercept.
- b_1, b_2, \dots, b_n are the coefficients (slopes) of the independent variables.

Linear regression is widely used for prediction and inference in various fields, including economics, finance, healthcare, and social sciences, due to its simplicity and interpretability.

Linear regression assumes that observations are independent of each other, which may not hold true for time series data.

Considerations and assumptions:

1. **Linearity:** The relationship between the independent and dependent variables should be approximately linear. This can be assessed using scatter plots or residual plots.
2. **Homoscedasticity:** The variance of the residuals should be constant across all levels of the independent variables. Residual plots can help check for homoscedasticity.
3. **Normality of Residuals:** The residuals should be normally distributed. This can be assessed using a histogram or a Q-Q plot of the residuals. Note: Linear regression does not make any assumptions about the distribution of the independent variables. They can follow any distribution, including normal, uniform, or skewed distributions. But it assumes that the dependent variable follows a normal distribution.
4. **Independence of Errors:** The residuals should be independent of each other. This assumption is often violated in time series data or spatial data, so careful consideration is needed.
5. **No Multicollinearity:** The independent variables should not be highly correlated with each other. Multicollinearity can lead to unstable estimates of the coefficients.
6. **Continuous Variables:** Linear regression assumes that the independent variables are continuous rather than categorical. Categorical variables can be encoded appropriately (e.g., through one-hot encoding) before fitting the model but it is not recommended since it might affect model interpretability.
7. **Outliers and Influential Points:** Outliers and influential points can disproportionately affect the estimated coefficients and should be investigated and, if necessary, addressed.
8. **Interactions and Non-linear Relationships:** Linear regression assumes a linear relationship between the independent and dependent variables. Interaction terms and polynomial terms can be included to capture non-linear relationships or interactions.
9. **Model Interpretability:** Linear regression provides interpretable coefficients that represent the change in the dependent variable for a one-unit change in the independent

variable. Ensure that the interpretation of coefficients makes sense in the context of the data and the research question.

- Cross-validation: Use cross-validation techniques to assess the generalization performance of the model and guard against overfitting.

Locally Weighted Regression:

LWR stands for Locally Weighted Regression. It's a non-parametric regression technique used for modeling the relationship between variables when the underlying relationship is not assumed to be linear.

In LWR, instead of fitting a single global model to the entire dataset, separate local models are fit to different regions of the data. The model gives more weight to data points closer to the point being predicted and less weight to points farther away. This means that the model focuses more on the local structure of the data, hence the name "locally weighted regression."

The key idea behind LWR is to adaptively fit a regression model to the data around each query point, rather than assuming a single global relationship for the entire dataset. This allows LWR to capture complex and nonlinear relationships between variables.

LWR is particularly useful when the relationship between variables varies across different regions of the data or when the data contains outliers or noise. It is commonly used in areas such as **time series analysis**, signal processing, and robotics.

Considerations and assumptions:

- Local Linearity:** LWR assumes that the relationship between the dependent variable and the independent variables is approximately linear within local regions of the data. It fits a separate regression model to each local region.
- Smoothness of Data:** LWR works best when the underlying data is smooth and continuous. It may not perform well with highly discontinuous or noisy data.
- Selection of Kernel Function:** LWR relies on a kernel function to assign weights to neighboring data points. The choice of kernel function can impact the performance of LWR, and different kernel functions may be more suitable for different datasets.
- Bandwidth Selection:** The bandwidth parameter controls the size of the local region over which LWR operates. Selecting an appropriate bandwidth is crucial, as a bandwidth that is too small may lead to overfitting, while a bandwidth that is too large may result in underfitting.
- Data Sparsity:** LWR may not perform well in regions of the data where there are few or no data points. In such cases, the estimates provided by LWR may be unreliable.
- Computational Complexity:** LWR involves fitting a separate regression model to each local region of the data, which can be computationally expensive, especially for large datasets. Efficient algorithms and implementations may be necessary for practical use.
- Robustness to Outliers:** LWR can be sensitive to outliers, especially if they are given high weights by the chosen kernel function. Outlier detection and robust kernel functions may be used to mitigate this issue.
- Interpretability:** While LWR provides flexible and locally adaptive predictions, the interpretation of the model can be challenging, especially if the relationship between the variables varies significantly across different regions of the data.

ARIMA:

ARIMA (AutoRegressive Integrated Moving Average) is a statistical model used for time series forecasting. It belongs to the class of linear models and is specifically designed to capture the temporal dependencies and patterns present in time series data. ARIMA models are widely used for modeling and predicting time series data with stationary or near-stationary properties.

The ARIMA model is composed of three main components:

- AutoRegressive (AR) Component:** This component models the linear relationship between an observation and a linear combination of past observations (autoregressive terms). It captures the effect of previous values in the time series on the current value.
- Integrated (I) Component:** The integrated component represents the differencing of the time series data to make it stationary. It involves subtracting the current observation from the previous observation to remove trends or seasonality.
- Moving Average (MA) Component:** The moving average component models the linear relationship between an observation and a linear combination of past white noise error terms. It captures the effect of past forecast errors on the current observation.

The ARIMA model is defined by three parameters: p, d, and q, which correspond to the order of the AR, I, and MA components, respectively. The ARIMA model can be extended to seasonal ARIMA (SARIMA or SARIMAX) to handle seasonal patterns in the data.

Overall, ARIMA models are versatile and widely used for time series forecasting in various domains such as finance, economics, weather forecasting, and more. They provide a powerful framework for modeling and predicting temporal dependencies in sequential data.

Considerations and assumptions:

- Stationarity:** ARIMA models assume that the time series data is stationary or can be transformed into a stationary series through differencing. Stationarity implies that the statistical properties of the series (such as mean, variance, and autocorrelation) remain constant over time.
- Linearity:** ARIMA models assume a linear relationship between the observations and their lagged values. While non-linear relationships can be captured by transforming the data, ARIMA may not be suitable for highly non-linear relationships.
- Independence of Errors:** ARIMA models assume that the errors (residuals) are independent and identically distributed (i.i.d.). Violations of this assumption may indicate that the model is not adequately capturing the underlying patterns in the data.
- Normality of Errors:** While not strictly required, ARIMA models often assume that the errors follow a normal distribution. Deviations from normality may affect the validity of statistical tests and confidence intervals derived from the model.
- Seasonality and Trends:** ARIMA models are designed to capture non-seasonal patterns in the data. For time series with seasonal or trend components, seasonal ARIMA (SARIMA) models are more appropriate. ARIMA models are designed to capture non-seasonal patterns in the data. If the time series exhibits seasonal patterns, it may require seasonal differencing or seasonal adjustment using methods like seasonal decomposition (e.g., seasonal-trend decomposition using LOESS, STL decomposition).
- Data Quality and Outliers:** ARIMA models are sensitive to outliers and anomalies in the data. Outliers can significantly impact model estimation and prediction accuracy, so it's essential to identify and handle them appropriately.
- Sufficient Data:** ARIMA models require a sufficient amount of data to estimate model parameters accurately. Insufficient data may lead to overfitting or unreliable parameter estimates.
- Model Selection:** Selecting the appropriate orders (p, d, q) for the ARIMA model requires careful consideration and may involve model diagnostics, such as ACF (autocorrelation function) and PACF (partial autocorrelation function) plots, and statistical tests (e.g., AIC, BIC).
- Validation and Testing:** After fitting the ARIMA model, it's crucial to validate its performance on unseen data using techniques like cross-validation or hold-out validation. This helps assess the model's generalization ability and identify potential overfitting.

- Binary Classification -

Logistic regression:

Logistic regression is a type of linear regression model used for binary classification tasks. Despite its name, logistic regression is actually a classification algorithm, not a regression algorithm.

In logistic regression, the algorithm models the probability that a given data point belongs to a particular class using a logistic (also called **sigmoid**) function. The logistic function maps any real-valued input to a value between 0 and 1, representing the probability of the positive class. The output of the logistic function can then be interpreted as the likelihood or probability that the data point belongs to the positive class.

Sigmoid function -> binary classification.

Mathematically, the logistic regression model can be represented as:

$$P(y=1|x) = \frac{1}{1+e^{-(\beta_0+\beta_1x_1+\beta_2x_2+\dots+\beta_nx_n)}}$$

Where:

- $P(y=1|x)$ is the probability that the output y (the dependent variable) is 1 given the input features x .
- $\beta_0, \beta_1, \dots, \beta_n$ are the coefficients (parameters) of the model.
- x_1, x_2, \dots, x_n are the input features.

During training, logistic regression learns the optimal values of the coefficients (β) that maximize the likelihood of the observed data given the model. This is typically

done using optimization algorithms like gradient descent or Newton's method.

Once trained, the logistic regression model can be used to predict the class labels (0 or 1) for new, unseen data points by evaluating the logistic function with the learned coefficients.

Considerations and assumptions:

1. **Input Features:** These features can be numerical, categorical, or a combination of both.
2. **Binary Outcome:** Logistic regression assumes that the dependent variable (the outcome or response variable) is binary, meaning it has only two possible outcomes or classes.
3. **Independence of Observations:** Logistic regression assumes that the observations in the dataset are independent of each other. In other words, the outcome of one observation does not influence the outcome of another observation.
4. **Linearity of Log Odds:** Logistic regression assumes that the relationship between the independent variables (predictors or features) and the log odds of the outcome is linear. This means that the log odds of the outcome is a linear function of the independent variables.
5. **No Multicollinearity:** Logistic regression assumes that there is little or no multicollinearity among the independent variables. Multicollinearity occurs when two or more independent variables are highly correlated with each other, which can lead to unstable estimates of the coefficients.
6. **Large Sample Size:** Logistic regression performs best with a relatively large sample size. Small sample sizes can lead to unstable estimates and unreliable inference.
7. **No Outliers:** Outliers in the data can disproportionately influence the estimated coefficients and affect the performance of logistic regression. It's essential to detect and handle outliers appropriately.
8. **No Perfect Separation:** Logistic regression assumes that there is no perfect separation of the outcome variable by the independent variables. Perfect separation occurs when the values of the independent variables perfectly predict the outcome, leading to infinite parameter estimates.
9. **Correct Specification of Model:** Logistic regression assumes that the model is correctly specified, meaning that all relevant independent variables are included in the model and that the functional form of the model is appropriate for the data.
10. **No Homoscedasticity Assumption:** Unlike linear regression, logistic regression does not assume homoscedasticity (constant variance of residuals) because it models probabilities rather than raw outcomes.
11. **Interpretation of Coefficients:** The coefficients estimated by logistic regression represent the change in the log odds of the outcome for a one-unit change in the corresponding independent variable, assuming all other variables are held constant.

- Multiclass Classification -

Softmax regression or Multiclass Logistic regression:

Softmax regression, also known as multinomial logistic regression, is a type of regression model used for multiclass classification tasks. It extends logistic regression, which is designed for binary classification, to handle multiple classes.

In softmax regression, the goal is to predict the probability that a given data point belongs to each of the possible classes. The model assigns a probability to each class, and the sum of these probabilities across all classes is equal to 1. This makes softmax regression suitable for scenarios where there are more than two mutually exclusive classes.

Considerations and assumptions are the same as Logistic Regression with the difference that the output is not binary but categorical.

1.2 Support Vector Machines:

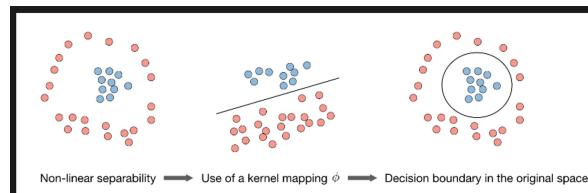
- Binary Classification -

SVM:

SVMs are particularly powerful for binary classification tasks, where the goal is to separate data points into two classes by finding the optimal hyperplane that maximizes the margin between the classes.

Here's a step-by-step explanation of how SVMs work for binary classification:

1. **Data Representation:** SVMs work with labeled training data, where each data point is represented by a feature vector x and a corresponding class label y indicating the class to which the data point belongs.
2. **Feature Space:** SVMs operate in a high-dimensional feature space, where each dimension corresponds to a feature of the data. The goal is to find a hyperplane that best separates the data points into two classes.
3. **Optimal Hyperplane:** The optimal hyperplane is the one that maximizes the margin, which is the distance between the hyperplane and the nearest data points (called support vectors) from each class. The margin represents the margin of confidence or separation between the classes.
4. **Maximizing Margin:** SVMs find the optimal hyperplane by solving an optimization problem that involves maximizing the margin while minimizing the classification error. This optimization problem is typically formulated as a quadratic programming problem.
5. **Kernel Trick:** In cases where the data is not linearly separable in the original feature space, SVMs can use the kernel trick to implicitly map the input features into a higher-dimensional space, where the data becomes linearly separable. Common kernel functions include linear, polynomial, radial basis function (RBF), and sigmoid kernels.
6. **Soft Margin SVM:** In situations where the data is not perfectly separable, SVMs can use a soft margin approach, allowing for some misclassification of data points. This is achieved by introducing a slack variable that penalizes misclassified points.
7. **Decision Boundary:** Once the optimal hyperplane is found, it serves as the decision boundary for classifying new, unseen data points. Data points on one side of the hyperplane are assigned to one class, while data points on the other side are assigned to the other class.



Considerations and assumptions:

1. **Linear Separability:** SVM assumes that the data is linearly separable, meaning that there exists at least one hyperplane that can perfectly separate the data points of different classes. When the data is not linearly separable, SVM may still work, but it may require additional techniques such as kernel functions or soft margin SVM.
2. **Margin Maximization:** SVM aims to find the hyperplane that maximizes the margin between the classes. The margin is the distance between the hyperplane and the closest data points (support vectors) from each class. This assumption implies that SVM seeks the most robust decision boundary that generalizes well to unseen data.
3. **Noisy Data:** SVM is sensitive to noise, outliers, and mislabeled data points, as these can affect the position and orientation of the optimal hyperplane. Preprocessing steps such as data cleaning and outlier detection may be necessary to improve the performance of SVM, especially in noisy datasets.
4. **Feature Scaling:** SVM performance can be sensitive to the scale of the features. It is generally recommended to scale or normalize the features to have zero mean and unit variance before training the SVM model. This ensures that all features contribute equally to the optimization process.
5. **Regularization Parameter (C):** SVM includes a regularization parameter C that controls the trade-off between maximizing the margin and minimizing the classification error. Higher values of C result in a narrower margin and may lead to overfitting, while lower values of C may lead to underfitting. Choosing an appropriate value of C is essential for optimizing the SVM model's performance.
6. **Kernel Selection:** When the data is not linearly separable, SVM can use kernel functions to implicitly map the input features into a higher-dimensional space where the data becomes separable. The choice of kernel function (e.g., linear, polynomial, radial basis function) and its parameters can significantly impact the performance of the SVM model.

- Multiclass Classification -

SVM:

Support Vector Machines (SVMs) can be extended to handle multiple-class classification tasks using two common approaches: one-vs-one (OvO) and one-vs-rest (OvR) strategies.

1. One-vs-One (OvO):
 - o In the one-vs-one approach, a separate binary SVM classifier is trained for each pair of classes.
 - o For K classes, $\binom{K}{2}$ binary classifiers are trained.
 - o During prediction, each classifier votes for one of the two classes in the binary classification problem.
 - o The class with the most votes across all classifiers is assigned as the predicted class.
2. One-vs-Rest (OvR):
 - o In the one-vs-rest approach, K binary classifiers are trained, each representing one class versus the rest.
 - o For each class, a binary classifier is trained to distinguish that class from all other classes combined.
 - o During prediction, each classifier predicts the probability of belonging to the class it represents.
 - o The class with the highest probability across all classifiers is assigned as the predicted class.

Both OvO and OvR strategies have their advantages and disadvantages. OvO requires training $\binom{K}{2}$ classifiers, which can be computationally expensive for large K . However, it tends to be more robust to imbalanced datasets. OvR, on the other hand, requires training only K classifiers, which can be more efficient for large K , but it may suffer from imbalanced class distributions.

In practice, the choice between OvO and OvR often depends on the specific characteristics of the dataset and the computational resources available. Both approaches are widely used and can effectively handle multiple-class classification tasks with SVMs.

Same assumptions and considerations as in binary classification.

- Regression -

Support Vector Regression (SVR):

Here's how SVR works:

1. **Data Representation:** Like in classification tasks, SVR also works with labeled training data, where each data point is represented by a feature vector x and a corresponding target value y .
2. **Linear Regression:** SVR aims to find a hyperplane that best fits the data points while maximizing the margin within a certain threshold (ϵ). This hyperplane is called the ϵ -insensitive tube.
3. **Loss Function:** SVR minimizes a loss function that penalizes errors outside the ϵ -insensitive tube. Data points within the tube or within a distance of ϵ from the predicted values incur no penalty, while points outside the tube incur a penalty proportional to their distance from the tube boundary.
4. **Kernel Trick:** As in classification tasks, SVR can use the kernel trick to implicitly map the input features into a higher-dimensional space. This allows SVR to capture nonlinear relationships between the features and the target variable.
5. **Regularization:** SVR includes a regularization term that penalizes large coefficients in the linear regression model. This helps prevent overfitting and improves the generalization performance of the model.
6. **Prediction:** Once the SVR model is trained, it can be used to predict the target values for new, unseen data points. The predicted values are obtained by evaluating the linear function learned by the model on the input features.
SVR is particularly useful for regression tasks where the relationship between the features and the target variable is nonlinear or where the data contains outliers. It provides a flexible and robust approach to regression modeling and can handle a wide range of datasets and problem domains.

Considerations and assumptions:

1. **Linear Relationship:** SVR assumes that there is a linear relationship between the input features and the target variable. However, SVR can capture nonlinear relationships through the use of kernel functions.
2. **ϵ -Insensitive Loss Function:** SVR minimizes a loss function that penalizes errors outside an ϵ -insensitive tube around the predicted values. Data points within the tube or within a distance of ϵ from the predicted values incur no penalty, while points outside the tube incur a penalty proportional to their distance from the tube boundary.
3. **Kernel Selection:** As in classification tasks, SVR can use kernel functions to implicitly map the input features into a higher-dimensional space. The choice of kernel function (e.g., linear, polynomial, radial basis function) and its parameters can significantly impact the performance of the SVR model.
4. **Regularization:** SVR includes a regularization parameter (often denoted as C) that controls the trade-off between maximizing the margin and minimizing the error. Higher values of C result in a narrower ϵ -insensitive tube and may lead to overfitting, while lower values of C may lead to underfitting. Choosing an appropriate value of C is essential for optimizing the SVR model's performance.
5. **Feature Scaling:** SVR performance can be sensitive to the scale of the features. It is generally recommended to scale or normalize the features to have zero mean and unit variance before training the SVR model. This ensures that all features contribute equally to the optimization process.
6. **Robustness to Outliers:** SVR is robust to outliers to some extent due to the ϵ -insensitive loss function. Outliers that fall within the ϵ -insensitive tube have minimal impact on the model's performance. However, outliers outside the tube can still affect the model's fit, and preprocessing steps such as outlier detection and removal may be necessary.

1.3 Tree-based Models:

Tree-based models can be used both for Regression and for Classification problems.

Classification and Regression Trees (CART):

CART is a flexible and interpretable algorithm that can handle both categorical and continuous features and can capture nonlinear relationships in the data. It is widely used in practice due to its simplicity, scalability, and ability to generate understandable decision rules.

CART, which stands for Classification and Regression Trees, is a type of tree-based model used for both classification and regression tasks. CART constructs a binary tree recursively by partitioning the feature space into regions, with each partition representing a decision rule based on the values of the input features.

Here's how CART works:

1. **Binary Tree Structure:** CART builds a binary tree structure, where each node in the tree represents a decision rule based on a feature and a threshold value. The tree starts with a root node that contains the entire dataset, and at each node, CART selects the best feature and threshold to split the data into two child nodes.
2. **Splitting Criterion:** CART uses a splitting criterion to determine the best feature and threshold for each split. For classification tasks, common splitting criteria include Gini impurity and entropy, which measure the purity or homogeneity of the classes within each partition. For regression tasks, the mean squared error (MSE) or mean absolute error (MAE) may be used as the splitting criterion.
3. **Recursive Partitioning:** CART recursively partitions the data into subsets based on the selected feature and threshold at each node. The process continues until a stopping criterion is met, such as reaching a maximum tree depth, having fewer than a minimum number of data points in a node, or when no further improvement in purity or error reduction can be achieved.
4. **Leaf Nodes:** Once the tree is fully grown, each terminal node (or leaf node) contains a subset of the data, and predictions are made by aggregating the target variable values within each leaf node. For classification tasks, the majority class within the leaf node is assigned as the predicted class label. For regression tasks, the mean (or median) value of the target variable within the leaf node is used as the predicted value.
5. **Pruning (Optional):** After the tree is built, it may be pruned to prevent overfitting and improve generalization performance. Pruning involves removing branches of the tree that do not contribute significantly to reducing impurity or error on a separate validation dataset.

Considerations and assumptions:

1. **Feature Importance:** CART provides a measure of feature importance based on how much each feature contributes to reducing impurity in the decision tree. Features that result in large reductions in impurity are considered more important.
2. **Nonlinear Relationships:** CART assumes that the decision boundary can be represented by axis-parallel splits in the feature space. While CART can capture nonlinear relationships through recursive partitioning, it may struggle with complex nonlinear relationships compared to other algorithms like neural networks or kernel-based methods.

3. **Handling Missing Values:** CART can handle missing values in the dataset by either imputing them or treating them as a separate category during the splitting process.
4. **Robustness to Outliers:** CART is generally robust to outliers, as it partitions the feature space based on relative rankings rather than absolute distances.

Random Forest:

Random Forest is a versatile and powerful machine learning algorithm that is used for both classification and regression tasks. It belongs to the family of ensemble learning methods, which combine the predictions of multiple individual models to improve overall performance.

Here's how Random Forest works:

1. **Ensemble of Decision Trees:** Random Forest is composed of a collection of decision trees, where each tree is trained independently on a random subset of the training data and features.
2. **Random Subsampling:** During training, Random Forest randomly selects a subset of the training data with replacement (bootstrapping) to create multiple bootstrap samples. Each decision tree is trained on one of these bootstrap samples.
3. **Feature Randomization:** In addition to subsampling the data, Random Forest also randomly selects a subset of features at each node of the decision tree. This process helps to decorrelate the trees and introduces diversity among them.
4. **Tree Building:** Each decision tree in the Random Forest is grown using a process similar to CART (Classification and Regression Trees). The tree is recursively built by selecting the best feature and threshold for splitting the data at each node based on a splitting criterion (such as Gini impurity or entropy).
5. **Voting for Classification:** For classification tasks, the predictions of individual trees are aggregated using a majority voting scheme. The class that receives the most votes among all trees is assigned as the final prediction.
6. **Averaging for Regression:** For regression tasks, the predictions of individual trees are averaged to obtain the final prediction. This averaging process helps to reduce variance and improve the stability of the predictions.
7. **Bagging and Bootstrap Aggregating:** Random Forest leverages the principles of bagging (bootstrap aggregating) to reduce overfitting and improve generalization performance. By training multiple trees on different subsets of data, Random Forest can capture different aspects of the underlying data distribution.
8. **Robustness to Overfitting:** Random Forest is known for its robustness to overfitting, especially when compared to individual decision trees. The ensemble nature of Random Forest helps to smooth out noise and reduce variance, leading to more reliable predictions.
9. **Feature Importance:** Random Forest provides a measure of feature importance based on how much each feature contributes to reducing impurity or error across all trees in the forest. Features that result in larger reductions in impurity or error are considered more important.

Considerations and assumptions:

1. **Ensemble of Decision Trees:** Random Forest is an ensemble learning method composed of multiple decision trees. As such, it inherits some considerations from decision trees, such as potential overfitting if the trees are allowed to grow too deep.
2. **Randomeess:** Random Forest introduces randomness by bootstrapping the data and randomly selecting subsets of features at each node. While this randomness helps to decorrelate the trees and improve generalization, it can also lead to variability in model performance.
3. **Model Interpretability:** While Random Forest provides feature importance scores, the ensemble nature of the algorithm can make it more challenging to interpret individual trees compared to standalone decision trees.
4. **Computational Cost:** Training a Random Forest can be computationally expensive, especially for large datasets or when using a large number of trees in the ensemble. However, it can often be parallelized and scaled to leverage modern computing resources effectively.
5. **Hyperparameter Tuning:** Random Forest has several hyperparameters that can impact its performance, such as the number of trees in the ensemble, the maximum depth of the trees, and the number of features considered at each split. Careful hyperparameter tuning is essential to optimize model performance.
6. **Handling Imbalanced Data:** While Random Forest can handle imbalanced datasets to some extent, it may still benefit from techniques such as class weighting or resampling to improve performance on minority classes.
7. **Feature Importance Interpretation:** While Random Forest provides feature importance scores, it's essential to interpret them with caution. The importance scores are relative and may not always reflect true causal relationships between features and the target variable.
8. **Nonlinear Relationships:** Random Forest is capable of capturing complex nonlinear relationships in the data. However, it may struggle with capturing certain types of nonlinearities, especially if they are highly intricate or require a large number of trees to model accurately.

- Boosting models -

Boosting is a specific type of ensemble learning technique where models are trained sequentially, with each model (or learner) focusing on the mistakes made by the previous models. In boosting, each model is trained to correct the errors of its predecessor, with more emphasis placed on the misclassified data points. The final prediction is typically made by combining the predictions of all models using a weighted voting scheme, where models with higher performance contribute more to the final prediction.

Adaboost:

AdaBoost, short for Adaptive Boosting, is a popular ensemble learning algorithm used for classification tasks. It operates by combining the predictions of multiple individual models, typically weak learners (e.g., decision trees with limited depth), to create a strong learner with improved predictive performance.

Here's how AdaBoost works:

1. **Initialization:** AdaBoost assigns equal weights to each sample in the training dataset initially.
2. **Sequential Training:** AdaBoost trains a series of weak learners sequentially. In each iteration:
 - a. **Weighted Training:** It fits a weak learner to the training data, with the sample weights adjusted to focus more on the samples that were misclassified in the previous iteration.
 - b. **Classifier Weight:** AdaBoost assigns a weight to each weak learner based on its performance (e.g., accuracy) on the training data. Better-performing weak learners receive higher weights, indicating their importance in the ensemble.
 - c. **Predictive Power:** Weak learners with higher weights contribute more to the final prediction, while those with lower weights have less influence.
3. **Combining Predictions:** AdaBoost combines the predictions of all weak learners using a weighted majority voting scheme. The final prediction is determined by the weighted sum of the predictions of individual weak learners.
4. **Boosting Iterations:** AdaBoost continues this process for a predefined number of iterations or until a specified level of performance is reached. Each iteration focuses on improving the classification accuracy of the previously misclassified samples, leading to iterative refinement of the model.

Asumptions and Considerations:

1. **Weak Learner Requirement:** AdaBoost works best when the weak learners used in the ensemble have performance better than random guessing, but they don't need to be overly complex or highly accurate models.
2. **Data Quality:** AdaBoost can be sensitive to noisy data and outliers since it tries to correct misclassifications in subsequent iterations. Outliers or noisy data points may receive undue focus in later iterations, leading to suboptimal performance.
3. **Performance:** The performance of AdaBoost heavily depends on the performance of the weak learners. If the weak learners consistently perform poorly, AdaBoost may not be able to achieve good results.
4. **Training Time:** AdaBoost can be computationally expensive, especially when using complex weak learners or training on large datasets. Each iteration of AdaBoost requires sequentially training a series of weak learners, which can increase training time.
5. **Number of Iterations:** The number of boosting iterations (or weak learners) in AdaBoost is a hyperparameter that needs to be tuned. Too few iterations may result in underfitting, while too many iterations may lead to overfitting.
6. **Imbalanced Data:** AdaBoost can struggle with imbalanced datasets, where one class is significantly more prevalent than others. In such cases, it may be necessary to adjust class weights or use sampling techniques to ensure balanced performance.
7. **Model Interpretability:** While the individual weak learners in AdaBoost are often simple models (e.g., decision stumps), the final ensemble model may be less interpretable, especially when using a large number of weak learners.
8. **Overfitting:** AdaBoost can be susceptible to overfitting, particularly if the weak learners are too complex or the number of iterations is too high. Regularization techniques, cross-validation, or early stopping can help mitigate overfitting.
9. **Non-Linear Relationships:** AdaBoost can capture non-linear relationships in the data, but its performance may degrade if the relationships are highly complex or require a large number of weak learners to model accurately.

10. **Robustness:** AdaBoost is generally robust to noise and outliers, as long as they don't dominate the training process. The adaptive nature of AdaBoost helps it focus on correcting misclassifications, rather than being overly influenced by outliers.

XGBoost:

XGBoost, short for eXtreme Gradient Boosting, is an advanced and highly efficient implementation of gradient boosting machines. It's an ensemble learning method that builds a strong predictive model by combining the predictions of multiple individual models, typically decision trees, in a sequential manner.

Here's how XGBoost works:

1. **Gradient Boosting Framework:** XGBoost follows the gradient boosting framework, where models are trained sequentially to correct the errors of the previous models.
2. **Objective Function:** XGBoost uses a customizable objective function to quantify the difference between the predicted and actual values. This objective function is optimized during training to minimize prediction errors.
3. **Gradient Descent Optimization:** XGBoost employs gradient descent optimization techniques to minimize the objective function. It iteratively updates the model parameters (e.g., tree structure and leaf values) in the direction that reduces the objective function's value the most.
4. **Regularization:** XGBoost includes built-in regularization techniques to prevent overfitting and improve generalization performance. Regularization terms are added to the objective function to penalize complex models or large parameter values.
5. **Tree Ensemble:** XGBoost builds an ensemble of decision trees during training, with each tree trained to correct the errors of the previous trees. The trees are typically shallow (with limited depth) to avoid overfitting and improve computational efficiency.
6. **Feature Importance:** XGBoost provides a measure of feature importance based on how much each feature contributes to reducing the objective function's value across all trees in the ensemble. Features that result in larger reductions in the objective function are considered more important.
7. **Parallel and Distributed Computing:** XGBoost is designed for scalability and efficiency, with support for parallel and distributed computing. It can leverage multiple CPU cores and distributed computing frameworks (e.g., Spark) to train models on large datasets efficiently.
8. **Wide Range of Applications:** XGBoost is widely used in various machine learning tasks, including classification, regression, ranking, and recommendation systems. It has proven to be highly effective in competitions and real-world applications across different domains.

Asumptions and Considerations:

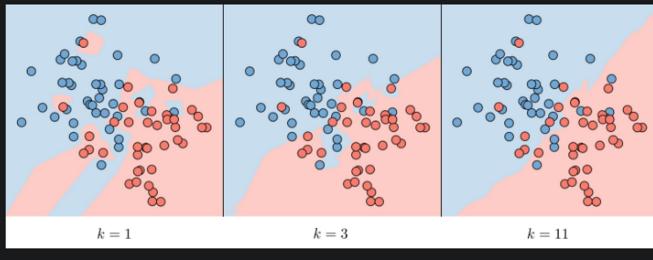
1. **Tree Structure:** XGBoost builds an ensemble of decision trees, so it assumes that the underlying data can be effectively partitioned into regions based on the feature space. This may not be suitable for data with complex or nonlinear relationships that cannot be adequately captured by decision trees.
2. **Gradient Optimization:** XGBoost relies on gradient-based optimization techniques to minimize the objective function. It assumes that the objective function is differentiable and continuous. Non-differentiable or discontinuous objective functions may lead to suboptimal results.
3. **Overfitting:** While XGBoost includes regularization techniques to mitigate overfitting, it's still possible to overfit the model, especially if the parameters are not properly tuned or if the dataset is noisy. Regularization parameters such as the learning rate and the depth of the trees should be carefully adjusted to prevent overfitting.
4. **Interpretability:** XGBoost builds complex ensemble models composed of multiple decision trees, which may be less interpretable compared to simpler models like linear regression. While XGBoost provides feature importance scores, interpreting the exact decision-making process of the ensemble can be challenging.
5. **Computational Resources:** XGBoost can be computationally intensive, especially when training large ensembles on large datasets. Users need to consider the computational resources available and potentially use distributed computing frameworks to train XGBoost models efficiently.
6. **Hyperparameter Tuning:** XGBoost has several hyperparameters that need to be tuned to achieve optimal performance. Finding the right combination of hyperparameters can require significant experimentation and computational resources.
7. **Imbalanced Data:** XGBoost may struggle with imbalanced datasets, where one class is significantly more prevalent than others. Techniques such as class weighting or resampling may be necessary to address class imbalance effectively.
8. **Feature Scaling:** While XGBoost is not as sensitive to feature scaling as some other algorithms, normalizing or standardizing features can still improve convergence speed and model performance, especially when using gradient-based optimization.
9. **Memory Usage:** XGBoost stores the entire ensemble of decision trees in memory during training, so memory usage can become a concern for large datasets or models with many trees. Users should monitor memory usage and potentially reduce model complexity or use memory-efficient settings if memory constraints are an issue.

1.3 non-parametric Models:

K-NN:

□ ***k*-nearest neighbors** — The *k*-nearest neighbors algorithm, commonly known as *k*-NN, is a non-parametric approach where the response of a data point is determined by the nature of its *k* neighbors from the training set. It can be used in both classification and regression settings.

*Remark: the higher the parameter *k*, the higher the bias, and the lower the parameter *k*, the higher the variance.*



Considerations and assumptions:

1. **Distance Metric:** k-NN relies on a distance metric (e.g., Euclidean distance) to measure the similarity between data points. The choice of distance metric can significantly impact the algorithm's performance, so it's crucial to select an appropriate metric based on the nature of the data.
2. **Feature Scaling:** Since k-NN calculates distances between data points, features with larger scales can dominate the distance calculations. It's often necessary to scale or normalize the features to ensure that all features contribute equally to the distance computations.
3. ***k*-Value Selection:** The choice of the *k*-value, the number of nearest neighbors considered, can influence the model's performance. A small *k*-value may lead to overfitting, while a large *k*-value may result in underfitting. Cross-validation or other validation techniques can help determine the optimal *k*-value for a given dataset.
4. **Curse of Dimensionality:** In high-dimensional spaces, the notion of distance becomes less meaningful, which can affect the performance of k-NN. As the number of dimensions increases, the distance between nearest neighbors may become less discriminative. Dimensionality reduction techniques or feature selection methods may be necessary to mitigate the curse of dimensionality.
5. **Computational Complexity:** The prediction time complexity of k-NN grows linearly with the size of the training data, as the algorithm needs to calculate distances to all training instances. For large datasets, this can result in high computational costs during inference, especially if the number of dimensions or the *k*-value is large.
6. **Imbalanced Data:** k-NN can be sensitive to imbalanced datasets, where one class is significantly more prevalent than others. In such cases, the majority class may dominate the prediction for new instances, leading to biased results. Techniques such as class weighting or resampling may be necessary to address class imbalance effectively.
7. **Local Structure:** k-NN assumes that data points that are close to each other in feature space have similar target values. Therefore, it works well when the underlying data has local structure or clusters. However, it may struggle with data that lacks clear local structure or has complex decision boundaries.
8. **Storage Requirements:** Unlike many other machine learning algorithms, k-NN requires storing the entire training dataset during inference, which can lead to high memory requirements, especially for large datasets.

1. Generative Learning models:

Generative models, on the other hand, explicitly model the joint probability distribution of the features and the classes. They learn the underlying data-generating process and can generate new samples from the learned distribution.

Gaussian Discriminant Analysis (GDA) and Naive Bayes are primarily used as classification (both binary and multiclass) algorithms and are not typically employed for generating new data directly. However, they can indirectly contribute to data generation through techniques like sampling from the learned probability distributions.

Gaussian Discriminant Analysis (GDA):

Gaussian Discriminant Analysis (GDA) is a probabilistic generative model used primarily for classification tasks. Here's a breakdown of how it works:

1. **Assumption of Gaussian Distribution:** GDA assumes that the features of each class come from a Gaussian (normal) distribution. This means that if you plot the distribution of each feature for each class, it will resemble a bell curve.
2. **Modeling Class Conditional Distributions:** GDA models the probability distribution of each feature given the class label. Mathematically, this is represented as $P(x|y)$, where x is the feature vector and y is the class label. For each class y , GDA estimates the parameters (mean and covariance) of the Gaussian distribution that best fits the feature vectors belonging to that class.
3. **Bayes' Theorem:** GDA uses Bayes' theorem to compute the posterior probability of each class given a feature vector. Bayes' theorem states that the posterior probability is proportional to the product of the prior probability of the class and the likelihood of the observed features given that class.
4. **Decision Rule:** To classify a new data point, GDA computes the posterior probability of each class given the observed features and selects the class with the highest posterior probability. In other words, it assigns the class label that maximizes the posterior probability.
5. **Parameter Estimation:** GDA requires estimating the parameters of the Gaussian distributions for each class. This involves calculating the mean vector and covariance matrix for each class based on the training data.
6. **Model Interpretability:** GDA provides insights into the distribution of features within each class, making it interpretable and allowing users to understand how different features contribute to classification decisions.
7. **Linear Decision Boundary:** When the covariance matrices of the Gaussian distributions are assumed to be equal for all classes (i.e., shared covariance), GDA results in a linear decision boundary. However, if the covariance matrices are different, the decision boundary can become quadratic or nonlinear.

Considerations and assumptions:

1. **Gaussian Distribution:** GDA assumes that the features within each class follow a Gaussian (normal) distribution. If the features don't exhibit a Gaussian distribution, GDA may not perform well.
2. **Homoscedasticity:** GDA assumes that the covariance matrix of each class is identical. This means that the spread of the data points around the mean is consistent across all classes. If the covariance matrices are significantly different, GDA may not accurately capture the underlying data distribution.
3. **Feature Independence:** GDA assumes that the features are statistically independent within each class. If features are highly correlated, GDA may not perform well as it doesn't take into account the interdependencies between features.
4. **Large Training Data:** GDA performs best with a sufficient amount of training data. Insufficient data can lead to inaccurate parameter estimates and affect the performance of the classifier.
5. **Linear Decision Boundaries:** When the covariance matrices are assumed to be equal for all classes (i.e., shared covariance), GDA results in linear decision boundaries. If the data has nonlinear decision boundaries, GDA may not be suitable without modifications.
6. **Imbalanced Classes:** GDA may not perform well when classes are heavily imbalanced, meaning that one class has significantly more samples than others. It can lead to biased parameter estimates and affect the classifier's ability to generalize.
7. **Sensitive to Outliers:** GDA is sensitive to outliers, as they can significantly affect the estimation of the mean and covariance parameters. Outliers should be handled appropriately to prevent them from skewing the model's performance.

Naive Bayes:

Naive Bayes is a simple yet powerful probabilistic classifier based on Bayes' theorem with the "naive" assumption of feature independence. Here's how it works:

1. **Bayes' Theorem:** At its core, Naive Bayes relies on Bayes' theorem, which calculates the probability of a hypothesis (class label) given the observed evidence (features). Mathematically, it is expressed as:
$$P(y|x) = \frac{P(x|y) \times P(y)}{P(x)}$$
where:
 - o $P(y|x)$ is the posterior probability of class y given features x ,
 - o $P(x|y)$ is the likelihood of observing features x given class y ,
 - o $P(y)$ is the prior probability of class y ,
 - o $P(x)$ is the probability of observing features x (also known as evidence).
2. **Naive Assumption:** Naive Bayes assumes that the features are conditionally independent given the class label. This means that the presence of one feature is independent of the presence of any other feature, given the class label. While this assumption is often violated in real-world data, Naive Bayes can still perform well in practice, especially with large datasets.
3. **Parameter Estimation:** Naive Bayes requires estimating the likelihood and prior probabilities from the training data. The likelihood $P(x|y)$ represents the probability distribution of features given each class, while the prior $P(y)$ represents the probability of each class in the dataset.
4. **Classification:** To classify a new instance, Naive Bayes calculates the posterior probability for each class using Bayes' theorem and selects the class with the highest probability as the predicted class label.
5. **Types of Naive Bayes:** There are different variants of Naive Bayes, including:
 - o Gaussian Naive Bayes: Assumes that features follow a Gaussian distribution.
 - o Multinomial Naive Bayes: Suitable for discrete features (e.g., word counts in text classification).
 - o Bernoulli Naive Bayes: Appropriate for binary features (e.g., presence or absence of a feature).
 - o Categorical Naive Bayes: Suitable for categorical features with more than two levels.
6. **Advantages:**
 - o Simple and easy to implement.
 - o Requires a small amount of training data.
 - o Performs well in many real-world scenarios, especially with text classification tasks.
7. **Limitations:**
 - o The "naive" assumption of feature independence may not hold true in practice.
 - o Requires careful preprocessing of data, especially with categorical features and handling of missing values.
 - o Can be sensitive to the presence of irrelevant features.

Considerations and limitations:

1. **Feature Independence:** The algorithm assumes that features are conditionally independent given the class label. This means that the presence of one feature is independent of the presence of any other feature, given the class. While this assumption simplifies the model, it may not hold true in all real-world scenarios.
2. **Gaussian Distribution (for Gaussian Naive Bayes):** If using Gaussian Naive Bayes, it assumes that the features within each class follow a Gaussian (normal) distribution. If the features do not exhibit a Gaussian distribution, the model's performance may be compromised.
3. **Discrete or Binary Features (for Multinomial and Bernoulli Naive Bayes):** For Multinomial and Bernoulli Naive Bayes, the algorithm assumes that the features are discrete or binary, respectively. These variants are commonly used in text classification tasks where features represent word counts or binary indicators of word presence.
4. **Handling Missing Values:** Naive Bayes does not explicitly handle missing values. Therefore, it's important to handle missing data appropriately before training the model, such as imputation or removal of instances with missing values.
5. **Categorical Features (for Categorical Naive Bayes):** Categorical Naive Bayes assumes that features are categorical with more than two levels. It's essential to ensure that categorical features are properly encoded before training the model.

- Data Sparsity:** Naive Bayes may perform poorly with rare or unseen feature combinations, especially in high-dimensional spaces. Techniques such as smoothing (e.g., Laplace smoothing) can help alleviate the issue of data sparsity.
- Balanced Classes:** While Naive Bayes can handle imbalanced datasets, it may produce biased predictions towards the majority class, especially when classes are highly imbalanced. Techniques such as class weighting or resampling may be necessary to address class imbalance effectively.
- Model Interpretability:** Naive Bayes provides interpretable results by estimating probabilities for each class given the observed features. This allows users to understand the model's decision-making process and gain insights into the relative importance of different features.

Evaluation metrics for Regression:

- Mean Absolute Error (MAE):** Represents the average absolute difference between the predicted and true values. It gives an idea of the magnitude of the errors without considering their direction.
- Mean Squared Error (MSE):** Measures the average of the squared differences between the predicted and true values. It penalizes larger errors more heavily than MAE and is sensitive to outliers.
- Root Mean Squared Error (RMSE):** Square root of the MSE, providing an interpretable measure in the same units as the target variable. RMSE is commonly used as it has the same units as the dependent variable, making it easier to interpret.
- R-squared (R^2):** Represents the proportion of variance in the dependent variable that is predictable from the independent variables. R-squared ranges from 0 to 1, where a higher value indicates a better fit of the model to the data.
- Mean Absolute Percentage Error (MAPE):** Measures the average percentage difference between the predicted and true values. It provides a relative measure of the accuracy of the model and is particularly useful when dealing with data of different scales.
- Median Absolute Error:** Computes the median of the absolute differences between the predicted and true values. It is less sensitive to outliers compared to mean-based metrics like MAE and MSE.
- Explained Variance Score:** Measures the proportion of variance explained by the model relative to the total variance in the data. A score of 1 indicates that the model perfectly predicts the target variable.

Evaluation metrics for Classification:

- Accuracy:** Measures the proportion of correctly classified instances out of total instances. It is a straightforward metric but may not be suitable for imbalanced datasets.
- Precision:** Measures the proportion of true positive predictions out of all positive predictions. It indicates the model's ability to avoid false positives.
- Recall (Sensitivity):** Measures the proportion of true positive predictions out of all actual positives. It indicates the model's ability to capture all positive instances.
- F1 Score:** Harmonic mean of precision and recall, providing a balance between them. It is useful when there is an imbalance between the classes.
- ROC-AUC (Receiver Operating Characteristic - Area Under Curve):** Measures the area under the ROC curve, which plots the true positive rate against the false positive rate. It evaluates the model's ability to distinguish between classes across different thresholds.
- Precision-Recall Curve:** Plots precision against recall, providing insights into the trade-off between them. It is particularly useful when dealing with imbalanced datasets.
- Confusion Matrix:** Tabulates the number of true positive, true negative, false positive, and false negative predictions. It provides a detailed breakdown of the model's performance across different classes.

Summary:

Algorithm	Type of Problems Solved	Discriminative /Generative	Type of Input	Need Normalization	Sensitivity to Outliers	Assumes Linear Relationship	Feature Independence Importance	Model Interpretability	Computational Cost	Capture Nonlinear Relationships	Sensitivity to Imbalanced Data	Complexity	Keywords
Linear Regression	Regression	Discriminative	Continuous	Yes	Yes	Yes	Yes + Homoscedasticity of residuals	High	Low	No	No	Low	Finds the straight line that models the relationship between x and y.
Locally Weighted Regression	Regression	Discriminative	Continuous	Yes	Yes	No (but expects Local Linearity)	Yes	Low	Medium	Yes	No	Low	Separate Linear Reg models for different data regions.
Logistic Regression	Binary Classification	Discriminative	Continuous/Categorical	Yes	Yes	No	Yes	Medium	Low	No	No	Low (best for large datasets)	Sigmoid function. With optimization (like gradient descent or newton's method)
Softmax Regression	Multiclass Classification	Discriminative	Continuous/Categorical	Yes	Yes	No	Yes	Medium	Low	No	No	Low (best for large datasets)	LogR but with softmax function.
SVM	Binary/Multiclass Classification (one vs one or one vs rest + voting)	Discriminative	Continuous	Yes	Yes	No	Yes	Medium	Medium	Yes	No	Medium	Finds the hyperplane that separates classes. high-dimensional feature space
SVR	Regression	Discriminative	Continuous	Yes	Yes	No	Yes	Medium	Medium	Yes	No	Medium	Finds the hyperplane that models the relationship between x and y. high dimension.
CART	Classification (majority class in leaf is pred val)/Regression (predicted value is the mean of the target variable at the leaf node) + Feature importance	Discriminative	Continuous/Categorical	No	No	No (but struggles with very complex nonlinear relationships)	Yes	High	Low	Yes	No	Medium	Binary tree structure. Rules and thresholds for each node decided by splitting criterion. Stopping splitting criterion. Pruning can prevent overfitting.
Random Forest	Classification (Voting) /Regression (Averaging) + Feature importance	Discriminative	Continuous/Categorical	No	Yes	No	Yes	Medium	Medium (can be parallelized)	Yes	No (but benefits if balanced)	High	Ensembles CARTs trained independently with subsets of training data. Robust to overfitting if trees don't grow too deep.
XGBoost	Classification/Regression (objective func like MSE or MAE)	Discriminative	Continuous/Categorical	No (but benefits if very different scales)	Yes	No (but struggles with very complex nonlinear relationships)	Yes	Medium	High	Yes	Yes	High	Ensembles sequentially (usually trees). Models focus on solving mistakes from previous models. Uses gradient descent optimization to minimize an objective function in each model. Final model produces final preds.
AdaBoost	Classification/Regression	Discriminative	Continuous/Categorical	No	Yes	No	Yes	Medium	High	Yes	Medium	High	Ensembles sequentially (usually trees). Models focus on solving mistakes from previous models. Weighted voting scheme. Weight is based on performance.
k-NN	Classification/Regression	Discriminative	Continuous/Categorical	Yes	No	No	No	High	High	Yes	Yes	High (avoid using many features, it increases dimensionality)	Uses distance metric (e.g. euclidean) to measure similarity between dt points. Pred is determined by 'k' neighbors from the training set.
Gaussian Discriminant Analysis (GDA)	Classification	Generative	Continuous	Yes	Yes	No	Yes + homoscedasticity between classes (covariance matrix for each class is identical) + Normal distributed features	High	Low	Yes	Yes	Low (Works best with large datasets)	Expect normal distributed features. Models the prob distribution for each class. Uses Bayes theorem to compute posterior probability. Pred is the class with highest posterior prob for given features sample.
Naive Bayes	Classification	Generative	Continuous (Gaussian naive bayes: assumes gaussian distrib of	Yes	Yes	No	Yes	High	Low	Yes	No (but might produce bias)	Low	Naive Bayes relies on Bayes' theorem, which calculates the

		<p>features)</p> <p>/Categorical -----</p> <p>Multinomial naive bayes for discrete features</p> <p>-----</p> <p>Bernoulli naive bayes for binary features</p> <p>-----</p> <p>Categorical naive bayes for categorical features with more than 2 lvs</p>						towards the majority class)	probability of a hypothesis (class label) given the observed evidence (features). Used to compute posterior prob and select the class with highest posterior prob.
--	--	---	--	--	--	--	--	-----------------------------	--

4 - Unsupervised Learning

jueves, 2 de mayo de 2024 12:25

The goal of unsupervised learning is to find hidden patterns in unlabeled data.

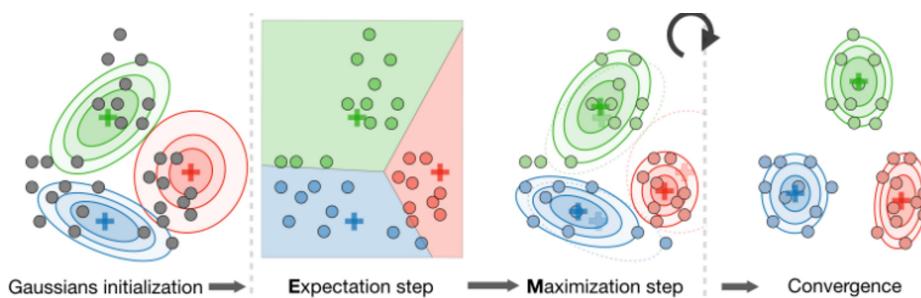
Clustering models:

Expectation-Maximization (EM):

Expectation-Maximization (EM) is an iterative algorithm used to estimate parameters in statistical models with latent variables, especially in situations where some data may be missing or unobserved.

Here's how it works:

1. **Expectation Step (E-step):**
 - In the E-step, the algorithm computes the expected value of the latent variables given the observed data and the current estimates of the model parameters.
 - It calculates the probabilities or likelihoods of the latent variables taking certain values based on the observed data and the current parameter estimates. This step involves filling in the missing or unobserved data.
2. **Maximization Step (M-step):**
 - In the M-step, the algorithm updates the parameters of the model to maximize the likelihood of the observed data, taking into account the expected values of the latent variables obtained from the E-step.
 - It adjusts the parameters of the model to better fit the observed data, incorporating the information provided by the expected values of the latent variables.
3. **Iterative Process:**
 - The E-step and M-step are repeated iteratively until the algorithm converges to a set of parameter estimates that maximize the likelihood of the observed data.
 - At each iteration, the algorithm typically improves the fit of the model to the data and converges to a local maximum of the likelihood function.
4. **Initialization:**
 - EM requires an initial guess of the model parameters to start the iterative process.
 - The algorithm may be sensitive to the choice of initial parameters, and multiple initializations with different starting points may be necessary to ensure convergence to a global maximum.
5. **Application:**
 - EM is widely used in various fields, including machine learning, computer vision, and natural language processing.
 - Common applications include clustering algorithms like Gaussian Mixture Models (GMMs), where EM is used to estimate the parameters of the mixture components.



Assumptions and considerations:

1. **Incomplete Data:** EM assumes that the observed data is incomplete or contains missing values, and that there are unobserved latent variables influencing the data generation process. It's designed to handle situations where some information about the underlying data generation process is missing or unobserved.
2. **Probabilistic Model:** EM assumes that the data generation process can be described by a probabilistic model with latent variables. The model typically involves joint probabilities or likelihoods of observed and latent variables.
3. **Local Convergence:** EM may converge to a local maximum of the likelihood function rather than a global maximum. The algorithm's performance can be sensitive to the choice of initial parameter values, and multiple initializations may be necessary to ensure convergence to a good solution.
4. **Convergence Criterion:** EM requires a convergence criterion to determine when to stop the iterative optimization process. Common criteria include reaching a maximum number of iterations, achieving a small change in parameter estimates between iterations, or when the log-likelihood of the observed data stabilizes.
5. **Model Complexity:** EM may struggle with high-dimensional or complex models, as the optimization process becomes computationally intensive and prone to local optima. Careful consideration of model complexity and regularization techniques may be necessary to ensure convergence and prevent overfitting.

6. **Dependence on Initialization:** The performance of EM can depend heavily on the choice of initial parameter values. Sensible initialization strategies, such as using results from simpler models or random initialization with multiple restarts, can help improve convergence and avoid getting stuck in local optima.
7. **Assumed Data Distribution:** EM may assume specific distributions for the latent and observed variables, such as Gaussian distributions for continuous variables or multinomial distributions for categorical variables. Violations of these distributional assumptions can lead to biased parameter estimates and affect the algorithm's performance.

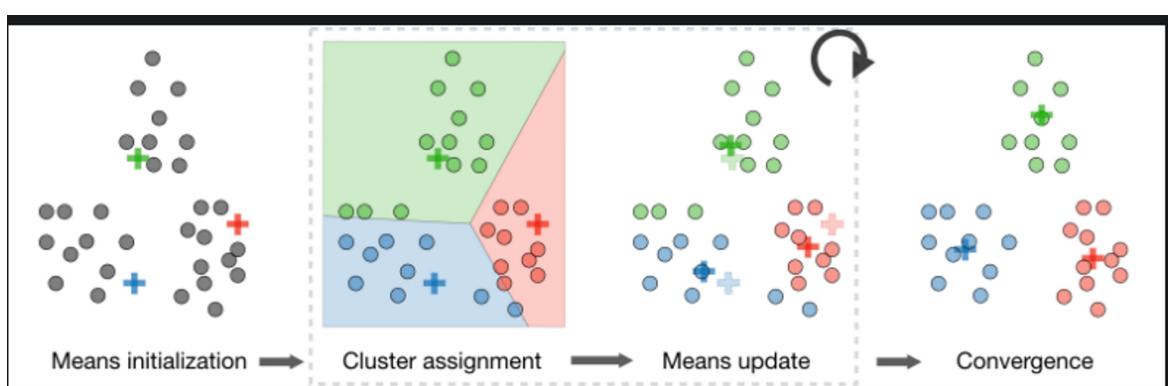
K-means clustering:

K-means clustering is a popular unsupervised machine learning algorithm used for partitioning a dataset into a pre-defined number of clusters.

Here's how it works:

1. **Initialization:** The algorithm starts by randomly initializing K cluster centroids (points in the feature space) within the data domain. These centroids represent the initial cluster centers.
2. **Assignment:** Each data point in the dataset is then assigned to the nearest cluster centroid based on some distance metric, typically the Euclidean distance. The data points are clustered based on their proximity to the centroids.
3. **Update:** After the assignment step, the algorithm recalculates the cluster centroids by taking the mean of all data points assigned to each cluster. This moves the centroids to the center of their respective clusters.
4. **Iteration:** Steps 2 and 3 are repeated iteratively until convergence criteria are met. Convergence may be reached when the cluster assignments no longer change or when a specified number of iterations is reached.
5. **Final Clustering:** Once convergence is achieved, the algorithm produces the final clustering, where each data point belongs to the cluster with the nearest centroid.
6. **Choosing K:** The number of clusters, K, is a hyperparameter that needs to be specified beforehand. Various methods, such as the elbow method or silhouette analysis, can be used to determine the optimal value of K based on clustering performance metrics.
7. **Initialization Sensitivity:** K-means is sensitive to the initial random selection of centroids. Different initializations can lead to different final cluster assignments and centroids. To mitigate this sensitivity, the algorithm is often run multiple times with different initializations, and the clustering with the lowest within-cluster variance is selected.
8. **Assumption of Clusters:** K-means assumes that the clusters are spherical and have similar sizes. It works well when the clusters are well-separated and have roughly equal variances, but may struggle with non-linearly separable clusters or clusters of different shapes and sizes.
9. **Scalability:** K-means is computationally efficient and scales well to large datasets. However, it may struggle with high-dimensional data or datasets with uneven cluster sizes.

In summary, K-means clustering is a simple yet effective algorithm for partitioning data into clusters. It is widely used in various domains such as customer segmentation, image segmentation, and anomaly detection, providing insights into the underlying structure of the data.



Assumptions and considerations:

1. **Number of Clusters (K):** K-means requires the number of clusters, K, to be specified beforehand. The algorithm partitions the data into exactly K clusters, so it's essential to choose an appropriate value of K based on domain knowledge or evaluation metrics.
2. **Initial Centroid Placement:** K-means is sensitive to the initial placement of centroids. Different initializations can lead to different final cluster assignments and centroids. To mitigate this sensitivity, the algorithm is often run multiple times with different initializations, and the clustering with the lowest within-cluster variance is selected.

3. **Cluster Shape and Size:** K-means assumes that clusters are spherical and have similar sizes. It works well when the clusters are well-separated and have roughly equal variances. However, it may struggle with non-linearly separable clusters or clusters of different shapes and sizes.
4. **Scalability:** K-means is computationally efficient and scales well to large datasets. However, its performance may degrade with high-dimensional data or datasets with uneven cluster sizes.
5. **Distance Metric:** K-means relies on a distance metric, typically the Euclidean distance, to measure the similarity between data points and centroids. It's important to choose an appropriate distance metric based on the nature of the data and the problem domain.
6. **Assumption of Equal Variance:** K-means assumes that all clusters have equal variance. This assumption may not hold true in practice, especially if the clusters have different shapes or sizes. In such cases, alternative clustering algorithms like Gaussian Mixture Models (GMMs) may be more appropriate.
7. **Outliers:** K-means can be sensitive to outliers, as they can disproportionately influence the position of centroids. Outliers can lead to suboptimal cluster assignments and centroids, affecting the overall clustering performance.
8. **Convergence:** K-means converges to a local minimum of the within-cluster variance objective function. It's important to monitor the convergence of the algorithm and ensure that it reaches a stable clustering solution.

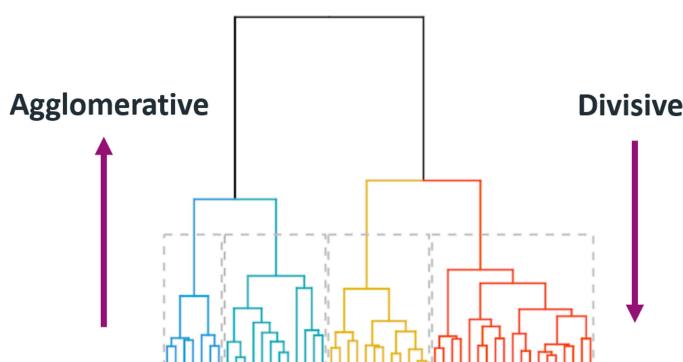
Hierarchical clustering:

Hierarchical clustering is a clustering algorithm used to group similar data points into clusters based on their proximity to each other. Unlike K-means, which requires specifying the number of clusters beforehand, hierarchical clustering produces a hierarchical decomposition of the dataset into a tree-like structure called a dendrogram.

Here's how it works:

- **Agglomerative or Divisive:** Hierarchical clustering can be agglomerative (bottom-up) or divisive (top-down).
 - Agglomerative: It starts with each data point as a separate cluster and then iteratively merges the closest clusters until only one cluster remains.
 - Divisive: It starts with all data points in a single cluster and then recursively divides the clusters into smaller clusters until each data point is in its own cluster.
- **Distance Metric:** Hierarchical clustering requires a distance metric to measure the dissimilarity or distance between data points. Common distance metrics include Euclidean distance, Manhattan distance, and cosine similarity. The choice of distance metric depends on the nature of the data and the problem domain.
- **Linkage Criteria:** In agglomerative hierarchical clustering, the choice of linkage criteria determines how the distance between clusters is calculated during the merging process. Common linkage criteria include:
 - Single Linkage: The distance between two clusters is defined as the shortest distance between any two points in the two clusters.
 - Complete Linkage: The distance between two clusters is defined as the maximum distance between any two points in the two clusters.
 - Average Linkage: The distance between two clusters is defined as the average distance between all pairs of points in the two clusters.
 - Ward's Linkage: It minimizes the variance when merging clusters and tends to produce compact, spherical clusters.
- **Dendrogram:** The output of hierarchical clustering is a dendrogram, which is a tree-like structure that illustrates the merging process. The dendrogram visualizes the hierarchical relationships between clusters and can be used to determine the optimal number of clusters.
- **Cutting the Dendrogram:** To obtain a specific number of clusters, practitioners can cut the dendrogram at a certain height or distance threshold. The height at which the dendrogram is cut corresponds to the desired number of clusters.
- **Scalability:** Hierarchical clustering can be computationally expensive, especially for large datasets, as it requires pairwise distance calculations between all data points. However, with efficient algorithms and data structures, hierarchical clustering can still be applied to datasets of moderate size.

In summary, hierarchical clustering is a flexible and interpretable clustering algorithm that produces a dendrogram to represent the hierarchical relationships between data points. It does not require specifying the number of clusters beforehand and can be useful for exploratory data analysis and visualizing cluster structures.



Assumptions and considerations:

- Assumption of Proximity:** Hierarchical clustering assumes that the proximity or similarity between data points can be measured using a distance metric. Common distance metrics include Euclidean distance, Manhattan distance, and cosine similarity. The choice of distance metric can impact the clustering results and should be chosen based on the nature of the data.
- Choice of Linkage Criteria:** Hierarchical clustering requires selecting a linkage criteria to determine how the distance between clusters is calculated during the merging process. Different linkage criteria, such as single linkage, complete linkage, average linkage, and Ward's linkage, can produce different clustering structures. The choice of linkage criteria should be made based on the clustering objectives and the characteristics of the data.
- Assumption of Cluster Structure:** Hierarchical clustering assumes that the data can be organized into a hierarchical structure, where clusters are nested within each other. This hierarchical structure can be visualized using a dendrogram, which illustrates the merging process and the hierarchical relationships between clusters.
- Sensitivity to Scaling:** Hierarchical clustering can be sensitive to the scaling of the data, as distance metrics may be affected by the scale of the features. It's important to standardize or normalize the features to ensure that all features contribute equally to the distance calculations.
- Computational Complexity:** Hierarchical clustering can be computationally expensive, especially for large datasets, as it requires pairwise distance calculations between all data points. The computational complexity of hierarchical clustering algorithms can be $O(n^2)$ or higher, where n is the number of data points. Efficient algorithms and data structures are needed to handle large datasets.
- Difficulty with High-Dimensional Data:** Hierarchical clustering may struggle with high-dimensional data, as distance metrics become less meaningful in high-dimensional spaces (the curse of dimensionality). Dimensionality reduction techniques may be necessary to reduce the dimensionality of the data before applying hierarchical clustering.
- Choice of Number of Clusters:** Hierarchical clustering produces a dendrogram, which does not directly provide the number of clusters. Determining the optimal number of clusters requires interpreting the dendrogram or cutting it at a certain height or distance threshold. Different cutting strategies can lead to different numbers of clusters, and the choice of the number of clusters should be guided by domain knowledge or clustering evaluation metrics.

DBSCAN:

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is a popular density-based clustering algorithm used for partitioning a dataset into clusters of varying shapes and sizes. Here's how it works:

1. Core Points, Border Points, and Noise:

- DBSCAN categorizes data points into three categories: core points, border points, and noise (or outlier) points.
- A core point is a data point that has at least a specified number of neighboring points (minPts) within a given distance (epsilon or ϵ).
- A border point is a data point that is within the epsilon distance of a core point but does not have enough neighbors to be considered a core point itself.
- A noise point (or outlier) is a data point that is neither a core point nor a border point.

2. Clustering Process:

- DBSCAN starts by randomly selecting an unvisited data point. If the selected point is a core point, DBSCAN forms a cluster around it by recursively adding all reachable points within the epsilon distance.
- If the selected point is a border point, it is assigned to the cluster of its nearest core point.
- The algorithm continues this process until all data points have been visited and classified into clusters or noise.

3. Parameter Selection:

- The key parameters in DBSCAN are epsilon (ϵ) and minPts .
- Epsilon defines the maximum distance between two points for them to be considered neighbors. It determines the size of the neighborhood around each data point.
- MinPts specifies the minimum number of neighboring points required for a data point to be considered a core point.
- The choice of epsilon and minPts depends on the characteristics of the dataset and the desired clustering outcome. Selecting appropriate values for these parameters is crucial for the effectiveness of DBSCAN.

4. Pros vs Cons:

- DBSCAN is robust to noise and can handle datasets with outliers effectively.
- It can discover clusters of arbitrary shapes and sizes, including clusters with irregular shapes or varying densities.
- DBSCAN does not require specifying the number of clusters beforehand, unlike many other clustering algorithms.
- DBSCAN may struggle with datasets of varying densities or with clusters of significantly different sizes.
- It is sensitive to the choice of epsilon and minPts parameters, and selecting appropriate values can be challenging.
- DBSCAN may not perform well on high-dimensional datasets due to the curse of dimensionality.

In summary, DBSCAN is a powerful density-based clustering algorithm that can effectively identify clusters in datasets with noise and outliers. By properly selecting the parameters and understanding its characteristics, DBSCAN can provide valuable insights into the underlying structure of

the data.

Assumptions and considerations:

1. **Density-Based Clusters:** DBSCAN assumes that clusters are dense regions in the feature space, separated by regions of lower density. It identifies clusters based on the density of data points rather than their proximity in the feature space.
2. **Parameter Sensitivity:** DBSCAN's performance is sensitive to the choice of its two key parameters: epsilon (ϵ) and minPts.
 - Epsilon (ϵ) determines the maximum distance between two points for them to be considered neighbors.
 - MinPts specifies the minimum number of neighboring points required for a data point to be considered a core point.
 - Selecting appropriate values for these parameters is crucial for the effectiveness of DBSCAN. Poor parameter choices can lead to either over-segmentation (too many small clusters) or under-segmentation (fewer or no clusters).
3. **Cluster Shape and Size:** DBSCAN can identify clusters of arbitrary shapes and sizes, including clusters with irregular shapes and varying densities. However, it may struggle with datasets containing clusters of significantly different sizes or non-convex shapes.
4. **Noise Handling:** DBSCAN is robust to noise and can effectively identify outliers as noise points. It classifies data points that do not belong to any cluster as noise (or outlier) points.
5. **Computational Complexity:** DBSCAN's computational complexity depends on the size of the dataset and the parameter settings. It requires calculating pairwise distances between data points and can be computationally expensive for large datasets or high-dimensional data.
6. **Non-Euclidean Distance Metrics:** While DBSCAN commonly uses the Euclidean distance metric, it can also work with other distance metrics appropriate for the dataset's domain. However, the choice of distance metric may affect the clustering results and should be carefully considered.
7. **Handling High-Dimensional Data:** DBSCAN may struggle with high-dimensional datasets due to the curse of dimensionality. Dimensionality reduction techniques or careful feature selection may be necessary to improve DBSCAN's performance on high-dimensional data.

OPTICS:

OPTICS (Ordering Points To Identify the Clustering Structure) is a density-based clustering algorithm that extends the concepts of DBSCAN (Density-Based Spatial Clustering of Applications with Noise). OPTICS generates a hierarchical clustering structure, called an OPTICS plot or reachability plot, which provides a more flexible and interpretable approach to density-based clustering. Here's how it works:

1. **Reachability Distance:**
 - OPTICS introduces the concept of reachability distance, which measures the distance at which a data point can be reached from another data point while taking into account the density of the data points in between.
 - The reachability distance of a point p from another point q is defined as the maximum of the core distance of q and the distance between p and q .
2. **Core Distance:**
 - The core distance of a data point p is the distance to its k -th nearest neighbor, where k is a user-defined parameter.
 - Core distance provides a measure of the local density around a data point. Points with lower core distances are considered to be in denser regions.
3. **OPTICS Plot:**
 - OPTICS constructs an ordering of the data points based on their reachability distances. This ordering forms a reachability plot or OPTICS plot, which represents the clustering structure of the dataset.
 - The OPTICS plot visualizes the hierarchical relationships between clusters and noise points. It reveals the density-based clustering structure of the dataset at different density levels.
4. **Clustering:**
 - To obtain actual clusters from the OPTICS plot, practitioners can use a parameter called the ϵ -neighborhood size to extract clusters at different density levels.
 - Clusters are formed by grouping together points that are within a specified ϵ -neighborhood distance. The choice of ϵ determines the granularity of the clustering.
5. **Advantages:**
 - OPTICS provides a more flexible and interpretable approach to density-based clustering compared to DBSCAN.
 - It produces a hierarchical clustering structure that allows users to identify clusters at different density levels and to adjust parameters such as the minimum cluster size and maximum cluster radius.
 - OPTICS is robust to noise and outliers and can handle datasets with varying densities and non-convex shapes.
6. **Limitations:**
 - OPTICS may require more computational resources compared to DBSCAN, especially for large datasets, due to the construction of the OPTICS plot.
 - Parameter selection, such as the k -nearest neighbors and ϵ -neighborhood size, can impact the clustering results and may require tuning.

Assumptions and considerations:

1. **Density-Based Clusters:** OPTICS assumes that clusters in the dataset are defined by regions of high data density separated by regions of low density. It identifies clusters based on the density of data points rather than their geometric proximity.
2. **Reachability Distance:** OPTICS relies on the concept of reachability distance to measure the connectivity between data points in the dataset. The reachability distance considers both the local density around a point and the distances between points in the dataset.
3. **Parameter Sensitivity:** OPTICS performance can be sensitive to the choice of parameters, such as the k -nearest neighbors and the ϵ -neighborhood size used for clustering. Proper parameter selection is essential to achieve meaningful clustering results.
4. **Hierarchical Clustering Structure:** OPTICS produces a hierarchical clustering structure in the form of an OPTICS plot, which represents the ordering of data points based on their reachability distances. Understanding and interpreting this hierarchical structure is crucial for identifying clusters at different density levels.
5. **Computation Time:** OPTICS may require more computational resources compared to other clustering algorithms, especially for large datasets, due to the construction of the OPTICS plot and the calculation of reachability distances.
6. **Handling Noise and Outliers:** OPTICS is robust to noise and outliers in the dataset and can effectively identify clusters while distinguishing noise points from dense regions. However, the choice of parameters and the definition of what constitutes noise may influence the clustering results.
7. **Cluster Shape and Size:** OPTICS can identify clusters of arbitrary shapes and sizes, including clusters with irregular shapes and varying densities. However, it may struggle with datasets containing clusters of significantly different sizes or non-convex shapes.
8. **Interpretation of the OPTICS Plot:** Interpreting the OPTICS plot and extracting meaningful clusters from it require understanding the hierarchical relationships between data points and adjusting clustering parameters accordingly.

Summary:

Algorithm	Type of Problems Solved	Type of Output	Assumptions/Considerations	Interpretability	Computational Cost	Sensitivity to Outliers	Sensitivity to Imbalanced Data	Complexity
EM (Expectation-Maximization)	Clustering	Probabilistic assignment of data points to clusters	- Assumes data is generated by a mixture of Gaussian distributions - Requires the number of clusters to be specified - May converge to local optima - Sensitive to initialization	Low	Medium	Sensitive	Not sensitive	Medium
k-means	Clustering	Hard assignment of data points to clusters	- Assumes clusters are spherical and of equal variance - Requires the number of clusters to be specified - Sensitive to initialization - Can converge to local optima	Low	Low	Sensitive	Not sensitive	Medium
Hierarchical Clustering	Clustering	Dendrogram representation of hierarchical clusters	- Does not require the number of clusters to be specified - May be agglomerative (bottom-up) or divisive (top-down) - Sensitive to the choice of linkage method and distance metric	High	High	Sensitive	Not sensitive	High
DBSCAN	Clustering	Density-based clustering with noise points	- Identifies clusters as dense regions separated by areas of lower density - Does not require the number of clusters to be specified - Sensitive to epsilon and minPts parameters - Handles irregularly shaped clusters and noise well	Low	Medium	Not sensitive	Sensitive	Medium
OPTICS	Clustering	Hierarchical clustering with reachability plot	- Orders data points based on reachability distances - Produces a hierarchical clustering structure - Sensitive to minPts parameter - Suitable for datasets with varying densities and non-convex shapes	Medium	High	Not sensitive	Sensitive	High

Clustering evaluation metrics:

Clustering assessment metrics are used to evaluate the quality of clustering algorithms by quantitatively measuring how well the algorithm has grouped similar data points together and separated dissimilar ones. Here are some commonly used clustering assessment metrics:

1. Silhouette Score:

- The silhouette score measures the cohesion and separation of clusters.
- For each data point, the silhouette score computes the mean distance between the point and all other points in its cluster (a) and the mean distance between the point and all points in the nearest neighboring cluster (b). The silhouette score for the point is then given by $(b - a) / \max(a, b)$.
- The silhouette score ranges from -1 to 1, where a high value indicates that the point is well-clustered, a value near zero indicates overlapping clusters, and a negative value indicates that the point may be assigned to the wrong cluster.

2. Davies-Bouldin Index (DB Index):

- The DB index measures the compactness and separation of clusters.
- It computes the average similarity between each cluster and its most similar cluster, normalized by the average within-cluster distance.
- A lower DB index indicates better clustering, with values closer to zero representing better-defined clusters.

3. Dunn Index:

- The Dunn index evaluates the compactness and separation of clusters.
- It computes the ratio of the minimum distance between clusters to the maximum diameter of clusters.
- A higher Dunn index indicates better clustering, with larger values representing better-defined and well-separated clusters.

4. Calinski-Harabasz Index:

- The Calinski-Harabasz index measures the ratio of between-cluster dispersion to within-cluster dispersion.
- It computes the ratio of the sum of between-cluster distances to the sum of within-cluster distances, multiplied by the ratio of the total number of points to the total number of clusters minus one.
- A higher Calinski-Harabasz index indicates better clustering, with larger values representing more compact and well-separated clusters.

5. Adjusted Rand Index (ARI):

- The ARI measures the similarity between the true clustering and the clustering produced by the algorithm, corrected for chance.
- It computes the fraction of pairs of data points that are assigned to the same or different clusters by both the true and predicted clusterings, adjusted for chance.
- The ARI ranges from -1 to 1, where a value of 1 indicates perfect clustering agreement, a value near zero indicates random clustering, and negative values indicate clustering disagreement.

6. Normalized Mutual Information (NMI):

- NMI measures the mutual information between the true clustering and the predicted clustering, normalized by the entropy of the two clusterings.
- It computes the information gain obtained by knowing the true clustering given the predicted clustering, normalized by the entropy of the true and predicted clusterings.
- The NMI ranges from 0 to 1, where a value of 1 indicates perfect clustering agreement, and a value of 0 indicates no mutual information between the clusterings.

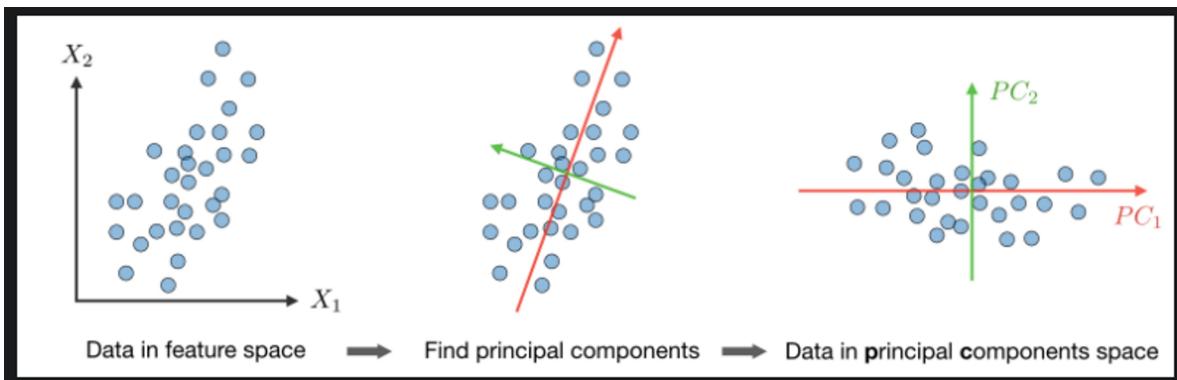
Dimension reduction:

Principal Component Analysis (PCA):

Principal Component Analysis (PCA) is a dimensionality reduction technique used to simplify complex datasets while preserving important information. It achieves this by transforming the original variables into a new set of orthogonal variables called principal components (PCs).

Here's how PCA works:

1. **Standardization:** PCA begins with the standardization of the dataset to ensure that all variables have a mean of zero and a standard deviation of one. This step is essential for PCA to give equal weight to all variables.
 2. **Covariance Matrix:** PCA computes the covariance matrix of the standardized dataset. The covariance matrix describes the relationships between pairs of variables, indicating how they vary together.
 3. **Eigenvalue Decomposition:** PCA then performs eigenvalue decomposition on the covariance matrix to extract its eigenvectors and eigenvalues. Eigenvectors represent the directions of maximum variance in the original feature space, while eigenvalues indicate the magnitude of variance along each eigenvector.
 4. **Selection of Principal Components:** PCA selects a subset of the eigenvectors, known as principal components, based on their corresponding eigenvalues. The number of principal components chosen typically depends on the amount of variance explained or a predefined threshold.
 5. **Projection:** Finally, PCA projects the original data onto the selected principal components, creating a new lower-dimensional space. This projection retains the maximum amount of variance possible while reducing the number of dimensions.
- PCA is widely used for various purposes, including data visualization, noise reduction, feature extraction, and speeding up machine learning algorithms. By reducing the dimensionality of the dataset, PCA can help uncover underlying patterns and relationships, making it a valuable tool in exploratory data analysis and model building.



Assumptions and considerations:

1. **Linearity:** PCA assumes that the relationships between variables are linear. If the relationships are highly nonlinear, PCA may not effectively capture the underlying structure of the data.
2. **Orthogonality:** PCA assumes that the principal components are orthogonal to each other, meaning they are uncorrelated. This assumption allows PCA to maximize the amount of variance explained by each component.
3. **Normality:** While not strictly required, PCA works best when the variables in the dataset are approximately normally distributed. Non-normal distributions may still work with PCA but could affect the interpretation of the principal components.
4. **Homoscedasticity:** PCA assumes that the variance of each variable is roughly constant across all levels of other variables. If there are significant differences in variances, PCA may give more weight to variables with larger variances.
5. **Large Variance Dominance:** PCA tends to emphasize variables with large variances. Therefore, variables with small variances may have less influence on the principal components and may be less accurately represented in the reduced-dimensional space.
6. **Scaling:** PCA is sensitive to the scale of the variables. It's important to standardize or normalize the variables before performing PCA to ensure that all variables contribute equally to the analysis.
7. **Linear Relationships:** PCA assumes that relationships between variables are linear. If the relationships are highly nonlinear, PCA may not capture the underlying structure effectively.
8. **Outliers:** PCA can be sensitive to outliers, especially when computing the covariance matrix. Outliers may disproportionately influence the principal components and should be handled appropriately before performing PCA.

Independent Component Analysis:

Independent Component Analysis (ICA) is a statistical technique used to separate a multivariate signal into additive, independent components. Unlike Principal Component Analysis (PCA), which finds orthogonal components that capture the maximum variance in the data, ICA aims to find components that are statistically independent of each other.

Here's how ICA works:

1. **Statistical Independence:** ICA assumes that the observed signals are generated by a linear combination of statistically independent source signals. These source signals are mixed together with some unknown mixing coefficients to produce the observed data.
2. **Unmixing Matrix:** The goal of ICA is to estimate an unmixing matrix that can separate the observed signals into their independent source components. This unmixing matrix is typically estimated using optimization algorithms such as gradient descent or maximum likelihood estimation.
3. **Independence Assumption:** ICA relies on the assumption that the source signals are statistically independent. This assumption allows ICA to identify the underlying structure of the data by finding components that are statistically unrelated to each other.
4. **Non-Gaussianity:** ICA works best when the source signals are non-Gaussian, as Gaussian signals are not typically independent. By exploiting the non-Gaussian nature of the source signals, ICA can separate them more effectively.
5. **Ambiguities:** ICA may encounter ambiguities in the estimation of the unmixing matrix, such as permutation and scaling ambiguities. Permutation ambiguity refers to the arbitrary ordering of the independent components, while scaling ambiguity refers to the arbitrary scaling of the components.
6. **Applications:** ICA has numerous applications in signal processing, neuroscience, image processing, and machine learning. It can be used for blind source separation, artifact removal from EEG signals, image denoising, and feature extraction.

Assumptions and considerations:

- Statistical Independence:** The primary assumption of ICA is that the observed signals are generated by a linear combination of statistically independent source signals. This assumption may not hold if the sources are not truly independent.
- Non-Gaussianity:** ICA works best when the source signals are non-Gaussian, as Gaussian signals are typically not independent. Non-Gaussianity enables ICA to differentiate between the sources and separate them effectively.
- Linear Mixing:** ICA assumes that the observed signals are generated by a linear mixing process, where each source signal is combined with some unknown mixing coefficients. Nonlinear mixing may violate this assumption and lead to suboptimal results.
- Number of Sources:** ICA requires an estimate of the number of independent sources in the data. If the number of sources is incorrectly specified, ICA may produce inaccurate results.
- Ambiguities:** ICA may encounter ambiguities in the estimation of the unmixing matrix, leading to permutation and scaling ambiguities. Permutation ambiguity refers to the arbitrary ordering of the independent components, while scaling ambiguity refers to the arbitrary scaling of the components.
- Preprocessing:** Proper preprocessing of the data is crucial for the success of ICA. This may include centering the data to have zero mean, whitening the data to have unit variance, and decorrelating the data to ensure that the covariance matrix is identity.
- Computational Complexity:** ICA can be computationally intensive, especially for high-dimensional data. Efficient algorithms and optimization techniques are needed to estimate the unmixing matrix accurately.
- Model Selection:** Choosing an appropriate algorithm and parameter settings for ICA can impact its performance. Different algorithms may be better suited for specific types of data or assumptions.

Summary:

Aspect	ICA	PCA
Objective	Separates mixed signals into statistically independent components	Reduces the dimensionality of data while preserving maximal variance
Assumption	Assumes observed signals are linear combinations of statistically independent sources	Assumes observed variables are linear combinations of orthogonal principal components
Non-Gaussianity	Works best when source signals are non-Gaussian	Not explicitly dependent on source signals' distribution, but Gaussianity is often assumed
Independence	Aims to find components that are statistically independent	Finds orthogonal components that capture maximal variance
Output	Independent components	Principal components
Ambiguities	May encounter permutation and scaling ambiguities in component ordering	No inherent ambiguities
Preprocessing	Often requires centering, whitening, and decorrelation of data	Typically starts with centering and may involve scaling
Computational Complexity	Can be computationally intensive, especially for high-dimensional data	Generally less computationally intensive than ICA
Applications	Signal processing, blind source separation, image processing	Dimensionality reduction, noise reduction, data visualization

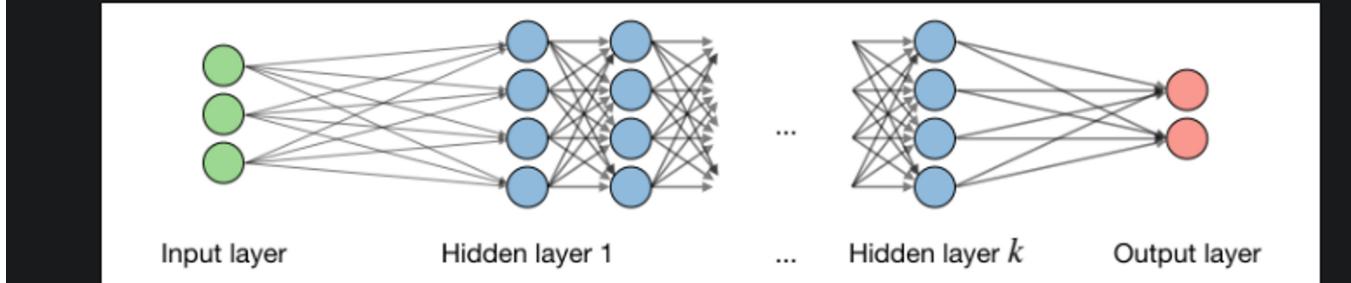
5 - Deep Learning

jueves, 2 de mayo de 2024 12:25

Key Concepts:

Neural networks are a class of models that are built with layers. Commonly used types of neural networks include convolutional and recurrent neural networks.

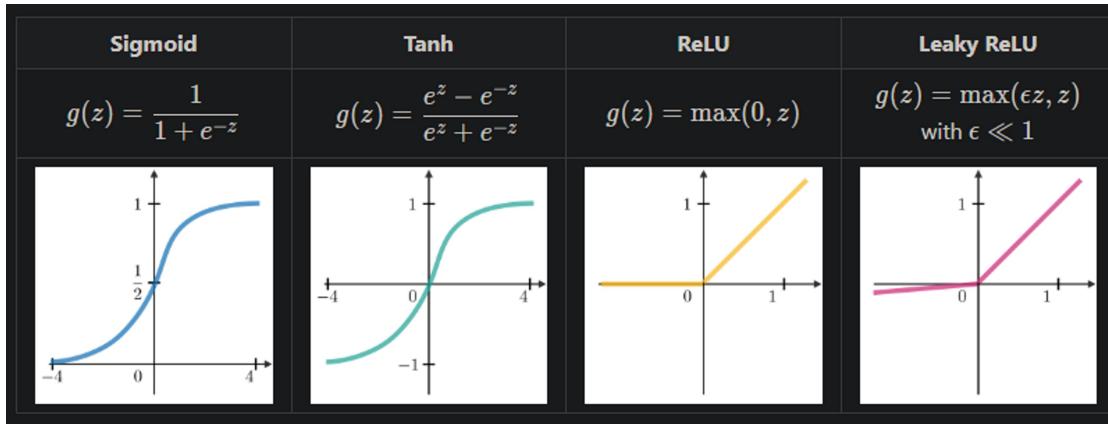
□ **Architecture** — The vocabulary around neural networks architectures is described in the figure below:



In each Hidden layer and in the output layer, we have an activation function.

Activation function:

Activation functions are used at the end of a hidden unit to introduce non-linear complexities to the model. Here are the most common ones:



Using a proper activation function is crucial for the effectiveness of a neural network. Here's how you can ensure you're using the right activation function:

- **Understand the Problem:** Different problems may benefit from different activation functions. For example, for binary classification tasks, the sigmoid activation function is commonly used in the output layer, while for multi-class classification tasks, the softmax activation function is preferred.
- **Consider Network Architecture:** The choice of activation function may depend on the network architecture and the specific layer where it is applied. For hidden layers, popular choices include ReLU (Rectified Linear Unit), Leaky ReLU, and variants like ELU (Exponential Linear Unit) or PReLU (Parametric ReLU).
- **Avoid Vanishing and Exploding Gradients:** Activation functions like sigmoid and tanh are prone to vanishing gradients, especially in deep networks. ReLU and its variants help mitigate this issue by avoiding saturation in the positive range, thereby preventing gradient vanishing.
- **Handle Non-linearity:** Activation functions introduce non-linearity into the network, allowing it to learn complex patterns and relationships in the data. Ensure that the chosen activation function can effectively capture the non-linearities

present in the problem domain.

- **Address Saturation:** Some activation functions, such as sigmoid and tanh, saturate at extreme input values, leading to slow convergence during training. ReLU variants address this by being non-saturating in the positive range, promoting faster training.
- **Consider Computational Efficiency:** Some activation functions are computationally more expensive than others. For large-scale models or deployment in resource-constrained environments, choosing computationally efficient activation functions can be beneficial.
- **Experiment and Validate:** Experiment with different activation functions and evaluate their performance using validation data. Compare the model's performance with various activation functions and select the one that yields the best results in terms of accuracy, convergence speed, and generalization.

□ **Learning rate** — The learning rate, often noted α or sometimes η , indicates at which pace the weights get updated. This can be fixed or adaptively changed. The current most popular method is called Adam, which is a method that adapts the learning rate.

□ **Backpropagation** — Backpropagation is a method to update the weights in the neural network by taking into account the actual output and the desired output. The derivative with respect to weight w is computed using chain rule and is of the following form:

$$\boxed{\frac{\partial L(z, y)}{\partial w} = \frac{\partial L(z, y)}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial w}}$$

As a result, the weight is updated as follows:

$$\boxed{w \leftarrow w - \alpha \frac{\partial L(z, y)}{\partial w}}$$

□ **Updating weights** — In a neural network, weights are updated as follows:

- Step 1: Take a batch of training data.
- Step 2: Perform forward propagation to obtain the corresponding loss.
- Step 3: Backpropagate the loss to get the gradients.
- Step 4: Use the gradients to update the weights of the network.

□ **Dropout** — Dropout is a technique meant to prevent overfitting the training data by dropping out units in a neural network. In practice, neurons are either dropped with probability p or kept with probability $1 - p$.

Loss functions:

Binary Classification:

- **Binary Cross-Entropy Loss (Log Loss)**: Used when the output has two classes (0 or 1).
- **Hinge Loss**: Particularly used with Support Vector Machines (SVMs) for binary classification.

Multi-Class Classification:

- **Categorical Cross-Entropy Loss (Softmax Loss)**: Used when the output has multiple classes and each sample belongs to exactly one class.
- **Sparse Categorical Cross-Entropy Loss**: Similar to categorical cross-entropy, but used when the true labels are integers instead of one-hot encoded vectors.

Regression:

- **Mean Squared Error (MSE) Loss:** Commonly used for regression tasks, where the model predicts continuous values.
- **Mean Absolute Error (MAE) Loss:** Another option for regression tasks, often preferred if the dataset contains outliers.

Ordinal Regression:

- **Ordinal Cross-Entropy Loss:** Used when the output represents ordered categories (e.g., low, medium, high) and the order matters.

Imbalanced Classification:

- **Weighted Cross-Entropy Loss:** Adjusts the loss for imbalanced datasets by assigning different weights to classes based on their frequency.
- **Focal Loss:** Focuses on hard-to-classify examples, reducing the loss contribution of well-classified examples.

Anomaly Detection:

- **Reconstruction Loss:** Used with autoencoders for anomaly detection tasks, where the loss measures how well the model can reconstruct the input data.

Semantic Segmentation:

- **Dice Loss:** Measures the overlap between predicted and true segmentation masks, often used for medical image segmentation.
- **Jaccard/IoU Loss:** Measures the intersection over union (IoU) between predicted and true segmentation masks.

Object Detection:

- **Focal Loss:** Particularly used with region-based convolutional neural networks (R-CNNs) for object detection tasks.

Dropout — Dropout is a technique meant to prevent overfitting the training data by dropping out units in a neural network. In practice, neurons are either dropped with probability p or kept with probability $1 - p$.

Exploding Gradients:

The term "gradient exploding" refers to a phenomenon that can occur during the training of neural networks, particularly in deep architectures with many layers. It occurs when the gradients of the model parameters (weights and biases) become extremely large during backpropagation, leading to numerical instability and making the optimization process difficult or impossible.

Gradient exploding is the counterpart of another common issue known as "gradient vanishing," where gradients become extremely small, causing the learning process to slow down significantly.

In the context of machine learning and neural networks, a gradient refers to a vector that represents the rate of change of a function with respect to its parameters.

Gradient exploding can manifest in several ways:

Loss Divergence: As gradients become excessively large, the loss function may diverge, leading to NaN (not a number) or infinity values.

Overflow: The numerical values of weights and activations may grow exponentially, exceeding the limits of numerical precision in computer arithmetic, resulting in overflow errors.

Unstable Training: The optimization algorithm (e.g., stochastic gradient descent) may fail to converge or oscillate wildly, making training unstable and unpredictable.

Gradient exploding can be caused by various factors, including:

- **Initialization:** Inadequate initialization of network weights, such as using large initial weights or inappropriate weight initialization techniques.
- **Activation Functions:** Certain activation functions, such as the sigmoid function, can exacerbate gradient exploding due to their saturating nature, leading to vanishing or exploding gradients.
- **Network Architecture:** Deep architectures with many layers are more prone to gradient exploding, especially if the gradients are not properly managed or regularized.

To mitigate gradient exploding, several techniques can be employed:

- **Gradient Clipping:** Limiting the magnitude of gradients during training by clipping them to a predefined threshold, preventing them from growing too large.

- **Proper Initialization:** Using appropriate weight initialization techniques, such as Xavier or He initialization, to ensure that gradients neither explode nor vanish during training.
- **Batch Normalization:** Normalizing the activations of each layer within mini-batches, which can help stabilize the gradients and accelerate training.
- **Learning Rate Scheduling:** Adjusting the learning rate dynamically during training, such as using techniques like learning rate decay or adaptive learning rate methods, to prevent abrupt changes that may lead to gradient exploding.

Gradient vanishing:

Gradient vanishing refers to a phenomenon that occurs during the training of deep neural networks, particularly in architectures with many layers. It happens when the gradients of the loss function with respect to the parameters (weights and biases) of the network become extremely small as they are backpropagated through multiple layers. Consequently, the gradients effectively "vanish," becoming close to zero or too small to effectively update the parameters during optimization.

Key points about gradient vanishing:

Deep Architectures: Gradient vanishing tends to be more pronounced in deep neural networks with many layers. As gradients are backpropagated through numerous layers, they can diminish significantly, especially in networks with sigmoid or hyperbolic tangent activation functions, which are prone to saturation.

Saturating Activation Functions: Activation functions like sigmoid and hyperbolic tangent have regions where their derivatives are close to zero, leading to vanishing gradients. When gradients are multiplied layer by layer during backpropagation, these small gradients compound, resulting in vanishing gradients deeper in the network.

Long-Term Dependencies: In recurrent neural networks (RNNs) and long short-term memory networks (LSTMs), gradient vanishing can hinder the learning of long-term dependencies. If important information from earlier time steps or layers is not effectively propagated due to vanishing gradients, the network may struggle to capture temporal patterns or context.

Implications for Training: Gradient vanishing can impede the convergence of the training process, slowing down or preventing learning altogether. It can also lead to suboptimal performance or poor generalization, as the network may fail to capture important features or relationships in the data.

To mitigate gradient vanishing, several techniques can be employed:

- **Initialization:** Proper initialization of network weights, such as using techniques like Xavier or He initialization, can help alleviate gradient vanishing by ensuring that gradients neither explode nor vanish during training.
- **Activation Functions:** Using activation functions like ReLU (Rectified Linear Unit) or variants (e.g., Leaky ReLU) that do not saturate in the positive region can mitigate gradient vanishing by maintaining non-zero gradients for positive inputs.
- **Normalization:** Techniques like batch normalization or layer normalization can stabilize training by normalizing activations within mini-batches or layers, reducing the likelihood of vanishing gradients.

Preprocessing needed:

Input data for neural networks needs to be numeric. Neural networks operate on numerical data and perform computations involving weights, biases, and activation functions, which require numeric inputs.

However, there are techniques available to handle non-numeric data:

- **Encoding Categorical Variables:** If your dataset contains categorical variables (e.g., color, gender, product category), you can encode them numerically using techniques like one-hot encoding or label encoding. This converts categorical variables into numerical representations that can be fed into neural networks.
- **Feature Engineering:** Sometimes, non-numeric data can be transformed or engineered into numeric features that capture relevant information. For example, text data can be converted into numerical representations using techniques like word embeddings or TF-IDF (Term Frequency-Inverse Document Frequency).
- **Preprocessing:** It's essential to preprocess your data before feeding it into neural networks. This may involve scaling numeric features to a similar range, handling missing values, and normalizing or standardizing the data to improve model performance.
- **Embedding Layers:** In the case of sequential or text data (e.g., sentences, documents), you can use embedding layers in

neural networks to automatically learn numerical representations (embeddings) of words or tokens from the input data.

Types of neural networks based on problem to solve:

Classification:

- **Convolutional Neural Networks (CNNs):** Used for image classification tasks such as identifying objects in images (e.g., CIFAR-10, ImageNet).
- **Recurrent Neural Networks (RNNs):** Applied to text classification tasks like sentiment analysis (e.g., IMDb movie reviews dataset).
- **Feedforward Neural Networks (FNNs):** Utilized for binary or multi-class classification problems with structured data (e.g., predicting customer churn).

Regression:

- **Feedforward Neural Networks (FNNs):** Used for regression tasks like predicting house prices based on features such as size, location, and number of bedrooms (e.g., Boston Housing dataset).
- **Recurrent Neural Networks (RNNs):** Applied to time series forecasting tasks such as predicting stock prices or weather conditions.

Clustering:

- **Self-Organizing Maps (SOMs):** Used to cluster similar data points in an unsupervised manner, often applied to tasks like customer segmentation.
- **Autoencoders:** Can be used for dimensionality reduction and clustering by learning a compressed representation of data.

Dimensionality Reduction:

- **Autoencoders:** Learn a low-dimensional representation of data by compressing and then reconstructing it. Useful for tasks like data visualization and feature extraction.

Sequence Modeling:

- **Long Short-Term Memory (LSTM) Networks:** Used for tasks involving sequential data with long-range dependencies, such as language translation, speech recognition, and time series prediction.
- **Gated Recurrent Units (GRUs):** Similar to LSTMs, but with fewer parameters and faster training times, commonly used for similar tasks.

Generative Modeling:

- **Generative Adversarial Networks (GANs):** Used to generate realistic images, music, or text by training two neural networks (generator and discriminator) in a competitive setting.
- **Variational Autoencoders (VAEs):** Can be used to generate new data samples by learning a probabilistic model of the data distribution.

Reinforcement Learning:

- **Deep Q-Networks (DQN):** Used for model-free reinforcement learning, particularly in tasks involving discrete action spaces like game playing (e.g., Atari games).
- **Policy Gradient Methods:** Used for policy optimization in tasks with continuous action spaces, such as robotic control or autonomous driving.

Description of mentioned architectures:

Feedforward Neural Networks (FNNs):

- Simplest form of neural network where information flows in one direction from input to output.
- Consists of input layer, hidden layers (if any), and output layer.
- Used for tasks like classification and regression.

Convolutional Neural Networks (CNNs):

- Designed for processing structured grid-like data such as images.
- Utilize convolutional layers with filters to extract features hierarchically.
- Commonly used in image recognition, object detection, and image segmentation tasks.

Recurrent Neural Networks (RNNs):

- Specialized for sequence data where the order of elements matters.
- Have connections that form directed cycles, allowing information to persist over time.
- Used in tasks like natural language processing, time series prediction, and speech recognition.

Long Short-Term Memory Networks (LSTMs):

- A type of RNN designed to address the vanishing gradient problem.
- Utilize gated cells to selectively remember or forget information over long sequences.
- Particularly effective for tasks requiring modeling of long-range dependencies.

Gated Recurrent Units (GRUs):

- Similar to LSTMs but with a simpler architecture, having fewer parameters.
- Comprise of update and reset gates to control information flow.
- Offer a good trade-off between performance and computational complexity.

Autoencoders:

- Neural networks trained to copy input data to output, typically through a bottleneck layer.
- Used for unsupervised learning, dimensionality reduction, and data denoising.
- Comprise an encoder to compress input and a decoder to reconstruct it.

Generative Adversarial Networks (GANs):

- Comprise of two neural networks, generator and discriminator, trained in a competitive setting.
- Generator creates realistic data samples, while discriminator tries to distinguish between real and generated data.
- Used for generating realistic images, videos, and other types of data.

Variational Autoencoders (VAEs):

- Combine elements of autoencoders and variational inference.
- Aim to learn a latent representation of data and generate new samples from the learned distribution.
- Often used for generating new images and performing data synthesis tasks.

Transformers:

- Utilize self-attention mechanisms to capture global dependencies in sequential data.
- Introduced in the context of natural language processing (NLP) for tasks like language translation and text generation.
- Have become the state-of-the-art architecture for various NLP tasks due to their parallelizability and ability to handle long-range dependencies.

6 - Time Series Forecasting

jueves, 2 de mayo de 2024 9:58

While some forecasting models may focus primarily on non-seasonal patterns, many models explicitly account for seasonal variations in the data. The choice of model depends on the specific characteristics of the time series and the nature of the patterns to be captured.

Models like **ARIMA** that focus on finding non-seasonal patterns, need to remove seasonality, trends and make them stationary (Stationarity implies that the statistical properties of the series (such as mean, variance, and autocorrelation) remain constant over time), during preprocessing.

LSTMs (Long Short-Term Memory networks) are a type of recurrent neural network (RNN) designed to handle sequences of data where dependencies exist across time steps. While LSTMs are powerful models for capturing temporal dependencies, they do not inherently require preprocessing steps such as handling seasonality, stationarity, and trends. However, these preprocessing steps might still be necessary depending on the characteristics of the data and the specific task at hand.

Here's a breakdown:

Seasonality:

- Seasonality refers to patterns that repeat at regular intervals, such as daily, weekly, or yearly cycles.
- If the data exhibits strong seasonal patterns, it may be beneficial to preprocess it to remove or adjust for seasonality. Techniques like seasonal differencing or seasonal decomposition can be used.
- Preprocessing for seasonality can help improve the LSTM's ability to capture long-term dependencies without being influenced by repetitive patterns.

Stationarity:

- Stationarity implies that the statistical properties of a time series (such as mean and variance) remain constant over time.
- While LSTMs can learn from non-stationary data, stationarity can aid in model training and interpretation.
- Preprocessing techniques such as differencing or transformation (e.g., log transformation) may be applied to make the data stationary before feeding it into the LSTM.

Trends:

- Trends refer to long-term changes or patterns in the data.
- If the data exhibits trends, detrending techniques may be applied to remove or reduce the trend component.
- Detrending can help focus the LSTM on capturing shorter-term dependencies and patterns that are not influenced by long-term trends.

Here are some common methods to treat them:

Seasonality:

- **Seasonal Differencing:** Subtracting the time series values from the corresponding values from the previous season. This helps remove the seasonal component.
- **Seasonal Decomposition:** Decomposing the time series into its seasonal, trend, and residual components using techniques like seasonal decomposition of time series (STL) or seasonal-trend decomposition using LOESS (STL).
- **Seasonal Adjustment:** Applying mathematical transformations or adjustments to remove or reduce the seasonal patterns in the data.

Stationarity:

- **Differencing:** Taking differences between consecutive observations to remove the trend component and make the data stationary.
- **Transformation:** Applying mathematical transformations such as logarithmic transformation or square root transformation to stabilize the variance and make the data more stationary.
- **Detrending:** Removing the trend component from the time series data using techniques like polynomial regression or moving averages.

Trends:

- **Detrending:** Removing the trend component from the time series using techniques like polynomial regression, moving averages, or exponential smoothing.
- **Transformation:** Applying mathematical transformations to stabilize the variance and remove trends from the data.
- **Seasonal-Trend Decomposition:** Decomposing the time series into its seasonal, trend, and residual components to isolate and remove the trend component.

7 - Intro Gen AI and LLMs

jueves, 2 de mayo de 2024 12:25

GenAI, or Generative Artificial Intelligence, refers to a category of AI algorithms and models that are designed to generate new data samples that are similar to, or indistinguishable from, existing data. Unlike traditional AI models that focus on tasks like classification or prediction, generative models aim to understand and mimic the underlying distribution of the data to produce novel samples. GenAI encompasses a variety of generative models, including Generative Adversarial Networks (GANs), Variational Autoencoders (VAEs), autoregressive models, and more. These models learn to generate new data samples by capturing the underlying patterns and structures present in the training data.

GenAI for NLP:

Transformers are a type of neural network architecture that has gained significant attention, particularly in the field of natural language processing (NLP). Introduced by Vaswani et al. in the paper "Attention is All You Need" in 2017, transformers have revolutionized NLP tasks by achieving state-of-the-art performance in various benchmarks.

Here's an elaboration on transformers:

Self-Attention Mechanism:

- At the core of transformers is the self-attention mechanism, which allows the model to weigh the importance of different words in a sequence when processing each word.
- Unlike recurrent neural networks (RNNs) and convolutional neural networks (CNNs), transformers can capture long-range dependencies in sequences efficiently, as each word can attend to all other words in the sequence directly.

Transformer Architecture:

- Transformers consist of an encoder and a decoder, each comprising multiple layers of self-attention mechanisms and feedforward neural networks.
- The encoder processes the input sequence, while the decoder generates the output sequence in sequence-to-sequence tasks like language translation.
- Each layer in the transformer architecture is composed of multiple self-attention heads, allowing the model to capture different aspects of the input sequence simultaneously.

Positional Encoding:

- Since transformers do not inherently maintain the order of elements in a sequence like RNNs, positional encoding is added to the input embeddings to provide information about the position of each word in the sequence.
- Positional encoding allows the model to differentiate between words based on their position in the sequence and learn contextual representations effectively.

Attention Mechanisms:

- Transformers employ different types of attention mechanisms, including self-attention (also known as intra-attention) and multi-head attention, which allow the model to focus on different parts of the input sequence simultaneously.
- Attention mechanisms enable transformers to capture dependencies between words that are distant from each other in the sequence, leading to improved performance in tasks requiring understanding of long-range context.

Applications:

- Transformers have been successfully applied to various NLP tasks, including language translation, text summarization, sentiment analysis, question answering, and named

- entity recognition.
- They have also been extended to other domains beyond NLP, such as image processing, time series analysis, and reinforcement learning.

Pre-trained Models:

- Large-scale pre-trained transformer models, such as BERT (Bidirectional Encoder Representations from Transformers) and GPT (Generative Pre-trained Transformer), have been released by organizations like Google and OpenAI.
- These pre-trained models, trained on massive text corpora, can be fine-tuned on specific downstream tasks with relatively small amounts of task-specific data, achieving state-of-the-art performance.

LLMs:

LLMs, or Large Language Models, are a type of artificial intelligence model that can understand and generate human-like text at a large scale. These models are typically based on deep learning architectures, such as transformers, and are trained on vast amounts of text data using techniques like unsupervised learning. LLMs have achieved remarkable success in various natural language processing (NLP) tasks, including language translation, text generation, sentiment analysis, question answering, and more.

Here's how applications based on LLMs are currently built:

Pretraining: LLMs are pretrained on large corpora of text data, often consisting of millions or even billions of sentences. During pretraining, the model learns to predict the next word in a sequence given the previous words. This process enables the model to capture semantic and syntactic patterns in the text.

Fine-Tuning: After pretraining, LLMs are fine-tuned on specific tasks or domains to adapt them to particular applications. Fine-tuning involves training the model on task-specific data with labeled examples to optimize its parameters for the target task. For example, a sentiment analysis application might fine-tune an LLM on a dataset of customer reviews labeled with sentiment labels (positive/negative).

Integration: Once the model is pretrained and fine-tuned, it can be integrated into various applications and platforms. This integration may involve developing APIs (Application Programming Interfaces) or SDKs (Software Development Kits) that allow developers to easily incorporate the model's capabilities into their software systems.

Deployment: Applications based on LLMs are deployed to production environments where they can interact with users in real-time. Deployment involves setting up infrastructure to host and serve the model, ensuring scalability, reliability, and low-latency performance.

Feedback Loop: As users interact with applications based on LLMs, feedback is collected and used to improve the model over time. This feedback loop helps refine the model's predictions, enhance its performance, and adapt it to evolving user needs and preferences.

Examples of applications built on LLMs include:

- Chatbots and virtual assistants that can engage in natural language conversations with users.
- Content generation tools that can automatically write articles, summaries, or product descriptions.
- Language translation services that can translate text between multiple languages accurately.
- Sentiment analysis tools that can analyze the sentiment of text data, such as social media posts or customer reviews.

Overall, LLMs have revolutionized the field of natural language processing and are powering a wide range of intelligent applications that leverage the power of human-like text generation and understanding.

Langchain:

Language Models Integration: LangChain integrates state-of-the-art language models, such as Large Language Models (LLMs) based on deep learning architectures like transformers. These models enable LangChain to offer advanced language-related services, including translation, interpretation, summarization, and sentiment analysis.

8 - GenAI and LLMs

martes, 30 de abril de 2024 10:08

What is Gen AI ?

GenAI is a term that typically refers to artificial intelligence (AI) systems or technologies designed to generate content, such as text, images, music, or even entire pieces of software, using algorithms and data inputs. It's basically smart computer programs that can create stuff on their own, often mimicking or augmenting human creativity and problem-solving abilities.

What is a Language Model ?

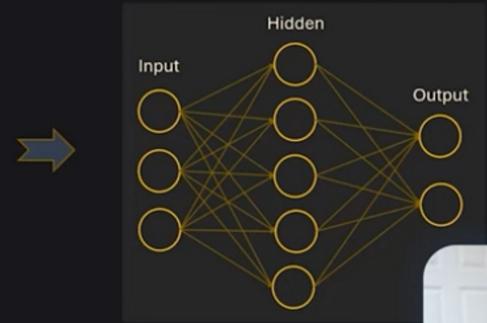
A Language Model is an AI model that can predict the next word (or set of words) for a given sequence of words.

Language Models are trained using a technique called self-supervised learning. It can be trained like so: Imagine we have a text like the one that follows. We will use this text to create a "problem", like try to fill the gaps in a given text. We will create this training pairs (problem, answer) and feed them to a neural network for training.

Self-supervised learning is like learning from hints you give yourself instead of having a teacher tell you the answers directly. The idea is to find patterns or connections in the data itself without needing someone to label it for you. It's a way for AI to learn on its own by making predictions about the data it's given.

Dadabhai Naoroji was the first Indian nationalist to embrace Swaraj as the destiny of the nation.^[55] Bal Gangadhar Tilak deeply opposed a British education system that ignored and defamed India's culture, history, and values. He resented the denial of freedom of expression for nationalists, and the lack of any voice or role for ordinary Indians in the affairs of their nation. For these reasons, he considered Swaraj as the natural and only solution. His popular sentence "Swaraj is my birthright, and I shall have it" became the source of inspiration for Indians.

first Indian _____ to embrace → nationalist
•
Indian nationalist to _____ Swaraj → embrace
Embrace Swaraj as the _____ of → destiny

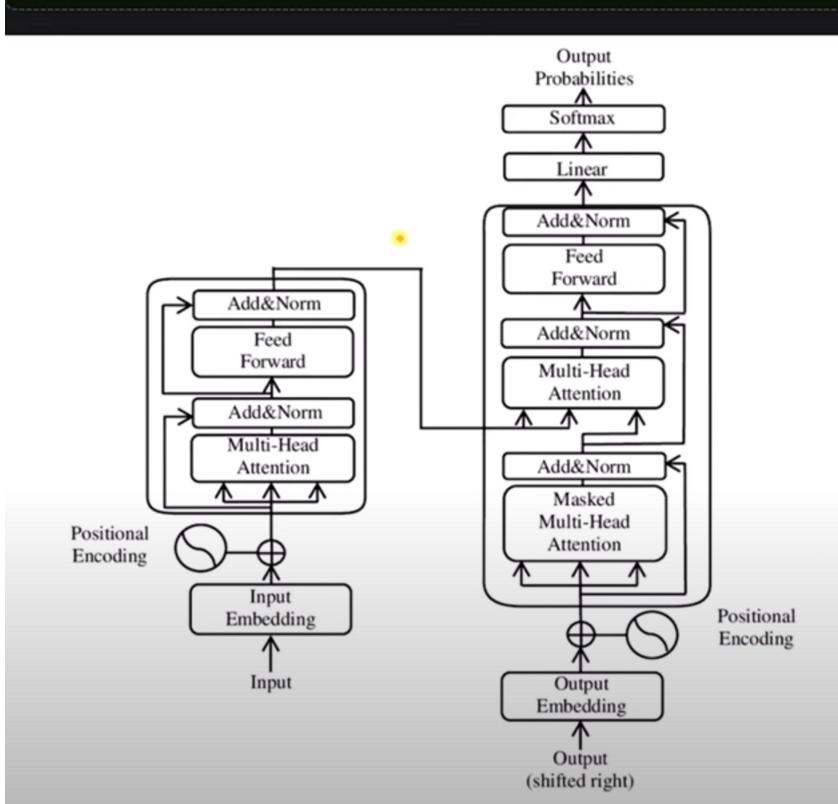


LLMs

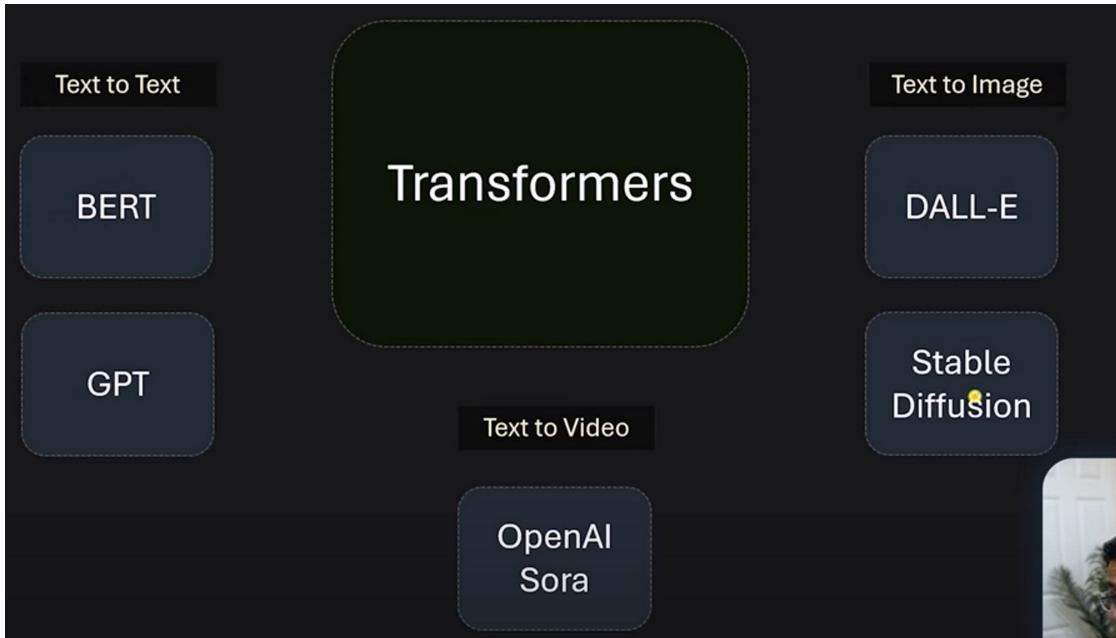
When we have very deep neural networks that receive many data from many different sources, we obtain a Large Language Model LLM. GPT4 is one of them, with 100 75 billion parameters.

LLMs use neural networks of type transformers .

Transformers



Current best transformer models.



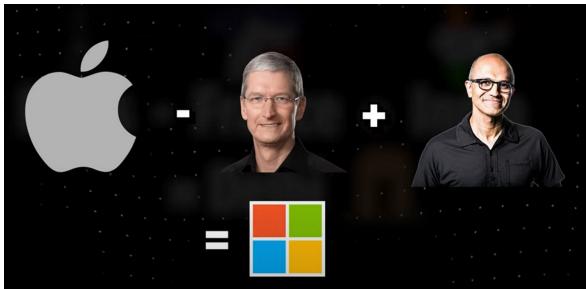
As well as GPT we have PaLM2 and LLaMA as state-of-the-art models atm.

On top of statistical predictions LLMs use Reinforcement Learning with Human Feedback. For the same input, an LLM can produce multiple answers, when we introduce human feedback, we can tell which answer is more accurate which is then reintroduced in the continuous learning of the algorithms to provide better answers.

LLMs use Embeddings.

What is an Embedding?

An Embedding is a numerical representation of text in the form of a vector, such as you can capture the meaning of that text. This also allows for operations like the following to be possible:



Embeddings are stored in a Vector Database that allows for efficient search of these embeddings to be able to provide the expected answers by the user.

Vector Databases use a technique called **Semantic Search** to retrieve the right embeddings. So search in these databases is not done by searching the exact 'keyword' matching but understanding the intent of the user query and using the context to perform the search.

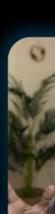
How do embeddings look like:

	Context: <i>Revenue of Apple</i>	Context: <i>Calories in Apple</i>	Context: <i>Nutrition in Orange</i>
Apple	Apple	Orange	
related_to_phones	$\begin{bmatrix} 1 \\ 0 \\ 1 \\ 82 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 95 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 50 \end{bmatrix}$
is_location			
has_stock			
revenue			
is_fruit			
calories			

It can be seen that apple and orange, which are both fruits, have similar embeddings. The technique used to pass from words to embeddings is called **Word2Vec**.

Once we have a set of words embedded, we can perform the following types of operations:

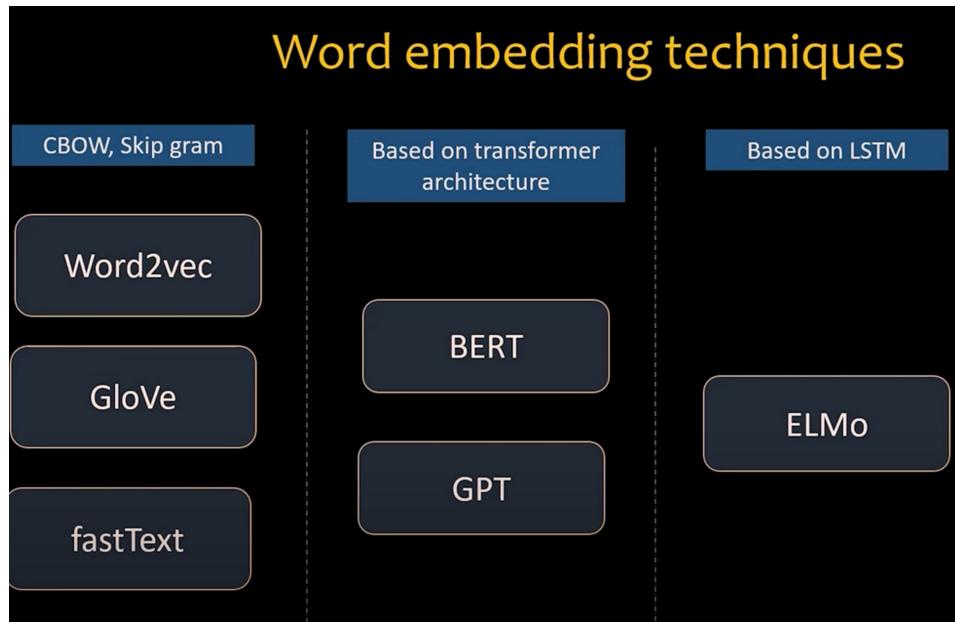
	battle	horse	king	man	queen	..	woman
authority	0	0.01	1	0.2	1	...	0.2
event	1	0	0	0	0	...	0
has tail?	0	1	0	0	0	...	0
rich	0	0.1	1	0.3	1	...	0.2
gender	0	1	-1	-1	1	...	1

King	- man	+ woman	Queen
$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ -1 \end{bmatrix}$	$\begin{bmatrix} 0.2 \\ 0 \\ 0 \\ 0.3 \\ -1 \end{bmatrix}$	$\begin{bmatrix} 0.2 \\ 0 \\ 0 \\ 0.2 \\ 1 \end{bmatrix}$	$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0.9 \\ 1 \end{bmatrix}$
=		\sim	

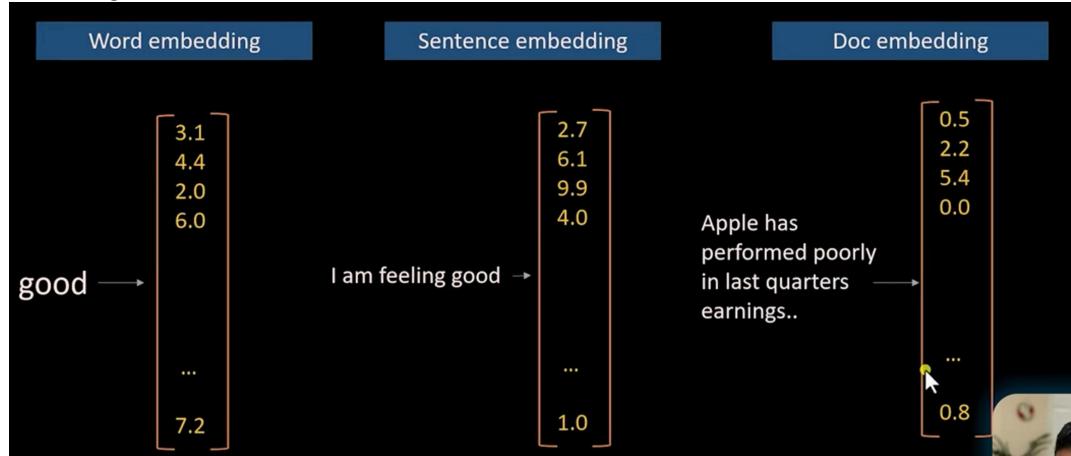
It's important to understand that this example uses hand-crafted features for better understanding, but that in reality, complex

statistical techniques are used to auto-derive the features of text.

Current state-of-the-art techniques for this purpose:



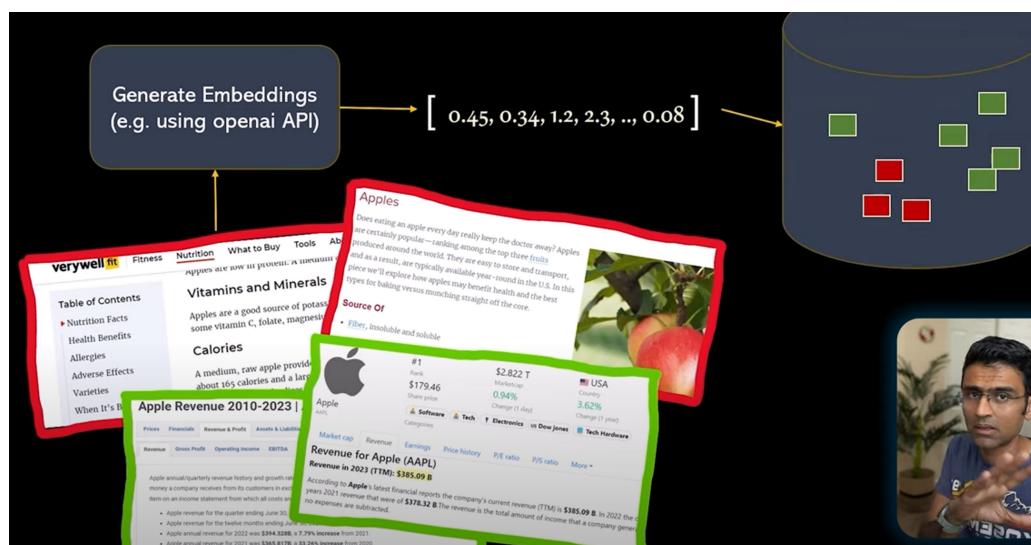
Embeddings are not limited to words, we can embedd sentences and docs too.



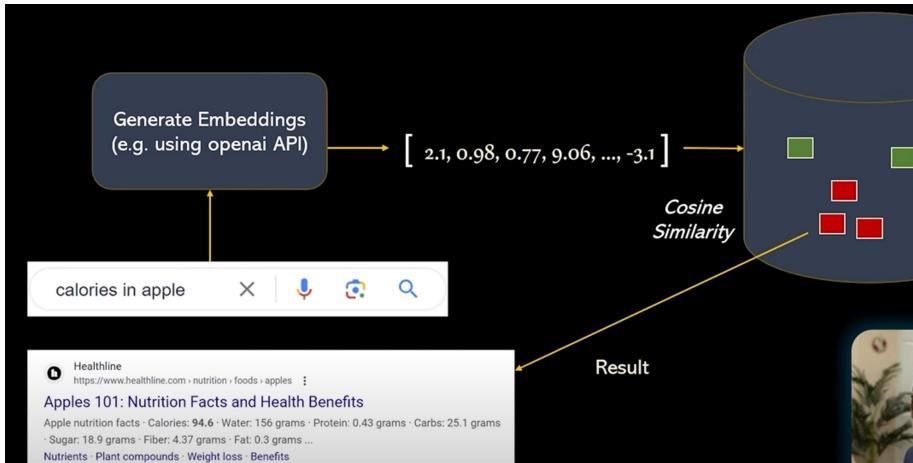
Vector databases:

Vector databases are used to store and retrieve the embeddings.

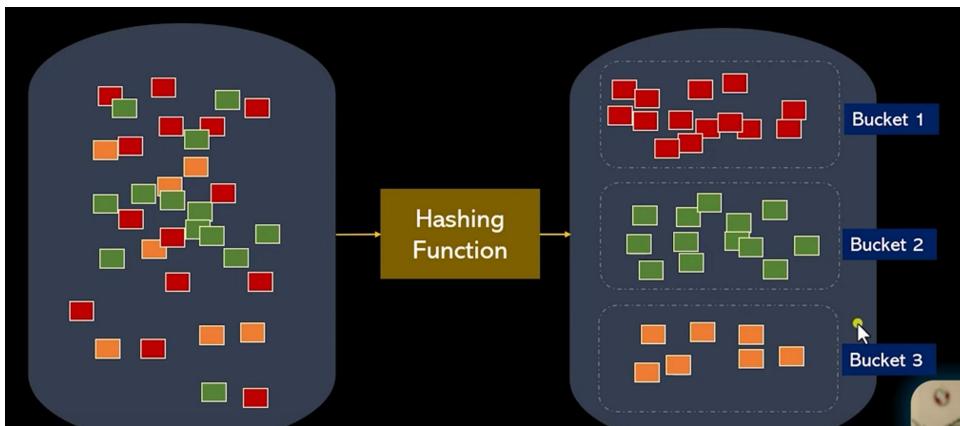
So in order to store some embeddings, we first have to generate the embeddings from the given data and store them, into a database.



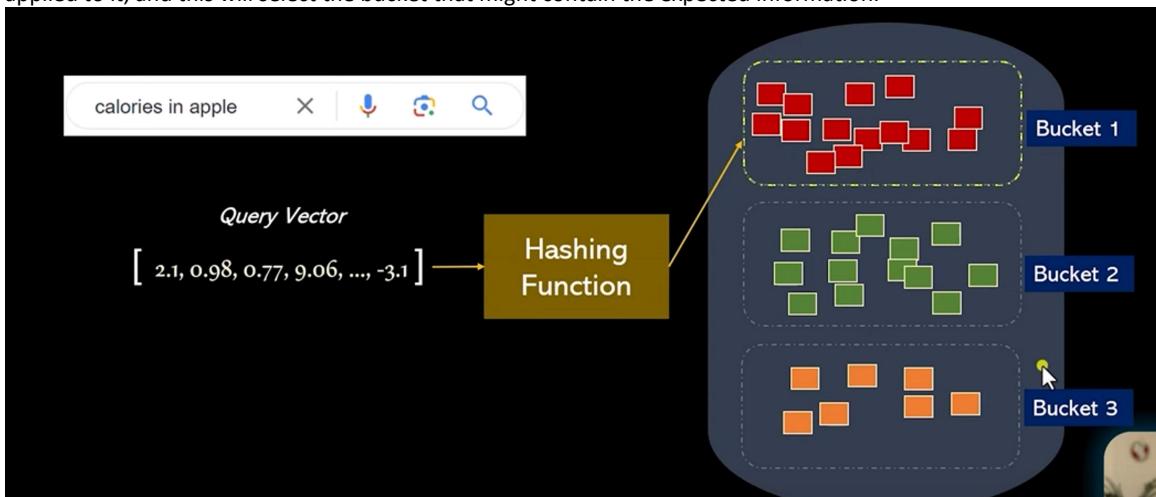
When we have a search query, we have to first generate an embedding for that query and we will compare this embedding with the ones stored into the database in order to retrieve the most similar from it.



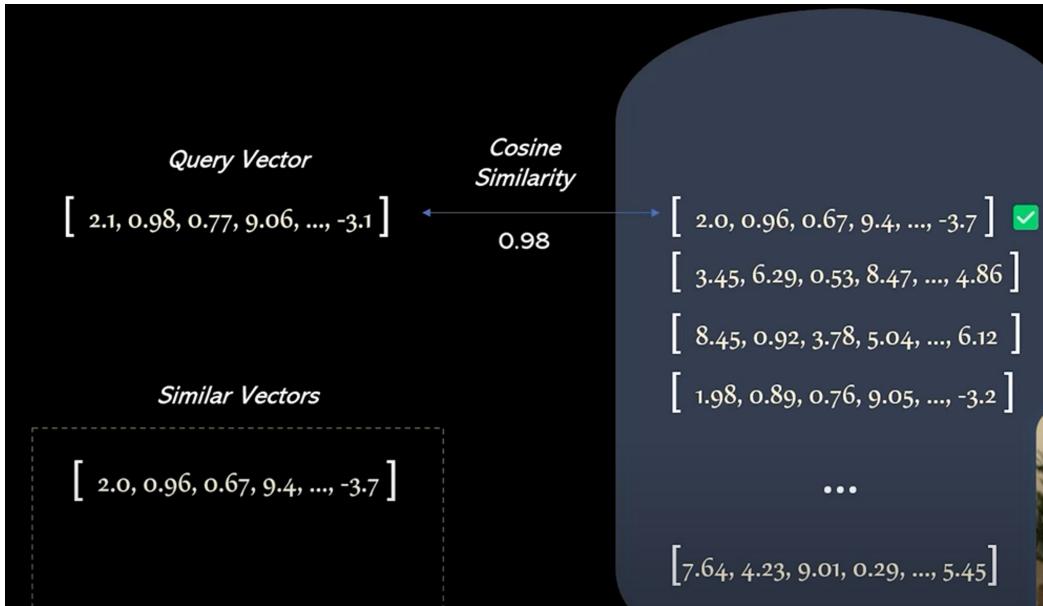
It is important to consider that these databases will have millions of stored embeddings, therefore linear search is not an option due to extremely long computation times. To improve this retrieval, embeddings are stored inside "buckets" of similar elements. This allocation is done by applying a Hashing function to the embeddings to place them in their corresponding bucket. This technique is called **Locality Sensitive Hashing LSH**.



Now, the same procedure will be done when we want to perform a query. The query will be embedded, a hashing function will be applied to it, and this will select the bucket that might contain the expected information.



Now that we only have to search inside a bucket, linear search becomes a viable option. We can use Cosine similarity to retrieve a close match and place it on the results.

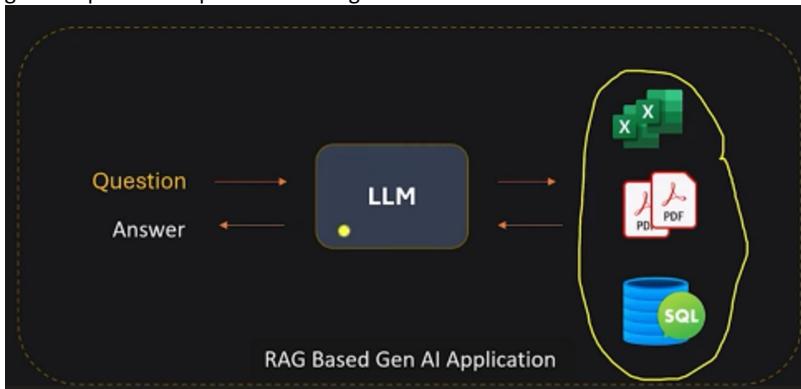


There are many search techniques at the moment LSH is just one of them.

Vector databases provides **Fast Search** and provides **Optimized Storage**. Vector databases often use specialized data structures and indexing techniques to support fast similarity searches, nearest neighbor queries, and other operations common in applications like recommendation systems, search engines, and natural language processing.

Retrieval Augmented Generation (RAG):

RAG is a framework used in Natural Language Generation (NLG) tasks, particularly in LLMs, for retrieving answers from owned data given a question or perform a task given an instruction.



RAG stands for "Retrieve, Aggregate, Generate." It's a framework used in natural language generation (NLG) tasks, particularly in large-scale language models like GPT (Generative Pre-trained Transformer). Here's what each step involves:

Retrieve: This step involves retrieving relevant information or context from a given knowledge source, such as a database, documents, or the internet. The retrieved information serves as the basis for generating the output.

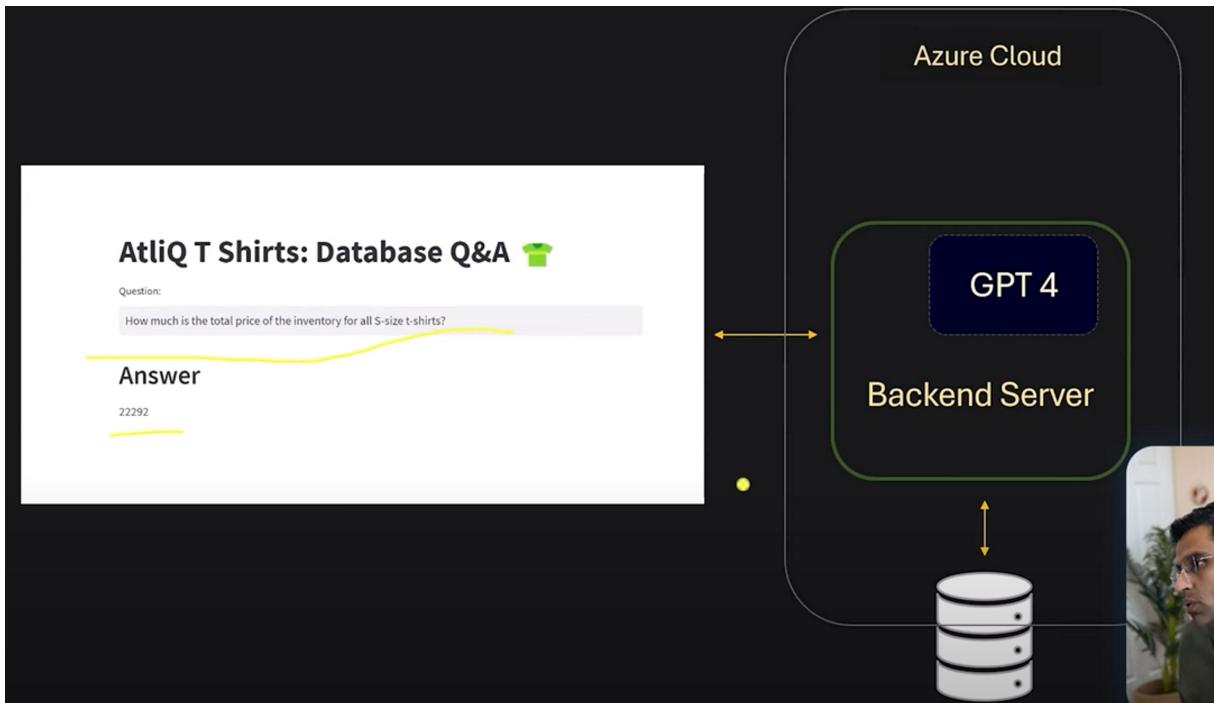
Aggregate: After retrieving relevant information, the next step is to aggregate or combine this information in a meaningful way. This may involve summarizing, condensing, or organizing the retrieved information to focus on the key points or concepts.

Generate: Finally, the aggregated information is used to generate the output text or response. This could involve generating a natural language response, answering a question, completing a task, or providing a summary based on the retrieved and aggregated information.

Overall, the RAG framework helps guide the process of generating natural language text by first retrieving relevant information, aggregating it into a coherent form, and then using it to generate the final output.

Currently RAG solutions tailored to specific companies, would normally finetune current LLMs like GPT4 or LLAMA to their specific datasets or data sources.

So imagine we have a retail company and we want to query our data like so: Here our GPT4 will be like the brain, it will transform the question into an SQL query, the backend server will perform this retrieval from the database and forwarding the info and all of this could be handled into a cloud for your private network use only so all your data is protected.



Langchain

So what happens if we want to switch and try different models without that affecting our Application or having to rewrite any code? We can use Langchain for this purpose. Langchain is a python framework that is used to build GenAI applications on top of LLMs.

