# Machine Learning Study Guide

January 10, 2024

# Contents

# 1 ML Algorithms' Categories

## 1.1 Parametric and Non-Parametric Algorithms

**Parametric Algorithms** assume that the data follows a specific distribution with a fixed number of parameters. Examples include linear regression and Gaussian Naive Bayes. Advantages include simplicity, speed, and less data requirement. However, they may suffer from limited flexibility and bias if the chosen model is incorrect.

**Non-Parametric Algorithms** do not assume any specific distribution for the data, learning the structure as they train. Examples include k-nearest neighbors and decision trees. They offer flexibility and fewer assumptions but may be complex, computationally intensive, and require more data.

Parametric algorithms are useful when the data distribution is clear and limited, while non-parametric algorithms are suitable when there's no prior knowledge about the data distribution or when the relationship between features and target variables is complex. Factors such as dataset size, computational resources, model interpretability, and problem nature should be considered when choosing between the two types of algorithms.

## 1.2 Supervised, Self-Supervised, and Unsupervised Learning

**Supervised Learning** involves training a model on labeled data, mapping inputs to target outputs. This approach excels in tasks such as image classification and speech recognition. However, it requires a substantial amount of labeled data, which can be costly and time-consuming to acquire.

**Self-Supervised Learning** leverages the structure of input data to create pseudo-labels, guiding the learning process. By solving a pretext task, the model learns representations transferable to various downstream tasks. While this approach capitalizes on unlabeled data, the quality of learned representations depends on the choice of pretext task, and it may be computationally expensive.

**Unsupervised Learning** focuses on capturing the underlying structure or distribution of data without relying on labels. Techniques such as clustering and dimensionality reduction are common in unsupervised learning. This paradigm also utilizes unlabeled data, but its applicability to specific tasks depends on the inherent structure of the data and the algorithm's ability to capture it.

# 2 Distance Measures in Machine Learning

Distance measures play a crucial role in many machine learning algorithms, such as clustering, nearest neighbor search, and dimensionality reduction. They quantify the similarity or dissimilarity between data points, which can be used to inform decisions made by algorithms.

## 2.1 L1 Distance (Manhattan Distance)

The L1 distance, also known as Manhattan distance or Taxicab distance, is the sum of the absolute differences between the coordinates of two points in a vector space. Given two vectors $\mathbf{x}$ and $\mathbf{y}$, the L1 distance is defined as:

$$d_{L1}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{n} |x_i - y_i|$$

## 2.2 L2 Distance (Euclidean Distance)

The L2 distance, also known as Euclidean distance, is the square root of the sum of the squared differences between the coordinates of two points in a vector space. Given two vectors $\mathbf{x}$ and $\mathbf{y}$, the L2 distance is defined as:

$$d_{L2}(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^{n} (x_i - y_i)^2}$$

## 2.3 Cosine Distance

Cosine distance is a measure of the angle between two vectors, and it is often used to compare the direction of the vectors rather than their magnitudes. Given two vectors $\mathbf{x}$ and $\mathbf{y}$, the cosine distance is defined as:

$$d_{\text{cosine}}(\mathbf{x}, \mathbf{y}) = 1 - \frac{\mathbf{x} \cdot \mathbf{y}}{||\mathbf{x}|| \cdot ||\mathbf{y}||}$$

## 2.4 Dot Product

The dot product is a measure of the projection of one vector onto another. It is not a distance measure in the strict sense, but it can be used to compute the similarity between two vectors. Given two vectors $\mathbf{x}$ and $\mathbf{y}$, the dot product is defined as:

$$\mathbf{x} \cdot \mathbf{y} = \sum_{i=1}^{n} x_i \cdot y_i$$

## 2.5 KL Divergence

Kullback-Leibler (KL) Divergence is a measure of the difference between two probability distributions. It is often used to compare the similarity between probability distributions in machine learning tasks. Given two probability distributions $P$ and $Q$, the KL divergence is defined as:

$$D_{KL}(P||Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

## 2.6 Comparison and Summary

- **L1 distance** is robust to outliers and less sensitive to small differences between coordinates. It is often used in sparse data settings and can be more interpretable than other distance measures.

- **L2 distance** is the most commonly used distance measure and is sensitive to differences in magnitude. It can be affected by outliers and may not be suitable for high-dimensional data due to the curse of dimensionality.

- **Cosine distance** focuses on the angle between vectors, making it less sensitive to differences in magnitude. It is often used in text mining and information retrieval tasks, where the direction of the vectors (e.g., term frequency) is more important than their magnitudes.

- **Dot Product** measures the projection of one vector onto another, and although it is not a distance measure itself, it can be used to assess similarity between vectors. It is commonly used in tasks involving high-dimensional data, such as collaborative filtering and image recognition.

- **KL Divergence** measures the difference between two probability distributions and is particularly useful for comparing the similarity between distributions in tasks such as language modeling and anomaly detection. However, it is not symmetric, which means that $D_{KL}(P||Q) \neq D_{KL}(Q||P)$.

The choice of distance measure depends on the specific problem and the nature of the data at hand. L1 and L2 distances are general-purpose distance measures, with L1 being more robust to outliers and L2 being more sensitive to differences in magnitude. Cosine distance is useful for comparing the direction of vectors, making it well-suited for text and high-dimensional data. Dot product can be used to assess similarity in high-dimensional spaces, while KL Divergence is more specialized for comparing probability distributions. It is essential to consider these characteristics when selecting a distance measure for a particular machine learning task.

# 3 Classic ML Algorithms

## 3.1 Linear Regression

Linear regression is a parametric supervised learning algorithm used for modeling the relationship between a dependent variable (target) and one or more independent variables (features). The goal is to find the best-fitting line (or hyperplane in the case of multiple features) that minimizes the sum of the squared differences between the actual and predicted values.

### 3.1.1 Formula

The formula for a simple linear regression (one feature) is:

$$y = \beta_0 + \beta_1 x$$

In multiple linear regression (multiple features), the formula is:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_n x_n$$

where:

- $y$ is the dependent variable (target)
- $x_i$ are the independent variables (features)
- $\beta_0$ is the intercept
- $\beta_i$ are the coefficients

### 3.1.2 Ordinary Least Squares (OLS)

OLS is a popular method for estimating the coefficients in linear regression. It aims to minimize the sum of squared residuals:

$$\min_{\beta_0,\beta_1,\dots,\beta_n} \sum_{i=1}^{m} (y_i - (\beta_0 + \beta_1 x_{i1} + \cdots + \beta_n x_{in}))^2$$

The closed-form solution for the coefficients can be expressed as:

$$\boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

where:

- $\mathbf{X}$ is the feature matrix with a column of ones for the intercept
- $\mathbf{y}$ is the target vector
- $\boldsymbol{\beta}$ is the vector of coefficients

### 3.1.3 Advantages

1. Simplicity: Linear regression is easy to understand and implement.
2. Interpretability: The coefficients provide insights into the relationship between features and the target.
3. Speed: It is computationally efficient.

### 3.1.4 Disadvantages

1. Limited complexity: Assumes a linear relationship between features and target, which may not always be true.

2. Sensitivity to outliers: Outliers can significantly impact the model's performance.

3. Multicollinearity: The presence of highly correlated features can lead to unstable estimates of coefficients.

### 3.1.5 Context

Linear regression is widely used for problems involving continuous target variables when the relationship between features and the target is assumed to be linear. It is a good starting point for understanding more complex models and serves as a baseline for evaluating other regression techniques.

## 3.2 Logistic Regression

Logistic regression is a parametric supervised learning algorithm used for binary classification problems. It models the probability of a data point belonging to a specific class using a logistic function (sigmoid function). The algorithm estimates the coefficients of the linear combination of features that best predict the log-odds of the target class.

### 3.2.1 Sigmoid Function

The logistic (sigmoid) function maps the input to a value between 0 and 1, representing the probability of the target class:

$$P(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \cdots + \beta_n x_n)}}$$

where:

- $P(y = 1|\mathbf{x})$ is the probability of the target class given the features

- $x_i$ are the independent variables (features)

- $\beta_0$ is the intercept

- $\beta_i$ are the coefficients

### 3.2.2 Cost Function

The cost function for logistic regression is derived from the cross-entropy loss. For a single data point, the cost function is:

$$J(\boldsymbol{\beta}) = -[y \log(P(y = 1|\mathbf{x})) + (1 - y) \log(1 - P(y = 1|\mathbf{x}))]$$

For a dataset with $m$ data points, the cost function is the average of the individual costs:

$$J(\boldsymbol{\beta}) = -\frac{1}{m} \sum_{i=1}^{m} [y_i \log(P(y_i = 1|\mathbf{x}_i)) + (1 - y_i) \log(1 - P(y_i = 1|\mathbf{x}_i))]$$

### 3.2.3   Gradient Descent

To minimize the cost function, gradient descent is used to update the coefficients iteratively. The update rule for gradient descent is:

$$\boldsymbol{\beta} := \boldsymbol{\beta} - \alpha \nabla J(\boldsymbol{\beta})$$

where:

- $\boldsymbol{\beta}$ is the vector of coefficients

- $\alpha$ is the learning rate

- $\nabla J(\boldsymbol{\beta})$ is the gradient of the cost function with respect to the coefficients

### 3.2.4   Advantages

1. Probabilistic interpretation: Logistic regression provides the probability of a data point belonging to a specific class.

2. Ease of implementation: It is relatively simple to implement and can be regularized to prevent overfitting.

3. Computationally efficient: The algorithm is not computationally intensive compared to more complex models.

### 3.2.5   Disadvantages

1. Linearity assumption: Assumes a linear relationship between features and the log-odds of the target class, which may not always be true.

2. Binary classification: It is primarily used for binary classification problems, although it can be extended to multi-class classification using techniques like one-vs-rest or one-vs-one.

3. Performance: More complex models like support vector machines or neural networks can outperform logistic regression in some cases.

### 3.2.6  Context

Logistic regression is widely used for binary classification problems where the relationship between features and the target class can be assumed to be linear. It is a good starting point for classification problems and serves as a baseline for evaluating other classification techniques.

## 3.3  Support Vector Machines

Support Vector Machines (SVM) is a non-parametric supervised learning algorithm used for binary and multiclass classification problems. SVM aims to find the optimal hyperplane that maximizes the margin between the two classes in the feature space. In cases where the data is not linearly separable, SVM can use kernel functions to transform the feature space to a higher-dimensional space where a separating hyperplane can be found.

### 3.3.1  Margin and Support Vectors

The margin is the distance between the closest data points of different classes and the separating hyperplane. The goal of SVM is to maximize this margin. The data points on the edge of the margin are called support vectors.

### 3.3.2  Optimization Problem

The optimization problem for SVM can be formulated as follows:

$$\min_{\mathbf{w},b} \frac{1}{2}||\mathbf{w}||^2$$

subject to:

$$y_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1, \quad i = 1, \ldots, m$$

where:

- $\mathbf{w}$ is the normal vector to the hyperplane
- $b$ is the bias term
- $\mathbf{x}_i$ are the feature vectors
- $y_i \in -1, 1$ are the target class labels
- $m$ is the number of data points

### 3.3.3  Kernel Functions

In cases where the data is not linearly separable, SVM can use kernel functions to implicitly map the data to a higher-dimensional space. Popular kernel functions include:

- Linear kernel: $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$

- Polynomial kernel: $K(\mathbf{x}_i, \mathbf{x}_j) = (\gamma \mathbf{x}_i^T \mathbf{x}_j + r)^d$

- Radial basis function (RBF) kernel: $K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma ||\mathbf{x}_i - \mathbf{x}_j||^2)$

- Sigmoid kernel: $K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\gamma \mathbf{x}_i^T \mathbf{x}_j + r)$

where $\gamma$, $r$, and $d$ are kernel-specific parameters.

### 3.3.4 Advantages

1. Effective in high-dimensional spaces: SVM can handle a large number of features.

2. Robustness: Only the support vectors affect the decision boundary, making the model robust to outliers.

3. Flexibility: Kernel functions allow SVM to model non-linear decision boundaries.

### 3.3.5 Disadvantages

1. Scalability: SVM is not well-suited for very large datasets due to its computational complexity.

2. Model interpretability: SVM models can be difficult to interpret, especially with non-linear kernel functions.

3. Sensitive to hyperparameters: The choice of kernel function and hyperparameters can have a significant impact on the model's performance.

### 3.3.6 Context

SVM is widely used for classification problems when the relationship between features and the target class is complex or non-linear. It is effective in high-dimensional spaces and can handle a large number of features.

## 3.4 Decision Trees

Decision trees are non-parametric supervised learning algorithms used for both classification and regression tasks. They model the relationship between features and the target variable by recursively partitioning the input space and learning simple decision rules based on the features. Decision trees can be visualized as a flowchart-like structure, with internal nodes representing decisions based on feature values and leaf nodes representing the predicted target values.

### 3.4.1   Nodes and Leaves

- **Internal nodes**: Represent decisions based on feature values. Each internal node tests a specific feature and has multiple branches, each corresponding to a possible value or range of values for that feature.

- **Leaf nodes**: Represent the predicted target value (class label for classification or continuous value for regression) associated with a specific combination of feature values.

### 3.4.2   Training Algorithms

There are various algorithms for constructing decision trees, such as ID3, C4.5, and CART. They generally use a greedy, top-down approach:

1. Select the best feature and split value based on a criterion (e.g., information gain, Gini impurity).

2. Partition the dataset into subsets based on the chosen feature and split value.

3. Recursively apply steps 1 and 2 on each subset until a stopping condition is met (e.g., maximum depth, minimum samples per leaf).

### 3.4.3   Inference

Inference in decision trees involves traversing the tree from the root to a leaf node. At each internal node, the decision is made based on the value of the corresponding feature in the input data. The final prediction is obtained from the leaf node reached at the end of the traversal.

### 3.4.4   Pruning

Pruning is a technique used to reduce the complexity of a decision tree and mitigate overfitting. There are two main types of pruning:

- **Pre-pruning**: Stop growing the tree before it perfectly fits the training data, based on a stopping condition (e.g., maximum depth, minimum samples per leaf).

- **Post-pruning**: First grow the tree to fit the training data, then iteratively remove branches that do not improve the tree's performance on a validation dataset.

### 3.4.5   Advantages

1. Interpretability: Decision trees are easy to understand and visualize, making them highly interpretable.

2. Minimal data preprocessing: Decision trees do not require feature scaling or normalization, and they can handle both numerical and categorical features.

3. Non-linear relationships: Decision trees can model complex, non-linear relationships between features and the target variable.

### 3.4.6 Disadvantages

1. Overfitting: Decision trees can easily overfit the training data, especially when they are deep or have few samples per leaf.

2. Instability: Small changes in the training data can lead to significantly different trees, making decision trees sensitive to noise and variance.

3. Greedy nature: The greedy, top-down approach used in tree construction may not always result in the optimal tree, as it does not consider all possible trees.

### 3.4.7 Context

Decision trees are widely used for problems where interpretability is important, and they can handle a mix of numerical and categorical features. They serve as the foundation for more advanced ensemble methods like random forests and gradient boosting machines.

## 3.5 K-Means Clustering

K-means clustering is a widely-used unsupervised learning algorithm for partitioning a dataset into a specified number of clusters ($K$). The algorithm aims to minimize the within-cluster sum of squares (inertia) by iteratively updating the cluster centroids and cluster assignments.

### 3.5.1 Algorithm

The K-means clustering algorithm proceeds as follows:

1. Initialize the $K$ cluster centroids randomly by selecting $K$ data points from the dataset.

2. Assign each data point to the nearest centroid.

3. Update the centroids by computing the mean of all data points assigned to each centroid.

4. Repeat steps 2 and 3 until convergence (i.e., the centroids do not change significantly or a maximum number of iterations is reached).

### 3.5.2 Advantages

Some advantages of K-means clustering include:

1. Simplicity and ease of implementation.

2. Efficient computation, making it suitable for large datasets.

3. Scalability to high-dimensional data when combined with dimensionality reduction techniques.

4. Wide applicability across various domains and tasks.

### 3.5.3  Disadvantages

However, K-means clustering also has several disadvantages:

1. The need to specify the number of clusters ($K$) beforehand, which may not be known a priori.

2. Sensitivity to initial centroid placement, potentially leading to different clustering results for different initializations.

3. Convergence to local optima, which can be mitigated by running the algorithm multiple times with different initializations.

4. The assumption of equal-sized and spherical clusters, which may not hold for real-world datasets.

5. Sensitivity to outliers and noise, which can be addressed by using robust variants of K-means or by preprocessing the data.

## 3.6  Hierarchical Clustering

Hierarchical clustering is an unsupervised learning algorithm that builds a tree-like structure (dendrogram) to represent the nested grouping of data points based on their similarity. There are two main approaches to hierarchical clustering: agglomerative (bottom-up) and divisive (top-down).

### 3.6.1  Algorithm

The agglomerative hierarchical clustering algorithm proceeds as follows:

1. Treat each data point as a singleton cluster.

2. Compute the pairwise distances between all clusters.

3. Merge the two clusters with the smallest distance, creating a new cluster.

4. Update the distance matrix by calculating the distances between the new cluster and the remaining clusters using a linkage criterion (e.g., single, complete, average, or Ward's linkage).

5. Repeat steps 2 to 4 until there is only one cluster left, which contains all the data points.

The divisive hierarchical clustering algorithm is less common and works by initially treating the entire dataset as a single cluster, then recursively splitting clusters into smaller clusters until each data point forms its own cluster.

### 3.6.2 Advantages

Some advantages of hierarchical clustering include:

1. No need to specify the number of clusters beforehand, as the dendrogram can be cut at different levels to obtain different clustering results.

2. The dendrogram provides a visual representation of the clustering process, aiding interpretation and understanding of the data structure.

3. Robustness to initialization, as the algorithm is deterministic and does not depend on random initial conditions.

4. Ability to handle non-spherical and non-convex clusters by using appropriate distance and linkage measures.

### 3.6.3 Disadvantages

However, hierarchical clustering also has several disadvantages:

1. Higher computational complexity compared to K-means, making it less suitable for large datasets.

2. Irreversibility of the clustering process, meaning that once clusters are merged (agglomerative) or split (divisive), they cannot be undone in subsequent steps.

3. Sensitivity to the choice of distance and linkage measures, which can affect the clustering results.

4. Difficulty in determining the optimal number of clusters and the appropriate cutting level in the dendrogram.

## 3.7 DBSCAN

Density-Based Spatial Clustering of Applications with Noise (DBSCAN) is an unsupervised clustering algorithm that groups data points based on their density in the feature space. Unlike K-means and hierarchical clustering, DBSCAN does not require the user to specify the number of clusters beforehand and is capable of handling noise and identifying clusters of arbitrary shapes.

### 3.7.1 Algorithm

1. For each unvisited data point, retrieve its neighborhood, which consists of all points within a specified distance $\epsilon$ (eps) from the point.

2. If the neighborhood contains at least a minimum number of points (minPts), consider the data point as a core point and form a cluster. Recursively expand the cluster by adding all directly density-reachable points and their respective density-reachable points.

3. If the data point is not a core point and not density-reachable from any core point, label it as noise.

4. Repeat steps 1 to 3 until all data points have been visited and assigned to a cluster or labeled as noise.

### 3.7.2 Advantages

1. No need to specify the number of clusters beforehand, as the algorithm automatically determines the number of clusters based on the density of the data points.

2. Ability to identify clusters of arbitrary shapes, as opposed to K-means, which is limited to spherical clusters.

3. Robustness to noise, as the algorithm explicitly labels and separates noise points from the clusters.

4. Deterministic clustering results, as long as the order of processing the data points is fixed.

### 3.7.3 Disadvantages

1. Sensitivity to the choice of $\epsilon$ and minPts parameters, which can significantly affect the clustering results.

2. Difficulty in determining the optimal values for $\epsilon$ and minPts, especially for datasets with varying densities.

3. Poor performance on high-dimensional data due to the curse of dimensionality, which affects the density estimation.

4. Inefficient for large datasets, as the algorithm has a quadratic time complexity in the worst case.

DBSCAN can find arbitrarily shaped clusters, is robust to noise, and does not require the number of clusters to be specified beforehand. However, it is sensitive to the choice of hyperparameters ($\varepsilon$ and MinPts) and may not perform well when clusters have varying densities.

## 3.8 PCA (Principal Component Analysis)

Principal Component Analysis (PCA) is a dimensionality reduction technique used to transform a high-dimensional dataset into a lower-dimensional space while preserving the maximum amount of variance in the data. PCA is particularly useful for visualization, data compression, and improving the efficiency of other machine learning algorithms.

The main steps involved in PCA are as follows:

1. Standardize the data by centering the features around their mean and scaling them to have unit variance.

2. Compute the covariance matrix of the standardized data.

3. Calculate the eigenvalues and eigenvectors of the covariance matrix.

4. Sort the eigenvectors in descending order of their corresponding eigenvalues.

5. Select the first $k$ eigenvectors with the highest eigenvalues to form a $d \times k$ dimensional projection matrix, where $d$ is the original dimensionality and $k$ is the target dimensionality.

6. Transform the original data by multiplying it with the projection matrix to obtain the lower-dimensional representation.

PCA assumes that the principal components are linear combinations of the original features and works well when the data has a linear structure. However, it may not capture complex, nonlinear relationships between the features.

## 3.9 t-SNE (t-Distributed Stochastic Neighbor Embedding)

t-SNE is a non-linear dimensionality reduction technique particularly suited for visualizing high-dimensional data in a 2D or 3D space. It is based on the concept of preserving the local distances between instances in the lower-dimensional space, making it effective for revealing clusters and structures in the data.

1. Define a probability distribution over pairs of instances in the high-dimensional space such that similar instances have a higher probability of being picked, while dissimilar instances have a lower probability. This is usually done by modeling the pairwise similarities using a Gaussian distribution.

2. Define a probability distribution over pairs of instances in the lower-dimensional space using a t-distribution with one degree of freedom (also known as the Student's t-distribution).

3. Minimize the Kullback-Leibler (KL) divergence between the two probability distributions with respect to the lower-dimensional embeddings using gradient descent. The KL divergence is a measure of how one probability distribution is different from another.

t-SNE has some advantages over PCA, as it can capture non-linear structures in the data and is more effective at visualizing clusters. However, it is computationally expensive and does not provide an explicit mapping function between the high-dimensional and low-dimensional spaces, making it unsuitable for out-of-sample data.

# 4 Ensemble Methods

Ensemble methods are techniques that combine multiple base learners to create a more powerful and accurate model. They can be used with various types of base learners, such as decision trees, logistic regression, or support vector machines. There are two main types of ensemble methods: bagging and boosting.

## 4.1 Bagging

Bagging (Bootstrap Aggregating) is an ensemble method that trains multiple base learners independently on different subsets of the training data, generated by random sampling with replacement (bootstrap samples). The final prediction is obtained by averaging (for regression) or voting (for classification) the predictions of all base learners.

### 4.1.1 Random Forests

Random Forests is an ensemble learning method that builds multiple decision trees and combines their predictions to improve the overall performance and robustness of the model. It is primarily used for both classification and regression tasks.

**Algorithm** The Random Forests algorithm works as follows:

1. Select a random subset of the data with replacement (bootstrap sample) to build each decision tree.

2. For each decision tree, at each split node, select a random subset of features and use the best split criterion (e.g., Gini impurity or information gain) to split the node.

3. Grow each decision tree to its maximum depth, without pruning.

4. Make a prediction for each decision tree in the ensemble.

5. Combine the predictions from all decision trees using majority voting (classification) or averaging (regression) to obtain the final prediction.

**Advantages** Some advantages of Random Forests include:

1. High predictive accuracy due to the combination of multiple decision trees, which helps to reduce overfitting and improve generalization.

2. Robustness to noise and outliers, as each decision tree is trained on a random subset of the data and features, making the overall model less sensitive to individual data points.

3. Ability to handle large datasets with many features and missing values.

4. Provides feature importance scores, which can be used for feature selection and understanding the underlying structure of the data.

**Disadvantages** However, Random Forests also has several disadvantages:

1. Increased computational complexity and memory requirements due to the construction and storage of multiple decision trees.

2. Slower prediction times compared to single decision trees, as predictions from all trees in the ensemble must be combined.

3. Reduced interpretability compared to single decision trees, as the ensemble model is more complex and harder to visualize.

## 4.2   Boosting

Boosting is an ensemble method that trains multiple base learners sequentially, with each learner focusing on the errors made by the previous one. The final prediction is obtained by combining the predictions of all base learners, typically using a weighted average or a weighted vote.

### 4.2.1   AdaBoost

AdaBoost (Adaptive Boosting) is a boosting algorithm that uses decision stumps (one-level decision trees) as base learners. AdaBoost adjusts the weights of the training samples based on the errors made by the previous learner and trains the next learner on the reweighted samples. The final prediction is a weighted vote of all base learners, with the weights determined by their classification accuracy.

### 4.2.2   Gradient Boosting Machine

Gradient Boosting Machine (GBM) is an ensemble learning method that builds a series of weak learners, often decision trees, and combines them to form a strong learner. The main idea behind GBM is to iteratively train weak learners to correct the errors made by the previous learners in the ensemble. It is mainly used for regression and classification tasks.

**Algorithm**   The Gradient Boosting Machine algorithm works as follows:

1. Initialize the model with a constant prediction value, typically the mean of the target variable for regression or the majority class for classification.

2. For each iteration, fit a weak learner (e.g., a shallow decision tree) to the residuals (the difference between the true target values and the current predictions) of the previous iteration.

3. Update the predictions by adding a fraction of the weak learner's predictions to the current predictions. This fraction, called the learning rate, controls the contribution of each weak learner to the ensemble.

4. Repeat steps 2-3 for a predefined number of iterations or until a convergence criterion is met.

**Advantages**   Some advantages of Gradient Boosting Machine include:

1. High predictive accuracy due to the combination of multiple weak learners, which helps to reduce overfitting and improve generalization.

2. Flexibility, as it can work with various loss functions and different types of weak learners, not limited to decision trees.

3. Robustness to noise and outliers, as each weak learner is trained on the residuals of the previous learners, making the overall model less sensitive to individual data points.

4. Provides feature importance scores, which can be used for feature selection and understanding the underlying structure of the data.

**Disadvantages**    However, Gradient Boosting Machine also has several disadvantages:

1. Increased computational complexity and longer training times due to the sequential nature of the algorithm, which makes it harder to parallelize compared to other ensemble methods like Random Forests.

2. Requires careful tuning of hyperparameters, such as the number of iterations, learning rate, and tree depth, to achieve optimal performance.

3. Reduced interpretability compared to single decision trees, as the ensemble model is more complex and harder to visualize.

**Intuition**    The intuition behind Gradient Boosting Machine is to iteratively improve the model's predictions by focusing on the errors made in previous iterations. By training weak learners on the residuals, GBM forces each new learner to focus on the areas where the current ensemble is performing poorly. This process can be thought of as a way to "correct" the mistakes made by the previous learners, allowing the ensemble to gradually improve its predictions over time. The learning rate parameter helps to control the pace at which the model learns, preventing overfitting and ensuring that each weak learner makes a modest contribution to the overall model.

### 4.2.3   XGBoost

XGBoost (eXtreme Gradient Boosting) is an optimized version of GBM that introduces several improvements, such as:

- Regularization: XGBoost includes L1 and L2 regularization terms in the loss function to prevent overfitting.

- Sparsity-aware: XGBoost can handle sparse data and missing values efficiently.

- Column block: XGBoost uses a compressed memory-efficient representation of the input data, which speeds up training and inference.

- Parallelization: XGBoost parallelizes the construction of decision trees across multiple CPU cores.

## 4.3   Advantages

1. Improved accuracy: Ensemble methods typically achieve better performance than single base learners.

2. Robustness: Ensembles are more robust to noise and outliers, as they combine the predictions of multiple base learners.

3. Flexibility: Ensemble methods can be used with various types of base learners and can be adapted for different tasks (classification, regression).

## 4.4 Disadvantages

1. Increased complexity: Ensemble methods can be more complex and harder to interpret than single base learners, especially with boosting algorithms.

2. Computational cost: Training and inference with ensemble methods can be computationally expensive, as they involve multiple base learners.

3. Risk of overfitting: Boosting algorithms, in particular, can be prone to overfitting if the base learners are too complex or the number of boosting iterations is too high.

# 5 Optimization

Optimization refers to the process of finding the best set of parameters for a machine learning model to minimize the objective function (e.g., the loss function). Gradient-based optimization algorithms use the gradient of the objective function to update the model parameters iteratively until a local minimum is reached.

## 5.1 Gradient Descent

Gradient descent is a first-order optimization algorithm that updates the model parameters based on the gradient of the objective function:

$$\boldsymbol{\theta} := \boldsymbol{\theta} - \alpha \nabla J(\boldsymbol{\theta})$$

where:

- $\boldsymbol{\theta}$ is the vector of model parameters
- $\alpha$ is the learning rate
- $\nabla J(\boldsymbol{\theta})$ is the gradient of the objective function with respect to the parameters

## 5.2 Stochastic Gradient Descent (SGD)

Stochastic Gradient Descent (SGD) is a variant of gradient descent that updates the model parameters based on the gradient of the objective function computed for a single data point or a small batch of data points:

$$\boldsymbol{\theta} := \boldsymbol{\theta} - \alpha \nabla J(\boldsymbol{\theta}; \mathbf{x}_i, y_i)$$

where $(\mathbf{x}_i, y_i)$ is a single data point or a small batch of data points.

### 5.2.1   Advantages of Stochastic Gradient Descent

- **Faster convergence:** Since the gradient is computed using only a small portion of the dataset, the computation is significantly faster, which allows for quicker updates to the model's parameters. This can lead to faster convergence to the optimal solution, especially when dealing with large datasets.

- **Noise can help escape local minima:** The inherent noise in the stochastic gradient estimate can actually help the algorithm escape local minima in non-convex optimization problems. This characteristic can be beneficial in deep learning models where the loss surface may have many local minima.

- **Memory efficiency:** SGD requires less memory compared to full-batch Gradient Descent since only a small portion of the dataset is loaded into memory at a time. This makes it possible to train models on large datasets that may not fit entirely into memory.

- **Tunable trade-off between computation and convergence:** By adjusting the batch size used in SGD, one can control the trade-off between computational efficiency and convergence rate. Smaller batch sizes will result in faster updates but may lead to a noisier convergence, while larger batch sizes can provide a more stable convergence but at a higher computational cost.

## 5.3   Momentum

Momentum is a technique that accelerates gradient descent and helps overcome local minima by introducing a velocity term, which is the exponentially decaying moving average of the gradients:

$$\mathbf{v} = \gamma \mathbf{v} + \alpha \nabla J(\boldsymbol{\theta})$$
$$\boldsymbol{\theta} = \boldsymbol{\theta} - \mathbf{v}$$

where $\gamma$ is the momentum coefficient (typically between 0.9 and 0.99).

## 5.4   RMSprop

RMSprop (Root Mean Square Propagation) is an optimization algorithm that adapts the learning rate for each parameter based on the moving average of the squared gradients:

$$\mathbf{s} := \rho\mathbf{s} + (1 - \rho)(\nabla J(\boldsymbol{\theta}))^2$$
$$\boldsymbol{\theta} := \boldsymbol{\theta} - \frac{\alpha}{\sqrt{\mathbf{s} + \epsilon}}\nabla$$
$$J(\boldsymbol{\theta})$$

where $\rho$ is the decay rate (typically between 0.9 and 0.99), and $\epsilon$ is a small constant to prevent division by zero.

## 5.5 ADAM

ADAM (Adaptive Moment Estimation) is an optimization algorithm for gradient-based optimization of neural networks and other machine learning models. It is an extension of stochastic gradient descent (SGD) that combines the benefits of adaptive learning rates and momentum. ADAM is widely used in deep learning due to its efficiency and robustness.

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1)\nabla J(\boldsymbol{\theta}t - 1)$$
$$v_t = \beta_2 vt - 1 + (1 - \beta_2)(\nabla J(\boldsymbol{\theta}_{t-1}))^2$$
$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$
$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$
$$\boldsymbol{\theta}t = \boldsymbol{\theta}t - 1 - \alpha\frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$$

where:

- $\boldsymbol{\theta}$ is the parameter vector
- $\nabla J(\boldsymbol{\theta})$ is the gradient of the objective function with respect to the parameters
- $m_t$ and $v_t$ are the first and second moment estimates, respectively
- $\beta_1$ and $\beta_2$ are the exponential decay rates for the first and second moment estimates (usually set to 0.9 and 0.999)
- $\alpha$ is the learning rate
- $\epsilon$ is a small constant to prevent division by zero (usually set to $10^{-8}$)
- $t$ is the iteration number

### 5.5.1 Intuition

ADAM combines the ideas of adaptive learning rates and momentum:

- **Adaptive learning rates**: ADAM computes individual learning rates for each parameter based on the square root of the second moment estimate ($\hat{v}_t$). This allows the algorithm to adapt the learning rate for each parameter based on the scale and sparsity of the gradients, which can help with convergence and stability.

- **Momentum**: ADAM maintains a moving average of the gradients (first moment estimate $m_t$) to incorporate information from past gradients. This helps to accelerate the convergence and overcome local minima or noisy gradients.

# 6  Loss Functions

Loss functions measure the difference between the predictions made by a machine learning model and the true target values. They are used during the training process to update the model's parameters to minimize the loss. Here, we discuss three commonly used loss functions: logistic loss, cross-entropy loss, and hinge loss.

## 6.1  Logistic Loss

The logistic loss, also known as the log loss or binary cross-entropy loss, is used in binary classification problems. It is defined as:

$$L(y, p) = -y \log(p) - (1 - y) \log(1 - p)$$

where:

- $y$ is the true class label (0 or 1)

- $p$ is the predicted probability of the positive class (1)

### 6.1.1  Intuition

The logistic loss penalizes wrong predictions with a logarithmic scale, where the penalty increases as the predicted probability deviates from the true class label. When the predicted probability is close to the true label, the loss is small, while the loss becomes large when the predicted probability is far from the true label. This encourages the model to assign high confidence to the correct class.

## 6.2  Cross-Entropy Loss

The cross-entropy loss, also known as categorical cross-entropy loss, is used in multiclass classification problems. It is defined as:

$$L(\boldsymbol{y}, \boldsymbol{p}) = -\sum_{i=1}^{C} y_i \log(p_i)$$

where:

- $y$ is the one-hot encoded true class label vector

- $p$ is the predicted probability vector for each class

- $C$ is the number of classes

- $y_i$ and $p_i$ are the true label and predicted probability, respectively, for class $i$

### 6.2.1 Intuition

The cross-entropy loss is an extension of the logistic loss for multiclass problems. It measures the dissimilarity between the true class label distribution and the predicted probability distribution. Like logistic loss, the cross-entropy loss penalizes wrong predictions with a logarithmic scale, encouraging the model to assign high confidence to the correct class.

## 6.3 Hinge Loss

The hinge loss is used in support vector machines (SVM) for binary classification problems. It is defined as:

$$L(y, f(\boldsymbol{x})) = \max(0, 1 - yf(\boldsymbol{x}))$$

where:

- $y$ is the true class label (-1 or 1)

- $f(\boldsymbol{x})$ is the raw output of the SVM, representing the signed distance from the input data point $\boldsymbol{x}$ to the decision boundary

### 6.3.1 Intuition

The hinge loss encourages a margin between the two classes in the feature space. It penalizes predictions that fall within the margin or on the wrong side of the decision boundary. The loss is zero for points that lie outside the margin and on the correct side of the boundary, which means that the model focuses on improving the predictions for the most difficult samples (those that are close to the decision boundary). This is consistent with the objective of SVMs to find the maximum-margin hyperplane that separates the classes.

## 6.4 Triplet Loss

Triplet loss is designed to optimize the relative distance between anchor, positive, and negative examples in the embedding space. Given a dataset with $N$ samples, each sample is represented as a triple $(a, p, n)$, where $a$ is the anchor, $p$ is a positive example with the same identity or class as the anchor, and $n$ is a negative example from a different identity or class. The triplet loss function is defined as:

$$\mathcal{L}(a, p, n) = \max \{d(a, p) - d(a, n) + \alpha, 0\} \tag{1}$$

where $d(a, b)$ denotes the distance metric, usually the Euclidean distance, between embeddings $a$ and $b$, and $\alpha$ is the margin. The goal is to minimize the total triplet loss:

$$\mathcal{L}_{total} = \sum_{i=1}^{N} \mathcal{L}(a_i, p_i, n_i) \tag{2}$$

### 6.4.1 Hyperparameters

- **Margin ($\alpha$):** The margin enforces a minimum separation between positive and negative pairs in the embedding space. A suitable value of $\alpha$ is critical for the convergence and performance of the model. A margin prevents the model from learning trivial solutions where all embeddings collapse to a single point or become very close to each other. This is because the model has to ensure that the anchor-positive distance is always smaller than the anchor-negative distance by at least $\alpha$. It also enforces the separation between different classes in the embedding space, leading to embeddings with higher discriminative power.

- **Batch size:** The batch size influences the number of triplets used in each iteration of the optimization process. A larger batch size may lead to more stable optimization, but can also increase computational requirements.

### 6.4.2 Hard Negative Mining

Hard negative mining is a strategy used to improve the convergence and performance of the model by selecting more challenging negative examples. In the context of triplet loss, hard negatives are those instances for which the difference in distances, $d(a, p) - d(a, n)$, is small. By focusing on hard negatives, the model is encouraged to learn more discriminative embeddings.

Two common approaches to hard negative mining are:

1. **Batch hard negative mining:** For each anchor-positive pair in a mini-batch, select the hardest negative example, i.e., the one with the smallest distance to the anchor.

2. **Online hard negative mining:** Continuously update the negative examples during training, selecting the hardest negatives for each anchor-positive pair at each iteration.

## 6.5 Mean Squared Error (MSE)

The mean squared error (MSE) is a widely used loss function for regression tasks. It is defined as:

$$L(\boldsymbol{y}, \hat{\boldsymbol{y}}) = \frac{1}{N} \sum_{i=1}^{N} (y_i - \hat{y}_i)^2$$

### 6.5.1 Intuition

The MSE measures the average squared difference between the true and predicted target values. It penalizes large errors more severely than small errors due to the squaring operation. This encourages the model to fit the data more closely and is sensitive to outliers.

## 6.6 Mean Absolute Error (MAE)

The mean absolute error (MAE) is another common loss function for regression tasks. It is defined as:

$$L(\boldsymbol{y}, \hat{\boldsymbol{y}}) = \frac{1}{N} \sum_{i=1}^{N} |y_i - \hat{y}_i|$$

### 6.6.1 Intuition

The MAE measures the average absolute difference between the true and predicted target values. Unlike MSE, it does not penalize large errors as severely due to the absence of the squaring operation. This makes the MAE more robust to outliers and can be a better choice for problems with noisy or skewed data.

# 7 Feature Selection

Feature selection is the process of selecting a subset of the most important features from the original set of features in a dataset. It is an essential step in machine learning as it can improve model performance, reduce overfitting, and decrease training time. There are three main categories of feature selection methods: filter methods, wrapper methods, and embedded methods.

## 7.1 Filter Methods

Filter methods are based on the intrinsic properties of the features and their relationship with the target variable. They are independent of any specific machine learning model and are computationally efficient. Some common filter methods are:

- **Variance Threshold**: Features with low variance are removed, as they do not contribute much information to the model.

- **Correlation Coefficient**: Features that have a high correlation with the target variable are selected, while features with a high correlation among themselves are removed to reduce multicollinearity.

- **Mutual Information**: Features that have a high mutual information with the target variable are selected, as they share more information about the target variable.

- **Chi-Squared Test**: In classification problems, the chi-squared test is used to measure the dependence between categorical features and the target variable. Features with high chi-squared values are more likely to be related to the target variable.

## 7.2 Wrapper Methods

Wrapper methods involve training a specific machine learning model with different feature subsets and selecting the subset that results in the best model performance. Although they can lead to better feature subsets, they are computationally expensive and prone to overfitting. Some common wrapper methods are:

- **Forward Selection**: Features are added one by one to the model, and the performance is evaluated at each step. The process stops when adding more features does not improve the performance.

- **Backward Elimination**: Features are removed one by one from the model, and the performance is evaluated at each step. The process stops when removing more features degrades the performance.

- **Recursive Feature Elimination**: A model is trained with all features, and the least important features are removed iteratively based on their importance. This process continues until the desired number of features is reached.

## 7.3 Embedded Methods

Embedded methods combine the advantages of both filter and wrapper methods by incorporating feature selection as part of the model training process. Some common embedded methods are:

- **Lasso Regression**: Lasso regression is a linear regression model with an L1 regularization term that encourages sparsity in the model parameters. It can be used to perform feature selection by setting the coefficients of irrelevant features to zero.

- **Ridge Regression**: Ridge regression is a linear regression model with an L2 regularization term that discourages large parameter values. Although it does not perform feature selection directly, it can be useful for ranking features based on their coefficients.

- **Elastic Net**: Elastic Net is a linear regression model that combines both L1 and L2 regularization terms. It can perform feature selection like Lasso while maintaining the stability of Ridge regression.

- **Decision Trees**: Decision tree algorithms, such as CART and C4.5, can inherently perform feature selection by selecting the most important features at the top of the tree based on information gain or Gini impurity.

- **Random Forests and Gradient Boosting Machines**: Ensemble methods like Random Forests and Gradient Boosting Machines can rank features based on their importance.

# 8 Feature Scaling

Feature scaling is an important preprocessing step that helps to bring the features of a dataset to a similar scale. This is crucial for many machine learning algorithms that are sensitive to the scale of input features, such as gradient descent-based methods, support vector machines, and k-means clustering. There are several techniques for feature scaling:

## 8.1 Normalization

Normalization, also known as min-max scaling, scales the features to a range of $[0, 1]$. The formula for normalization is as follows:

$$x_{normalized} = \frac{x - x_{min}}{x_{max} - x_{min}} \tag{3}$$

where $x_{normalized}$ is the normalized value of feature $x$, $x_{min}$ and $x_{max}$ are the minimum and maximum values of the feature, respectively.

## 8.2 Standardization

Standardization scales the features so that they have a mean of 0 and a standard deviation of 1. This technique is useful when the features follow a Gaussian distribution. The formula for standardization is:

$$x_{standardized} = \frac{x - \mu}{\sigma} \tag{4}$$

where $x_{standardized}$ is the standardized value of feature $x$, $\mu$ is the mean of the feature, and $\sigma$ is the standard deviation of the feature.

## 8.3 Robust Scaling

Robust scaling is a method that scales the features based on the median and interquartile range (IQR), making it robust to outliers. The formula for robust scaling is:

$$x_{robust} = \frac{x - Q_1}{Q_3 - Q_1} \tag{5}$$

where $x_{robust}$ is the robust-scaled value of feature $x$, $Q_1$ is the first quartile (25th percentile) of the feature, and $Q_3$ is the third quartile (75th percentile) of the feature.

## 8.4 When to Use Feature Scaling

Feature scaling is essential for algorithms that rely on the distance between data points, like k-means clustering and k-nearest neighbors, as well as for models that use gradient descent for optimization, such as linear regression and neural networks. It is less important for tree-based models like decision trees and random forests, as they are less sensitive to the scale of input features.

# 9 Handling Categorical Variables

Categorical variables are non-numeric data types that represent discrete categories or labels. Many machine learning algorithms require numeric input, so it is essential to convert categorical variables into a suitable numeric representation. There are several techniques for handling categorical variables:

## 9.1 Label Encoding

Label encoding assigns an integer value to each unique category. This method is useful when dealing with ordinal categorical variables, where the order of the categories is meaningful. However, it might not be suitable for nominal categorical variables, as the integer values can introduce an artificial order that does not exist in the data.

## 9.2 One-Hot Encoding

One-hot encoding creates binary features for each category of a categorical variable. For a variable with $k$ categories, one-hot encoding generates $k$ new binary features, with each feature representing the presence (1) or absence (0) of a specific category. This method is suitable for nominal categorical variables, as it does not introduce an artificial order. However, it can significantly increase the dimensionality of the dataset, especially if the categorical variables have many unique categories.

## 9.3 Dummy Encoding

Dummy encoding is similar to one-hot encoding but generates $k - 1$ binary features for a variable with $k$ categories. This method helps to avoid the dummy variable trap, a situation where multicollinearity occurs due to the inclusion of redundant information in the dataset.

## 9.4 Target Encoding

Target encoding, also known as mean encoding, replaces the categorical variable with the mean of the target variable for each category. This method can capture the relationship between the categorical variable and the target variable while reducing the dimensionality of the dataset. However, it might introduce leakage if not done correctly, so it is essential to perform target encoding separately for the training and validation/test sets.

## 9.5 Embeddings

Embeddings are another way to represent categorical variables, especially when dealing with high cardinality categorical variables. Embeddings map categorical variables to low-dimensional vectors that can capture complex relationships between categories. This method is commonly used with deep learning models, such as neural networks and transformers.

# 10    Model Architectures

## 10.1    Two-Tower Architecture for Search

The two-tower architecture is a neural network design pattern commonly used in information retrieval tasks, such as search and recommendation systems. This architecture consists of two separate neural networks, or "towers," which are trained to generate embeddings for queries and documents (or items). The main advantage of this architecture is its ability to efficiently compute similarity scores between queries and documents in the embedding space, enabling fast and accurate retrieval of relevant items.

### 10.1.1    Architecture Overview

In a two-tower architecture, the query tower and the document tower have separate neural network structures that are tailored to their specific input types (e.g., text, images, or a combination of features). These towers can be constructed using various neural network components, such as fully connected layers, convolutional layers, or transformer layers, depending on the input data and the desired complexity of the model.

The query tower processes the user's query and generates a query embedding, while the document tower processes the documents or items and generates document embeddings. Once the embeddings are generated, they can be compared using a similarity metric, such as the dot product, cosine similarity, or Euclidean distance, to rank the documents according to their relevance to the query.

### 10.1.2    Training the Two-Tower Model

To train a two-tower model, the objective is to minimize a loss function that encourages the query and document embeddings to be similar when they are relevant to each other and dissimilar otherwise. One common approach is to use a ranking loss function, such as the triplet loss or the contrastive loss, which operates on pairs or triplets of query-document combinations with varying degrees of relevance.

During training, the query and document towers learn to generate embeddings that satisfy the desired similarity constraints. This is achieved by updating the model parameters using backprop-agation and gradient descent, typically in a mini-batch setting. The trained model can then be used to generate query and document embeddings for information retrieval tasks, such as search and recommendation.

### 10.1.3    Advantages for Search

The two-tower architecture offers several benefits for search applications:

- **Efficiency**: By precomputing the document embeddings and storing them in an efficient data structure, such as an approximate nearest neighbors (ANN) index, the search engine can

quickly retrieve relevant documents by comparing their embeddings with the query embedding. This significantly reduces the computational cost of search, as the expensive document processing steps are performed offline.

- **Semantic matching**: The two-tower model learns to generate embeddings that capture the underlying semantics of the queries and documents, enabling the search engine to return more relevant and accurate results based on the meaning of the query, rather than just keyword matching.

- **Flexibility**: The two-tower architecture is highly adaptable to different types of input data and can be easily extended or combined with other techniques, such as feature engineering, text classification, or clustering, to further improve the search performance.

- **Scalability**: The two-tower architecture scales well with large datasets and high-dimensional feature spaces, as the complexity of the model can be adjusted by changing the depth and width of the neural network towers.

## 10.2 Autoencoders

Autoencoders are a type of unsupervised neural network used for tasks such as dimensionality reduction, denoising, and representation learning. They consist of two main components: an encoder that compresses the input data into a lower-dimensional representation, and a decoder that reconstructs the input data from the compressed representation. Autoencoders are trained by minimizing the reconstruction error, which is the difference between the original input and the reconstructed output.

### 10.2.1 Architecture

The architecture of an autoencoder typically consists of a symmetrical arrangement of layers, with the encoder and decoder having a mirrored structure. The encoder maps the input data to a lower-dimensional latent space, also known as the bottleneck layer, while the decoder maps the latent space back to the original input space.

The encoder and decoder can be composed of fully connected layers, convolutional layers, or recurrent layers, depending on the nature of the input data and the desired complexity of the model.

### 10.2.2 Training

Autoencoders are trained by minimizing the reconstruction error, which measures the dissimilarity between the original input and the reconstructed output. Common loss functions used for this purpose include the mean squared error (MSE) for continuous data and the cross-entropy loss for categorical data.

During training, the encoder learns to compress the input data into a compact representation, while the decoder learns to reconstruct the original data from the compressed representation.

The model parameters are updated using backpropagation and gradient descent, typically in a mini-batch setting.

### 10.2.3 Variational Autoencoders (VAEs)

Variational Autoencoders (VAEs) are a type of autoencoder that introduces a probabilistic perspective to the encoding process. Instead of learning a deterministic mapping from the input data to the latent space, VAEs learn the parameters of a probability distribution in the latent space.

The architecture of a VAE is similar to that of a standard autoencoder, with the key difference being that the encoder outputs the mean and variance of a Gaussian distribution instead of a single point in the latent space. The decoder takes a random sample from this distribution and reconstructs the input data.

Training a VAE involves minimizing a combination of the reconstruction error and a regularization term called the Kullback-Leibler (KL) divergence. The KL divergence measures the difference between the learned distribution in the latent space and a prior distribution, usually chosen to be a standard Gaussian. This regularization term encourages the model to learn a smooth and continuous mapping between the input data and the latent space.

VAEs have been used for various tasks, such as image synthesis, data augmentation, and semi-supervised learning. They offer several advantages over standard autoencoders, including the ability to generate diverse and realistic samples from the learned distribution in the latent space, as well as providing a more robust and expressive representation of the input data.

AEs and Variational Autoencoders VAEs were quite popular and had numerous applications in the past. However, their popularity has declined somewhat over the past few years due to the following reasons:

**Advancements in other deep learning architectures**: The rapid development of other deep learning architectures, such as Generative Adversarial Networks (GANs) and Transformers, has somewhat overshadowed AEs and VAEs. GANs, in particular, have demonstrated superior performance in generating high-quality and realistic samples, which has made them a popular choice for tasks like image synthesis and style transfer.

**Limited capacity**: AEs and VAEs have a limited capacity to model complex data distributions due to their bottleneck architecture, which constrains the information flow through the network. This limitation can make it challenging to learn meaningful representations for high-dimensional and diverse data, such as images and text.

**Mode collapse**: VAEs are known to suffer from a phenomenon called mode collapse, where the generated samples tend to be concentrated around a few modes of the data distribution, leading to a lack of diversity in the output. While there have been efforts to address this issue, it remains a challenge for VAEs.

**Difficulty in training**: Training AEs and VAEs can be challenging, particularly when it comes to balancing the reconstruction loss and the regularization term in VAEs. If the balance is not properly maintained, the model may either fail to learn meaningful representations or overfit to the training data.

**Competition from unsupervised and self-supervised learning techniques**: Recent advancements in unsupervised and self-supervised learning techniques, such as contrastive learning and masked language modeling, have shown promising results in learning rich and expressive representations from data without explicit supervision. These techniques have gained popularity in domains like natural language processing and computer vision, reducing the reliance on AEs and VAEs for representation learning.

# 11 Model Evaluation

## 11.1 Classification Model Evaluation

Model evaluation is the process of assessing the performance of a machine learning model based on specific metrics. Model selection involves choosing the best model from a set of candidate models based on their performance on a validation set. Below, we discuss various evaluation metrics and concepts related to model evaluation, particularly for classification tasks.

### 11.1.1 True Positive (TP), False Positive (FP), True Negative (TN), and False Negative (FN)

- **True Positive (TP)**: The number of positive instances correctly classified as positive.
- **False Positive (FP)**: The number of negative instances incorrectly classified as positive.
- **True Negative (TN)**: The number of negative instances correctly classified as negative.
- **False Negative (FN)**: The number of positive instances incorrectly classified as negative.

### 11.1.2 Confusion Matrix

A confusion matrix is a table that summarizes the performance of a classification model by comparing the predicted class labels with the true class labels. It contains the counts of TP, FP, TN, and FN for each class. The confusion matrix helps in understanding the strengths and weaknesses of a model, particularly in terms of its performance on different classes.

### 11.1.3 Accuracy

Accuracy is the proportion of correctly classified instances out of the total instances. It is calculated as:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Accuracy is a common metric, but it may not be suitable for imbalanced datasets, where the majority class can dominate the performance measure.

### 11.1.4 Precision

Precision is the proportion of true positive instances among the instances predicted as positive. It is calculated as:

$$\text{Precision} = \frac{TP}{TP + FP}$$

Precision measures the ability of the classifier to correctly identify positive instances and avoid false positives.

### 11.1.5 Recall

Recall, also known as sensitivity or true positive rate, is the proportion of true positive instances among the actual positive instances. It is calculated as:

$$\text{Recall} = \frac{TP}{TP + FN}$$

Recall measures the ability of the classifier to identify all the positive instances, even at the expense of including some false positives.

### 11.1.6 F-score

The F-score, also known as the F1-score, is the harmonic mean of precision and recall. It is calculated as:

$$F\text{-score} = 2\frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

The F-score balances precision and recall and is particularly useful when dealing with imbalanced datasets, as it takes both false positives and false negatives into account.

### 11.1.7 Imbalanced Datasets

In imbalanced datasets, one class has a significantly larger number of instances than the other classes. In such cases, accuracy may not be an appropriate evaluation metric, as the majority class can dominate the performance measure. Metrics like precision, recall, and F-score are more suitable for evaluating models on imbalanced datasets.

### 11.1.8 ROC Curve and AUC

The Receiver Operating Characteristic (ROC) curve is a graphical representation used to evaluate the performance of binary classification models. It plots the true positive rate (sensitivity or recall) against the false positive rate (1 - specificity) at various classification threshold levels.

The ROC curve is created by varying the classification threshold and calculating the TPR and FPR at each threshold. A model with perfect classification performance would have an ROC curve that goes directly from (0, 0) to (0, 1) and then to (1, 1). A model with random classification performance would have an ROC curve that follows the diagonal line from (0, 0) to (1, 1).

The Area Under the ROC Curve (AUC) is a single scalar value that measures the overall performance of a binary classifier. It is the area under the ROC curve and ranges from 0 to 1, where a higher value indicates better classification performance. The AUC can be interpreted as the probability that a randomly chosen positive instance will have a higher predicted probability than a randomly chosen negative instance.

## 11.2 Regression Evaluation Metrics

Regression evaluation metrics are used to assess the performance of regression models by comparing the predicted values to the true values. These metrics help quantify the errors made by the model, providing insights into the model's overall accuracy and effectiveness.

### 11.2.1 Mean Absolute Error (MAE)

Mean Absolute Error (MAE) measures the average magnitude of errors made by the model, without considering their direction. It is computed as the average of the absolute differences between the predicted values and the true values:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i|$$

where $y_i$ represents the true value, $\hat{y}_i$ represents the predicted value, and $n$ is the number of observations.

MAE is easy to interpret and is measured in the same unit as the target variable. Lower values of MAE indicate better model performance.

### 11.2.2 Mean Squared Error (MSE)

Mean Squared Error (MSE) measures the average squared difference between the predicted and true values. It emphasizes larger errors by squaring the differences:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

MSE is sensitive to outliers, as it penalizes larger errors more heavily. Lower values of MSE indicate better model performance.

### 11.2.3   Root Mean Squared Error (RMSE)

Root Mean Squared Error (RMSE) is the square root of MSE. It measures the average magnitude of errors, similar to MAE, but with more emphasis on larger errors due to the squaring operation:

$$\text{RMSE} = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}$$

RMSE is measured in the same unit as the target variable and is more sensitive to outliers than MAE. Lower values of RMSE indicate better model performance.

### 11.2.4   R-squared ($R^2$)

R-squared, also known as the coefficient of determination, measures the proportion of the variance in the true values that can be explained by the model. It is computed as the square of the correlation between the predicted and true values:

$$R^2 = 1 - \frac{\text{MSE}}{\text{Var}(y)}$$

where MSE is the mean squared error and $\text{Var}(y)$ is the variance of the true values.

$R^2$ ranges from 0 to 1, with higher values indicating better model performance. A value of 1 indicates that the model explains 100% of the variance in the true values, whereas a value of 0 indicates that the model does not explain any variance.

### 11.2.5   Summary

Different regression evaluation metrics provide different insights into the performance of a model. MAE and RMSE focus on the average magnitude of errors, with RMSE being more sensitive to outliers. MSE emphasizes larger errors due to the squaring operation. $R^2$ provides a proportion of the variance in the true values that can be explained by the model. It is important to consider multiple evaluation metrics when assessing the performance of a regression model, as each metric has its advantages and limitations.

## 11.3   Clustering Model Evaluation Metrics

### 11.3.1   Silhouette Score

The Silhouette Score is an internal evaluation metric that measures the cohesion and separation of clusters. It is calculated for each data point in the dataset and then averaged to obtain the overall score. The Silhouette Score ranges from -1 to 1, with higher values indicating better clustering quality. A high Silhouette Score suggests that data points within a cluster are close to each other and far from points in other clusters.

### 11.3.2 Normalized Mutual Information

Normalized Mutual Information (NMI) is an external evaluation metric that compares the clustering results with a given ground truth. It is based on the mutual information between the clustering assignment and the ground truth, normalized by the entropies of both the clustering assignment and the ground truth. NMI ranges from 0 to 1, with higher values indicating better clustering quality. A value of 1 indicates perfect agreement between the clustering results and the ground truth, while a value of 0 indicates no mutual information between them.

### 11.3.3 Summary

The choice of clustering evaluation metrics depends on the specific problem and the availability of ground truth information. The Silhouette Score is useful when there is no ground truth available, while NMI provides a more objective assessment of clustering quality when ground truth is available. In practice, it is often useful to consider multiple evaluation metrics to obtain a more comprehensive understanding of the clustering quality and the specific context of the problem.

# 12 Model Selection

## 12.1 Hyperparameter Tuning and Model Selection

In machine learning, hyperparameters are parameters that are set before the training process begins, and they influence the behavior of the algorithm during training. Some common hyperparameters include learning rate, regularization strength, and model complexity (e.g., the number of layers in a neural network). Tuning these hyperparameters is essential for achieving optimal model performance. Some popular techniques for hyperparameter tuning include:

1. **Grid Search**: In grid search, a predefined range of values is selected for each hyperparameter, and the model is trained and evaluated on all possible combinations of these values. Grid search is computationally expensive.

2. **Random Search**: Random search involves sampling hyperparameter values from a specified distribution, rather than iterating through all possible combinations. This method can be more efficient than grid search.

3. **Bayesian Optimization**: Bayesian optimization is a more sophisticated technique that models the relationship between hyperparameters and the objective function (e.g., model performance) using a probabilistic model. The algorithm then uses this model to decide which hyperparameter values to try next, focusing on the most promising regions of the hyperparameter space.

4. **Evolutionary Algorithms**: Evolutionary algorithms use principles from evolutionary biology, such as mutation, crossover, and selection, to explore the hyperparameter space. These algorithms maintain a population of candidate hyperparameter settings, and iteratively refine the population based on their performance.

## 12.2    Feature Engineering

1. **Feature Transformation**: Feature transformation involves applying mathematical functions to the features to change their scale, distribution, or relationship with the target variable. Common transformations include normalization, which scales features to a standard range using the formula $x' = \frac{x - x_{min}}{x_{max} - x_{min}}$, and standardization, which scales features to have zero mean and unit variance using the formula $x' = \frac{x - \mu}{\sigma}$. Power transformations, such as logarithmic ($x' = \log(x)$), square root ($x' = \sqrt{x}$), or Box-Cox transformations, can also be used. Transforming features can improve model performance by making the data more amenable to the assumptions of the learning algorithm.

2. **Feature Construction and Crossings**: Feature construction is the process of creating new features from existing ones, typically by combining or aggregating information from multiple features. This can be done using domain knowledge, unsupervised learning methods (e.g., clustering, dimensionality reduction), or automated feature learning techniques (e.g., deep learning). Constructing new features can capture complex patterns or relationships in the data that may be difficult for the model to learn directly from the raw features. Feature crossings, a type of feature construction, involve creating higher-order interaction terms by multiplying or otherwise combining two or more input features, e.g., $x_1 x_2$. These crossed features can help capture complex, nonlinear relationships within the data and improve model performance, particularly for linear models.

3. **Handling Categorical Features**: Many machine learning algorithms require numerical input, so categorical features (i.e., features with discrete, non-numeric values) must be transformed into numerical representations. Common techniques for handling categorical features include one-hot encoding (creating binary features for each unique category), ordinal encoding (assigning integer values to categories based on their rank), and target encoding (using the mean of the target variable for each category as the encoding). The choice of encoding method depends on the nature of the categorical feature and the assumptions of the learning algorithm.

4. **Feature Discretization**: Converting continuous features into discrete, categorical representations by dividing their value range into intervals or bins. This can simplify the learning problem, reduce dimensionality and noise, enhance interpretability, and enable the use of feature crossings with continuous variables. One common method of discretization is equal-width binning, which divides the range of values into equal-sized intervals, or equal-frequency binning, which divides the values into intervals containing an equal number of data points.

Feature engineering is an iterative process that requires domain knowledge, intuition, and experimentation. The choice of features and their representation can significantly impact the performance of the machine learning model, so it is important to evaluate different feature engineering strategies and continuously refine the feature set based on model performance and validation results.

# 13    Cross Validation

Cross-validation is a model validation technique used to assess the performance of a machine learning model on unseen data. It helps to minimize overfitting and provides a more accurate

estimate of the model's performance. In cross-validation, the dataset is divided into $k$ smaller subsets or folds, and the model is trained and tested $k$ times, with each fold used as the test set exactly once.

## 13.1  k-Fold Cross Validation

In k-fold cross-validation, the dataset is randomly partitioned into $k$ equally sized folds. For each iteration, one of the folds is used as the test set, and the remaining $(k-1)$ folds are combined to form the training set. The model is trained on the training set and evaluated on the test set. This process is repeated $k$ times, with each fold being used as the test set once. The final performance estimate is the average of the performance metrics (e.g., accuracy, F1-score, etc.) obtained from each iteration.

## 13.2  Choosing a Good Value for k

Choosing a good value for $k$ in k-fold cross-validation depends on the size of the dataset and the trade-offs between computation time, variance, and bias. Here are some guidelines for selecting an appropriate value for $k$:

- **k = 5 or 10**: These values are commonly used in practice, as they have been found to provide a good balance between computation time and performance estimation. With a larger dataset, a smaller value of $k$ (e.g., 5) can be chosen to reduce computation time, while with a smaller dataset, a larger value of $k$ (e.g., 10) can provide a more accurate performance estimate.

- **Leave-One-Out Cross Validation (LOOCV)**: If the dataset is very small, setting $k$ to the number of instances in the dataset results in LOOCV, where each instance is used as a test set exactly once. LOOCV provides the least biased performance estimate but has a high computational cost, as the model must be trained and tested for every instance in the dataset.

- **Stratified k-Fold Cross Validation**: For imbalanced datasets, it is important to maintain the class distribution in each fold. Stratified k-fold cross-validation ensures that each fold contains approximately the same proportion of instances for each class. This helps to obtain a more accurate performance estimate, especially for imbalanced datasets.

In general, it is recommended to perform a sensitivity analysis to determine the best value of $k$ for your specific dataset and problem. This involves conducting cross-validation with different values of $k$ and comparing the resulting performance estimates.

# 14  Bias / Variance Trade-off

In machine learning, the bias/variance trade-off is a fundamental concept that helps to understand the performance of a model in terms of its ability to generalize to unseen data. The prediction error of a model can be decomposed into three components: bias, variance, and irreducible error.

- **Bias**: Bias is the error introduced by approximating a real-world problem with a simplified model. High bias models make strong assumptions about the underlying data structure and can result in underfitting, where the model does not capture the underlying patterns in the data.

- **Variance**: Variance is the error introduced by the model's sensitivity to small fluctuations in the training set. High variance models capture the noise in the training data and can result in overfitting, where the model performs well on the training set but poorly on unseen data.

- **Irreducible error**: This is the error that cannot be reduced by improving the model, as it is caused by inherent noise or randomness in the data.

The bias/variance trade-off is the balance between the model's complexity and its ability to generalize to new data. A model that is too simple will have high bias and low variance, while a model that is too complex will have low bias and high variance. The goal is to find a balance between bias and variance that results in the lowest overall prediction error.

# 15 Regularization Techniques

Regularization techniques are used in machine learning to reduce the complexity of a model and prevent overfitting by adding a penalty term to the objective function. This penalty term discourages the model from assigning large weights to the features, which can lead to overfitting. Some common regularization techniques are:

## 15.1 L1 Regularization (Lasso)

L1 regularization, also known as Lasso, adds an L1 penalty term to the objective function, which is the sum of the absolute values of the model parameters:

$$\text{Objective}(\boldsymbol{w}) = \text{Loss}(\boldsymbol{w}) + \alpha \sum_{i=1}^{d} |w_i|$$

$$= \frac{1}{2n} \sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_i)^2 + \alpha \sum_{i} |w_i|$$

where $\boldsymbol{w}$ is the vector of model parameters, $d$ is the number of features, and $\alpha$ is a hyperparameter that controls the strength of the regularization.

L1 regularization encourages sparsity in the model parameters, effectively setting some of the less important features to zero. This can help with feature selection and result in a more interpretable model.

## 15.2    L2 Regularization (Ridge)

L2 regularization, also known as Ridge, adds an L2 penalty term to the objective function, which is the sum of the squares of the model parameters:

$$\text{Objective}(\boldsymbol{w}) = \text{Loss}(\boldsymbol{w}) + \alpha \sum_{i=1}^{d} w_i^2$$

$$= \frac{1}{2n} \sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_i)^2 + \alpha \sum_{i} w_i^2$$

L2 regularization discourages large parameter values but does not force them to zero like L1 regularization. This can result in a more stable model, especially when there are correlations between features.

## 15.3    Dropout

Dropout is a regularization technique used to prevent overfitting in neural networks. During the training process, dropout randomly sets a fraction of neurons' outputs to zero at each training iteration, effectively "dropping out" those neurons. This helps prevent the model from relying too heavily on any single neuron and encourages the model to learn more robust and generalizable features.

Dropout is applied with a probability $p$, called the dropout rate, which determines the fraction of neurons that are dropped out at each iteration. A common value for $p$ is between 0.5 and 0.8. Dropout is only applied during training, and during inference, all neurons are used, with their outputs scaled by the dropout rate to maintain consistency.

## 15.4    Covariate Shift

Covariate shift is a phenomenon that occurs when the input feature distribution changes between the training and test data, leading to a mismatch between the data used to train the model and the data encountered during testing. In the context of neural networks, the internal covariate shift refers to the changes in the distribution of hidden layer activations during training. This shift can slow down the training process and negatively impact the model's generalization capability.

The internal covariate shift problem can be formally defined as follows. Let $P(x)$ and $P(y|x)$ represent the input feature distribution and the conditional distribution of the output label given the input, respectively. Covariate shift occurs when the input distribution in the training data, $P_{train}(x)$, differs from the input distribution in the test data, $P_{test}(x)$, while the conditional distribution $P(y|x)$ remains unchanged. Mathematically, this can be expressed as:

$$P_{train}(x) \neq P_{test}(x) \quad \text{and} \quad P(y|x) = P_{train}(y|x) = P_{test}(y|x) \tag{6}$$

To mitigate the effects of covariate shift, various techniques have been proposed, including domain adaptation and normalization methods such as batch normalization and layer normalization. These

techniques aim to reduce the internal covariate shift by normalizing activations and stabilizing the distributions throughout the training process. By minimizing the effects of covariate shift, these methods can improve convergence, training stability, and generalization performance of neural networks.

## 15.5    Layer Normalization

Layer Normalization is a technique used to normalize the activations of a neural network layer to have zero mean and unit variance. It helps to address the issue of internal covariate shift, where the distribution of inputs to a given layer changes during training due to updates in the previous layers' parameters. This can slow down training and make it harder for the network to converge.

In Layer Normalization, the normalization is performed independently for each instance and across all activations of a single layer. The mean and variance are computed for each instance and used to normalize the activations, followed by a scale and shift operation using learnable parameters. This helps to stabilize the training process and allows the use of higher learning rates, leading to faster convergence.

## 15.6    Batch Normalization

Batch Normalization is another normalization technique used in neural networks to address the internal covariate shift issue. In Batch Normalization, the normalization is performed independently for each feature and across instances in a mini-batch.

The mean and variance are computed for each feature in the mini-batch and used to normalize the activations. These normalized activations are then scaled and shifted using learnable parameters. Like Layer Normalization, Batch Normalization helps stabilize the training process and allows for the use of higher learning rates.

In addition to improving convergence, Batch Normalization has been shown to have a regularizing effect, reducing the need for other regularization techniques like dropout. However, it is dependent on the batch size and can be sensitive to the choice of the mini-batch.

Both Layer Normalization and Batch Normalization are widely used in deep learning, and the choice between them depends on the specific architecture and problem. Layer Normalization is more suited for recurrent neural networks (RNNs) and networks with variable-sized inputs, while Batch Normalization is commonly used in convolutional neural networks (CNNs) and feedforward networks.

# 16    Handling Data

In real-world scenarios, data can often be incomplete, imbalanced, or have distribution shifts. Effective handling of these data issues is crucial for building robust and accurate machine learning models.

## 16.1 Missing Data

Missing data occurs when some instances in the dataset have incomplete feature values. Handling missing data is essential to avoid biases and inaccuracies in the model. Common strategies for handling missing data include:

- **Deletion**: Remove instances with missing values. This is only advisable when the amount of missing data is small, and the deletion does not introduce biases in the dataset.

- **Imputation**: Fill in missing values with estimated values. Common imputation methods include:

  - Mean or median imputation: Replace missing values with the mean (for continuous features) or median (for ordinal features) of the non-missing values.

  - Mode imputation: Replace missing values with the mode (for categorical features) of the non-missing values.

  - k-Nearest Neighbors (k-NN) imputation: Replace missing values with the average (for continuous features) or mode (for categorical features) of the k nearest neighbors in the feature space.

  - Model-based imputation: Train a regression or classification model (e.g., linear regression, random forest, etc.) to predict the missing values using the other features.

- **Indicator variables**: Create a binary indicator variable for each feature with missing values, indicating whether the value is missing or not. This allows the model to learn the relationship between the missingness and the target variable.

## 16.2 Imbalanced Data

Imbalanced data occurs when the distribution of the target variable classes is uneven. This can lead to biased models that favor the majority class. Strategies for handling imbalanced data include:

- **Resampling**:

  - *Oversampling*: Increase the number of instances in the minority class by replicating them or generating synthetic instances using techniques like the Synthetic Minority Over-sampling Technique (SMOTE).

  - *Undersampling*: Reduce the number of instances in the majority class by randomly removing instances or using techniques like Tomek links or neighborhood cleaning rule.

- **Cost-sensitive learning**: Assign different misclassification costs to the majority and minority classes, making the model more sensitive to the minority class.

- **Ensemble methods**: Use ensemble methods like bagging and boosting with balanced or stratified sampling to improve the performance on the minority class.

- **Evaluation metrics**: Use appropriate evaluation metrics like precision, recall, F1-score, or area under the ROC curve (AUC-ROC) that take into account both false positives and false negatives, instead of relying solely on accuracy.

## 16.3 Data Distribution Shifts

Data distribution shifts occur when the distribution of the data changes between the training and test sets or over time. This can lead to poor generalization and model performance degradation. Strategies for handling data distribution shifts include:

- **Domain adaptation**: Train a model on a source domain and adapt it to a target domain with different data distributions, using techniques like transfer learning, feature space alignment, or adversarial training.

- **Covariate shift adaptation**: If the input feature distribution changes but the conditional distribution of the target variable given the features remains the same.

# 17 Feedforward Neural Networks

Feedforward Neural Networks (FFNNs) are a type of artificial neural network where information flows in one direction from the input layer to the output layer without looping back. They consist of multiple layers of interconnected neurons, with each layer performing a transformation on the input data. The main components of an FFNN are:

- **Input layer**: The input layer receives the input features and passes them to the next layer.

- **Hidden layers**: One or more hidden layers perform non-linear transformations on the input data using a combination of weights, biases, and activation functions. Each neuron in a hidden layer is connected to all neurons in the previous and subsequent layers.

- **Output layer**: The output layer produces the final predictions or classifications. The number of neurons in the output layer corresponds to the number of target variables or classes.

An essential aspect of FFNNs is the activation function, which introduces non-linearity into the network, allowing it to learn complex patterns and relationships in the data.

## 17.1 Forward Pass

The forward pass of a feedforward neural network involves computing the weighted sum of the inputs and a bias term for each neuron in a layer, passing the result through an activation function, and then feeding the output as input to the next layer. This process is repeated for all layers until the output layer is reached. The forward pass can be represented by the following formula:

$$z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)} \tag{7}$$

$$a^{(l)} = f(z^{(l)}) \tag{8}$$

where $W^{(l)}$ and $b^{(l)}$ are the weight matrix and bias vector for layer $l$, respectively, $a^{(l-1)}$ is the activation of the previous layer, $z^{(l)}$ is the weighted sum for layer $l$, $a^{(l)}$ is the activation of layer $l$, and $f$ is the activation function.

## 17.2   Backward Pass

The backward pass, also known as backpropagation, is the process of updating the network's weights and biases to minimize the error between the predicted output and the target output. This is done by computing the gradient of the loss function with respect to each weight and bias, and then updating the weights and biases using gradient descent. The gradient computation involves applying the chain rule of calculus to compute the derivatives of the loss function with respect to the weights and biases. The backward pass can be represented by the following formulas:

$$\delta^{(L)} = \nabla_a C \odot f'(z^{(L)}) \tag{9}$$

$$\delta^{(l)} = ((W^{(l+1)})^T \delta^{(l+1)}) \odot f'(z^{(l)}) \tag{10}$$

$$\frac{\partial C}{\partial b^{(l)}} = \delta^{(l)} \tag{11}$$

$$\frac{\partial C}{\partial W^{(l)}} = \delta^{(l)} (a^{(l-1)})^T \tag{12}$$

where $C$ is the loss function, $\delta^{(l)}$ is the error term for layer $l$, and $f'(z^{(l)})$ is the derivative of the activation function with respect to the weighted sum $z^{(l)}$. The weights and biases are then updated using gradient descent:

$$W^{(l)} = W^{(l)} - \eta \frac{\partial C}{\partial W^{(l)}} \tag{13}$$

$$b^{(l)} = b^{(l)} - \eta \frac{\partial C}{\partial b^{(l)}} \tag{14}$$

where $\eta$ is the learning rate, a hyperparameter that controls the step size of the gradient descent update.

## 17.3   Activation Functions

Activation functions are applied to the output of each neuron in the hidden layers and determine whether the neuron should be "activated" based on its input. Some common activation functions and their intuitions are:

- **Sigmoid**: The sigmoid activation function is defined as

$$f(x) = \frac{1}{1 + e^{-x}}$$

  It maps the input to a value between 0 and 1, making it suitable for binary classification or when the output should represent a probability. However, the sigmoid function suffers from the vanishing gradient problem, where the gradients become very small for large input values, slowing down the training process.

- **Tanh**: The hyperbolic tangent (tanh) activation function is defined as

$$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

It maps the input to a value between -1 and 1, providing a zero-centered output that can be advantageous during training. Like the sigmoid function, tanh also suffers from the vanishing gradient problem.

- **ReLU**: The Rectified Linear Unit (ReLU) activation function is defined as

$$f(x) = \max(0, x)$$

It is computationally efficient and addresses the vanishing gradient problem by having a constant gradient for positive input values. However, ReLU units can suffer from the dying ReLU problem, where some neurons become inactive and produce zero outputs during training, effectively "dying" and not contributing to the learning process.

- **Leaky ReLU**: The Leaky ReLU activation function is defined as

$$f(x) = \max(\alpha x, x)$$

where $\alpha$ is a small constant (e.g., 0.01). It addresses the dying ReLU problem by allowing a small non-zero gradient for negative input values, keeping the neurons "alive" and contributing to the learning process.

- **Softmax**: The softmax activation function is used in the output layer for multi-class classification problems. It converts the input into a probability distribution over the classes, ensuring that the sum of the output probabilities is 1. The softmax function is defined as

$$f_i(x) = \frac{e^{x_i}}{\sum_{j=1}^{K} e^{x_j}}$$

where $x_i$ is the input value for class $i$ and $K$ is the total number of classes.

These activation functions introduce non-linearity into the network

## 17.4   Pooling

Pooling is an operation commonly used in neural networks, especially in Convolutional Neural Networks (CNNs), to reduce the spatial dimensions of feature maps while preserving important features. Pooling helps in reducing the number of parameters and computational complexity, making the network more efficient and less prone to overfitting. There are several types of pooling techniques, with max-pooling being the most common. Other types include average pooling and global average pooling.

### 17.4.1   Max-Pooling

Max-pooling is a downsampling operation that selects the maximum value from a fixed-size window (also called a kernel) within the input feature map and outputs it as the pooled value. This

operation is applied with a certain stride, which determines the step size between consecutive pooling operations. Max-pooling helps retain the most prominent features while discarding less significant ones, making the network more robust to small variations and translations in the input data.

### 17.4.2 Average Pooling

Average pooling, as the name suggests, calculates the average value within the fixed-size window instead of selecting the maximum value. This technique provides a smoother representation of the input feature map, reducing the risk of overfitting. However, it might not be as effective in preserving the most important features compared to max-pooling.

### 17.4.3 Global Average Pooling

Global average pooling is a pooling operation that takes the average of the entire input feature map, outputting a single value. This technique is often used as an alternative to fully connected layers in deep learning models, particularly in image classification tasks. Global average pooling reduces the number of parameters, thus decreasing the risk of overfitting and reducing computational complexity.

# 18 Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM)

Recurrent Neural Networks (RNNs) are a type of neural network designed to handle sequential data. They can process variable-length input sequences by maintaining a hidden state that can capture information from previous time steps. RNNs are particularly useful for tasks involving time series data, natural language processing, and other sequential data.

## 18.1 Backpropagation Through Time (BPTT)

Backpropagation Through Time (BPTT) is the training algorithm used for RNNs. BPTT is an extension of the standard backpropagation algorithm, modified to handle the recurrent connections and temporal dependencies in RNNs. The main idea of BPTT is to "unfold" the RNN in time, creating a feedforward network with multiple copies of the RNN, one for each time step in the sequence. The unfolded network is then trained using the standard backpropagation algorithm.

## 18.2 Vanishing and Exploding Gradient Problem

One of the main challenges in training deep neural networks, particularly recurrent neural networks (RNNs), is the vanishing and exploding gradient problem. This issue arises during backpropagation through time (BPTT) in RNNs, when gradients are backpropagated through the network's time steps.

### 18.2.1  Vanishing and Exploding Gradients in RNNs

RNNs are particularly susceptible to the vanishing and exploding gradient problem due to their recurrent nature. During BPTT, the gradients are multiplied by the recurrent weight matrix at each time step. If the eigenvalues of this matrix are less than 1, the gradients can become very small (vanish) as they are backpropagated. Conversely, if the eigenvalues are greater than 1, the gradients can grow exponentially (explode), leading to unstable training.

Vanishing gradients make it difficult for the network to learn long-range dependencies in the data, as the influence of earlier time steps is diminished. Exploding gradients can cause the model's parameters to diverge, resulting in erratic behavior and poor generalization.

### 18.2.2  Vanishing and Exploding Gradients in Non-RNN Architectures

Vanishing and exploding gradients can also occur in deep feedforward neural networks, particularly those with many layers. As the gradients are backpropagated through the layers, they can be multiplied by small or large weights, leading to vanishing or exploding gradients, respectively. This can result in poor convergence and training instability.

### 18.2.3  Mitigating the Vanishing and Exploding Gradient Problem

Several techniques can help mitigate the vanishing and exploding gradient problem in neural networks:

**Weight Initialization:** Proper weight initialization, such as Glorot (Xavier) initialization or He initialization, can help prevent gradients from vanishing or exploding during training. These methods scale the initial weights based on the size of the input and output layers, helping to maintain a stable signal throughout the network.

**Gradient Clipping:** Gradient clipping is a technique used to prevent exploding gradients by limiting the maximum magnitude of the gradients during training. If the gradient's norm exceeds a predefined threshold, the gradient is rescaled to keep its norm below the threshold. This can help stabilize training and prevent divergence.

**Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU) Cells:** LSTMs and GRUs are advanced RNN architectures designed to address the vanishing gradient problem. They incorporate gating mechanisms that allow the network to selectively remember or forget information, making it easier to learn long-range dependencies in the data. These architectures have been shown to alleviate the vanishing gradient problem in many sequence-to-sequence learning tasks.

**Batch Normalization:** Batch normalization is a technique that normalizes the activations of each layer using the mean and variance of the current mini-batch. By normalizing the activations, the technique helps maintain a stable distribution and reduces the internal covariate shift, which can help prevent vanishing and exploding gradients.

**Layer Normalization:** Layer normalization is another normalization technique that normalizes the activations within each layer based on the layer's mean and variance. Unlike batch normal-

ization, layer normalization does not depend on the mini-batch, making it particularly useful for RNNs, where the sequence length may vary. Layer normalization can help stabilize the training process and alleviate the vanishing and exploding gradient problem in deep networks.

# 19 Recurrent Neural Networks (RNNs)

Recurrent Neural Networks (RNNs) are a type of neural network designed to handle sequential data. They process variable-length input sequences by maintaining a hidden state that captures information from previous time steps. RNNs are particularly useful for tasks involving time series data, natural language processing, and other sequential data.

## 19.1 Long Short-Term Memory (LSTM)

Long Short-Term Memory (LSTM) is an advanced type of RNN architecture that addresses the vanishing and exploding gradient problem that occurs during the training of standard RNNs. LSTM networks introduce memory cells, which can store and update information over long sequences, allowing the network to learn and remember long-range dependencies.

The key component of an LSTM is the memory cell, which consists of several gates that control the flow of information in and out of the cell:

- **Input gate**: Determines how much of the new input should be added to the memory cell. The input gate is updated using its own weight matrix and an activation function, typically the sigmoid function.

- **Forget gate**: Determines how much of the existing memory cell content should be retained. The forget gate is also updated using its own weight matrix and an activation function, usually the sigmoid function.

- **Output gate**: Determines how much of the memory cell content should be used to compute the output at the current time step. The output gate has its own weight matrix and employs an activation function, typically the sigmoid function.

These gates allow the LSTM to selectively update, retain, or output information from the memory cell, enabling it to capture long-range dependencies and mitigate the vanishing and exploding gradient problem. LSTMs have been widely used for various tasks involving sequential data, such as natural language processing, speech recognition, and time series forecasting.

In addition to the gates, the LSTM cell also has a mechanism to update the memory cell state itself, which typically involves the use of a pointwise multiplication operation with the input and forget gate outputs, as well as a pointwise addition operation with the new input. The memory cell state is updated using its own weight matrix and an activation function, often the hyperbolic tangent (tanh) function.

## 19.2 Gated Recurrent Units (GRU)

Gated Recurrent Units (GRU) are another advanced type of RNN architecture that addresses the vanishing and exploding gradient problem, similar to LSTMs. GRUs simplify the LSTM architecture by combining the forget and input gates into a single "update gate." As a result, GRUs have fewer parameters and can be computationally more efficient than LSTMs, while still being able to capture long-range dependencies.

The key component of a GRU is the gating mechanism, which consists of two gates that control the flow of information:

- **Update gate**: Determines how much of the previous hidden state should be retained and how much of the new input should be used to update the hidden state. The update gate effectively combines the functionality of the input and forget gates in an LSTM.

- **Reset gate**: Determines how much of the previous hidden state should be used to compute the candidate hidden state, which is then used to update the actual hidden state.

These gates allow the GRU to selectively update and retain information from the hidden state, enabling it to capture long-range dependencies and mitigate the vanishing and exploding gradient problem. GRUs have been widely used for various tasks involving sequential data, such as natural language processing, speech recognition, and time series forecasting, and often achieve similar performance to LSTMs with reduced computational complexity.

# 20 Sequence-to-Sequence (Seq2Seq) Models

Sequence-to-Sequence (Seq2Seq) models are a type of neural network architecture designed for tasks that involve mapping input sequences to output sequences. They are particularly useful for natural language processing tasks, such as machine translation, text summarization, and dialogue systems, where both the input and output are variable-length sequences.

Seq2Seq models consist of two main components:

- **Encoder**: The encoder is an RNN (typically an LSTM or GRU) that processes the input sequence and encodes it into a fixed-size vector representation, called the context vector. The encoder captures the information from the input sequence and passes it to the decoder.

- **Decoder**: The decoder is another RNN (also typically an LSTM or GRU) that takes the context vector from the encoder and generates the output sequence. The decoder is trained to predict the next token in the output sequence, given the context vector and the previously generated tokens.

During training, the Seq2Seq model learns to map input sequences to output sequences by minimizing the difference between the predicted output tokens and the ground-truth output tokens. The model is typically trained using teacher forcing, where the decoder is fed the ground-truth output tokens as input at each time step, rather than its own previous predictions. This helps the model to converge faster and learn more effectively.

One of the main challenges in Seq2Seq models is handling variable-length input and output sequences. This is usually addressed by using techniques such as padding and masking, where input

and output sequences are padded to a fixed length, and a mask is used to ignore the padded positions during training and evaluation.

Another challenge is that the context vector may not be able to capture all the information from long input sequences, leading to a loss of information and reduced performance. To address this issue, attention mechanisms can be added to the Seq2Seq model, allowing the decoder to selectively focus on different parts of the input sequence while generating the output sequence. Attention mechanisms have been shown to significantly improve the performance of Seq2Seq models on various tasks, such as machine translation and text summarization.

# 21 Attention Mechanism

Attention mechanisms were introduced to improve the performance of sequence-to-sequence (Seq2Seq) models by allowing the model to selectively focus on different parts of the input sequence when generating the output sequence. This alleviates the issue of having to compress all the information from the input sequence into a fixed-size context vector, which can be especially challenging for long input sequences.

The attention mechanism computes a set of attention weights for each time step in the input sequence, which are used to create a weighted sum of the encoder hidden states. This weighted sum, called the context vector, is then combined with the current hidden state of the decoder to generate the output token at each time step.

The attention mechanism can be formulated as follows:

1. Compute an attention score for each encoder hidden state at each time step of the decoder. The attention score can be computed using various methods, such as dot product, scaled dot product, or a small feedforward neural network.

2. Apply a softmax function to the attention scores to obtain the attention weights, which sum to 1.

3. Compute the context vector as the weighted sum of the encoder hidden states using the attention weights.

4. Combine the context vector with the decoder's hidden state to generate the output token at the current time step.

In this way, the attention mechanism allows the model to focus on different parts of the input sequence based on the current decoding step, making it more flexible and capable of capturing long-range dependencies.

There are several types of attention mechanisms, including:

- **Global attention**: The model computes attention scores for all time steps in the input sequence, allowing it to attend to any part of the input sequence.

- **Local attention**: The model focuses on a small window of the input sequence, reducing the computational complexity compared to global attention.

- **Self-attention**: Also known as intra-attention, this mechanism computes attention scores within the same sequence, allowing the model to capture relationships between different parts of the sequence. Self-attention is a key component of the Transformer architecture, which has achieved state-of-the-art performance on various natural language processing tasks.

- **Cross Attention**: Is used only in the decoder part of the Transformer model. It allows each token in the target sequence (during decoding) to attend to all tokens in the source sequence (encoder output). This mechanism helps the decoder to capture relevant information from the encoder output, which is useful for tasks like machine translation or summarization. Cross-attention is applied after self-attention in the decoder and follows the same Query, Key, Value approach. However, the Query vectors come from the decoder's previous layer, while the Key and Value vectors come from the encoder output.

- **Masked Attention** is used when we want to prevent the model from attending to certain positions in the input sequence. This is particularly useful in the decoding phase of sequence-to-sequence models, where the model should not have access to future tokens. In masked attention, a mask is applied to the attention scores before the softmax function. The mask is typically a matrix with the same dimensions as the attention scores, where masked positions have a very large negative value (e.g., -1e9) that effectively zeros out the softmax output for those positions. This ensures that the model does not attend to masked positions while calculating the attention-weighted output.

Attention mechanisms have been widely adopted in deep learning models for various tasks, such as machine translation, text summarization, image captioning, and speech recognition, leading to significant improvements in performance and the ability to handle long input sequences.

# 22 Teacher Forced Learning

Teacher forcing is a training technique used in sequence-to-sequence models, where during training, the true target tokens are provided as input to the decoder at each time step, rather than using the model's own predictions from previous time steps. This approach has both advantages and disadvantages.

## 22.1 Advantages

- Faster convergence: As the model receives the correct target tokens during training, it learns the correct dependencies more effectively, leading to faster convergence.

- Better performance: The model is less prone to compounding errors, as it is not affected by its own mistakes during training.

## 22.2 Disadvantages

- Exposure bias: During inference, the model does not have access to the true target tokens and must rely on its own predictions. This can lead to a discrepancy between the training

and inference scenarios, known as exposure bias.

- Error propagation: If the model generates an incorrect token during inference, the error can propagate and affect the entire generated sequence, as it continues to use its own predictions as input.

To mitigate the issues caused by teacher forcing, researchers have proposed techniques such as scheduled sampling, which combines the use of true target tokens and model predictions during training to make the training and inference scenarios more similar.

# 23 Transformer Architecture

The Transformer architecture, introduced by Vaswani et al. in the paper "Attention is All You Need," is a groundbreaking neural network architecture designed for sequence-to-sequence tasks. It has achieved state-of-the-art performance on a wide range of natural language processing tasks, such as machine translation, text summarization, and sentiment analysis. The key innovation of the Transformer architecture is the use of self-attention mechanisms, which enable the model to capture dependencies between all input tokens, regardless of their distance in the sequence.

## 23.1 Architecture Overview

The Transformer consists of an encoder and a decoder, similar to other sequence-to-sequence models. However, instead of using RNNs, the Transformer relies entirely on self-attention mechanisms and feedforward neural networks. The encoder and decoder are both composed of multiple identical layers, each containing two main components:

### 23.1.1 Multi-Head Self-Attention

The input of the multi-head self-attention mechanism is a set of token embeddings. These embeddings are linearly transformed into three different sets of vectors: Query (Q), Key (K), and Value (V) vectors. Each set is obtained by applying a separate learned weight matrix to the input embeddings.

The self-attention mechanism computes attention scores for each token in the input sequence with respect to all other tokens in the sequence. This is done by calculating the dot product between the Query vector of the current token and the Key vectors of all other tokens. The dot product effectively measures the similarity between the current token and the other tokens in the sequence.

Next, a softmax function is applied to the attention scores, which normalizes them and converts them into attention weights that sum to 1. The attention weights represent the relative importance of each token in the context of the current token.

Finally, the attention weights are used to create a weighted sum of the Value vectors. This weighted sum forms the output of the self-attention mechanism for the current token, which is then passed through the rest of the Transformer layers.

In the multi-head self-attention mechanism, multiple sets of Query, Key, and Value vectors are generated using different weight matrices. This allows the model to attend to different aspects of the input sequence simultaneously, capturing various relationships between the tokens. The outputs from all attention heads are then concatenated and linearly transformed to produce the final output of the multi-head self-attention layer.

$$Output = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

### 23.1.2    Position-Wise Feedforward Neural Networks

In the Transformer architecture, the Position-wise Feedforward Neural Network (FFNN) is an essential component of both the encoder and decoder layers. It is applied to each token's output from the multi-head self-attention mechanism independently and identically, allowing the model to learn non-linear interactions between the token features.

The Position-wise Feedforward Neural Network consists of two linear layers (fully connected layers) with a non-linear activation function in between, such as ReLU or GELU. The FFNN can be represented by the following formula:

$$FFNN(x) = W_2 \cdot \text{ReLU}(W_1 \cdot x + b_1) + b_2$$

In the context of the Transformer, the Position-wise FFNN operates on each token's output from the multi-head self-attention independently. This means that the same FFNN is applied to every position in the sequence, which allows the model to efficiently capture local patterns and interactions between the token features.

After applying the Position-wise FFNN, the outputs are combined with the input of the layer using residual connections, followed by layer normalization. This helps stabilize the training process and improves model performance.

Residual connections and layer normalization are important components of the Transformer architecture, as they help stabilize the training process and improve model performance.

1. **Residual connections**: After the output of the Position-wise Feedforward Neural Network (FFNN) is computed, it is combined with the input of the FFNN using residual connections (also known as skip connections). This is done by adding the input to the output of the FFNN:

$$y = x + FFNN(x)$$

where $x$ is the input of the FFNN and $y$ is the resulting output after the residual connection. The intuition behind residual connections is that they allow the model to learn to approximate the residual (or difference) between the input and the desired output, rather than learning the entire transformation. This can help mitigate the vanishing gradient problem, as the gradients can flow more easily through the network, enabling the model to learn deeper representations.

2. **Layer normalization**: After the residual connection, layer normalization is applied to the output. Layer normalization is a technique that normalizes the outputs of a layer across the feature dimension, ensuring that the mean and variance of the layer's output remain

consistent during training. This is done by computing the mean and variance of the output across the feature dimension, and then normalizing the output using these statistics:

$$z = \frac{y - \text{mean}(y)}{\sqrt{\text{var}(y) + \epsilon}}$$

where $y$ is the output after the residual connection, $z$ is the normalized output, and $\varepsilon$ is a small constant added for numerical stability. After normalization, the output is scaled by a learnable gain parameter and shifted by a learnable bias parameter.

Layer normalization helps stabilize the training process by mitigating the covariate shift problem, which arises when the distribution of the inputs to a layer changes during training. By normalizing the outputs, layer normalization ensures that the subsequent layers receive inputs with a consistent distribution, allowing the model to learn more effectively.

## 23.2   Positional Encoding

Since the Transformer does not use recurrent or convolutional layers, it lacks the ability to capture the relative positions of tokens in the input sequence. To address this issue, positional encoding is added to the input token embeddings, which provides information about the token positions.

The positional encoding can be computed using various methods, such as learned embeddings or fixed sinusoidal functions. In the original Transformer paper, the authors used sinusoidal functions with different frequencies, which allows the model to capture both short- and long-range dependencies:

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

where $pos$ is the position of the token in the sequence, $i$ is the dimension of the positional encoding, and $d_{model}$ is the model's hidden size.

## 23.3   Encoder

The encoder of the Transformer architecture is responsible for processing the input sequence and generating a continuous representation that captures the relationships between tokens in the sequence. The encoder consists of a stack of identical layers, each containing a multi-head self-attention mechanism and a position-wise feedforward neural network, connected through residual connections and followed by layer normalization.

The input to the encoder is a sequence of token embeddings, which are combined with positional encodings to provide information about the token positions. The input sequence is then passed through the multi-head self-attention mechanism, allowing the model to capture dependencies between all input tokens. The output of the self-attention mechanism is then fed into the position-wise feedforward neural network, which captures local patterns and non-linear interactions between token features. This process is repeated for all layers in the encoder stack, resulting in a continuous representation of the input sequence that can be used as input to the decoder.

## 23.4 Decoder

The decoder of the Transformer architecture is responsible for generating the output sequence, given the continuous representation produced by the encoder. Similar to the encoder, the decoder consists of a stack of identical layers, each containing two multi-head self-attention mechanisms and a position-wise feedforward neural network, connected through residual connections and followed by layer normalization.

The first multi-head self-attention mechanism in each decoder layer processes the output sequence produced so far, allowing the model to capture dependencies between output tokens. This self-attention mechanism is masked to prevent the model from attending to future tokens in the output sequence, ensuring that the predictions are based solely on the available information.

The second multi-head self-attention mechanism in each decoder layer attends to the continuous representation generated by the encoder. This allows the model to align the output tokens with the relevant input tokens, incorporating the context from the input sequence into the output sequence.

The output of the second self-attention mechanism is then fed into the position-wise feedforward neural network, which captures local patterns and non-linear interactions between token features. This process is repeated for all layers in the decoder stack, resulting in a continuous representation of the output sequence.

Finally, the continuous representation of the output sequence is transformed into a probability distribution over the target vocabulary using a linear layer followed by a softmax activation function. The model is trained to maximize the likelihood of the correct target token at each position in the output sequence, given the input sequence and the preceding output tokens.

During inference, the decoder generates the output sequence one token at a time, conditioning on the input sequence and the previously generated output tokens. This can be done using greedy decoding, where the model selects the token with the highest probability at each step, or more advanced decoding strategies such as beam search or top-k sampling, which explore multiple candidate sequences to improve the quality of the generated output.

## 23.5 Memory Constraints

The main memory constraint in the Transformer architecture comes from the computation of the self-attention mechanism, which requires storing attention scores for all pairs of tokens in the input sequence. This leads to a quadratic memory complexity with respect to the sequence length, which can be problematic for very long sequences.

Several techniques have been proposed to address this issue, such as:

- **Local attention**: Restricting the attention mechanism to a small window around each token, reducing the memory complexity to linear with respect to the sequence length.

- **Sparse attention**: We can decompose the attention matrix into a product of two sparse matrices, each with only $O(N \log(N))$ non-zero elements. The sparsity pattern of these matrices can be determined by clustering technique.

- **Longformer**: Adapting the self-attention mechanism to focus on local windows and a small number of global tokens, which allows the model to scale to longer sequences while still capturing important global context.

- **Transformer-XL**: introduces the notion of recurrence to the deep self-attention network. Instead of computing the hidden states from scratch for each new segment, Transformer-XL reuses the hidden states obtained in previous segments. The reused hidden states serve as memory for the current segment, which builds up a recurrent connection between the segments. As a result, modeling very long-term dependency becomes possible because information can be propagated through the recurrent connections. As an additional contribution, the Transformer-XL uses a new relative positional encoding formulation that generalizes to attention lengths longer than the one observed during training.

# 24 Knowledge Distillation, Model Size Reduction, and Pruning

Deep neural networks, while being powerful models capable of achieving state-of-the-art performance, often suffer from high computational requirements and large memory footprints. To mitigate these issues and enable the deployment of such models on resource-constrained devices, several techniques have been developed, such as knowledge distillation, model size reduction, and pruning.

## 24.1 Knowledge Distillation

Knowledge distillation (KD) is a technique used to compress a large, pre-trained model, known as the teacher model, into a smaller model, known as the student model. The idea is to transfer the knowledge embedded in the teacher model to the student model, preserving as much of the teacher's performance as possible while reducing the model's size and computational requirements.

The distillation process generally involves training the student model to mimic the output probabilities of the teacher model, rather than the ground-truth labels. This is achieved by minimizing a loss function that combines the traditional cross-entropy loss with the Kullback-Leibler (KL) divergence between the teacher's and student's output probabilities. The KL divergence term encourages the student model to learn the teacher's soft labels, which contain more information about the relationships between classes than the hard labels.

The training process for knowledge distillation can be represented as:

$$\mathcal{L}KD = (1 - \alpha)\mathcal{L}CE(\mathbf{y}, \mathbf{z}S) + \alpha T^2 \mathcal{L}KL(\mathbf{z}_T/T, \mathbf{z}_S/T) \tag{15}$$

where $\mathcal{L}CE$ is the cross-entropy loss, $\mathcal{L}KL$ is the KL divergence, $\mathbf{y}$ are the ground-truth labels, $\mathbf{z}_T$ and $\mathbf{z}_S$ are the logits from the teacher and student models, respectively, $T$ is a temperature parameter, and $\alpha$ is a hyperparameter that balances the contribution of the two loss terms.

## 24.2 Model Size Reduction

Model size reduction aims to reduce the number of parameters in a neural network without significantly compromising its performance. One popular method for model size reduction is network architecture search (NAS), which involves searching for optimal network architectures with fewer parameters and lower computational requirements. Techniques like weight sharing and quantization can also be used to compress the weights of a model, reducing its memory footprint.

### 24.2.1 Network Architecture Search

Network architecture search (NAS) is an automated process that explores different neural network architectures to find an optimal configuration with fewer parameters and lower computational requirements while maintaining performance. NAS algorithms can be guided by various search strategies, such as evolutionary algorithms, reinforcement learning, or gradient-based optimization.

### 24.2.2 Weight Sharing and Quantization

Weight sharing and quantization are techniques used to compress the weights of a neural network. Weight sharing involves clustering similar weights together and representing them with a single value, while quantization maps the continuous weights to a smaller set of discrete values. These techniques can significantly reduce the memory footprint of a model, making it more suitable for deployment on resource-constrained devices.

## 24.3 Pruning

Pruning is a technique used to remove redundant or unimportant neurons or connections from a neural network, thereby reducing its size and computational requirements. Pruning can be performed during training or after training the model (post-training pruning). There are several pruning strategies, including weight pruning, neuron pruning, and structured pruning.

### 24.3.1 Weight Pruning

Weight pruning involves removing individual connections in the network based on the magnitude of their weights. Connections with smaller weights are considered less important and are pruned, while connections with larger weights are retained. Weight pruning can be performed iteratively, with a certain percentage of the smallest weights removed at each step. After pruning, the network is fine-tuned to recover any lost performance.

$$\mathbf{W}_{pruned} = \mathbf{W} \odot \mathbb{K}(|\mathbf{W}| > \tau) \tag{16}$$

where $\mathbf{W}$ represents the weight matrix, $\mathbb{K}(\cdot)$ is an indicator function, and $\tau$ is a threshold for pruning.

### 24.3.2 Neuron Pruning

Neuron pruning involves removing entire neurons and their corresponding connections from the network based on a specific criterion, such as the sum of the absolute values of the incoming weights or the activation values. Neurons with lower criterion values are considered less important and are pruned.

$$\mathbf{W}_{\text{pruned}} = \mathbf{W} \odot \mathbb{K}(\sum_{i,j} |w_{ij}| > \tau) \tag{17}$$

where $\mathbf{W}$ represents the weight matrix, $\mathbb{K}(\cdot)$ is an indicator function, and $\tau$ is a threshold for pruning.

### 24.3.3 Structured Pruning

Structured pruning involves removing entire structures within the network, such as channels, filters, or layers. This type of pruning not only reduces the number of parameters and computational requirements but also simplifies the network architecture, making it easier to implement on hardware accelerators. Structured pruning is typically guided by a criterion such as the sum of the absolute values of the weights or the average percentage of zeros in the activations.

$$\mathbf{W}_{pruned} = \mathbf{W} \odot \mathbb{K}(\sum i,j|w_{ij}| > \tau) \tag{18}$$

where $\mathbf{W}$ represents the weight matrix, $\mathbb{K}(\cdot)$ is an indicator function, and $\tau$ is a threshold for pruning.

These techniques can be combined in various ways to create efficient and compact neural networks that are suitable for deployment on resource-constrained devices while maintaining high performance.

# 25 Embeddings

Embeddings are a powerful technique for representing discrete objects, such as words or tokens, in continuous vector spaces. By mapping these objects to continuous vectors, models can learn meaningful relationships and similarities between them. In the context of natural language processing, various embedding methods have been developed to represent words and their meanings.

## 25.1 One-hot Encoding

One-hot encoding is a simple technique for representing categorical variables, including words, as binary vectors. In the context of text data, one-hot encoding involves creating a vocabulary of unique words present in the dataset, and then representing each word as a vector with the same

length as the vocabulary. The vector contains a "1" in the position corresponding to the word in the vocabulary and "0" everywhere else.

The main advantage of one-hot encoding is its simplicity and ease of implementation. However, it has several disadvantages, such as:

- High dimensionality: The size of the one-hot vectors grows with the size of the vocabulary, which can be large for real-world text datasets.

- Sparse representation: One-hot vectors are sparse, as they mostly contain zeros. This can be inefficient in terms of memory usage and computation.

- No semantic information: One-hot encoding does not capture any semantic relationships between words, as all one-hot vectors are orthogonal to each other.

## 25.2   TF-IDF

Term Frequency-Inverse Document Frequency (TF-IDF) is a frequency-based method for representing text data that takes into account the importance of words in the context of a collection of documents. The TF-IDF score for a word in a document is calculated as the product of two components:

- **Term Frequency (TF)**: The number of times the word appears in the document, usually normalized by the total number of words in the document, or normalized by the frequency of the most common word in the document.

- **Inverse Document Frequency (IDF)**: A measure of how rare or unique the word is across the entire collection of documents. It is calculated as the logarithm of the total number of documents divided by the number of documents containing the word.

$$tf(t, d) = \frac{f_{t,d}}{max\left\{f_{w,d} : w \in d\right\}}$$
$$df(t) = \#\left\{d : tf(t; d) > 0\right\}$$
$$idf(t) = \log\left(\frac{N}{df(t)}\right)$$
$$TFIDF(t, d) = tf(t, d) \cdot idf(t)$$

The resulting TF-IDF score assigns a higher weight to words that are more important or informative in the document, while down-weighting common words that appear frequently across multiple documents. The main advantage of TF-IDF is that it captures the context-specific importance of words and can lead to better performance in certain NLP tasks, such as text classification and information retrieval.

However, TF-IDF also has some limitations, including:

- High dimensionality: Similar to one-hot encoding, the size of the TF-IDF vectors is determined by the vocabulary size, which can be large for real-world text datasets.

- No semantic information: Although TF-IDF captures the importance of words within a document, it does not capture semantic relationships between words.

## 25.3  BM-25

$$\text{score}(D, Q) = \sum_{i=1}^{n} \text{IDF}(q_i) \cdot \frac{f(q_i, D) \cdot (k_1 + 1)}{f(q_i, D) + k_1 \cdot (1 - b + b \cdot \frac{|D|}{\text{avgdl}})} \tag{19}$$

BM-25 is an iteration on TF-IDF. It introduces two concepts.

### 25.3.1  Term Saturation

If a document contains 100 occurrences of the term "computer," it's not twice as relevant as a document that contains computers 50 times. The importance of a term plateaus. BM25 solves this issue by introducing a term, k1 that controls the shape of this saturation curve.

### 25.3.2  Document Length Normalization

Instead of looking at the frequency of a word in a document, we should look at the proportional frequency of a word relative to the length of a document. If a document is very long and contains the word computer once, it's not very relevant. If the document is only length 2, and the computer occurs once, then it's probably very important.

## 25.4  Word2Vec

Word2Vec, introduced by Mikolov et al., is a widely-used word embedding technique that learns continuous vector representations for words based on their co-occurrence patterns in large text corpora. Word2Vec consists of two main architectures: Continuous Bag-of-Words (CBOW) and Skip-Gram.

- **CBOW**: In the CBOW architecture, the model learns to predict a target word based on the context words surrounding it. This approach helps the model capture the semantic meaning of words by learning from the words that frequently appear together.

- **Skip-Gram**: The Skip-Gram architecture is the inverse of CBOW. It learns to predict the context words surrounding a target word. This approach allows the model to learn more fine-grained representations, as it focuses on predicting individual context words rather than aggregating them.

### 25.4.1  Negative Sampling for Word2Vec Models

Training Word2Vec models using the full softmax over the entire vocabulary can be computationally expensive, especially for large corpora and vocabularies. Negative sampling is an efficient optimization technique introduced to address the computational complexity of softmax in Word2Vec

models. It is a form of noise contrastive estimation (NCE), where the model learns to distinguish true context-target word pairs from artificially generated "negative" pairs. Instead of updating all the word vectors in the vocabulary for each training sample, negative sampling updates only a small subset of "negative" samples and the target word vector.

Given a true context-target word pair $(c, t)$, the objective of negative sampling is to maximize the probability of observing the true pair while minimizing the probability of observing $k$ negative samples $(c, t_i)$, where $t_i$ is a randomly sampled word from the vocabulary. The loss function for negative sampling can be formulated as follows:

$$\mathcal{L} = \log \sigma(\mathbf{v}_t^T \mathbf{v}_c) + \sum_{i=1}^{k} \log \sigma(-\mathbf{v}_{t_i}^T \mathbf{v}_c) \tag{20}$$

where $\sigma(x) = \frac{1}{1+e^{-x}}$ is the sigmoid function, $\mathbf{v}_t$ and $\mathbf{v}_c$ are the target and context word vectors, respectively, and $\mathbf{v}_{t_i}$ are the negative sample word vectors.

By updating only a small number of negative samples and the target word vector, negative sampling significantly reduces the computational complexity of training Word2Vec models. The choice of $k$ is a trade-off between computational efficiency and model performance, with smaller values of $k$ leading to faster training but potentially reduced accuracy.

### 25.4.2 Advantages

*Distributed representation*: Word2Vec learns continuous, dense vector representations for words, which can effectively capture semantic and syntactic relationships between them.

*Scalability*: Word2Vec can be trained on large corpora efficiently, thanks to the use of techniques like negative sampling and hierarchical softmax.

*Transfer learning*: Pre-trained Word2Vec embeddings can be used as input features for downstream tasks, improving performance by leveraging the learned semantic relationships.

### 25.4.3 Disadvantages of Word2Vec:

*Static embeddings*: Word2Vec learns a single, context-independent embedding for each word, which cannot capture the different meanings of a word depending on its context.

*Suboptimal handling of rare words*: Due to the nature of the training objective, Word2Vec may not learn high-quality embeddings for rare words, as it relies on co-occurrence patterns that may be sparse for infrequent words.

*Lack of interpretability*: Although Word2Vec embeddings can capture semantic relationships, it is difficult to interpret the dimensions of the learned embeddings directly.

## 25.5 GloVe Embeddings

GloVe (Global Vectors for Word Representation) is a popular unsupervised learning algorithm for obtaining dense vector representations of words. These embeddings are designed to capture both

semantic and syntactic information in a continuous, multi-dimensional space. GloVe embeddings were introduced by Pennington et al. (2014) as an alternative to the widely-used Word2Vec embeddings, aiming to address some of the limitations of the latter.

### 25.5.1 Concept and Motivation

The main idea behind GloVe is to leverage global co-occurrence information from the entire corpus to create more accurate word embeddings. Word2Vec, on the other hand, relies on local context windows to generate word representations, which can limit its ability to capture global semantic relationships between words. GloVe addresses this issue by constructing a global co-occurrence matrix that captures how frequently words co-occur across the entire corpus and then factorizing this matrix to obtain the word embeddings.

### 25.5.2 Co-occurrence Matrix and Factorization

GloVe starts by constructing a co-occurrence matrix $X$, where each entry $X_{ij}$ represents the number of times word $j$ appears in the context of word $i$ in the corpus. The context is typically defined by a fixed-size window around each word. This matrix captures the global co-occurrence statistics of the words in the corpus.

Next, GloVe factorizes the co-occurrence matrix using a weighted least squares optimization objective to obtain the word embeddings. The objective function aims to minimize the difference between the dot product of the word vectors and the logarithm of their co-occurrence count, with a weighting function to handle rare and common co-occurrences differently:

$$J = \sum_{i,j=1}^{V} f(X_{ij}) \left( w_i^T \tilde{w}_j + b_i + \tilde{b}j - \log Xij \right)^2 \tag{21}$$

Here, $w_i$ and $\tilde{w}_j$ are the word vectors for words $i$ and $j$, $b_i$ and $\tilde{b}j$ are the respective bias terms, $V$ is the size of the vocabulary, and $f(X_{ij})$ is the weighting function that assigns different weights to different co-occurrence counts.

## 25.6 ELMo

ELMo (Embeddings from Language Models) is a contextual word embedding technique introduced by Peters et al. Unlike Word2Vec and GloVe, which learn static word embeddings, ELMo learns dynamic, context-dependent word representations using a bidirectional LSTM language model. This allows ELMo to capture the different meanings of a word depending on its context, resulting in more accurate and expressive embeddings.

## 25.7 Bidirectional Encoder Representations from Transformers (BERT)

The Bidirectional Encoder Representations from Transformers (BERT) is a pre-trained language model based on the transformer architecture, which utilizes only the encoder portion of the original transformer. BERT is designed to generate bidirectional contextualized word representations, making it suitable for various natural language understanding tasks.

- **Encoder-only architecture:** BERT employs only the encoder blocks of the original transformer architecture. Each block consists of multi-head self-attention, followed by layer normalization, a position-wise feed-forward network (FFNN), and another layer normalization. These encoder blocks are stacked on top of each other to form the BERT architecture. Using only the encoder portion allows BERT to efficiently process and generate contextualized embeddings for each token in parallel, as there is no autoregressive constraint.

- **Bidirectional context:** BERT is designed to generate bidirectional contextualized word representations by processing input text in both forward and backward directions. This is achieved through the use of self-attention mechanisms in the encoder layers, allowing the model to capture dependencies and relationships between words in both directions. This bidirectional context enables BERT to perform well in various natural language understanding tasks, such as question answering, sentiment analysis, and named entity recognition.

- **Pre-training tasks:** BERT is pre-trained on two unsupervised tasks: masked language modeling (MLM) and next sentence prediction (NSP). The MLM task involves randomly masking a percentage of input tokens and training the model to predict the masked tokens based on the context provided by the remaining unmasked tokens. The NSP task trains the model to predict whether two input sentences are consecutive or not, helping BERT learn sentence-level relationships.

After pretraining, BERT can be fine-tuned on specific tasks by adding a task-specific output layer on top of the pretrained model. During fine-tuning, the entire model, including the transformer layers, is updated to adapt to the specific task.

The key advantage of BERT over previous models is its ability to effectively model bidirectional context, which helps it better understand the relationships between words and their meanings in different contexts.

## 25.8 RoBERTa

RoBERTa (Robustly Optimized BERT Approach): RoBERTa is an optimized version of BERT that improves upon its training process and hyperparameters.

- Trains the model for longer, with bigger batches, over more data

- Removes the NSP task

- Trains on longer sequences

## 25.9 T5: Text-to-Text Transfer Transformer

T5 is a transformer model introduced by that unifies various NLP tasks into a single text-to-text format. The main idea behind T5 is to cast all NLP tasks as a sequence-to-sequence problem, where both input and output are text sequences. This allows for consistent pretraining and fine-tuning across different tasks.

The T5 model is pretrained on a denoising autoencoder objective. During pretraining, some of the input tokens are masked, and the model learns to reconstruct the original sequence by predicting the masked tokens. Once the pretraining is completed, T5 can be fine-tuned on specific tasks by providing task-specific input and output pairs.

One of the key innovations of T5 is the use of a task prefix to indicate the specific task to be performed during fine-tuning. For example, for a sentiment analysis task, the input sequence could be prepended with the prefix *"sentiment: "* to inform the model that it should classify the sentiment of the subsequent text.

$$task\_prefix + task\_specific\_input \Rightarrow task\_specific\_output \tag{22}$$

## 25.10 Generative Pre-trained Transformer (GPT)

The Generative Pre-trained Transformer (GPT) is a language model based on the transformer architecture, which utilizes only the decoder portion of the original transformer. One of the key features of GPT is its autoregressive nature achieved by using masked self-attention in the decoder layers.

- **Decoder-only architecture:** GPT employs only the decoder blocks of the original transformer architecture. Each block consists of masked multi-head self-attention, followed by layer normalization, a position-wise feed-forward network (FFNN), and another layer normalization. These decoder blocks are stacked on top of each other to form the GPT architecture.

- **Positional embeddings:** GPT incorporates positional embeddings to provide the model with information about the position of each word in the input sequence. These embeddings are added to the input token embeddings before being fed into the model.

- **Autoregressive modeling and masked self-attention:** GPT generates text in an autoregressive manner, predicting the next word based on the previously generated words. To ensure autoregressive behavior, GPT uses masked self-attention in the decoder layers. The attention mechanism is modified so that each word can only attend to the words before it, preventing the model from looking into the future words during training and inference. This is accomplished by applying a mask to the self-attention scores, effectively setting the future words' scores to $-\infty$, and ensuring that they have no influence on the current word prediction.

# 26 Tokenization

Tokenization is a fundamental pre-processing step in natural language processing, which involves splitting a given text into smaller units called tokens. Tokens can be words, phrases, sentences, or subword units, depending on the requirements of the task and the modeling approach. Tokenization helps in simplifying text data and facilitates further text processing and analysis. Some common tokenization techniques are:

1. **Word Tokenization**: Word tokenization involves breaking the text into individual words. It is the most basic form of tokenization and is commonly used in many NLP tasks. Word tokenization typically splits text on whitespace and punctuation, but it may not handle contractions, possessives, or language-specific features (e.g., hyphenated words) optimally.

2. **Sentence Tokenization**: Sentence tokenization, also known as sentence segmentation, splits the text into sentences. It is often used when the goal is to analyze or process text at the sentence level, such as in sentiment analysis or machine translation. Sentence tokenization typically relies on punctuation marks and capitalization patterns to identify sentence boundaries, but may need additional rules or heuristics to handle abbreviations, quotations, or other language-specific features.

3. **Subword Tokenization**: Subword tokenization breaks the text into smaller units that are more flexible and expressive than individual words. Subword tokenization is particularly useful for handling out-of-vocabulary words, morphologically rich languages, and noisy or informal text (e.g., social media posts). Some popular subword tokenization methods include:

   - **Byte Pair Encoding (BPE)**: BPE is a data compression algorithm that is adapted for subword tokenization. It iteratively merges the most frequent character or subword pairs in the text until a predefined vocabulary size is reached.

When selecting a tokenization method, it is important to consider the characteristics of the text data, the requirements of the NLP task, and the assumptions of the modeling approach. The choice of tokenization can have a significant impact on the performance and efficiency of the NLP model, so it is crucial to experiment with different tokenization techniques and evaluate their effect on the model's performance.

## 26.1 Creating the Embedding Lookup Table

The embedding lookup table is typically created by training an unsupervised algorithm, such as Word2Vec or GloVe, on a large corpus of text. The training process aims to capture semantic and syntactic information about words in a lower-dimensional continuous space. The steps to create an embedding lookup table are as follows:

1. **Define the vocabulary:** The vocabulary is a collection of unique tokens (words) found in the training corpus. It is common to limit the vocabulary size to a fixed number of the most frequent tokens, while rare tokens are replaced with a special out-of-vocabulary (OOV) token.

2. **Initialize the embedding matrix:** Create a matrix of shape $|V| \times d$, where $|V|$ is the size of the vocabulary and $d$ is the desired dimension of the embeddings. The matrix is usually

initialized with small random values.

3. **Train the embedding algorithm:** Train an unsupervised algorithm, such as Word2Vec or GloVe, on the corpus. The goal is to learn embeddings that capture meaningful relationships between words. During training, the embedding matrix is updated to minimize a specific loss function.

4. **Normalize the embeddings:** After training, it is common to normalize the embeddings to have unit length. This step can help improve the performance of certain tasks, as it ensures that the embeddings lie on a unit hypersphere, making distance-based comparisons more meaningful.

## 26.2 Retrieving Dense Embeddings from Token IDs

Once the embedding lookup table has been created, we can retrieve dense embeddings for tokens using their corresponding IDs. Here's how this process works:

1. **Tokenize the input text:** The input text is first tokenized into a sequence of tokens. Tokenization can be performed using various techniques, such as splitting by whitespace, subword tokenization, or byte-pair encoding (BPE).

2. **Convert tokens to IDs:** Each token in the sequence is mapped to its corresponding ID based on the vocabulary. If a token is not found in the vocabulary, it is replaced with the ID of the OOV token.

3. **Retrieve embeddings:** For each token ID, the corresponding row in the embedding matrix is retrieved. This row represents the dense embedding for the given token. The result is a sequence of dense vectors that can be used as input to a machine learning model.

# 27 Learning To Rank

## 27.1 Problem Formulation

Learning to Rank (LTR) problems can be formulated in various ways, depending on the desired output and the granularity of the relevance scores. The main problem formulations are:

1. **Binary classification**: In this formulation, the LTR problem is treated as a binary classification task, where each item is assigned a label indicating whether it is relevant or not relevant. This approach is suitable when relevance judgments are binary, and the goal is to simply distinguish relevant items from non-relevant ones.

2. **Multiclass classification**: In the multiclass classification formulation, the LTR problem is treated as a task where each item is assigned one of $K$ discrete relevance labels. This approach allows for more granularity in the relevance judgments and can be useful when there is a natural ordering among the relevance levels.

3. **Regression**: In the regression formulation, the LTR problem is treated as a task where each item is assigned a continuous relevance score. This approach provides the highest level of

granularity in the relevance judgments and can be used when the goal is to produce a precise ranking of the items based on their relevance.

4. **Ordinal regression**: In the ordinal regression formulation, the LTR problem is treated as a task where each item is assigned a relevance score from an ordered set of discrete levels. This approach is a hybrid of multiclass classification and regression, as it combines the discrete nature of classification with the ordinal structure of regression. It is suitable when there is a natural ordering among the relevance levels, but the precise differences between levels are not important.

## 27.2 Ranking Approaches

There are three primary approaches to Learning to Rank (LTR), which differ in how they treat the ranking problem and the structure of the input data. The main ranking approaches are:

1. **Pointwise approach**: In the pointwise approach, each item is assigned a relevance score independently, and the items are then ranked based on these scores. This approach treats the LTR problem as a regression, binary classification, or multiclass classification task, depending on the problem formulation. It assumes that the relevance of an item can be determined without considering the other items in the list, and it focuses on learning a scoring function that maps the input features to the relevance scores.

   Some popular pointwise algorithms include linear regression, logistic regression, and support vector machines for ranking.

2. **Pairwise approach**: The pairwise approach focuses on learning the relative order between pairs of items. It treats the LTR problem as a binary classification task, where the goal is to determine which item is more relevant in each pair. The pairwise approach assumes that the relative order between items is more important than their absolute relevance scores, and it focuses on learning a scoring function that can correctly rank the pairs of items.

   Some popular pairwise algorithms include RankNet, RankBoost, and RankSVM.

3. **Listwise approach**: In the listwise approach, the entire list of items is considered together to learn the optimal ranking. This approach typically involves minimizing a loss function that directly measures the difference between the predicted ranking and the ground truth ranking, such as the ListNet loss, ListMLE loss, or the LambdaMART loss. The listwise approach assumes that the ranking of the items is dependent on the overall structure of the list, and it focuses on learning a scoring function that can generate the correct ranking for the entire list.

   Some popular listwise algorithms include ListNet, ListMLE, LambdaMART, and AdaRank.

Each ranking approach has its own advantages and drawbacks, and the choice of approach depends on the specific requirements of the LTR task and the nature of the data. Pointwise approaches are often simpler and easier to implement, but they may not capture the complex relationships between items that are important for ranking. Pairwise approaches can capture some of these relationships by focusing on the relative order between items, but they may still miss some aspects of the overall list structure. Listwise approaches are designed to capture the full list structure, but they can be more complex and computationally expensive than the other approaches.

## 27.3 Collaborative Filtering

Collaborative filtering is a technique used in recommendation systems to provide personalized recommendations based on the preferences or behavior of users. It leverages the assumption that users who have agreed in the past will continue to have similar preferences in the future. Collaborative filtering can be divided into two main approaches:

1. **User-based collaborative filtering:** In this approach, the similarity between users is calculated based on their historical preferences or behavior. Recommendations for a target user are generated based on the preferences of similar users.

2. **Item-based collaborative filtering:** This approach focuses on calculating the similarity between items based on users' interactions with them. Recommendations are made by suggesting items similar to those the target user has liked or interacted with in the past.

One limitation of collaborative filtering is the cold-start problem, which arises when the system has little or no information about new users or items.

## 27.4 Content-Based Filtering

Content-based filtering, on the other hand, focuses on the features of items to provide recommendations. It recommends items to users based on the similarity between the items they have liked or interacted with in the past and the features of other items. The main idea is that if a user has shown interest in a particular item, they are likely to be interested in other items with similar characteristics.

Content-based filtering can be applied to various types of data, such as text, images, or video, by extracting relevant features and comparing them to generate recommendations.

The main advantage of content-based filtering is that it does not suffer from the cold-start problem for new items, as recommendations can be generated based on the item features. However, it may be less effective in discovering novel or serendipitous recommendations since it only relies on the features of items users have already interacted with.

## 27.5 Multi-Armed Bandit in Learning to Rank

The Multi-Armed Bandit (MAB) problem is a well-known reinforcement learning problem, wherein an agent must make a series of decisions to maximize the cumulative reward over time. In the context of Learning to Rank, the MAB approach can be applied to explore and exploit different ranking strategies in an online setting, with the objective of optimizing user engagement, satisfaction, or other metrics.

### 27.5.1 MAB Formulation for Learning to Rank

In the Learning to Rank setting, each "arm" of the multi-armed bandit corresponds to a ranking strategy or model. When a user submits a query, the MAB algorithm selects one of the ranking strategies and presents the ranked results to the user. The user's interaction with the search results

(e.g., clicks, dwell time, etc.) is considered as the reward signal. The MAB algorithm learns from these reward signals to decide which ranking strategy to use in the future.

### 27.5.2 MAB Algorithms for Learning to Rank

Several MAB algorithms have been adapted to Learning to Rank, including:

**Epsilon-Greedy:** A simple algorithm that selects the best-known ranking strategy with probability $1 - \epsilon$ and chooses a random strategy with probability $\epsilon$. This approach balances exploration and exploitation, with the $\epsilon$ parameter controlling the trade-off.

**Upper Confidence Bound (UCB):** An algorithm that selects the ranking strategy with the highest upper confidence bound on the expected reward. This balances exploration and exploitation by accounting for both the observed rewards and uncertainty in the estimates.

**Thompson Sampling:** A Bayesian approach that maintains a probability distribution over the expected reward of each ranking strategy. At each step, a sample is drawn from each distribution, and the strategy with the highest sampled reward is selected. This algorithm balances exploration and exploitation by implicitly considering both the observed rewards and uncertainty in the estimates.

### 27.5.3 Advantages of MAB for Learning to Rank

**Adaptability:** MAB algorithms can adapt to changing user preferences or new ranking strategies, as they continuously learn from user feedback.

**Efficient Exploration:** MAB algorithms efficiently balance exploration and exploitation, ensuring that the most promising ranking strategies are used more often while still exploring alternative strategies.

**Online Learning:** MAB algorithms operate in an online setting, making them suitable for real-time ranking optimization in response to user interactions.

### 27.5.4 Challenges and Limitations

**Reward Attribution:** Determining the appropriate reward signal from user interactions can be challenging, as users may interact with search results for various reasons, and some interactions may be noisy or biased.

**Delayed Feedback:** In some cases, the reward signal may be delayed or only partially observable, making it challenging for the MAB algorithm to learn effectively.

**Contextual Information:** Incorporating contextual information, such as user features or query characteristics, into the MAB framework can be complex but is often necessary for improving the quality of the ranking.

Despite these challenges, the multi-armed bandit approach has shown promise in Learning to Rank applications, providing an efficient and adaptive way to optimize ranking strategies based on user feedback.

## 27.6 LambdaRank

LambdaRank is an extension of RankNet, a pairwise learning-to-rank algorithm. RankNet is based on neural networks and utilizes the gradient descent optimization technique. LambdaRank aims to improve RankNet by directly optimizing the ranking metric, which is usually non-differentiable. In this section, we will discuss the key concepts of LambdaRank, its algorithm, and its advantages over RankNet.

### 27.6.1 Background

Before diving into LambdaRank, let us briefly review RankNet. RankNet is a pairwise learning-to-rank method that learns the optimal ranking function by minimizing the average number of inversions in the ranked list. It uses a neural network to model the ranking function and employs a cross-entropy loss function for optimization.

However, RankNet has certain limitations. It does not optimize ranking metrics directly, such as Normalized Discounted Cumulative Gain (NDCG) or Mean Average Precision (MAP). LambdaRank addresses this issue by introducing a new loss function that directly optimizes the ranking metric.

### 27.6.2 LambdaRank Algorithm

LambdaRank modifies the gradient descent optimization of RankNet to incorporate the ranking metric. The main idea is to compute the Lambda, which represents the amount of change in the ranking metric caused by swapping two documents. The algorithm can be summarized as follows:

1. Initialize a neural network with random weights.

2. For each query, obtain the ranked list of documents by computing their scores using the neural network.

3. Calculate the pairwise Lambda values for all pairs of documents with different relevance levels.

4. Update the neural network weights using the Lambda values and a chosen learning rate.

5. Iterate steps 2 to 4 until the algorithm converges or a stopping criterion is met.

The Lambda value, denoted as $\lambda_{ij}$, is the amount of change in the ranking metric when the positions of documents $i$ and $j$ are swapped. It can be computed as follows:

$$\lambda_{ij} = \frac{\partial C_{ij}}{\partial s_i} = \frac{1}{1 + e^{s_i - s_j}} (\Delta \text{metric}_{ij}) \tag{23}$$

where $C_{ij}$ is the cost function, $s_i$ and $s_j$ are the scores of documents $i$ and $j$, and $\Delta \text{metric}_{ij}$ is the change in the ranking metric due to swapping the positions of documents $i$ and $j$. The Lambda values are used to update the neural network weights.

### 27.6.3   Advantages of LambdaRank

LambdaRank has several advantages over RankNet:

- Direct optimization of ranking metrics: LambdaRank directly optimizes the chosen ranking metric (e.g., NDCG or MAP), which leads to better ranking performance.

- Efficient computation: LambdaRank computes Lambda values efficiently, making it suitable for large-scale ranking problems.

- Flexibility: LambdaRank can be easily adapted to optimize different ranking metrics by simply changing the computation of $\Delta\text{metric}_{ij}$.

### 27.6.4   Example

Suppose we have a query $q$ and its corresponding document set $D = d_1, d_2, d_3$ with relevance labels $y = 2, 3, 1$.

## 27.7   Example: LambdaRank

1. Initialize the neural network with random weights $w$.

2. Compute the scores for all documents in $D$ using the neural network

$$s_1 = f(d_1, w)$$
$$s_2 = f(d_2, w)$$
$$s_3 = f(d_3, w)$$

   Rank the documents based on their scores. Assume the current ranking order is $(d_2, d_1, d_3)$.

3. Calculate the pairwise Lambda values for all pairs of documents with different relevance levels. For example, compute $\lambda_{12}$, $\lambda_{13}$, and $\lambda_{23}$:

$$\lambda_{12} = \frac{\partial C_{12}}{\partial s_1} = \frac{1}{1 + e^{s_1 - s_2}}(\Delta\text{metric}_{12})$$
$$\lambda_{13} = \frac{\partial C_{13}}{\partial s_1} = \frac{1}{1 + e^{s_1 - s_3}}(\Delta\text{metric}_{13})$$
$$\lambda_{23} = \frac{\partial C_{23}}{\partial s_2} = \frac{1}{1 + e^{s_2 - s_3}}(\Delta\text{metric}_{23})$$

   Here, $\Delta\text{metric}_{ij}$ is the change in the ranking metric due to swapping the positions of documents $i$ and $j$. For example, if we are optimizing NDCG, we can compute $\Delta\text{metric}_{ij}$ as:

$$\Delta\text{metric}_{ij} = \text{NDCG}(y_{swap(i,j)}) - \text{NDCG}(y)$$

   where $y_{swap(i,j)}$ is the relevance vector $y$ after swapping the positions of documents $i$ and $j$.

4. Update the neural network weights $w$ using the Lambda values and a chosen learning rate $\eta$:

$$w_{new} = w_{old} - \eta \sum_{i=1}^{3} \frac{\partial f(d_i, w)}{\partial w} \sum_{j:y_j \neq y_i} \lambda_{ij}$$

5. Iterate steps 2 to 4 until the algorithm converges or a stopping criterion is met.

This example demonstrates the process of training a LambdaRank model on a simple dataset. In practice, the algorithm is applied to much larger datasets with more complex neural network architectures.

## 27.8   Listwise Ranking Losses

Listwise ranking methods aim to optimize a ranking function based on the entire ranked list of documents. Different listwise ranking methods use different cost functions to measure the quality of the ranking. Here are some popular cost functions used in listwise ranking methods:

1. **ListMLE** (Listwise Maximum Likelihood Estimation) uses a likelihood-based cost function. It models the probability of the observed ranked list and maximizes the likelihood of the true list. The cost function for ListMLE can be defined as:

$$L(\theta) = - \sum_{i=1}^{n} \log P(y_i | \mathbf{x}_i, \theta) \tag{24}$$

where $n$ is the number of queries, $y_i$ is the true ranked list for query $i$, $\mathbf{x}_i$ is the feature vector for the documents, and $\theta$ represents the model parameters.

2. **ListNet** is another listwise ranking method that uses a cross-entropy loss function. ListNet compares the predicted ranked list with the true ranked list by transforming them into probability distributions. The cost function for ListNet is defined as:

$$L(\theta) = - \sum_{i=1}^{n} \sum_{j=1}^{m} P_{ij}^{(g)} \log P_{ij}^{(f)} \tag{25}$$

where $P_{ij}^{(g)}$ is the ground-truth probability distribution, $P_{ij}^{(f)}$ is the predicted probability distribution, $n$ is the number of queries, and $m$ is the number of documents for each query.

3. **SoftRank** uses a smooth approximation of the ranking metric, such as NDCG, as its cost function. SoftRank introduces a temperature parameter to smooth out the original ranking metric and make it differentiable. The cost function for SoftRank can be defined as:

$$L(\theta) = \frac{1}{n} \sum_{i=1}^{n} -\text{SmoothedNDCG}(f(\mathbf{x}_i, \theta), y_i) \tag{26}$$

where $n$ is the number of queries, $\mathbf{x}_i$ is the feature vector for the documents, $y_i$ is the true ranked list, and $\theta$ represents the model parameters.

These cost functions are specifically designed for listwise ranking methods, which focus on optimizing the ranking model by considering the entire ranked list of documents.

## 27.9 Evaluation Metrics

Evaluation metrics are crucial for assessing the quality of Learning to Rank (LTR) algorithms, as they provide a quantitative measure of how well the predicted ranking aligns with the ground truth ranking. There are several evaluation metrics commonly used in LTR, each with its own strengths and weaknesses. Some of the most popular metrics are:

### 27.9.1 Mean Average Precision (MAP)

MAP is a widely used metric in information retrieval that measures the average precision at different recall levels across all queries. It is especially useful when dealing with binary relevance judgments, as it focuses on the precision of the retrieved relevant items. MAP is sensitive to the order of the ranked items and penalizes false positives.

Let's consider the same small dataset with five documents and their binary relevance scores (0 for non-relevant and 1 for relevant) for a given query:

| Document | True Binary Relevance | Predicted Ranking |
|----------|-----------------------|-------------------|
| A        | 1                     | 2                 |
| B        | 0                     | 4                 |
| C        | 1                     | 3                 |
| D        | 1                     | 1                 |
| E        | 1                     | 5                 |

Precision at rank $k$ is defined as the number of relevant documents among the top $k$ documents divided by $k$. To compute the MAP, we first calculate the precision values at each rank where a relevant document is retrieved:

$$P(1) = \frac{1}{1} = 1$$
$$P(2) = \frac{2}{2} = 1$$
$$P(3) = \frac{3}{3} = 1$$
$$P(5) = \frac{4}{5} = 0.8$$

Next, we compute the average precision (AP) by taking the average of the precision values at each relevant rank:

$$AP = \frac{1 + 1 + 1 + 0.8}{4} = 0.95$$

In this example, the dataset contains only one query, so the MAP is equal to the AP of that query:

$$MAP = 0.95$$

The MAP value of 0.95 indicates that the predicted ranking is relatively close to the ideal ranking, as the average precision across the relevant documents is high. Similar to NDCG, MAP takes into account both the relevance of the documents and their positions in the ranking, making it a suitable metric for evaluating the performance of ranking algorithms.

### 27.9.2 NDCG (Normalized Discounted Cumulative Gain) & DCG

is a popular evaluation metric for ranking algorithms, but it is not directly differentiable as it involves discrete ranking positions. To use NDCG in a differentiable manner, researchers often employ surrogate loss functions that approximate NDCG and are differentiable, such as RankNet, LambdaMART, or SoftRank.

NDCG is calculated based on DCG (Discounted Cumulative Gain) and IDCG (Ideal Discounted Cumulative Gain). The formula for DCG is as follows:

$$DCG_p = \sum_{i=1}^{p} \frac{rel_i}{\log_2{(i+1)}}$$

where $rel_i$ is the relevance score of the document at position $i$ and $p$ is the number of positions considered.

The Ideal Discounted Cumulative Gain (IDCG) is the maximum possible DCG for a given query, obtained by sorting the documents in descending order of their relevance scores.

$$IDCG_p = \sum_{i=1}^{p} \frac{rel_i^*}{\log_2{(i+1)}}$$

where $rel_i^*$ is the relevance score of the document at position $i$ in the ideal ranking.

Finally, NDCG is computed by normalizing DCG with respect to IDCG:

$$NDCG_p = \frac{DCG_p}{IDCG_p}$$

NDCG ranges between 0 and 1, where 1 indicates a perfect ranking and 0 indicates the worst possible ranking. This metric is useful for evaluating ranking algorithms as it takes into account both the relevance of the documents and their positions in the ranking, with higher weights assigned to more relevant documents ranked higher in the list.

| Document | True Relevance | Predicted Ranking |
|----------|----------------|-------------------|
| A | 3 | 2 |
| B | 1 | 4 |
| C | 2 | 3 |
| D | 4 | 1 |
| E | 2 | 5 |

**Example**

First, we calculate the DCG for the predicted ranking:

$$DCG_5 = \frac{rel_1}{\log_2(1+1)} + \frac{rel_2}{\log_2(2+1)} + \frac{rel_3}{\log_2(3+1)} + \frac{rel_4}{\log_2(4+1)} + \frac{rel_5}{\log_2(5+1)}$$

where $rel_i$ corresponds to the relevance score of the document at position $i$ in the predicted ranking. Based on the predicted ranking, we have:

$$DCG_5 = \frac{4}{\log_2(1+1)} + \frac{3}{\log_2(2+1)} + \frac{2}{\log_2(3+1)} + \frac{1}{\log_2(4+1)} + \frac{2}{\log_2(5+1)} \approx 8.10$$

Next, we compute the IDCG by sorting the documents in descending order of their true relevance scores:

$$IDCG_5 = \frac{rel_1}{\log_2(1+1)} + \frac{rel_2}{\log_2(2+1)} + \frac{rel_3}{\log_2(3+1)} + \frac{rel_4}{\log_2(4+1)} + \frac{rel^*_5}{\log_2(5+1)}$$

In the ideal ranking, the documents are sorted as D, A, C, E, B. Thus, we have:

$$IDCG_5 = \frac{4}{\log_2(1+1)} + \frac{3}{\log_2(2+1)} + \frac{2}{\log_2(3+1)} + \frac{2}{\log_2(4+1)} + \frac{1}{\log_2(5+1)} \approx 8.14$$

Finally, we compute the NDCG by dividing the DCG by the IDCG:

$$NDCG_5 = \frac{DCG_5}{IDCG_5} \approx \frac{8.10}{8.14} \approx 0.995$$

The NDCG for this small dataset is approximately 0.953, indicating that the predicted ranking is relatively close to the ideal ranking.

### 27.9.3 Precision at K (P@K)

P@K is a metric that measures the proportion of relevant items among the top K items in the ranking list. It is a simple and intuitive metric that focuses on the quality of the top-ranked items, making it suitable for scenarios where only the top results are important, such as search engines or recommendation systems. However, P@K does not consider the order of the items within the top K results.

### 27.9.4 Recall at K (R@K)

R@K is a metric that measures the proportion of relevant items retrieved in the top K results relative to the total number of relevant items. It is useful for scenarios where it is important to retrieve as many relevant items as possible within the top K results, but it does not consider the order of the items within the top K results.

### 27.9.5 F-score at K (F@K)

F@K is a metric that combines precision and recall at K using the harmonic mean, providing a balanced measure of both aspects. It is useful when both precision and recall are important, but like P@K and R@K, it does not consider the order of the items within the top K results.

### 27.9.6 Mean Reciprocal Rank (MRR)

MRR is a metric that measures the average of the reciprocal ranks of the first relevant item in each ranking list. It is useful for scenarios where the goal is to retrieve at least one relevant item as early as possible in the ranking list, but it is not sensitive to the ranking of the other relevant items.

# 28 Search

## 28.1 Inverted Indices

An inverted index is a fundamental data structure used in information retrieval systems to enable efficient querying and searching of documents. Given a collection of documents, the main goal of an inverted index is to quickly retrieve a list of documents that contain a specific query term or a set of terms.

### 28.1.1 Construction of an Inverted Index

The construction of an inverted index involves several steps:

1. **Tokenization:** Split the text in each document into words or terms.

2. **Normalization:** Perform case-folding, stemming, and stopword removal to normalize the terms.

3. **Posting Lists:** Create a posting list for each unique term, which contains document IDs of the documents containing the term.

4. **Inverted Index:** Combine all posting lists into a single data structure, typically a hashmap, where the keys are the terms and the values are the corresponding posting lists.

### 28.1.2 Querying an Inverted Index

To process a query using an inverted index, follow these steps:

1. **Parse and Normalize Query:** Tokenize the query and normalize its terms in the same manner as the document collection.

2. **Retrieve Posting Lists:** For each query term, retrieve its corresponding posting list from the inverted index.

3. **Merge Posting Lists:** Depending on the query operator (e.g., AND, OR), merge the posting lists to produce a final list of matching documents.

4. **Rank Documents:** Apply a retrieval model (e.g., TF-IDF, BM25) to rank the matching documents according to their relevance to the query.

Inverted indices significantly reduce the search space and improve the efficiency of information retrieval systems. However, they may require large amounts of memory to store the posting lists, especially for extensive document collections or vocabularies.

## 28.2 Query Processing and Optimization

Query processing and optimization are crucial steps in the search pipeline, aiming to improve search efficiency and the relevance of search results. This subsection discusses key aspects of query processing and optimization.

### 28.2.1 Query Tokenization and Normalization

Query tokenization is the process of converting a user query into a sequence of tokens, usually words or phrases. Normalization includes transforming tokens into a canonical form by applying techniques such as lowercasing, stemming, and lemmatization.

- **Lowercasing:** Convert all tokens to lowercase to ensure consistent matching with indexed documents.

- **Stemming:** Reduce words to their root form, e.g., "running" to "run", to match different forms of the same word.

- **Lemmatization:** Convert words to their base form using morphological analysis, e.g., "went" to "go", to account for inflectional variations.

### 28.2.2 Query Expansion

Query expansion aims to improve search recall by including additional related terms in the query. Common query expansion techniques are:

- **Synonyms:** Add synonyms of query terms to retrieve documents containing semantically similar words.

- **Stemming and Lemmatization:** Include stemmed or lemmatized forms of query terms to account for morphological variations.

### 28.2.3 Relevance Feedback and Pseudo-Relevance Feedback

Relevance feedback involves using user feedback on the relevance of initial search results to refine the query and retrieve more relevant documents. Pseudo-relevance feedback assumes that top-ranked documents in the initial search results are relevant and uses them to expand the query automatically.

- **Rocchio Algorithm:** Update the query vector by adding a linear combination of relevant and non-relevant document vectors.

- **Query Expansion using Top-k Documents:** Select top-ranked terms from the top-k search results and add them to the query.

# 29 Tail Query Expansion

Tail query expansion focuses on improving the performance of information retrieval systems for long-tail queries, which are less frequent and often more specific than popular queries. These queries usually have lower recall due to the limited number of relevant documents and the sparse use of specific terms. Tail query expansion aims to enhance search results by identifying and incorporating semantically related terms from the long-tail distribution.

## 29.1 Tail Query Expansion Using Query Logs and Click-through Data

Utilizing query logs and click-through data is another approach for tail query expansion that leverages user behavior analysis to identify patterns and relationships in long-tail queries. This method extracts additional terms that are frequently associated with the original query terms, thereby improving recall by incorporating user-generated knowledge into the query expansion process.

Query logs contain records of users' past queries, while click-through data provide information on which search results were clicked by users in response to their queries. Analyzing these data sources can reveal relationships between terms that may not be evident from the document collection alone. There are several techniques to exploit this information for query expansion:

1. **Query reformulation**: By examining the query logs, we can identify common query reformulations where users modify their initial query to obtain better search results. If a particular reformulation is frequently observed, we can infer that the terms added or modified in the reformulated query may be relevant for expanding the original query.

2. **Term co-occurrence analysis**: Analyzing the co-occurrence of terms in query logs or click-through data can help identify terms that are semantically related to the original query terms. For instance, if two terms frequently appear together in user queries or in clicked search results, they may be related and could be used for query expansion.

3. **Session-based expansion**: In this approach, we consider the entire search session of a user as a context for query expansion. If a user submits multiple queries within a single session and clicks on related search results, we can assume that the terms from those queries and results are related to the user's information need. These terms can be used to expand the original query.

## 29.2   Noise in Query Expansion

Query expansion aims to improve the recall of search results by incorporating related terms into the original query. However, this process may introduce noise, which can negatively impact the precision of search results. In this section, we discuss the sources of noise in query expansion and various techniques to mitigate their effects.

### 29.2.1   Sources of Noise

Noise in query expansion can arise from several sources:

1. **Ambiguity:** Expanding a query with ambiguous terms, such as homonyms or polysemous words, can lead to irrelevant search results if the expanded terms do not match the user's intended meaning.

2. **Irrelevant expansions:** The inclusion of terms that are semantically related to the query terms but not relevant to the user's information need can introduce noise. For example, expanding a query for "apple" with "orange" may be semantically related but may not be relevant to a user searching for information about the technology company Apple Inc.

3. **Over-expansion:** Including too many expanded terms can dilute the focus of the original query and result in the retrieval of less relevant documents.

### 29.2.2   Mitigating Noise

To mitigate the effects of noise in query expansion, several techniques can be employed:

1. **Context-aware expansion:** Utilize context-aware methods like word embeddings or context-sensitive language models (e.g., BERT, ELMo) to ensure that expanded terms are semantically relevant to the user's query. Additionally, leveraging user behavior data, such as click-through data or past queries, can help provide context and filter out ambiguous or irrelevant terms.

2. **Term weighting:** Apply term weighting techniques to assign lower weights to expanded terms compared to the original query terms. This approach helps maintain the focus on the original query while considering the expanded terms as supplementary information.

3. **Expansion term selection:** Limit the number of expanded terms and set a threshold for term similarity to ensure only highly related terms are included in the expansion. This helps prevent over-expansion and reduces the risk of including irrelevant terms.

4. **User feedback:** Incorporate user feedback, such as relevance judgments or clicks on search results, to refine the expanded terms or adjust term weights based on the user's perceived relevance. This enables a more personalized and targeted query expansion process.

5. **Evaluation and tuning:** Regularly evaluate the performance of the query expansion process using appropriate metrics (e.g., precision, recall, F1-score, MAP, NDCG) and methodologies (e.g., cross-validation, A/B testing, user studies). Based on the evaluation results, tune the expansion parameters and techniques to optimize the balance between recall improvement and noise mitigation.

By employing these strategies, it is possible to reduce the impact of noise in query expansion and enhance the overall effectiveness of information retrieval systems.

## 29.3 Language Models in Information Retrieval

Language models are widely used in information retrieval to estimate the likelihood of a query given a document.

### 29.3.1 Probabilistic Classic Models

Probabilistic classic models, such as the Binary Independence Model (BIM) and the Probabilistic Information Retrieval (PIR) Model, use probability theory to rank documents based on their relevance to a query.

- **Binary Independence Model (BIM):** BIM assumes that terms in a document are independent of each other and models the probability of relevance given the presence or absence of terms in the document. BIM computes a relevance score based on term weights, which are typically computed using term frequency and inverse document frequency.

- **Probabilistic Information Retrieval (PIR) Model:** PIR extends the BIM by introducing term dependencies and considering document length normalization. It calculates the probability of relevance by comparing the document term frequencies to the term frequencies in the entire collection.

### 29.3.2 Language Models

Language models in information retrieval estimate the probability of generating a query given a document. The documents are ranked based on their likelihood of generating the query.

**Pros:**

- Language models provide a natural way of incorporating query term frequency in the ranking process.

- Smoothing techniques can be easily incorporated to handle data sparsity and unseen terms.

- Language models can be extended to consider term dependencies and phrase queries.

### 29.3.3 Smoothing

Smoothing is required in language models for information retrieval to address the problem of data sparsity and avoid zero probabilities for unseen terms in a document. When estimating the probability of a term in a document using the maximum likelihood estimation (MLE), if a term does not appear in a document, its probability will be zero. This can lead to the entire probability of the query given the document being zero, which is undesirable for ranking purposes. Smoothing helps to assign a non-zero probability to unseen terms, preventing the issue of zero probabilities and providing more robust and reliable probability estimates.

Here are some commonly used smoothing techniques and their formulas:

**Dirichlet smoothing** - This method assumes a Dirichlet prior on the document language model.

$$P(t|D) = \frac{c(t, D) + \mu \cdot P(t|C)}{|D| + \mu}$$

where $P(t|D)$ is the smoothed probability of term $t$ in document $D$, $c(t, D)$ is the count of term $t$ in document $D$, $P(t|C)$ is the probability of term $t$ in the collection $C$, —D— is the length of the document, and $\mu$ is a hyperparameter that controls the strength of the prior.

**Bayesian smoothing with Gaussian priors** - This method assumes a Gaussian prior on the term probabilities in the document language model.

$$P(t|D) = \frac{c(t, D) + a}{|D| + b}$$

where $P(t|D)$ is the smoothed probability of term $t$ in document $D$, $c(t, D)$ is the count of term $t$ in document $D$, $|D|$ is the length of the document, and $a$ and $b$ are hyperparameters that control the strength and shape of the prior.

By applying these smoothing techniques, the language models become more robust and better at estimating term probabilities for unseen terms, ultimately improving the performance of information retrieval systems.

## 29.4 ANNs

ANNs can improve search performance by enabling efficient and scalable retrieval of semantically similar documents in high-dimensional spaces. They are particularly useful for tail queries as they can identify relevant documents even if the exact terms do not match, by exploiting the semantic relationships captured in the embeddings.

However, there are certain scenarios where ANNs might face challenges:

- **Similar embeddings:** If the search space contains a large number of very similar embeddings, the ANN algorithm may struggle to differentiate between them, leading to suboptimal results.

- **Live updating:** Incorporating new data or updates into the ANN index can be difficult, as many ANN algorithms are not designed for efficient incremental updates. This can be a challenge in dynamic environments where the underlying data frequently changes.

- **Lack of diversity:** ANN algorithms often prioritize speed and efficiency over diversity in search results. This could lead to the retrieval of a large number of similar items, reducing the overall usefulness of the search results for users.

Despite these challenges, combining embeddings and ANNs can offer significant benefits for search, including improved semantic understanding and scalability, making them a powerful tool

# 30 Scaling Large Language Models

## 30.1 Scaling Laws for Large Language Models

Scaling laws for large language models describe the relationship between model performance, model size, training data size, and computational resources. These laws help guide the development of more efficient and effective models by providing insights into how improvements can be achieved by increasing the scale of these factors. Key aspects of scaling laws include:

1. **Model size:** As the number of parameters in a language model increases, its performance typically improves, up to a certain point. Larger models are generally capable of capturing more complex patterns and relationships in the data, leading to better performance on various NLP tasks. The loss can be predicted using

$$L \sim N^{-\alpha}$$

   Where N is the model size, measured in the number of parameters. $\alpha$ generally falls in the range of 0.005 to 0.2 depending on the architecture and task.

2. **Training data size:** Increasing the amount of training data available to a language model can lead to significant improvements in its performance. This is particularly true for larger models, which have the capacity to effectively utilize more data to learn more nuanced representations. The loss can be predicted using

$$L \sim D^{-\beta}$$

   Where D is the dataset size measured in tokens, and $\beta$ is another constant that depends on the task and model architecture. Typically it's between 0.1 and 0.2

3. **Computational resources:** The amount of computational resources dedicated to training a model, such as the number of GPUs or TPUs, affects the model's performance. Increased computational resources allow for faster training, larger batch sizes, and more advanced optimization techniques, all of which can contribute to better model performance. The loss can be predicted using

$$L \sim C^{-\lambda}$$

   Where C is the compute measured in PF-days, and $\lambda$ is another constant that depends on the task and model architecture. Typically it's between 0.05 and 0.15.

Using simple algebra we can relate all of these relationships together, and determine how much we need to scale model size and computational resources if we scale training size.

## 30.2 Batch Size for Training Large Language Models

Training large language models requires careful consideration of batch size, as it has a significant impact on computational efficiency, gradient noise scale, and generalization. In this subsection, we discuss the importance of scaling the batch size when increasing data, compute, and model size, as well as the rationale behind using smaller batch sizes during fine-tuning.

1. **Computational Efficiency:** Utilizing larger batch sizes helps maintain computational efficiency by maximizing parallelism and throughput on modern hardware accelerators, such as GPUs and TPUs. This is particularly important when dealing with larger models and datasets, as efficient use of computational resources is crucial for reducing training time and cost. For pre-training large language models, batch sizes typically range from thousands to tens of thousands of examples, depending on the specific model and hardware setup.

2. **Gradient Noise Scale:** The gradient noise scale, defined as the ratio of the standard deviation of the stochastic gradient to the actual gradient, plays a crucial role in convergence and generalization. According to the power-law scaling relationship, increasing the batch size leads to a decrease in gradient noise scale, which in turn can lead to better convergence and generalization. As the data and model size increase, it is often necessary to scale the batch size to maintain an optimal level of gradient noise.

3. **Generalization:** As the data and model size increase, so does the capacity of the model. To effectively learn from larger datasets, it is often beneficial to use larger batch sizes to maintain a balance between learning the underlying patterns and avoiding overfitting to the training data. This helps ensure that the model generalizes well to new, unseen data.

4. **Fine-tuning:** During fine-tuning, it is common to use smaller batch sizes compared to pre-training. The primary reason for this is that fine-tuning often involves adapting the model to a specific task or domain with a limited amount of labeled data. Using a smaller batch size helps prevent overfitting by introducing more gradient noise, which acts as a form of regularization. Additionally, a smaller batch size can lead to faster convergence when fine-tuning, as updates are made more frequently.

However, it is important to note that scaling the batch size indefinitely may not always lead to improved performance. There is a trade-off between the benefits of larger batch sizes (such as reduced gradient noise and better convergence) and the diminishing returns on parallelism as the batch size grows. In practice, researchers typically experiment with different batch sizes to find the optimal value that leads to the best performance for a given model and dataset.

## 30.3 Upper-Bounds and Scaling Relationships

While scaling laws provide insights into how increasing various factors can improve language model performance, there are often practical limits and diminishing returns. As model size and training data size increase, the improvements in model performance tend to become smaller. This suggests that there may be an optimal point beyond which further increases in model size or training data size yield diminishing returns on performance improvement.

### 30.3.1   Distributed Training

Distributed training is essential for scaling large language models, as it enables the use of multiple devices to work on different parts of the model or data simultaneously. There are three primary approaches to distributed training:

1. **Data Parallelism:** In data parallelism, multiple devices process different mini-batches of data simultaneously while using the same model. After each iteration, the gradients calculated by each device are averaged and applied to the model's weights. This approach is relatively simple to implement and provides good speedup when training on multiple devices. However, it may suffer from communication overhead, especially when using a large number of devices or when the model's size is large.

2. **Model Parallelism:** Model parallelism partitions the model itself across multiple devices, so that each device is responsible for computing a portion of the model's forward and backward pass. This approach is particularly useful when the model is too large to fit into a single device's memory. The challenge in model parallelism is balancing the workload across devices and minimizing communication overhead, as intermediate results need to be exchanged between devices during both the forward and backward passes.

3. **Pipeline Parallelism:** In pipeline parallelism, the model is divided into multiple stages, each of which is processed by a separate device. The input data is passed through the stages in a pipelined fashion, with each device computing its portion of the forward and backward passes before passing the result to the next device. This approach can be effective in reducing communication overhead and allows for concurrent processing of multiple mini-batches. However, it requires careful scheduling and balancing of the workload between devices to avoid idle time and maintain high utilization.

### 30.3.2   Memory Optimizations

As model and dataset sizes grow, memory optimization techniques become increasingly important to efficiently utilize the available hardware resources. Several methods for memory optimization include:

1. **Mixed-Precision Training:** Mixed-precision training involves using lower-precision data types, such as float16 or bfloat16, for some parts of the model's computation, while still maintaining higher precision (e.g., float32) for critical operations. This reduces memory usage and can speed up training, as lower-precision operations are typically faster on specialized hardware. To avoid loss of numerical precision, techniques such as loss scaling and dynamic loss scaling can be employed during gradient computation and weight updates.

2. **Gradient Checkpointing:** Gradient checkpointing is a technique that trades off computation for memory savings during the backward pass of model training. Instead of storing all intermediate activations in memory, only a subset of them are stored, and the rest are recomputed on-the-fly during the backward pass. This reduces memory usage at the cost of additional computation. The choice of which activations to store and which to recompute can be made based on heuristics or more sophisticated algorithms to minimize the overall computational cost.

3. **Memory-Efficient Optimizers:** Some optimization algorithms, such as Adam and its variants, require additional memory to store per-parameter state, such as first and second moment estimates. Memory-efficient optimizers, such as LAMB and AdaBelief, use techniques to reduce memory overhead while still providing competitive optimization performance. These optimizers can be particularly useful when training large models with limited memory resources.

By leveraging distributed training techniques and memory optimization methods, it is possible to train large language models more efficiently, enabling researchers and practitioners to tackle more complex NLP tasks and push the boundaries of what is possible with current hardware.

# 31 Case Studies

## 31.1 Named Entity Recognition (NER)

Named Entity Recognition (NER) is a subtask of information extraction that involves identifying and classifying named entities, such as person names, organizations, and locations, in text. NER is essential for various NLP applications, including question answering, sentiment analysis, and relationship extraction. The following steps outline how NER is typically performed:

1. **Data Collection and Annotation:** Obtain a labeled dataset with annotated named entities. Often, this dataset consists of sentences or documents with entities tagged using the IOB (Inside, Outside, Beginning) or BIOES (Beginning, Inside, Outside, End, Single) schemes. Standard NER datasets include CoNLL-2003, OntoNotes, and MUC-7.

2. **Preprocessing and Tokenization:** Clean and preprocess the text by removing unnecessary characters, correcting misspellings, and converting the text to lowercase, if necessary. Tokenize the text into words or subwords, depending on the chosen model architecture. Tokenization can be performed using pre-trained tokenizers or custom tokenization methods, depending on the problem domain and language.

3. **Feature Extraction:** Convert the tokenized text into feature vectors that can be fed into the NER model. This step can involve using word embeddings (e.g., Word2Vec, GloVe, or FastText), character embeddings, or contextualized word representations generated by pre-trained language models (e.g., BERT, ELMo, or RoBERTa). Additionally, linguistic features such as part-of-speech tags, dependency parse trees, and capitalization patterns can be used to enhance the model's input representation.

4. **Model Selection and Training:** Choose an appropriate model architecture for NER. Deep learning models, such as recurrent neural networks (RNNs) with long short-term memory (LSTM) or gated recurrent units (GRU), bidirectional LSTMs (BiLSTMs), and transformer-based models (e.g., BERT or RoBERTa), have achieved state-of-the-art performance on NER tasks. Fine-tune the selected model on the training dataset using an appropriate loss function (e.g., cross-entropy loss) and optimization algorithm (e.g., Adam or AdamW).

5. **Evaluation and Hyperparameter Tuning:** Evaluate the performance of the NER model on a validation set during training to monitor overfitting and select the best model. Common evaluation metrics for NER tasks include precision, recall, F1-score, and entity-level F1-score.

Fine-tune the model hyperparameters, such as learning rate, batch size, and regularization, based on the validation set performance.

6. **Inference and Post-processing:** Apply the trained NER model to unseen text data to identify and classify named entities. Use the model's output, typically in the form of entity class probabilities, to generate the most likely sequence of entity labels for a given input. Post-processing techniques, such as the Viterbi algorithm, can be used to enforce label consistency and correct potential label sequence inconsistencies.

7. **Deployment and Continuous Improvement:** Integrate the trained NER model into downstream NLP applications and systems. Continuously monitor and update the model as new data becomes available or as the underlying language and named entity distribution evolve over time.

## 31.2 Text Summarization with Transformer-based Systems

Text summarization is the process of creating a shorter version of a given text that conveys the most important information from the original document. There are two primary types of summarization: extractive, which selects important sentences or phrases from the source text, and abstractive, which generates a new summary that captures the essence of the original text. In this subsection, we focus on transformer-based systems for text summarization, emphasizing context-length issues and potential solutions.

1. **Preprocessing and Tokenization:** Clean the input text by removing unnecessary characters, correcting misspellings, and handling special symbols. Tokenize the text into subwords or words using a pre-trained tokenizer associated with the chosen transformer model, such as BERT, GPT, or T5.

2. **Model Selection and Fine-tuning:** Choose a suitable pre-trained transformer-based model for text summarization. Models such as BART, T5, and PEGASUS have been specifically designed for sequence-to-sequence tasks like summarization. Fine-tune the selected model on a labeled summarization dataset, such as the CNN/Daily Mail dataset or Gigaword, using appropriate loss functions and optimization algorithms.

3. **Context-Length Issues:** Transformer-based models have a maximum sequence length constraint (e.g., 512 tokens for BERT or 1024 tokens for GPT-3), which can pose challenges when summarizing long documents. There are several techniques to address this issue:

   - *Text Segmentation:* Divide the input document into smaller, non-overlapping segments, and generate a summary for each segment separately. The individual summaries can then be combined to form the final summary. This approach might lead to the loss of coherence and important context information between segments.

   - *Sliding Window:* Split the input text into overlapping segments, and apply the summarization model on each segment. Fuse the generated summaries using heuristic methods or trainable fusion techniques. This method can help retain context information but may introduce redundancy in the final summary.

   - *Hierarchical Approaches:* Apply a two-step summarization process, where an initial extractive summarization step selects salient sections of the document, and a subsequent

abstractive summarization step generates a coherent summary. This approach can help reduce the input length while maintaining important context information.

- *Long-range Transformer Architectures:* Utilize transformer models specifically designed to handle long sequences, such as Longformer or BigBird. These models employ sparse attention mechanisms or other techniques to scale the computational complexity and accommodate longer input texts.

4. **Evaluation and Hyperparameter Tuning:** Assess the performance of the summarization model using evaluation metrics like ROUGE (Recall-Oriented Understudy for Gisting Evaluation), which measures the similarity between the generated summary and human-written reference summaries. Fine-tune model hyperparameters, such as learning rate, batch size, and number of training epochs, based on the validation set performance.

5. **Inference and Post-processing:** Apply the fine-tuned model to unseen text data to generate summaries. Post-process the generated summaries to remove inconsistencies, correct grammar errors, or handle other issues related to text quality and readability.

6. **Deployment and Continuous Improvement:** Integrate the trained text summarization model into downstream NLP applications and systems. Continuously monitor and update the model as new data becomes available, or as the underlying language and summary requirements evolve over time. Collect user feedback to refine the model and improve its performance in real-world scenarios.

## 31.3    Query Auto-Completion

Query auto-completion is a feature in search engines that predicts and suggests possible query completions as users type their search queries. It enhances the search experience by saving time, reducing keystrokes, and helping users formulate their queries more effectively. This case study differs from NER or summarization as it focuses on predicting user intent in real-time based on partial query input. In the following sections, we discuss various approaches, including trie-based methods and transformer-based models, for query auto-completion.

1. **Data Collection and Preprocessing:** Collect a dataset of historical search queries, which may include user search logs, query frequencies, and click-through data. Clean and preprocess the data by removing duplicates, correcting misspellings, and filtering out inappropriate queries.

2. **Baseline Model - Trie:** A trie (prefix tree) is a tree-like data structure that stores a dynamic set of strings, where each node represents a character. Tries can be used as a simple yet effective baseline model for query auto-completion. To implement a trie-based auto-completion system:

   - Build a trie structure with the collected search queries, where each path from the root node to a leaf node represents a complete query.

   - As the user types their query, traverse the trie based on the input characters to find the subtree containing all possible completions.

   - Retrieve the top-k most frequent or relevant query completions from the subtree.

3. **Language Models:** N-gram language models can be employed to predict the next word or character in the user's query based on the preceding n-1 words or characters. Train a language model on the historical query dataset and use it to generate query completions based on the user's input. However, n-gram models have limitations, such as sparsity and inability to capture long-range dependencies.

4. **Transformer-based Approaches:** Transformer models, such as GPT or BERT, can be fine-tuned for the query auto-completion task to overcome the limitations of n-gram models and better capture context information. To implement a transformer-based auto-completion system:

   - Fine-tune a pre-trained transformer model, such as GPT, on the historical query dataset. Use a masked language modeling (MLM) or causal language modeling (CLM) objective to predict the next token in the query sequence.

   - Given the user's partial query input, use the fine-tuned model to generate a list of possible next tokens, along with their probabilities.

   - Select the top-k most probable tokens to construct the suggested query completions. This process can be performed iteratively to generate longer completions.

5. **Evaluation and Hyperparameter Tuning:** Evaluate the performance of the query auto-completion model using metrics such as Mean Reciprocal Rank (MRR), precision@k, and keystroke savings ratio. Fine-tune model hyperparameters, such as learning rate, batch size, and number of training epochs, based on the validation set performance.

6. **Inference and Post-processing:** Deploy the trained model to predict query completions in real-time as users type their search queries. Post-process the generated completions to handle special characters, casing, or other formatting requirements.

7. **Deployment and Continuous Improvement:** Integrate the query auto-completion model into the search engine and continuously monitor its performance. Update the model with new query data to account for evolving user behavior, trends, and language patterns. Collect user feedback to refine the model and improve its relevance and effectiveness in real-world scenarios.

## 31.4  Search Result Ranking for Amazon

Designing a search result ranking system for Amazon requires considering the unique characteristics of an e-commerce platform.

### 31.4.1  Approach

For an e-commerce platform like Amazon, a listwise learning-to-rank approach is suitable, as it models the entire ranked list of products for a query. This approach captures the overall user preferences and optimizes the ranking model based on list-level metrics, such as NDCG or MAP, which aligns better with the platform's goal of providing a comprehensive and relevant search experience.

the trade-off between computational complexity and ranking quality should be considered. Listwise methods typically provide better ranking performance because they directly optimize the metrics relevant to search ranking, such as NDCG or MAP, capturing the overall user preferences. Pointwise methods, while more computationally efficient, might not deliver the same level of ranking quality since they don't model the entire ranked list.

To address the computational challenges of listwise methods at Amazon's scale, several techniques can be employed:

**Approximation Techniques:** Employ approximate ranking algorithms that reduce the complexity of listwise methods without significantly sacrificing ranking quality. These techniques can be based on sampling, clustering, or other efficient approximation strategies.

**Parallelization and Distributed Computing:** Leverage parallel and distributed computing frameworks to efficiently process large-scale data and perform model training and inference in parallel.

**Model Compression and Pruning:** Apply model compression techniques, such as quantization, pruning, or knowledge distillation, to reduce the size and complexity of the ranking models without significant performance loss.

**Two-stage Ranking:** Implement a two-stage ranking process, where a more efficient pointwise or pairwise model is first used to filter out irrelevant products, followed by a listwise model to fine-tune the ranking of the remaining products.

By employing these strategies, it is possible to mitigate the computational challenges of listwise learning-to-rank methods and achieve a balance between ranking quality and scalability for a large-scale platform like Amazon.

### 31.4.2   Non-Personalization Features

In addition to traditional information retrieval features (e.g., TF-IDF, BM25) and query-document similarity based on embeddings, the following features are particularly relevant for an e-commerce search ranking system:

- *Query Document Embeddings:* [Fill out]

- *Product Popularity:* The number of sales, views, or wishlist additions for a product can indicate its popularity and relevance to a user's query.

- *Product Ratings and Reviews:* Aggregate user ratings and the number of reviews can provide insights into product quality and user satisfaction.

- *Price and Discounts:* Product price and available discounts can influence user preferences and purchase decisions.

- *Shipping Information:* Shipping speed and cost can impact user preferences, especially for time-sensitive purchases.

- *Category Information:* The product category and subcategory can help improve the ranking by ensuring that products relevant to the user's query context are ranked higher.

### 31.4.3 Personalization Features

- *User Behavior Analysis:* Analyze user browsing and purchase history to identify preferences and interests. Use this information to boost the ranking of products that align with the user's past behavior.

- *Collaborative Filtering:* Implement user-based or item-based collaborative filtering techniques to recommend products based on the preferences of similar users or the similarity between products.

- *Demographic Information:* Leverage user demographic information, such as age, gender, and location, to customize search results based on regional preferences or demographic-specific trends.

- *Contextual Information:* Incorporate contextual information, such as time of day, day of the week, or season, to adjust the ranking based on time-sensitive trends or user behavior patterns.

### 31.4.4 Addressing the Cold-start Problem

One of the challenges faced by e-commerce search ranking systems is the cold-start problem, where new products or users with limited interaction history pose difficulties in estimating product relevance or personalizing rankings. To address this issue, we propose the following strategies:

**Content-based Filtering for New Products:** For new products with limited interaction data, utilize content-based filtering techniques to recommend items based on their similarity to other products in terms of features or descriptions. This can be achieved by computing similarity scores using techniques such as cosine similarity on textual features or deep learning-based similarity measures for images.

**User Onboarding and Preference Elicitation:** For new users with limited interaction history, the system can employ onboarding techniques, such as asking users to rate or express their preferences for a few selected products or categories during account creation. This initial information can be used to bootstrap the personalization process and refine recommendations over time as the user interacts with the platform.

**Fallback to Non-personalized Ranking:** In situations where there is insufficient data to personalize search results, the system can fall back to non-personalized ranking strategies, such as using general popularity, category-specific trends, or seasonal factors to rank products.

**Exploiting Social and Collaborative Information:** When data is sparse for a user, consider incorporating social network information or leveraging collaborative filtering techniques to make use of the preferences and behaviors of similar users or items.

### 31.4.5   Ensuring Result Diversity

Diversity in search results is essential to avoid over-personalization and filter bubbles that can lead to a narrow range of products being recommended, thus reducing user exposure to new or diverse items. Ensuring result diversity can also improve user satisfaction by presenting a broader set of relevant products that cater to different aspects of a user's query or preferences. To achieve result diversity in the search ranking system, we propose the following strategies:

**Maximum Marginal Relevance (MMR):**   MMR is a technique that balances the trade-off between relevance and diversity by iteratively selecting items with the highest marginal relevance, considering both their similarity to the query and their dissimilarity to already-selected items in the result list. By incorporating MMR into the ranking model, we can promote diverse results while maintaining relevance to the user's query.

**Subtopic Coverage:**   Identify various subtopics or aspects related to a user's query and ensure that the search results cover a diverse range of these subtopics. This can be achieved by clustering products based on their features, descriptions, or categories, and selecting representative items from different clusters to include in the search results.

**Diversity-Aware Ranking Metrics:**   Incorporate diversity-aware metrics, such as intent-aware metrics or diversity metrics like $\alpha$-NDCG, into the ranking optimization process. These metrics consider both the relevance and diversity aspects of the results, encouraging the ranking model to produce a diverse set of search results while maintaining high relevance.

**Personalization with Diversity Constraints:**   Integrate diversity constraints into the personalization process to ensure a diverse set of recommendations. This can be achieved by setting limits on the number of similar items or categories recommended to a user, ensuring exposure to a broader range of products that cater to different user preferences or interests.

By employing these strategies, the search ranking system can ensure a diverse set of search results that cater to various user preferences and interests, reducing the risk of over-personalization and promoting a more satisfying user experience.

## 31.5   Latent Diffusion Models for Image Generation

Latent diffusion models (LDMs) are a class of generative models that have gained attention for their ability to produce high-quality images while avoiding some of the limitations associated with other generative models like variational autoencoders (VAEs). In this section, we provide an in-depth overview of the LDM architecture and discuss its advantages, including a comparison with VAEs.

### 31.5.1 Latent Diffusion Model Architecture

The core idea behind LDMs is to model the generation process as a diffusion process, starting from a simple noise distribution and progressively transforming it into a complex data distribution through a series of denoising steps. The diffusion process can be viewed as a continuous-time stochastic process, where the image is gradually corrupted by injecting noise at each time step. The reverse process, known as the denoising process, aims to learn a denoising function that can remove the injected noise and recover the original image.

The denoising function in LDMs is typically parameterized by a deep neural network, which is trained to minimize the difference between the denoised image and the ground truth image at each time step. The architecture of this network can vary, but it often includes components such as convolutional layers, residual connections, and normalization layers to facilitate efficient learning.

Once the denoising function is learned, it can be used to generate new images by sampling from the noise distribution and applying the denoising process in reverse. This is done by iteratively applying the learned denoising function to the noisy samples, progressively refining the generated images until they resemble samples from the target data distribution.

### 31.5.2 Stable Diffusion

A recent development in LDMs is the introduction of the stable diffusion framework, which improves the training and generation dynamics of LDMs by incorporating a stability-inducing regularization term into the training objective. This regularization term encourages the denoising function to be Lipschitz continuous, which helps prevent unstable behaviors during training and generation. As a result, stable diffusion models can generate even higher-quality images while maintaining stable training dynamics.

### 31.5.3 Advantages of Latent Diffusion Models

LDMs offer several advantages over other generative models, such as VAEs:

**Reconstruction Quality:** LDMs can generate high-quality images with sharp details, in contrast to VAEs, which often produce over-smoothed images due to the pixel-wise reconstruction loss term in their objective function.

**Posterior Collapse:** LDMs do not suffer from posterior collapse, a common issue in VAEs where the learned latent variables become less expressive and do not contribute significantly to the generative process.

**Stability and Training Dynamics:** LDMs benefit from a more direct and stable optimization problem, which results in improved training dynamics and stable convergence, leading to better image generation performance.

**Flexibility:** LDMs offer a flexible framework that can be adapted to a variety of tasks, such as conditional image generation and domain adaptation, by conditioning the denoising function on additional information.

In conclusion, latent diffusion models provide a powerful and flexible approach to image generation, with several advantages over alternative generative models like VAEs. The stable diffusion framework further enhances the performance and stability of LDMs, making them a promising choice for various image generation tasks.

# 32 Questions

1. Explain the difference between parametric and non-parametric algorithms. Provide examples of each type.

2. What is the cost function for linear regression? Derive it and explain its significance.

3. Explain the role of the sigmoid function in logistic regression and how it helps with binary classification.

4. Describe the concept of kernel functions in Support Vector Machines. How do they enable SVMs to handle non-linearly separable data?

5. What is the role of entropy in decision tree algorithms? How does it help with feature selection at each node?

6. Explain the intuition behind ensemble methods. How do bagging and boosting improve the performance of base models?

7. Describe the Adam optimization algorithm. How does it combine the ideas of momentum and RMSprop?

8. Compare and contrast various loss functions used in machine learning, such as mean squared error (MSE), mean absolute error (MAE), logistic loss, and hinge loss. Explain their applications and intuitions.

9. Explain the concept of feature selection and its importance in machine learning. Describe some common techniques used for feature selection.

10. Describe the confusion matrix and its elements. How is it used to evaluate the performance of classification models?

11. Explain the ROC curve and the AUC metric. How do they help in evaluating classification models, particularly when dealing with imbalanced datasets?

12. What is cross-validation? Explain the k-fold cross-validation technique and how to choose a good value for k.

13. Compare and contrast various clustering algorithms, such as K-means, hierarchical clustering, and DBSCAN. Explain their advantages and disadvantages.

14. Explain the concept of dimensionality reduction and the role of techniques like PCA and t-SNE.

15. Describe the bias-variance trade-off in machine learning. How do regularization techniques like L1 and L2 help in managing this trade-off?

16. Explain the concept of dropout in neural networks. How does it help with regularization and preventing overfitting?

17. Describe the architecture of feedforward neural networks. Explain the role of activation functions and their intuitions.

18. Explain the backpropagation through time (BPTT) algorithm used in training recurrent neural networks (RNNs). How does it address the vanishing/exploding gradient problem?

19. Compare and contrast RNNs, LSTMs, and GRUs. Explain their architectures and applications.

20. Describe the sequence-to-sequence (seq2seq) model architecture and its applications.

21. Explain the attention mechanism and its role in improving the performance of seq2seq models.

22. Describe the Transformer architecture, including its key components like multi-head self-attention, positional encoding, and position-wise feedforward neural networks. Explain the intuition behind each component.

23. Explain the concept of word embeddings and their significance in natural language processing. Compare and contrast various word embedding techniques like Word2Vec, GloVe, ELMo, BERT, and GPT.

24. Describe the CBOW and Skip-gram models in Word2Vec. How do they learn word embeddings?

25. Explain how natural language processing techniques can be applied to search, particularly in the context of approximate nearest neighbor (ANN) search and semantic search.

26. Describe the triplet loss function and its applications in machine learning, particularly in the context of similarity learning.

27. Explain the two-tower architecture and its potential applications in search.

28. Describe the architecture of autoencoders, including the role of the encoder and decoder components.

29. Explain the role of batch normalization in neural networks. How does it help with regularization and improving the training process?

30. Describe the concept of weight tying in neural network architectures. How does it help with regularization and reducing the number of trainable parameters?

31. What are the challenges in handling missing data, imbalanced data, and data distribution shifts in machine learning? Describe some techniques to address these issues.

32. Explain the concept of tokenization in natural language processing. What are its challenges, and how do different tokenization techniques address them?

33. Describe the one-hot encoding and TF-IDF techniques for representing text data. Compare their strengths and weaknesses.

34. In the context of NLP, what are some challenges faced by word embeddings like Word2Vec and GloVe, and how do contextual embeddings like ELMo, BERT, and GPT address these challenges?

35. Explain the key differences between BERT and GPT architecture and training objectives. How do these differences affect their applicability to various NLP tasks?

36. How can subword tokenization, such as the Byte-Pair Encoding (BPE) or WordPiece, help improve the performance of NLP models and address issues like handling out-of-vocabulary words?

37. How does the Transformer-XL architecture improve upon the original Transformer model, and what are the advantages of the Transformer-XL in dealing with long-range dependencies in sequence data?

38. Explain the concept of knowledge distillation and its application in compressing large models like BERT and GPT into smaller, more efficient models for deployment in resource-constrained environments.

39. What is the role of pre-training and fine-tuning in transfer learning for NLP tasks? How do models like BERT leverage these concepts to achieve state-of-the-art performance on a wide range of NLP tasks?

40. How does the concept of zero-shot and few-shot learning apply to NLP models like GPT-3, and what are the implications of these capabilities for building more efficient and adaptable NLP systems?

41. Explain the differences between seq2seq models with attention and the original Transformer architecture. How do these architectural differences affect the performance and applicability of these models to various NLP tasks?

42. What are some challenges faced by transformer-based models in terms of computational resources, and what are some possible approaches to address these challenges while maintaining high performance?

43. In the context of NLP, explain the differences between unsupervised, semi-supervised, and supervised learning, and provide examples of tasks or models that fall under each category.

44. Explain the concept of self-supervised learning in NLP and its significance in recent advancements in transformer-based models. Provide examples of self-supervised pre-training objectives.

45. Explain the concept of positional encoding in transformers and its importance in handling sequence data. How do absolute and relative positional encodings differ, and what are the advantages of each approach?

46. Compare and contrast the self-attention mechanism in transformers with the attention mechanism in RNN-based seq2seq models. How do their architectures and applications differ?

47. Describe the concept of layer normalization in neural networks. How does it differ from batch normalization, and what are the advantages and disadvantages of each technique?

48. Explain the role of masked language modeling (MLM) in pre-training transformers like BERT. How does MLM differ from traditional language modeling, and what are its advantages and limitations?

49. Explain the concept of negative sampling in the context of learning word embeddings. How does it help improve the efficiency of training models like Word2Vec and GloVe?

50. Describe the architecture of seq2seq models with attention and the Transformer-XL. How do these models address the problem of long-range dependencies in sequence data, and what are the advantages and disadvantages of each approach?

51. Compare and contrast different tokenization strategies, such as character-level, word-level, and subword-level tokenization. What are the advantages and limitations of each strategy in handling different languages and domains?

52. Explain the concept of universal language models, such as multilingual BERT and XLM-R. How do they differ from single-language models, and what are the challenges and benefits of training and using universal language models?

53. Describe the concept of continuous learning in NLP models. How can models like GPT-3 and BERT be adapted for continuous learning, and what are the potential benefits and challenges of this approach?

54. Explain the role of contrastive learning in unsupervised and self-supervised learning tasks. How can contrastive learning be applied to NLP tasks, such as learning sentence representations or document embeddings?

55. Compare and contrast the use of reinforcement learning and supervised learning in NLP tasks. Provide examples of tasks where reinforcement learning has been successfully applied in NLP and discuss the challenges and benefits of using reinforcement learning in these tasks.

56. Describe the concept of prompt engineering in the context of large-scale language models like GPT-3. How does prompt engineering help in adapting these models to various tasks, and what are the challenges and limitations of this approach?

57. Explain the concept of adversarial training in NLP and its potential applications. How can adversarial training be used to improve the robustness and generalization of NLP models?

58. Compare and contrast the architectures and objectives of OpenAI's GPT-3 and Google's T5 (Text-to-Text Transfer Transformer) models. How do these differences affect their performance and applicability to various NLP tasks?