# Week 8

1. **Which of the following statements accurately describes the functionality of a Parameter Server in the context of distributed machine learning?**

A) The Parameter Server handles data preprocessing by scaling features and normalizing values before training.

B) The Parameter Server distributes a model over multiple machines and provides two main operations: Pull (to query parts of the model) and Push (to update parts of the model).

C) The Parameter Server exclusively supports model training using (Stochastic) gradient descent and does not handle other machine learning algorithms.

D) The Parameter Server uses the Collapsed Gibbs Sampling method to update model parameters by aggregating push updates via subtraction.

**Answer:**

**A) Incorrect - The Parameter Server handles data preprocessing by scaling features and normalizing values before training.**
**Explanation:** Data preprocessing is generally performed by data processing pipelines or workers before the model is distributed to the Parameter Server, not by the Parameter Server itself.

**B) Correct - The Parameter Server distributes a model over multiple machines and provides two main operations: Pull (to query parts of the model) and Push (to update parts of the model).**
**Explanation:** A Parameter Server is a distributed system architecture that manages model parameters in a machine learning context. It enables multiple worker nodes to access and update the model parameters efficiently. The Pull operation allows workers to fetch the latest parameters, while the Push operation lets them send back the updates, such as gradients, to the server for model synchronization.

**C) Incorrect - The Parameter Server exclusively supports model training using (Stochastic) gradient descent and does not handle other machine learning algorithms.**

**Explanation:** Parameter Servers are versatile and can support various machine learning algorithms beyond gradient descent.

**D) Incorrect - The Parameter Server uses the Collapsed Gibbs Sampling method to update model parameters.**
**Explanation:** Parameter Servers typically rely on gradient-based optimization methods, not specifically on Gibbs Sampling.

2. **Why is PageRank considered important in the context of information retrieval on the World Wide Web?**

A) It helps to categorize web pages based on the quality of their content, thus improving the accuracy of search results.

B) It provides an objective and mechanical method for rating the importance of web pages based on the link structure of the web, addressing challenges of page relevance amidst a large number of web pages.

C) It ensures that all web pages are indexed equally, regardless of their content or link structure.

D) It automatically filters out irrelevant web pages by analyzing their content and metadata.

**Answer:**

**A) Incorrect - It helps to categorize web pages based on the quality of their content.**
**Explanation:** PageRank does not categorize pages based on content quality; instead, it focuses on their link structure.

**B) Correct - It provides an objective and mechanical method for rating the importance of web pages based on the link structure of the web, addressing challenges of page relevance amidst a large number of web pages.**
**Explanation:** PageRank assesses the importance of web pages by analyzing the structure of links between them. It ranks pages based on the quantity and quality of links they receive, helping search engines deliver relevant results in a vast and complex web landscape.

**C) Incorrect - It ensures that all web pages are indexed equally, regardless of their content or link structure.**
**Explanation:** PageRank prioritizes pages linked by other significant pages, meaning not all pages are treated equally.

**D) Incorrect - It automatically filters out irrelevant web pages by analyzing their content and metadata.**
**Explanation:** While content and metadata are considered in some search algorithms, PageRank primarily operates on link structure.


3.What role does the outerJoinVertices() operator serve in Apache Spark's GraphX?

A) It removes all vertices that are not present in the input RDD.

B) It returns a new graph with only the vertices from the input RDD.

C) It joins the input RDD data with vertices and includes all vertices, whether present in the input RDD or not.

D) It creates a subgraph from the input RDD and vertices.

**Answer:**

**A) Incorrect - It removes all vertices that are not present in the input RDD.**
**Explanation:** This operator does not exclude vertices; it retains them regardless of whether they have matching data.

**B) Incorrect - It does not return a new graph with only the vertices from the input RDD.**
**Explanation:** It includes all vertices from the original graph.

**C) Correct - It joins the input RDD data with vertices and includes all vertices, whether present in the input RDD or not.**
**Explanation:** The outerJoinVertices() operator allows for a join between vertex properties and RDD data, including all vertices in the graph. This is useful when you want to retain vertices that may not have corresponding data in the input RDD.

**D) Incorrect - It does not create a subgraph from the input RDD and vertices.**
**Explanation:** It creates a new graph that retains all vertices from the original graph and merges them with input data.

4. **Which of the following statements accurately describes a key feature of GraphX, a component built on top of Apache Spark Core?**

A) GraphX focuses exclusively on performing machine learning tasks and does not support graph processing.

B) GraphX allows for efficient graph processing and analysis, supports high-level graph measures like triangle counting, and integrates the Pregel API for graph traversal.

C) GraphX is primarily used for data ingestion and preprocessing and does not provide functionalities for graph algorithms or analytics.

D) GraphX provides only basic graph visualization capabilities and does not include algorithms like PageRank or triangle counting.

**Answer:**

**A) Incorrect - GraphX focuses exclusively on performing machine learning tasks and does not support graph processing.**
**Explanation:** GraphX is primarily focused on graph processing and not solely on machine learning tasks.

**B) Correct - GraphX allows for efficient graph processing and analysis, supports high-level graph measures like triangle counting, and integrates the Pregel API for graph traversal.**
**Explanation:** GraphX is designed for processing graphs at scale and includes functionalities for both analytical and algorithmic operations, such as triangle counting and custom graph traversals using the Pregel API.

**C) Incorrect - GraphX is primarily used for data ingestion and preprocessing and does not provide functionalities for graph algorithms or analytics.**
**Explanation:** While it can handle data ingestion, its main purpose is graph processing and analysis.

**D) Incorrect - GraphX provides only basic graph visualization capabilities and does not include algorithms like PageRank or triangle counting.**
**Explanation:** GraphX includes a variety of advanced graph algorithms, not just basic visualization capabilities.

5. **Why are substantial indexes and data reuse important in graph processing?**

A) To create decorative elements within graphs.

C) To add redundancy to graphs for fault tolerance.

D) To increase the file size of graphs for better storage.

**Answer:**

**A) Incorrect - To create decorative elements within graphs.**
**Explanation:** Indexes and data reuse are focused on computational efficiency, not aesthetics.

**B) Correct - To save memory and processing resources by reusing routing tables and edge adjacency information.**
**Explanation:** Efficient graph processing relies on reducing redundancy and optimizing memory usage. Substantial indexes allow for quicker access to data without requiring excessive memory or processing power, leading to more efficient computations.

**C) Incorrect - To add redundancy to graphs for fault tolerance.**
**Explanation:** While redundancy can provide fault tolerance, substantial indexes are more about optimizing resource usage.

**D) Incorrect - To increase the file size of graphs for better storage.**
**Explanation:** Increasing file size is not a goal of using indexes or data reuse; rather, the aim is to reduce unnecessary data duplication.

**6. Which of the following statement(s) accurately describe the functionality of operators in Apache Spark's GraphX?**

A) Join operators add data to graphs and produce new graphs.

B) Structural operators operate on the structure of an input graph and produce a new graph.

C) Property operators modify the vertex or edge properties using a user-defined map function and produce a new graph.

D) All of the above

**Answer:**


**A)** Join operators do indeed add data to graphs and produce new graphs.
**B)** Structural operators operate on the graph's structure and can create new graphs.
**C)** Property operators modify vertex or edge properties using user-defined functions, producing new graphs.

**D) Correct - All of the above statements are true.**

**Explanation:**

Each type of operator serves to extend the capabilities of graph processing in GraphX, enabling various transformations and manipulations of graph data.


**7.Which RDD operator would you use to combine two RDDs by aligning their keys and producing a new RDD with tuples of corresponding values?**

A) union

B) join

C) sample

D) partitionBy


**Answer:**


**A) Incorrect - union.**
**Explanation:** The union operator combines two RDDs without regard to key relationships; it simply appends the elements of both RDDs.

**B) Correct - join.**
**Explanation:** The join operator is specifically designed for combining two RDDs based on their keys, resulting in an RDD of key-value pairs where the keys are aligned.

**C) Incorrect - sample.**
**Explanation:** The sample operator creates a new RDD by taking a random sample from an existing RDD, without combining two RDDs.

**D) Incorrect - partitionBy.**
**Explanation:** The partitionBy operator is used to control how data is partitioned across nodes, not for combining RDDs.

**8. Which of the following is a primary benefit of using graph-based methods in data mining and machine learning?**

A) Reducing the dimensionality of the data

B) Identifying influential people and information, and finding communities

C) Improving the speed of data retrieval from databases

D) Enhancing the accuracy of linear regression models

**Answer:**

**A) Incorrect - Reducing the dimensionality of the data.**
**Explanation:** While graph methods can assist in dimensionality reduction, techniques like PCA are more directly aimed at this task.

**B) Correct - Identifying influential people and information, and finding communities.**
**Explanation:** Graph-based methods excel in analyzing relationships and interactions within data. They help identify key players in a network and can reveal clusters or communities based on connectivity.

**C) Incorrect - Improving the speed of data retrieval from databases.**
**Explanation:** Graph methods are not primarily focused on data retrieval speed; database indexing is more relevant for that purpose.

**D) Incorrect - Enhancing the accuracy of linear regression models.**
**Explanation:** Graph-based methods are not designed to specifically improve the accuracy of linear regression; they focus on relationships in the data.

**9.Which of the following accurately describes a strategy used to optimize graph computations in distributed systems?**

A) Recasting graph systems optimizations as distributed join optimization and incremental materialized maintenance

B) Encoding graphs as simple arrays and using linear algebra operations

C) Expressing graph computation in sequential algorithms and optimizing with single-node processing

D) Implementing graph algorithms using recursive function calls and minimizing parallelism

**Answer:**
**A) Correct - Recasting graph systems optimizations as distributed join optimization and incremental materialized maintenance.**
**Explanation:** Treating graph operations similarly to joins allows for better optimization strategies, which can reduce computational overhead and improve efficiency.

**B) Incorrect - Encoding graphs as simple arrays and using linear algebra operations.**
**Explanation:** While some algorithms may leverage linear algebra, this is not a universal strategy for optimizing graph computations.

**C) Incorrect - Expressing graph computation in sequential algorithms and optimizing with single-node processing.**
**Explanation:** Distributed systems benefit from parallel processing; sequential approaches do not effectively utilize the capabilities of distributed systems.

**D) Incorrect - Implementing graph algorithms using recursive function calls and minimizing parallelism.**
**Explanation:** Recursive calls can be inefficient for large datasets; parallelism is crucial for effective distributed graph processing.

**10. What are the defining traits of a Parameter Server in distributed machine learning?**

S1: Distributes a model over multiple machines.

S2:  It offers two operations:

 (i) Pull for query parts of the model

 (ii) Push for update parts of the model.

A) Only S1 is true.

B) Only S2 is true.

D) NeitherS1 nor S2 are true.

**The correct answer is: C) Both S1 and S2 are true.**

A Parameter Server in distributed machine learning:

- **Distributes a model over multiple machines:** This allows for efficient training of large models on clusters of machines.
- **Offers two operations:**
  - **Pull:** Workers can query parts of the model from the Parameter Server.
  - **Push:** Workers can update parts of the model by pushing their computed gradients to the Parameter Server.

These two operations are essential for the coordination and synchronization of distributed machine learning algorithms.