

PGA Question Bank

SQL and Databases (100 Questions)

1. Q: What is SQL?

 A: SQL (Structured Query Language) is a tool for managing databases. It lets you retrieve, insert, update, or delete data in tables. For example, SELECT * FROM users fetches all user data. It's essential for data analysis and reporting.

2. Q: What's the difference between INNER JOIN and LEFT JOIN?

 A: INNER JOIN returns only rows with matching values in both tables. LEFT JOIN includes all rows from the left table, with NULLs for non-matching rows from the right. Use INNER for strict matches, LEFT for complete left-table data.

3. Q: What is a primary key?

 A: A primary key is a unique column (like an ID) that identifies each row in a table. It can't be NULL and ensures no duplicates. Only one primary key exists per table for data integrity.

4. Q: What are aggregate functions in SQL?

 A: Aggregate functions calculate a single value from multiple rows. Examples include SUM(), AVG(), COUNT(), MAX(), and MIN(). For instance, SELECT COUNT(*) FROM orders counts total orders.

5. Q: What is a foreign key?

 A: A foreign key links a column in one table to a primary key in another, ensuring valid relationships. For example, order.customer_id links to customers.id. It maintains referential integrity.

6. Q: What's the GROUP BY clause used for?

 A: GROUP BY groups rows with the same values for summary calculations. For example, SELECT department, SUM(salary) FROM employees GROUP BY department totals salaries by department. It's used with aggregates.

7. Q: What's the difference between WHERE and HAVING?

 A: WHERE filters individual rows before grouping, like WHERE age > 20. HAVING filters groups after GROUP BY, like HAVING COUNT(*) > 5. Use WHERE for row-level, HAVING for group-level.

8. Q: What is a subquery?

 A: A subquery is a query inside another query, providing data for the outer

query. For example, `SELECT name FROM employees WHERE salary > (SELECT AVG(salary) FROM employees)` finds above-average earners.

9. Q: What is normalization?

 **A:** Normalization organizes a database to reduce redundancy and improve efficiency. It splits data into related tables using keys. For example, separating customers and orders avoids duplicate customer details.

10. Q: What is a view in SQL?

 **A:** A view is a virtual table created by a query, like `CREATE VIEW top_sales AS SELECT * FROM orders WHERE amount > 1000`. It simplifies queries and restricts data access for security.

11. Q: What is a database index?

 **A:** An index speeds up data retrieval by creating a lookup for columns. For example, indexing `customer_id` makes searches faster. However, it slows down inserts and updates.

12. Q: What's the difference between DELETE and TRUNCATE?

 **A:** `DELETE` removes specific rows based on a condition, like `DELETE FROM users WHERE age < 18`. `TRUNCATE` deletes all rows without conditions, resetting the table. `DELETE` is slower but reversible.

13. Q: What is a stored procedure?

 **A:** A stored procedure is a pre-written SQL script saved in the database, like a function. For example, `CALL calculate_bonus()` can automate salary updates. It improves reusability and security.

14. Q: What is a transaction in SQL?

 **A:** A transaction is a group of SQL operations executed as one unit. For example, transferring money requires debiting one account and crediting another. Use `COMMIT` to save or `ROLLBACK` to undo.

15. Q: What are ACID properties?

 **A:** ACID (Atomicity, Consistency, Isolation, Durability) ensures reliable database transactions. Atomicity treats operations as one unit, Consistency maintains data rules, Isolation prevents interference, and Durability saves changes permanently.

16. Q: What's the difference between UNION and UNION ALL?

 **A:** `UNION` combines results from two queries, removing duplicates. `UNION ALL` includes all rows, keeping duplicates. Use `UNION ALL` for faster performance when duplicates are okay.

17. Q: What is a self-join?

 **A:** A self-join joins a table with itself to compare rows. For example, `SELECT a.name, b.name FROM employees a, employees b WHERE a.manager_id = b.id` lists employees and their managers.

18. Q: What is the CASE statement in SQL?

 **A:** The CASE statement adds conditional logic to queries. For example, `SELECT name, CASE WHEN salary > 50000 THEN 'High' ELSE 'Low' END AS status FROM employees` categorizes salaries.

19. Q: What's the difference between a table and a view?

 **A:** A table stores actual data, while a view is a virtual table based on a query. Views don't store data but simplify access. For example, a view can show only high-value orders.

20. Q: What is denormalization?

 **A:** Denormalization combines tables to improve read performance, accepting some redundancy. For example, storing customer names in an orders table avoids joins. It's used in data warehouses for speed.

21. Q: What is a composite key?

 **A:** A composite key uses multiple columns to uniquely identify rows, like order_id and product_id in an order details table. It's used when no single column is unique.

22. Q: What's the HAVING clause used for?

 **A:** HAVING filters grouped data after GROUP BY. For example, `SELECT department, COUNT(*) FROM employees GROUP BY department HAVING COUNT(*) > 10` shows departments with over 10 employees.

23. Q: What is a clustered index?

 **A:** A clustered index determines the physical order of data in a table, with only one per table. For example, indexing id sorts the table by id, speeding up queries but slowing updates.

24. Q: What is a non-clustered index?

 **A:** A non-clustered index is a separate structure pointing to table data, allowing multiple per table. It speeds up searches, like on email, without affecting the table's physical order.

25. Q: What is the DISTINCT keyword?

 **A:** DISTINCT removes duplicate rows from query results. For example, `SELECT DISTINCT city FROM customers` lists unique cities. It's useful for summarizing data but can slow queries.

26. Q: What is a schema in a database?

 **A:** A schema is a logical container for database objects like tables and views. For example, a sales schema might group all sales-related tables. It organizes and secures data access.

27. Q: What's the difference between DROP and TRUNCATE?

 **A:** DROP deletes an entire table or database, including its structure. TRUNCATE removes all rows but keeps the table structure. DROP is permanent, while TRUNCATE is faster for clearing data.

28. Q: What is a trigger in SQL?

 **A:** A trigger is a special stored procedure that runs automatically when specific events occur, like INSERT or UPDATE. For example, a trigger can log changes to an audit table.

29. Q: What is the LIKE operator used for?

 **A:** The LIKE operator searches for patterns in text. For example, SELECT name FROM customers WHERE name LIKE 'A%' finds names starting with 'A'. Use % for wildcards.

30. Q: What is a temporary table?

 **A:** A temporary table stores data for a session or transaction, deleted automatically when done. For example, CREATE TEMPORARY TABLE temp_users holds intermediate results. It's useful for complex queries.

31. Q: What is the COALESCE function?

 **A:** COALESCE returns the first non-NULL value from a list. For example, SELECT COALESCE(phone, email, 'Unknown') FROM contacts ensures a contact method is shown, avoiding NULLs.

32. Q: What's the difference between CHAR and VARCHAR?

 **A:** CHAR stores fixed-length strings, padding with spaces, while VARCHAR stores variable-length strings. Use CHAR for consistent lengths (like codes) and VARCHAR for varying text (like names).

33. Q: What is a database transaction log?

 **A:** A transaction log records all database changes, enabling recovery after crashes. It tracks operations like inserts or updates, ensuring data consistency and supporting rollbacks.

34. Q: What is the IN operator used for?

 **A:** The IN operator checks if a value matches a list. For example, SELECT name FROM employees WHERE dept IN ('HR', 'IT') finds employees in HR or IT departments.

35. Q: What is a correlated subquery?

 **A:** A correlated subquery depends on the outer query's rows, running for each row. For example, `SELECT name FROM employees e WHERE salary > (SELECT AVG(salary) FROM employees WHERE dept = e.dept)` compares salaries within departments.

36. Q: What is the BETWEEN operator?

 **A:** BETWEEN selects values within a range, inclusive. For example, `SELECT * FROM orders WHERE amount BETWEEN 100 AND 500` finds orders with amounts from 100 to 500.

37. Q: What is a materialized view?

 **A:** A materialized view stores query results physically, unlike regular views. It's refreshed periodically, like `CREATE MATERIALIZED VIEW sales_summary AS SELECT`. It's used for performance in data warehouses.

38. Q: What is the MERGE statement?

 **A:** The MERGE statement (or UPSERT) combines INSERT, UPDATE, and DELETE based on conditions. For example, it updates existing records or inserts new ones, streamlining data synchronization.

39. Q: What is a database constraint?

 **A:** A constraint enforces rules, like NOT NULL, UNIQUE, or CHECK. For example, `CHECK (age > 18)` ensures only adults are added. Constraints maintain data quality.

40. Q: What is the LIMIT clause?

 **A:** LIMIT restricts the number of rows returned, like `SELECT * FROM products LIMIT 5` for the top 5 products. It's useful for pagination or sampling data.

41. Q: What is a window function?

 **A:** Window functions perform calculations across a set of rows without grouping, like `ROW_NUMBER() OVER (PARTITION BY dept ORDER BY salary)`. They're used for rankings or running totals.

42. Q: What is the RANK() function?

 **A:** RANK() assigns rankings to rows within a partition, with ties getting the same rank. For example, `SELECT name, RANK() OVER (ORDER BY score)` ranks students by scores.

43. Q: What's the difference between RANK() and DENSE_RANK()?

 **A:** RANK() skips numbers after ties (e.g., 1, 1, 3), while DENSE_RANK() doesn't (e.g., 1, 1, 2). Use DENSE_RANK() for continuous rankings, like in leaderboards.

44. Q: What is a CTE (Common Table Expression)?

 A: A CTE is a temporary result set defined with WITH, like WITH sales AS (SELECT ...) SELECT * FROM sales. It improves query readability and supports recursion.

45. Q: What is SQL injection?

 A: SQL injection is a security flaw where attackers insert malicious SQL code into inputs, like user_id = '1; DROP TABLE users;'. Prevent it with parameterized queries or input validation.

46. Q: What is the NULL value in SQL?

 A: NULL represents missing or unknown data, not zero or empty. For example, SELECT * FROM employees WHERE email IS NULL finds employees without emails. It requires special handling.

47. Q: What is the ORDER BY clause?

 A: ORDER BY sorts query results, like SELECT name FROM customers ORDER BY name ASC. Use ASC for ascending or DESC for descending order to organize output.

48. Q: What is a cross join?

 A: A cross join combines every row from two tables, creating a Cartesian product. For example, SELECT * FROM colors CROSS JOIN sizes pairs all colors with all sizes.

49. Q: What is the EXISTS operator?

 A: EXISTS checks if a subquery returns any rows. For example, SELECT name FROM customers WHERE EXISTS (SELECT 1 FROM orders WHERE customer_id = customers.id) finds customers with orders.

50. Q: What is a database schema?

 A: A schema organizes database objects (tables, views) into logical groups. For example, a hr schema might hold employee-related tables, simplifying management and access control.

51. Q: What is the CAST function?

 A: CAST converts data types, like SELECT CAST('123' AS INTEGER) to turn a string into a number. It's useful for ensuring compatibility in calculations or comparisons.

52. Q: What is a partitioned table?

 A: A partitioned table splits large data into smaller, manageable chunks based on a key, like date ranges. It improves query performance and maintenance for big datasets.

53. Q: What is the WITH clause?

 **A:** The WITH clause defines a CTE for cleaner queries. For example, WITH temp AS (SELECT ...) SELECT * FROM temp organizes complex logic and supports reuse.

54. Q: What is a full-text search?

 **A:** Full-text search finds text patterns in columns, like SELECT * FROM articles WHERE MATCH(content) AGAINST('data'). It's faster than LIKE for searching large text data.

55. Q: What is a database trigger used for?

 **A:** Triggers automatically execute actions on events like INSERT. For example, a trigger can update a stock table when a sale is recorded, ensuring data consistency.

56. Q: What is the NOW() function?

 **A:** NOW() returns the current date and time, like SELECT NOW(). It's used for logging timestamps, such as recording when a record was created.

57. Q: What is the IFNULL function?

 **A:** IFNULL replaces NULL with a specified value, like SELECT IFNULL(salary, 0) FROM employees. It's similar to COALESCE but simpler for single-column checks.

58. Q: What is a recursive CTE?

 **A:** A recursive CTE builds hierarchical or iterative results, like WITH RECURSIVE org AS (...) to list employee hierarchies. It's used for tree-like data structures.

59. Q: What is the TOP clause?

 **A:** TOP limits rows returned, like SELECT TOP 10 * FROM sales for the top 10 sales. It's used in SQL Server, similar to LIMIT in MySQL.

60. Q: What is a database role?

 **A:** A database role groups permissions for users, like read_only for viewing data. For example, GRANT SELECT TO read_only simplifies access control for multiple users.

61. Q: What is the UPDATE statement?

 **A:** UPDATE modifies existing rows, like UPDATE employees SET salary = salary + 1000 WHERE dept = 'IT'. It's used to change data based on conditions.

62. Q: What is a deadlock in databases?

 **A:** A deadlock occurs when two transactions lock resources each other

needs, causing a stall. Databases resolve it by terminating one transaction. Proper design minimizes deadlocks.

63. Q: What is the ALTER TABLE statement?

 **A:** ALTER TABLE modifies a table's structure, like ALTER TABLE employees ADD email VARCHAR(50) to add a column. It's used for schema changes.

64. Q: What is a check constraint?

 **A:** A check constraint enforces a condition, like CHECK (age >= 18) on an employees table. It ensures only valid data is inserted or updated.

65. Q: What is the CONCAT function?

 **A:** CONCAT combines strings, like SELECT CONCAT(first_name, ' ', last_name) AS full_name FROM users. It's used to format text output cleanly.

66. Q: What is a database backup?

 **A:** A database backup copies data to recover it after loss or corruption. For example, BACKUP DATABASE sales TO DISK = 'sales.bak' saves the database state.

67. Q: What is the COUNT function?

 **A:** COUNT returns the number of rows, like SELECT COUNT(*) FROM orders WHERE status = 'shipped'. Use it to summarize data, like tracking order volume.

68. Q: What is a pivot operation in SQL?

 **A:** Pivoting transforms rows into columns, like turning sales by month into columns per month. For example, PIVOT in SQL Server summarizes data for reporting.

69. Q: What is the IS NULL operator?

 **A:** IS NULL checks for NULL values, like SELECT * FROM employees WHERE email IS NULL. It's used to find missing data in columns.

70. Q: What is a sequence in SQL?

 **A:** A sequence generates unique numbers, like CREATE SEQUENCE order_id START WITH 1. It's used for auto-incrementing IDs, similar to IDENTITY but more flexible.

71. Q: What is the MIN function?

 **A:** MIN finds the smallest value in a column, like SELECT MIN(price) FROM products. It's an aggregate function used to identify lowest values.

72. Q: What is the MAX function?

 **A:** MAX finds the largest value, like SELECT MAX(salary) FROM employees. It's used to identify highest values in data analysis tasks.

73. Q: What is a database synonym?

 **A:** A synonym is an alias for a database object, like CREATE SYNONYM emp FOR hr.employees. It simplifies queries and hides object locations.

74. Q: What is the ROUND function?

 **A:** ROUND rounds numbers to specified decimals, like SELECT ROUND(123.456, 2) to get 123.46. It's used for formatting numerical outputs.

75. Q: What is the DATEDIFF function?

 **A:** DATEDIFF calculates the difference between dates, like SELECT DATEDIFF(day, order_date, ship_date) FROM orders. It's used for time-based analysis.

76. Q: What is a database cursor?

 **A:** A cursor processes query results row by row, like in stored procedures. For example, it can loop through employee records to apply bonuses. Use sparingly for performance.

77. Q: What is the SUM function?

 **A:** SUM adds up values in a column, like SELECT SUM(amount) FROM sales. It's an aggregate function for totaling numerical data, like revenue.

78. Q: What is a schema-level trigger?

 **A:** A schema-level trigger runs on events like user logins, not tied to a table. For example, it can log all CREATE TABLE actions for auditing.

79. Q: What is the AVG function?

 **A:** AVG calculates the average of a column, like SELECT AVG(score) FROM tests. It's used to summarize data, ignoring NULL values.

80. Q: What is a database link?

 **A:** A database link connects to another database, like CREATE DATABASE LINK remote_db It allows queries across databases, such as fetching remote sales data.

81. Q: What is the ROW_NUMBER() function?

 **A:** ROW_NUMBER() assigns unique numbers to rows in a result, like SELECT ROW_NUMBER() OVER (ORDER BY salary) AS rank FROM employees. It's used for pagination or ranking.

82. Q: What is a partitioned index?

 **A:** A partitioned index splits an index across table partitions, like by date ranges. It improves query performance for large, partitioned tables.

83. Q: What is the SUBSTRING function?

 **A:** SUBSTRING extracts part of a string, like `SELECT SUBSTRING(name, 1, 3)` FROM users to get the first three characters. It's used for text manipulation.

84. Q: What is a database audit?

 **A:** A database audit tracks actions like logins or data changes for security. For example, enabling audit logs ensures compliance with regulations like GDPR.

85. Q: What is the UNIQUE constraint?

 **A:** The UNIQUE constraint ensures no duplicate values in a column, like email in a users table. Unlike a primary key, it allows NULLs.

86. Q: What is a database snapshot?

 **A:** A snapshot captures a database's state at a moment, like `CREATE DATABASE snapshot_sales AS` It's used for reporting or testing without affecting the main database.

87. Q: What is the CURRENT_DATE function?

 **A:** CURRENT_DATE returns today's date, like `SELECT CURRENT_DATE`. It's used for date-based queries, such as filtering recent orders.

88. Q: What is a global temporary table?

 **A:** A global temporary table is shared across sessions but holds session-specific data. For example, it's used for temporary calculations, deleted when the session ends.

89. Q: What is the LENGTH function?

 **A:** LENGTH returns the number of characters in a string, like `SELECT LENGTH(name) FROM users`. It's used for validating or analyzing text data.

90. Q: What is a database role-based access control?

 **A:** Role-based access control assigns permissions to roles, not users. For example, `GRANT SELECT ON sales TO analyst_role` simplifies managing user access.

91. Q: What is the TRUNCATE TABLE statement?

 **A:** TRUNCATE TABLE removes all rows but keeps the table structure, like `TRUNCATE TABLE logs`. It's faster than `DELETE` and resets auto-increment counters.

92. Q: What is a database synonym used for?

 **A:** A synonym simplifies object references, like `CREATE SYNONYM sales FOR db1.sales_table`. It hides complexity and improves query readability across schemas.

93. Q: What is the EXTRACT function?

 A: EXTRACT pulls parts of a date, like SELECT EXTRACT(YEAR FROM order_date) FROM orders. It's used for date-based analysis, like yearly sales.

94. Q: What is a database view used for?

 A: A view simplifies complex queries or restricts data access, like CREATE VIEW active_users AS SELECT * FROM users WHERE active = 1. It's virtual and doesn't store data.

95. Q: What is the UPPER function?

 A: UPPER converts text to uppercase, like SELECT UPPER(name) FROM users. It's used for standardizing or comparing text data case-insensitively.

96. Q: What is a database partition?

 A: A partition divides a table into smaller pieces, like by year. It improves query speed and maintenance for large datasets, like PARTITION BY RANGE (order_date).

97. Q: What is the LOWER function?

 A: LOWER converts text to lowercase, like SELECT LOWER(email) FROM users. It's used for case-insensitive searches or data standardization.

98. Q: What is a database audit trail?

 A: An audit trail logs database actions, like updates or logins, for security and compliance. For example, it tracks who modified a sales table.

99. Q: What is the NVL function?

 A: NVL replaces NULL with a value, like SELECT NVL(salary, 0) FROM employees (Oracle). It's similar to IFNULL or COALESCE for handling missing data.

100. Q: What is a database cluster?

 A: A database cluster is a group of servers working together for high availability or load balancing. For example, it ensures uptime if one server fails.

Data Visualization (80 Questions)

1. Q: What is Tableau used for?

 A: Tableau creates interactive visualizations and dashboards without coding. It connects to data sources like Excel or SQL databases, turning raw data into charts for decision-making. It's popular in business intelligence.

2. Q: How do you create a calculated field in Tableau?

 A: Right-click in the Data pane, select "Create Calculated Field," and write a

formula, like `SUM(Sales) / COUNT(Orders)`. Name it and save. Use it for custom metrics in visuals.

3. Q: What's the difference between a dashboard and a worksheet in Tableau?

 **A:** A worksheet holds a single visualization, like a bar chart. A dashboard combines multiple worksheets, filters, and text into one interactive view for a complete data story.

4. Q: What is Power BI?

 **A:** Power BI is Microsoft's tool for data visualization and analytics. It builds interactive reports and dashboards, connecting to Excel, databases, or cloud data. It's user-friendly for business users.

5. Q: What are slicers in Power BI?

 **A:** Slicers are interactive filters letting users select data, like choosing a year to view sales. Add them via the Visualizations pane, enhancing report interactivity for users.

6. Q: How do you create a pivot table in Excel?

 **A:** Select data, go to Insert > PivotTable, and choose a location. Drag fields to Rows, Columns, or Values to summarize, like sales by region. It's great for quick analysis.

7. Q: What's the difference between a bar chart and a histogram?

 **A:** Bar charts compare categories with separate bars, like sales by product. Histograms show numerical data distribution, like test scores, with adjacent bars for ranges.

8. Q: What is a live connection in Tableau?

 **A:** A live connection queries the data source in real-time, reflecting updates instantly. It's ideal for dynamic data but slower for large datasets compared to extracts.

9. Q: How do you add filters in Power BI?

 **A:** Drag a field to the Filters pane and set conditions, like selecting specific regions. Filters apply to visuals, pages, or reports, focusing data for analysis.

10. Q: What is a KPI visual in Power BI?

 **A:** A KPI visual shows a metric, like sales vs. target, with a trend line. It uses colors to indicate performance, making it easy to monitor goals in dashboards.

11. Q: What is a calculated column in Power BI?

 **A:** A calculated column adds a new column using DAX, like `TotalCost = [Price] * [Quantity]`. It's computed row by row and used in visuals or calculations.

12. Q: What is a measure in Power BI?

 **A:** A measure is a dynamic calculation, like $\text{TotalSales} = \text{SUM}(\text{Sales}[\text{Amount}])$, created with DAX. Unlike calculated columns, measures aggregate data for visuals, like summing sales.

13. Q: What is a parameter in Tableau?

 **A:** A parameter is a user-defined input, like a sales threshold. Create it via the Data pane, then use it in calculations or filters, like $\text{Sales} > [\text{Threshold}]$.

14. Q: What's the difference between a line chart and an area chart?

 **A:** A line chart shows data points connected by lines, like stock prices over time. An area chart fills the space below, emphasizing volume, like cumulative sales.

15. Q: How do you create a heat map in Tableau?

 **A:** Drag measures (like sales) to Color and Size, and dimensions (like regions) to Rows/ [Ideal Response]: In Tableau, drag a measure (e.g., Sales) to Color and Size, and a dimension (e.g., Region) to Rows or Columns. Use a square mark type to create a heat map showing intensity by color and size.

16. Q: What is DAX in Power BI?

 **A:** DAX (Data Analysis Expressions) is a formula language for creating calculations in Power BI. It's used for measures and calculated columns, like $\text{Profit} = \text{SUM}(\text{Sales}[\text{Revenue}]) - \text{SUM}(\text{Sales}[\text{Cost}])$, enhancing data analysis.

17. Q: What is a data source in Tableau?

 **A:** A data source in Tableau is the database, file, or server (like Excel or SQL) connected for analysis. It defines the data structure and relationships used in visualizations.

18. Q: How do you publish a dashboard in Tableau?

 **A:** Save your dashboard, then select Server > Publish Workbook. Log into Tableau Server or Online, choose a project, and publish. Users can now view it online.

19. Q: What is a treemap in visualization?

 **A:** A treemap displays hierarchical data as nested rectangles, with size and color representing values. For example, it shows sales by category, where larger rectangles mean higher sales.

20. Q: What is a drill-down feature in Power BI?

 **A:** Drill-down lets users click a visual to see detailed data, like clicking a region to view city-level sales. Enable it via the visual's settings for interactive exploration.

21. Q: How do you create a dual-axis chart in Tableau?

 **A:** Drag two measures to Rows, right-click the second measure's axis, and select "Dual Axis." Synchronize axes and adjust mark types for combined visuals, like sales and profit.

22. Q: What is a filter shelf in Tableau?

 **A:** The filter shelf in Tableau limits data shown in a visualization. Drag fields (like Date or Category) to the shelf and set conditions, like "Year = 2023."

23. Q: What is a card visual in Power BI?

 **A:** A card visual displays a single value, like total sales, prominently. Add it via the Visualizations pane, select a measure, and format for clarity in dashboards.

24. Q: How do you blend data in Tableau?

 **A:** Data blending combines multiple data sources in Tableau without joining. Link fields (like CustomerID) across sources, and Tableau merges data for analysis, useful for unrelated datasets.

25. Q: What is a waterfall chart?

 **A:** A waterfall chart shows cumulative changes, like revenue from start to end, with bars for increases and decreases. It's ideal for financial analysis in Power BI or Excel.

26. Q: How do you create a pie chart in Excel?

 **A:** Select data, go to Insert > Pie Chart, and choose a style. Customize labels or explode slices for clarity. Use pie charts for showing proportions, like market share.

27. Q: What is a sparkline in Excel?

 **A:** A sparkline is a tiny chart in a cell showing trends, like monthly sales. Select data, go to Insert > Sparkline, and choose Line or Column type.

28. Q: What is a data model in Power BI?

 **A:** A data model defines relationships between tables in Power BI, like linking Orders to Customers via CustomerID. It enables accurate calculations and visuals across datasets.

29. Q: How do you export a Tableau dashboard?

 **A:** Go to File > Export As, choose PDF, Image, or PowerPoint, and save. Alternatively, publish to Tableau Server for online sharing with stakeholders.

30. Q: What is a gauge chart in Power BI?

 **A:** A gauge chart shows a value within a range, like sales vs. target, with a needle. Add it via custom visuals to highlight performance metrics visually.

31. Q: What's the difference between a live connection and an extract in Tableau?

 **A:** A live connection queries the data source in real-time, while an extract is a static data snapshot. Extracts are faster for large data but need refreshing.

32. Q: How do you create a bar chart in Power BI?

 **A:** Select the Bar Chart visual, drag a category (like Product) to Axis, and a measure (like Sales) to Value. Customize colors and labels for clarity.

33. Q: What is a calculated table in Power BI?

 **A:** A calculated table is created using DAX, like NewTable = FILTER(Sales, Sales[Year] = 2023). It generates a new table for specific analysis within the data model.

34. Q: How do you add a title to a Tableau dashboard?

 **A:** In the dashboard, drag a Text object from the Objects pane, type your title, and format it (font, size, color). It clarifies the dashboard's purpose.

35. Q: What is a donut chart?

 **A:** A donut chart is a pie chart with a hollow center, showing proportions, like sales by region. Create it in Power BI or Excel for visual appeal.

36. Q: How do you create a combo chart in Excel?

 **A:** Select data, go to Insert > Combo Chart, and choose types (like Column + Line). Assign measures to primary/secondary axes, like sales and profit trends.

37. Q: What is a tooltip in Tableau?

 **A:** A tooltip shows details when hovering over a visual, like sales amount. Customize it via Worksheet > Tooltip to include fields or formatting for clarity.

38. Q: How do you schedule a data refresh in Power BI?

 **A:** In Power BI Service, go to the dataset's Settings, select Refresh, and set a schedule (e.g., daily). Ensure credentials are updated for automatic data updates.

39. Q: What is a funnel chart?

 **A:** A funnel chart shows stages in a process, like sales pipeline stages, with decreasing sizes. Create it in Power BI to analyze conversion rates.

40. Q: How do you create a hierarchy in Tableau?

 **A:** Drag a field (like State) onto another (like Country) in the Data pane to create a hierarchy. Use it for drill-down analysis, like Country > State.

41. Q: What is conditional formatting in Excel?

 **A:** Conditional formatting highlights cells based on rules, like coloring sales

above 1000 green. Go to Home > Conditional Formatting and set criteria for visual insights.

42. Q: How do you create a map visual in Power BI?

 **A:** Select the Map visual, drag a location field (like City) to Location, and a measure (like Sales) to Size or Color. It plots data geographically.

43. Q: What is a LOD calculation in Tableau?

 **A:** Level of Detail (LOD) calculations control aggregation, like {FIXED [Region]: SUM(Sales)} for region-level totals. Use FIXED, INCLUDE, or EXCLUDE for precise analysis.

44. Q: How do you create a table visual in Power BI?

 **A:** Select the Table visual, drag fields to Values, like Customer and Sales. It displays raw data in rows, ideal for detailed reports or summaries.

45. Q: What is a box plot in Tableau?

 **A:** A box plot shows data distribution (median, quartiles, outliers), like test scores. Drag a measure to Rows, set Mark to Box-and-Whisker, and add dimensions.

46. Q: How do you import data into Power BI?

 **A:** Click Get Data, choose a source (Excel, SQL, Web), and connect. Follow prompts to load or transform data, then use it for visuals or models.

47. Q: What is a reference line in Tableau?

 **A:** A reference line adds a benchmark, like an average, to a chart. Go to Analytics > Drag Reference Line, set a value (e.g., AVG(Sales)), and format.

48. Q: How do you create a matrix visual in Power BI?

 **A:** Select the Matrix visual, drag fields to Rows, Columns, and Values, like Region, Year, and Sales. It's like a pivot table for cross-tabulated data.

49. Q: What is a bullet chart?

 **A:** A bullet chart compares a value to a target, like sales vs. goal, with ranges for performance. Use it in Power BI for concise KPI tracking.

50. Q: How do you create a filter action in Tableau?

 **A:** In a dashboard, go to Worksheet > Actions > Add Action > Filter. Set source and target sheets, so clicking one visual filters another, enhancing interactivity.

51. Q: What is Power Query in Power BI?

 **A:** Power Query transforms data before loading, like cleaning or merging tables. Access it via Transform Data to filter, pivot, or combine datasets easily.

52. Q: How do you create a Gantt chart in Tableau?

 **A:** Drag start and end dates to Columns, use a calculated field for duration, and set Mark to Gantt Bar. It shows task timelines, like project schedules.

53. Q: What is a scatter plot used for?

 **A:** A scatter plot shows relationships between two variables, like sales vs. profit. Create it in Tableau or Power BI to identify correlations or outliers.

54. Q: How do you share a Power BI report?

 **A:** Publish to Power BI Service via File > Publish, then share via links, dashboards, or apps. Set permissions to control access for stakeholders.

55. Q: What is a Pareto chart?

 **A:** A Pareto chart combines bars and a line to show cumulative impact, like 80/20 rule for defects. Create it in Excel or Power BI for prioritization.

56. Q: How do you create a story in Tableau?

 **A:** Create worksheets, then go to Story > New Story. Add sheets as points, annotate, and navigate to present a data-driven narrative to stakeholders.

57. Q: What is a stacked bar chart?

 **A:** A stacked bar chart splits bars into segments, like sales by product within regions. Create it in Tableau or Power BI to show part-to-whole relationships.

58. Q: How do you create a custom visual in Power BI?

 **A:** Use the Power BI Developer Tools or import from AppSource. Code visuals in TypeScript, test in Power BI, and publish for unique dashboard needs.

59. Q: What is a trend line in Tableau?

 **A:** A trend line shows data patterns, like linear or exponential. Drag Trend Line from Analytics, apply to a scatter plot, and view statistical details.

60. Q: How do you create a KPI card in Excel?

 **A:** Use a cell with a formula (like SUM(Sales)), apply conditional formatting for color-coding, and add a sparkline nearby. It mimics Power BI KPI visuals.

61. Q: What is a bubble chart?

 **A:** A bubble chart uses bubbles to show three variables: x-axis, y-axis, and size. For example, sales (x), profit (y), and volume (size) in Tableau.

62. Q: How do you create a dashboard layout in Power BI?

 **A:** Add visuals to the canvas, resize and align using the View > Snap to Grid option. Group related visuals and add text for a clean, intuitive layout.

63. Q: What is a set in Tableau?

 **A:** A set groups data based on conditions, like top 10 customers by sales. Create via Data pane > Create Set, then use in filters or visuals.

64. Q: How do you create a gauge in Excel?

 **A:** Use a donut chart for the gauge background, overlay a pie chart for the needle, and calculate angles with formulas. It shows metrics like progress.

65. Q: What is a word cloud in Power BI?

 **A:** A word cloud visualizes text frequency, with larger words for higher counts. Import from AppSource, add a text field, and use for sentiment analysis.

66. Q: How do you create a calculated field in Excel?

 **A:** In a pivot table, go to Analyze > Fields, Items & Sets > Calculated Field. Add a formula, like Profit = Sales - Cost, for custom metrics.

67. Q: What is a highlight table in Tableau?

 **A:** A highlight table uses color to show values, like sales by region. Drag a measure to Color and dimensions to Rows/Columns for a visual table.

68. Q: How do you create a bookmark in Power BI?

 **A:** Go to View > Bookmarks, capture the current report state, and name it. Use bookmarks to save filter states or create presentation flows.

69. Q: What is a motion chart?

 **A:** A motion chart animates data over time, like sales by region. Create in Tableau with Pages shelf (e.g., Year) to show trends dynamically.

70. Q: How do you create a radar chart in Excel?

 **A:** Select data, go to Insert > Other Charts > Radar. It compares multiple variables, like skills ratings, with lines connecting data points.

71. Q: What is a calculated item in Excel?

 **A:** In a pivot table, go to Analyze > Fields, Items & Sets > Calculated Item. Add a formula, like Q1 = Jan + Feb + Mar, for custom groupings.

72. Q: How do you create a dashboard action in Tableau?

 **A:** Go to Dashboard > Actions > Add Action (Filter, Highlight, or URL). Link visuals, like clicking a bar to filter a related chart, for interactivity.

73. Q: What is a decomposition tree in Power BI?

 **A:** A decomposition tree breaks down a metric, like sales, by dimensions (e.g., Region, Product). Add via Visualizations and drill down interactively.

74. Q: How do you create a sunburst chart in Excel?

 A: Use a hierarchy dataset, go to Insert > Sunburst. It shows nested categories, like sales by Region > Category, with concentric rings.

75. Q: What is a parameter action in Tableau?

 A: A parameter action updates a parameter when clicking a visual, like setting a threshold. Go to Worksheet > Actions > Change Parameter for dynamic controls.

76. Q: How do you create a thermometer chart in Excel?

 A: Use a stacked column chart with two series (goal and actual), hide the goal series, and format the actual as a thermometer shape for progress tracking.

77. Q: What is a chiclet slicer in Power BI?

 A: A chiclet slicer displays filter options as buttons, like product categories. Import from AppSource and add a field for a visually appealing filter.

78. Q: How do you create a reference band in Tableau?

 A: Drag Reference Band from Analytics to the axis, set a range (e.g., confidence interval), and format. It highlights a value range in visuals.

79. Q: What is a small multiple in Power BI?

 A: Small multiples repeat a chart for each category, like sales by region. Add a field to Small Multiples in the Visualizations pane for comparison.

80. Q: How do you create a dashboard in Excel?

 A: Use pivot tables, charts, and slicers on a sheet. Link slicers to pivot tables, add visuals like bar charts, and format for an interactive dashboard.

Data Analysis and Statistics (80 Questions)

1. Q: What's the difference between mean and median?

 A: Mean is the average, calculated by summing values and dividing by count. Median is the middle value when sorted. Mean is sensitive to outliers, while median is robust, like for skewed salary data.

2. Q: What is a p-value in statistics?

 A: A p-value shows the probability of results occurring by chance under the null hypothesis. A low p-value (< 0.05) suggests statistical significance, rejecting the null. It's used in hypothesis testing.

3. Q: What is standard deviation?

 A: Standard deviation measures data spread around the mean. A low value means data is close to the mean, like consistent test scores. Calculate it as the square root of variance.

4. Q: What is a confidence interval?

 **A:** A confidence interval estimates a population parameter's range, like 95% confidence that the true mean lies between 50 and 60. It's based on sample data and margin of error.

5. Q: What is correlation?

 **A:** Correlation measures how two variables move together, ranging from -1 to 1. For example, height and weight may have a positive correlation ($r = 0.7$). It doesn't imply causation.

6. Q: What is the difference between Type I and Type II errors?

 **A:** Type I error (false positive) rejects a true null hypothesis, like convicting an innocent person. Type II error (false negative) fails to reject a false null, like missing a guilty person.

7. Q: What is a normal distribution?

 **A:** A normal distribution is a bell-shaped curve where most data clusters around the mean. It's symmetric, with 68% of data within one standard deviation. Many statistical tests assume normality.

8. Q: What is regression analysis?

 **A:** Regression analysis models relationships between variables, like predicting sales from advertising spend. Linear regression fits a line, while logistic regression predicts categories, like customer churn.

9. Q: What is an outlier?

 **A:** An outlier is a data point far from others, like a \$1M sale in \$100 average sales. Detect using box plots or z-scores and handle by investigating or removing if erroneous.

10. Q: What is the central limit theorem?

 **A:** The central limit theorem states that the mean of many random samples approaches a normal distribution, regardless of the population's shape, given a large sample size. It's key for inferential statistics.

11. Q: What is variance?

 **A:** Variance measures how spread out data is, calculated as the average squared deviation from the mean. A high variance, like varied test scores, indicates diverse data.

12. Q: What is a t-test?

 **A:** A t-test compares means between two groups, like test scores of two classes. It checks if differences are significant, assuming normality and equal variances.

13. Q: What is ANOVA?

 A: ANOVA (Analysis of Variance) compares means across three or more groups, like sales by region. It tests if group differences are significant, avoiding multiple t-tests.

14. Q: What is a chi-square test?

 A: A chi-square test checks if categorical data fits an expected distribution, like comparing observed vs. expected survey responses. It's used for independence or goodness-of-fit tests.

15. Q: What is causation vs. correlation?

 A: Correlation shows variables move together, like ice cream sales and temperature. Causation means one causes the other, proven by experiments, not just correlation.

16. Q: What is a z-score?

 A: A z-score measures how many standard deviations a value is from the mean, like $(x - \text{mean}) / \text{SD}$. It's used to compare data points across distributions.

17. Q: What is a hypothesis test?

 A: A hypothesis test evaluates claims using data, like testing if a drug improves health. Set a null hypothesis, calculate a p-value, and decide to reject or accept it.

18. Q: What is a sample in statistics?

 A: A sample is a subset of a population used for analysis, like surveying 100 customers from 10,000. It's chosen randomly to represent the population accurately.

19. Q: What is skewness?

 A: Skewness measures a distribution's asymmetry. Positive skew has a longer right tail (e.g., income), negative skew a longer left tail (e.g., age at death).

20. Q: What is kurtosis?

 A: Kurtosis measures a distribution's tails and peakedness. High kurtosis has heavy tails and a sharp peak, like financial returns. Normal distributions have zero excess kurtosis.

21. Q: What is a binomial distribution?

 A: A binomial distribution models successes in fixed trials with two outcomes, like heads in 10 coin flips. It uses probability p and number of trials n .

22. Q: What is the law of large numbers?

 A: The law of large numbers states that as sample size grows, the sample

mean approaches the population mean. It ensures reliable estimates with enough data.

23. Q: What is a Poisson distribution?

 **A:** A Poisson distribution models rare events in a fixed interval, like customer arrivals per hour. It uses a rate parameter λ for probability calculations.

24. Q: What is a population in statistics?

 **A:** A population is the entire group of interest, like all customers of a store. It's studied via samples due to size or cost constraints.

25. Q: What is a confidence level?

 **A:** A confidence level (e.g., 95%) shows the probability a confidence interval contains the true parameter. Higher levels widen the interval, increasing certainty.

26. Q: What is a box plot?

 **A:** A box plot shows data distribution with median, quartiles, and outliers, like test scores. The box spans Q1 to Q3, with whiskers and dots for extremes.

27. Q: What is multicollinearity?

 **A:** Multicollinearity occurs when independent variables in regression are highly correlated, skewing results. Detect with VIF (Variance Inflation Factor) and fix by removing variables.

28. Q: What is a paired t-test?

 **A:** A paired t-test compares means of the same group under two conditions, like pre- and post-training scores. It accounts for individual differences in paired data.

29. Q: What is a one-sample t-test?

 **A:** A one-sample t-test compares a sample mean to a known value, like testing if average height equals 170 cm. It checks for significant differences.

30. Q: What is a two-sample t-test?

 **A:** A two-sample t-test compares means of two independent groups, like male vs. female test scores. It assumes normality and tests for significant differences.

31. Q: What is a null hypothesis?

 **A:** A null hypothesis assumes no effect or difference, like "a drug has no impact." It's tested against an alternative hypothesis using statistical methods.

32. Q: What is the alternative hypothesis?

 **A:** The alternative hypothesis claims an effect or difference, like “a drug improves health.” It’s accepted if the null hypothesis is rejected based on data.

33. Q: What is a probability distribution?

 **A:** A probability distribution shows the likelihood of outcomes, like rolling a die. Examples include normal, binomial, and Poisson, used for modeling random variables.

34. Q: What is a histogram?

 **A:** A histogram visualizes numerical data distribution with adjacent bars, like test scores in ranges. It shows frequency and helps identify patterns or skewness.

35. Q: What is a percentile?

 **A:** A percentile shows the percentage of data below a value, like the 75th percentile score. It’s used to rank data, like test performance.

36. Q: What is a QQ plot?

 **A:** A QQ (Quantile-Quantile) plot checks if data follows a distribution, like normal. Plot sample quantiles against theoretical ones; a straight line suggests conformity.

37. Q: What is a residual in regression?

 **A:** A residual is the difference between observed and predicted values in regression, like actual vs. predicted sales. Small residuals indicate a good model fit.

38. Q: What is a standard error?

 **A:** Standard error measures the variability of a sample statistic, like the sample mean. It’s calculated as standard deviation divided by the square root of sample size.

39. Q: What is a dependent variable?

 **A:** A dependent variable is the outcome predicted in analysis, like sales in a regression model. It depends on independent variables, like advertising spend.

40. Q: What is an independent variable?

 **A:** An independent variable influences the dependent variable, like hours studied affecting test scores. It’s manipulated or controlled in experiments or models.

41. Q: What is a scatterplot?

 **A:** A scatterplot plots two variables as points, like height vs. weight, to show relationships. It helps identify correlations, trends, or outliers in data.

42. Q: What is a time series?

 **A:** A time series is data collected over time, like daily sales. It's analyzed for trends, seasonality, or forecasts using methods like moving averages.

43. Q: What is seasonality in data?

 **A:** Seasonality is a recurring pattern in time series data, like higher retail sales in December. It's modeled to improve forecasts and understand cycles.

44. Q: What is a moving average?

 **A:** A moving average smooths time series data by averaging a fixed window, like 3-month sales. It reduces noise and highlights trends.

45. Q: What is exponential smoothing?

 **A:** Exponential smoothing forecasts time series by weighting recent data more, like predicting next month's sales. It's simple and effective for short-term trends.

46. Q: What is a confidence interval for a proportion?

 **A:** A confidence interval for a proportion estimates a population percentage, like 95% confidence that 60-70% of customers are satisfied, based on sample data.

47. Q: What is a one-way ANOVA?

 **A:** One-way ANOVA compares means across multiple groups for one factor, like test scores by school. It tests if differences are significant.

48. Q: What is a two-way ANOVA?

 **A:** Two-way ANOVA tests two factors' effects, like teaching method and gender on scores. It checks main effects and interactions between factors.

49. Q: What is a power of a test?

 **A:** The power of a test is the probability of rejecting a false null hypothesis, avoiding Type II errors. Higher power (e.g., 0.8) means better detection.

50. Q: What is a sample size calculation?

 **A:** Sample size calculation determines how many observations are needed for reliable results, considering effect size, power, and significance level, like for surveys.

51. Q: What is a statistical model?

 **A:** A statistical model describes relationships in data, like a regression predicting sales from advertising. It simplifies reality for analysis and prediction.

52. Q: What is a goodness-of-fit test?

 **A:** A goodness-of-fit test checks if data matches an expected distribution, like a chi-square test for survey responses. It validates model assumptions.

53. Q: What is a paired sample?

 **A:** A paired sample involves related observations, like before-and-after scores for the same students. Use paired t-tests to compare paired data.

54. Q: What is a test statistic?

 **A:** A test statistic summarizes data for hypothesis testing, like a t-value in a t-test. It's compared to critical values to decide significance.

55. Q: What is a degrees of freedom?

 **A:** Degrees of freedom (df) is the number of values free to vary in a calculation, like $n-1$ in a t-test. It affects test statistic distributions.

56. Q: What is a probability density function?

 **A:** A probability density function (PDF) describes probabilities for continuous variables, like heights. The area under the curve gives probability over a range.

57. Q: What is a cumulative distribution function?

 **A:** A cumulative distribution function (CDF) gives the probability a variable is less than or equal to a value, like $P(X \leq 5)$ for test scores.

58. Q: What is a moment in statistics?

 **A:** A moment describes a distribution's shape, like mean (first moment) or variance (second). Higher moments include skewness and kurtosis for detailed analysis.

59. Q: What is a likelihood function?

 **A:** A likelihood function shows how likely parameters explain data, used in maximum likelihood estimation. For example, it estimates a coin's bias from flips.

60. Q: What is a non-parametric test?

 **A:** A non-parametric test analyzes data without assuming a distribution, like the Mann-Whitney U test for comparing medians. It's used for non-normal data.

61. Q: What is a Mann-Whitney U test?

 **A:** The Mann-Whitney U test compares two independent groups'

distributions, like test scores between schools. It's non-parametric, ideal for non-normal data.

62. Q: What is a Kruskal-Wallis test?

 **A:** The Kruskal-Wallis test compares medians across three or more groups, like customer satisfaction by store. It's a non-parametric alternative to ANOVA.

63. Q: What is a Wilcoxon signed-rank test?

 **A:** The Wilcoxon signed-rank test compares paired data, like pre- and post-test scores, non-parametrically. It's used when data isn't normally distributed.

64. Q: What is a bootstrap method?

 **A:** Bootstrapping estimates statistics by resampling data with replacement, like calculating a mean's confidence interval. It's flexible for complex or small datasets.

65. Q: What is a Monte Carlo simulation?

 **A:** A Monte Carlo simulation uses random sampling to model uncertainty, like predicting project costs. It runs many scenarios to estimate probabilities.

66. Q: What is a Bayesian approach?

 **A:** The Bayesian approach updates probabilities with new data, using prior beliefs. For example, it refines disease risk estimates as test results arrive.

67. Q: What is a prior probability?

 **A:** Prior probability is the initial likelihood of an event before new data, like assuming a 50% chance of rain. It's used in Bayesian analysis.

68. Q: What is a posterior probability?

 **A:** Posterior probability is the updated likelihood after new data, like 70% rain chance after clouds appear. It combines prior and observed data in Bayesian methods.

69. Q: What is a likelihood ratio?

 **A:** A likelihood ratio compares the probability of data under two hypotheses, like disease vs. no disease. It's used in diagnostic tests for decision-making.

70. Q: What is a survival analysis?

 **A:** Survival analysis studies time-to-event data, like patient recovery times. It uses methods like Kaplan-Meier curves to estimate survival probabilities.

71. Q: What is a Kaplan-Meier estimator?

 **A:** The Kaplan-Meier estimator plots survival probabilities over time, like patient survival rates. It accounts for censored data, common in medical studies.

72. Q: What is a hazard function?

⚠️ A: A hazard function shows the instantaneous risk of an event, like failure at time t. It's used in survival analysis to model event rates.

73. Q: What is a Cox proportional hazards model?

📊 A: The Cox model predicts hazard rates based on covariates, like age affecting survival. It assumes proportional hazards and is widely used in medical research.

74. Q: What is a log-rank test?

📈 A: The log-rank test compares survival distributions between groups, like treatment vs. control. It's non-parametric and used in survival analysis.

75. Q: What is a contingency table?

📋 A: A contingency table cross-tabulates categorical data, like gender vs. preference. It's analyzed with chi-square tests to check for associations.

76. Q: What is a correlation coefficient?

🔗 A: A correlation coefficient (r) measures linear relationship strength, from -1 to 1. For example, $r = 0.8$ indicates a strong positive correlation between variables.

77. Q: What is a partial correlation?

📈 A: Partial correlation measures the relationship between two variables, controlling for others. For example, it checks income vs. education, adjusting for age.

78. Q: What is a factor analysis?

📊 A: Factor analysis identifies underlying factors in data, like grouping survey questions into themes (e.g., satisfaction). It reduces dimensions for analysis.

79. Q: What is a principal component analysis (PCA)?

📈 A: PCA reduces data dimensions by creating new variables (components) capturing maximum variance. It's used for visualization or simplifying complex datasets.

80. Q: What is a discriminant analysis?

📈 A: Discriminant analysis classifies data into groups, like predicting customer churn based on behavior. It maximizes group separation using linear combinations.

Python for Data (80 Questions)

1. Q: What is Pandas in Python?

⌚ A: Pandas is a Python library for data analysis, providing DataFrames and

Series. It handles tasks like reading CSVs, filtering, or grouping data. For example, `df['sales'].mean()` calculates average sales.

2. Q: How do you read a CSV file using Pandas?

 **A:** Use `pd.read_csv('file.csv')`, like `import pandas as pd; df = pd.read_csv('data.csv')`. It loads the CSV into a DataFrame for analysis, with options for delimiters or headers.

3. Q: What is a DataFrame in Pandas?

 **A:** A DataFrame is a table-like structure with rows and columns, like a spreadsheet. For example, `df = pd.DataFrame({'name': ['A', 'B'], 'age': [25, 30]})` creates a simple DataFrame.

4. Q: How do you filter rows in Pandas?

 **A:** Use boolean indexing, like `df[df['age'] > 25]` to get rows where age exceeds 25. Combine conditions with & or |, like `df[(df['age'] > 25) & (df['city'] == 'NY')]`.

5. Q: What is the groupby function in Pandas?

 **A:** groupby groups data by a column for aggregation, like `df.groupby('region')['sales'].sum()`. It summarizes data, such as total sales per region.

6. Q: How do you merge DataFrames in Pandas?

 **A:** Use `merge`, like `df_merged = pd.merge(df1, df2, on='id', how='inner')`. Options include inner, left, right, or outer joins, similar to SQL.

7. Q: What is NumPy used for?

 **A:** NumPy is a Python library for numerical computations, offering arrays and math functions. For example, `np.array([1, 2, 3]).mean()` calculates the average. It's fast and foundational for data science.

8. Q: How do you handle missing values in Pandas?

 **A:** Use `df.dropna()` to remove rows with missing values or `df.fillna(0)` to replace them with a value, like 0. Check with `df.isnull().sum()` to identify missing data.

9. Q: What is a list comprehension in Python?

 **A:** A list comprehension creates lists concisely, like `[x*2 for x in [1, 2, 3]]` to get `[2, 4, 6]`. It's faster than loops for simple transformations.

10. Q: How do you plot data using Matplotlib?

 **A:** Use `plt.plot()`, like `import matplotlib.pyplot as plt; plt.plot(x, y); plt.show()`. It creates line plots, customizable with labels, titles, or styles for data visualization.

11. Q: What is the difference between a list and a NumPy array?

 A: A list is a flexible Python structure with mixed types, while a NumPy array is fixed-size, numeric, and supports vectorized operations, like `np.array([1, 2]) + 3`.

12. Q: How do you sort a DataFrame in Pandas?

 A: Use `df.sort_values('column')`, like `df.sort_values('sales', ascending=False)` for descending order. Add `inplace=True` to modify the DataFrame directly.

13. Q: What is a lambda function in Python?

 A: A lambda function is an anonymous, one-line function, like `lambda x: x*2`. Use it for short tasks, like `df.apply(lambda x: x*2)` to double values.

14. Q: How do you pivot a DataFrame in Pandas?

 A: Use `df.pivot(index='row', columns='column', values='value')`, like `df.pivot(index='date', columns='product', values='sales')` to reshape data into a wide format.

15. Q: What is Seaborn in Python?

 A: Seaborn is a visualization library built on Matplotlib, offering stylish plots. For example, `sns.boxplot(x='region', y='sales', data=df)` creates a box plot with ease.

16. Q: How do you drop a column in Pandas?

 A: Use `df.drop('column', axis=1)`, like `df.drop('age', axis=1, inplace=True)` to remove the age column. `axis=1` specifies columns, not rows.

17. Q: What is a dictionary in Python?

 A: A dictionary stores key-value pairs, like `d = {'name': 'Alice', 'age': 25}`. Access values with `d['name']`. It's used for mapping data, like configs.

18. Q: How do you calculate correlations in Pandas?

 A: Use `df.corr()`, like `df[['sales', 'profit']].corr()` to get a correlation matrix. Values range from -1 to 1, showing variable relationships.

19. Q: What is the apply function in Pandas?

 A: `apply` runs a function on DataFrame rows or columns, like `df['price'].apply(lambda x: x*1.1)` to increase prices by 10%. It's flexible for custom operations.

20. Q: How do you create a histogram in Seaborn?

 A: Use `sns.histplot(data=df, x='sales')` to plot a histogram of sales. Customize bins or add a kernel density estimate (KDE) for distribution insights.

21. Q: What is a tuple in Python?

💡 A: A tuple is an immutable sequence, like `t = (1, 2, 3)`. Access with `t[0]`. Use for fixed data, like coordinates, since it can't be changed.

22. Q: How do you rename columns in Pandas?

👉 A: Use `df.rename(columns={'old': 'new'})`, like `df.rename(columns={'sales': 'revenue'}, inplace=True)`. It updates column names for clarity or consistency.

23. Q: What is a set in Python?

💡 A: A set is an unordered collection of unique items, like `s = {1, 2, 2}` resulting in `{1, 2}`. Use for deduplication or operations like union/intersection.

24. Q: How do you handle duplicates in Pandas?

✍️ A: Use `df.drop_duplicates()`, like `df.drop_duplicates(subset='id', keep='first')` to keep the first occurrence of each id. It cleans redundant data.

25. Q: What is broadcasting in NumPy?

💻 A: Broadcasting applies operations to arrays of different shapes, like `np.array([1, 2]) + 3` adding 3 to each element. It simplifies vectorized calculations.

26. Q: How do you create a scatter plot in Matplotlib?

💡 A: Use `plt.scatter(x, y)`, like `plt.scatter(df['sales'], df['profit'])`; `plt.show()`. Customize with colors or sizes to visualize relationships.

27. Q: What is the loc accessor in Pandas?

🔍 A: `loc` selects rows/columns by labels, like `df.loc[0:2, 'name']` for rows 0-2 of the name column. It's explicit and flexible for indexing.

28. Q: What is the iloc accessor in Pandas?

💡 A: `iloc` selects rows/columns by integer positions, like `df.iloc[0:2, 0]` for the first two rows of the first column. It's for numerical indexing.

29. Q: How do you create a box plot in Seaborn?

📊 A: Use `sns.boxplot(x='category', y='value', data=df)` to show distribution, like sales by region. It highlights medians, quartiles, and outliers.

30. Q: What is a Series in Pandas?

📋 A: A Series is a one-dimensional

A: labeled array in Pandas, like a column in a DataFrame. For example, `pd.Series([1, 2, 3], index=['a', 'b', 'c'])` creates a Series with custom indices. It's used for single-column data operations.

31. Q: How do you concatenate DataFrames in Pandas?

 A: Use `pd.concat([df1, df2])` to stack DataFrames vertically or horizontally. For example, `pd.concat([df1, df2], axis=0)` appends rows, while `axis=1` merges columns. Ensure matching indices or columns.

32. Q: What is a pivot table in Pandas?

 A: A pivot table summarizes data, like `df.pivot_table(index='region', columns='year', values='sales', aggfunc='sum')`. It aggregates sales by region and year, similar to Excel pivot tables.

33. Q: How do you save a DataFrame to CSV?

 A: Use `df.to_csv('file.csv')`, like `df.to_csv('output.csv', index=False)` to exclude the index. It exports the DataFrame to a CSV file for sharing or storage.

34. Q: What is the describe method in Pandas?

 A: `describe()` provides summary statistics, like `df.describe()` for count, mean, min, max, and quartiles of numeric columns. It's quick for understanding data distribution.

35. Q: How do you create a bar plot in Matplotlib?

 A: Use `plt.bar(x, height)`, like `plt.bar(df['category'], df['sales'])`; `plt.show()`. Customize with colors or labels to compare categorical data visually.

36. Q: What is a virtual environment in Python?

 A: A virtual environment isolates Python packages for projects, avoiding conflicts. Create with `python -m venv env`, activate with `source env/bin/activate`, and install packages like Pandas.

37. Q: How do you read Excel files in Pandas?

 A: Use `pd.read_excel('file.xlsx')`, like `df = pd.read_excel('data.xlsx', sheet_name='Sheet1')`. Install openpyxl or xlrd for Excel file support.

38. Q: What is the value_counts method in Pandas?

 A: `value_counts()` counts unique values in a Series, like `df['city'].value_counts()`. It's useful for summarizing categorical data, like customer locations.

39. Q: How do you create a heatmap in Seaborn?

 A: Use `sns.heatmap(df.corr())` to visualize correlations, with colors showing strength. Customize with `cmap='coolwarm'` or `annot=True` to display values.

40. Q: What is the astype method in Pandas?

 A: `astype` changes a column's data type, like `df['age'].astype(float)` to convert integers to floats. It ensures compatibility for calculations or analysis.

41. Q: How do you create a line plot in Seaborn?

 **A:** Use `sns.lineplot(x='date', y='sales', data=df)` to plot trends over time. It's smoother than Matplotlib and supports hue for multiple lines.

42. Q: What is a generator in Python?

 **A:** A generator yields values one at a time, saving memory, like `def my_gen(): yield 1; yield 2`. Use in loops or with `next()` for large datasets.

43. Q: How do you reset a DataFrame index?

 **A:** Use `df.reset_index()`, like `df.reset_index(drop=True)` to create a new integer index and discard the old one. It's useful after filtering or grouping.

44. Q: What is the numpy.random module?

 **A:** `numpy.random` generates random numbers, like `np.random.rand(5)` for 5 random floats. It's used for simulations, testing, or sampling in data analysis.

45. Q: How do you apply a function to a DataFrame?

 **A:** Use `df.apply(func)`, like `df['price'].apply(lambda x: x * 1.1)` to increase prices. For entire rows, use `df.apply(func, axis=1)`.

46. Q: What is the melt function in Pandas?

 **A:** `melt` reshapes wide data to long format, like `pd.melt(df, id_vars='id', value_vars=['sales', 'profit'])`. It's used for tidying data for analysis.

47. Q: How do you create a violin plot in Seaborn?

 **A:** Use `sns.violinplot(x='category', y='value', data=df)` to show data distribution with density curves. It combines box plot and kernel density for insights.

48. Q: What is a try-except block in Python?

 **A:** A try-except block handles errors, like `try: x = 1/0 except ZeroDivisionError: print('Error')`. It prevents crashes and logs issues in scripts.

49. Q: How do you join DataFrames on index?

 **A:** Use `df1.join(df2)`, like `df1.join(df2, how='inner')` to merge on indices. Ensure aligned indices for correct matching, similar to `merge`.

50. Q: What is the numpy.where function?

 **A:** `np.where` applies conditions, like `np.where(df['sales'] > 100, 'High', 'Low')`. It's a vectorized alternative to loops for conditional assignments.

51. Q: How do you create a pair plot in Seaborn?

 **A:** Use `sns.pairplot(df)` to plot pairwise relationships for all numeric columns. Add `hue='category'` to color by a categorical variable, like regions.

52. Q: What is a list in Python?

 A: A list is a mutable sequence, like `l = [1, 2, 3]`. Access with `l[0]`, append with `l.append(4)`. It's versatile for storing ordered data.

53. Q: How do you slice a DataFrame in Pandas?

 A: Use `df[start:end]` for rows, like `df[0:5]`. For columns, use `df[['col1', 'col2']]`. Combine with `loc` or `iloc` for precise slicing.

54. Q: What is the cumsum method in Pandas?

 A: `cumsum` calculates cumulative sums, like `df['sales'].cumsum()` for running total sales. It's used for time series or progressive calculations.

55. Q: How do you create a pie chart in Matplotlib?

 A: Use `plt.pie(values, labels=labels)`, like `plt.pie(df['sales'], labels=df['category'])`; `plt.show()`. Customize with colors or `explode` for emphasis.

56. Q: What is a context manager in Python?

 A: A context manager handles setup/teardown, like with `open('file.txt')` as `f`. It ensures resources, like files, are closed properly after use.

57. Q: How do you check DataFrame data types?

 A: Use `df.dtypes` to list column data types, like `int64` or `object`. It helps verify data for analysis or debugging type-related errors.

58. Q: What is the rolling method in Pandas?

 A: `rolling` computes moving window calculations, like `df['sales'].rolling(window=3).mean()` for a 3-period moving average. It's used for smoothing time series.

59. Q: How do you create a stacked bar plot in Seaborn?

 A: Use `sns.catplot(x='year', y='sales', hue='product', kind='bar', data=df)` with `stacked=True`. It shows sales by product, stacked within years.

60. Q: What is the idxmax method in Pandas?

 A: `idxmax()` returns the index of the maximum value, like `df['sales'].idxmax()` for the row with highest sales. It's useful for identifying peaks.

61. Q: How do you create a density plot in Seaborn?

 A: Use `sns.kdeplot(data=df['sales'])` to plot a kernel density estimate. It shows the probability density of continuous data, like sales distribution.

62. Q: What is a decorator in Python?

 A: A decorator modifies a function's behavior, like `@timer` to time execution. Define with `@decorator` above a function to add functionality, like logging.

63. Q: How do you sample data in Pandas?

 **A:** Use `df.sample(n=5)` to randomly select 5 rows. Add `random_state=42` for reproducibility or `frac=0.1` for 10% of the data.

64. Q: What is the numpy.linspace function?

 **A:** `np.linspace(start, stop, num)` creates an array of evenly spaced numbers, like `np.linspace(0, 10, 5)` for [0, 2.5, 5, 7.5, 10]. It's used for plotting or simulations.

65. Q: How do you create a box plot in Matplotlib?

 **A:** Use `plt.boxplot(df['sales'])` to show distribution with median and outliers. Customize with `labels=['Sales']` or `patch_artist=True` for styling.

66. Q: What is a frozen set in Python?

 **A:** A frozen set is an immutable set, like `fs = frozenset([1, 2, 3])`. Use for unchangeable collections, like dictionary keys, unlike regular sets.

67. Q: How do you encode categorical variables in Pandas?

 **A:** Use `pd.get_dummies(df['category'])` for one-hot encoding, creating binary columns. Alternatively, map values, like `df['gender'].map({'M': 0, 'F': 1})`.

68. Q: What is the numpy.concatenate function?

 **A:** `np.concatenate([arr1, arr2])` joins arrays along an axis, like stacking rows (`axis=0`) or columns (`axis=1`). It's used for combining datasets.

69. Q: How do you create a regression plot in Seaborn?

 **A:** Use `sns.regplot(x='sales', y='profit', data=df)` to plot data with a fitted regression line. It shows linear relationships and confidence intervals.

70. Q: What is the shift method in Pandas?

 **A:** `shift` moves data by periods, like `df['sales'].shift(1)` to get previous period's sales. It's used for time series lags or differences.

71. Q: How do you create a swarm plot in Seaborn?

 **A:** Use `sns.swarmplot(x='category', y='value', data=df)` to plot individual points without overlap. It's great for small datasets to show distribution.

72. Q: What is a module in Python?

 **A:** A module is a Python file with functions, classes, or variables, like `math`. Import with `import math` to use `math.sqrt()` for reusable code.

73. Q: How do you calculate moving averages in Pandas?

 **A:** Use `df['sales'].rolling(window=3).mean()` for a 3-period moving average. It smooths data for trend analysis in time series.

74. Q: What is the numpy.zeros function?

💡 A: np.zeros(shape) creates an array of zeros, like np.zeros((2, 3)) for a 2x3 matrix. It's used for initializing arrays in computations.

75. Q: How do you create a count plot in Seaborn?

📊 A: Use sns.countplot(x='category', data=df) to plot the frequency of categories, like customer types. Add hue='region' for subgroup counts.

76. Q: What is a comprehension in Python?

📝 A: A comprehension creates collections concisely, like [x**2 for x in range(5)] for [0, 1, 4, 9, 16]. It includes list, dict, and set comprehensions.

77. Q: How do you fill missing values with the mean in Pandas?

🛠️ A: Use df['sales'].fillna(df['sales'].mean()) to replace missing values with the column's mean. Apply inplace=True to modify the DataFrame.

78. Q: What is the numpy.dot function?

➡️ A: np.dot(a, b) computes the dot product of arrays, like matrix multiplication. For example, np.dot([1, 2], [3, 4]) yields 11 (13 + 24).

79. Q: How do you create a strip plot in Seaborn?

📍 A: Use sns.stripplot(x='category', y='value', data=df) to plot points along an axis, like sales by region. It's similar to swarm but allows overlap.

80. Q: What is the to_datetime function in Pandas?

📅 A: to_datetime converts strings to datetime, like pd.to_datetime(df['date']). It enables date-based operations, like extracting year with df['date'].dt.year.

Business Analysis (40 Questions)

1. Q: What does a business analyst do?

💼 A: A business analyst bridges business and tech, gathering requirements, analyzing processes, and proposing solutions. For example, they document user needs for a new app to improve efficiency.

2. Q: How do you handle difficult stakeholders?

🤝 A: Listen actively, clarify concerns, and use data to support solutions. Communicate regularly, align with their goals, and build trust to resolve conflicts and ensure collaboration.

3. Q: What is a requirement elicitation?

📋 A: Requirement elicitation gathers stakeholder needs through interviews, surveys, or workshops. For example, discussing with users to define features for a CRM system ensures project alignment.

4. Q: What is a use case diagram?

 **A:** A use case diagram shows system interactions, with actors (users) and use cases (functions). For example, a “Login” use case connects a “Customer” to an app.

5. Q: What is SWOT analysis?

 **A:** SWOT analyzes Strengths, Weaknesses, Opportunities, and Threats. For example, a company’s strong brand (S) faces new competitors (T), guiding strategic planning.

6. Q: What is a process flow diagram?

 **A:** A process flow diagram maps steps in a process, like order fulfillment. It uses shapes (e.g., ovals for start/end) to visualize workflows and identify inefficiencies.

7. Q: What is gap analysis?

 **A:** Gap analysis compares current and desired states, like current sales vs. target. It identifies gaps and suggests actions, like training to boost performance.

8. Q: What is a stakeholder matrix?

 **A:** A stakeholder matrix maps stakeholders by influence and interest, like high-influence executives vs. low-interest staff. It guides communication and engagement strategies.

9. Q: What is a BRD?

 **A:** A Business Requirements Document (BRD) outlines business needs, goals, and solutions. For example, it details features, scope, and stakeholders for a new system.

10. Q: What is a FRD?

 **A:** A Functional Requirements Document (FRD) specifies system functions, like “users can reset passwords.” It details technical requirements based on the BRD.

11. Q: What is a user story?

 **A:** A user story describes a feature from the user’s perspective, like “As a customer, I want to track orders.” It follows the format: As a [user], I want [function], so [benefit].

12. Q: What is MoSCoW prioritization?

 **A:** MoSCoW prioritizes requirements: Must have, Should have, Could have, Won’t have. For example, “login” is Must, while “dark mode” is Could, guiding project focus.

13. Q: What is a wireframe?

 **A:** A wireframe is a low-fidelity sketch of a user interface, like a webpage layout. It shows structure and functionality, aiding design discussions with stakeholders.

14. Q: What is a decision tree in business analysis?

 **A:** A decision tree maps choices and outcomes, like choosing a vendor based on cost and quality. It visualizes decisions to support strategic planning.

15. Q: What is BPMN?

 **A:** BPMN (Business Process Model and Notation) standardizes process diagrams, using symbols like arrows for flow. It models workflows, like order processing, for clarity.

16. Q: What is a RACI chart?

 **A:** A RACI chart assigns roles: Responsible, Accountable, Consulted, Informed. For example, a developer is Responsible for coding, while a manager is Accountable.

17. Q: What is root cause analysis?

 **A:** Root cause analysis identifies the source of a problem, like low sales. Use techniques like the 5 Whys to find issues, such as poor marketing, and fix them.

18. Q: What is a data dictionary?

 **A:** A data dictionary defines data elements, like “customer_id: unique integer.” It ensures consistency in databases or reports, aiding developers and analysts.

19. Q: What is a feasibility study?

 **A:** A feasibility study assesses a project’s viability, like launching a new app. It evaluates technical, financial, and operational aspects to guide decisions.

20. Q: What is a persona in business analysis?

 **A:** A persona is a fictional user profile, like “Sarah, 30, tech-savvy.” It represents target users to guide product design and meet user needs.

21. Q: What is a flowchart?

 **A:** A flowchart visualizes a process with symbols, like diamonds for decisions. For example, it maps customer support steps to identify bottlenecks.

22. Q: What is a cost-benefit analysis?

 **A:** Cost-benefit analysis compares project costs to benefits, like \$10K for a system saving \$20K yearly. It justifies investments by quantifying returns.

23. Q: What is a traceability matrix?

 **A:** A traceability matrix links requirements to deliverables, ensuring all needs are met. For example, it tracks “login feature” from BRD to code.

24. Q: What is agile methodology?

 **A:** Agile is an iterative approach to projects, using sprints to deliver small features. For example, a team builds an app incrementally, adapting to feedback.

25. Q: What is a sprint in agile?

 **A:** A sprint is a short, fixed period (e.g., 2 weeks) in agile to complete tasks, like coding a feature. It ends with a review and planning.

26. Q: What is a product backlog?

 **A:** A product backlog lists prioritized tasks for a project, like features or bugs. It's managed by the product owner and refined in agile sprints.

27. Q: What is a burndown chart?

 **A:** A burndown chart tracks work completed in a sprint, plotting remaining tasks vs. time. It helps teams monitor progress and meet deadlines.

28. Q: What is a Kanban board?

 **A:** A Kanban board visualizes workflow with columns like “To Do,” “In Progress,” “Done.” Tasks move across, improving transparency and efficiency.

29. Q: What is a stakeholder interview?

 **A:** A stakeholder interview gathers insights on needs or issues, like asking managers about system pain points. It informs requirements and builds consensus.

30. Q: What is a business case?

 **A:** A business case justifies a project, detailing benefits, costs, and risks. For example, it argues for a new CRM to boost sales efficiency.

31. Q: What is a workshop in business analysis?

 **A:** A workshop is a collaborative session with stakeholders to define requirements or solve problems, like brainstorming app features. It fosters alignment and creativity.

32. Q: What is a context diagram?

 **A:** A context diagram shows a system’s interactions with external entities, like users or databases. It defines scope, like an app’s data flows.

33. Q: What is a PESTLE analysis?

 **A:** PESTLE analyzes Political, Economic, Social, Technological, Legal, and

Environmental factors. For example, it assesses market entry risks for a new product.

34. Q: What is a stakeholder analysis?

 **A:** Stakeholder analysis identifies and prioritizes stakeholders by influence and interest. For example, it ensures key executives are engaged early in projects.

35. Q: What is a risk register?

 **A:** A risk register lists project risks, their likelihood, and mitigation plans. For example, “server downtime” has a backup plan to ensure uptime.

36. Q: What is a change request?

 **A:** A change request proposes project modifications, like adding a feature. It's reviewed for impact on scope, cost, and timeline before approval.

37. Q: What is a value stream map?

 **A:** A value stream map visualizes the flow of value, like order processing steps. It identifies waste, like delays, to optimize processes.

38. Q: What is a decision table?

 **A:** A decision table lists conditions and actions, like “if payment late, send reminder.” It clarifies complex logic for business rules.

39. Q: What is a user acceptance test (UAT)?

 **A:** UAT verifies a system meets user needs, like testing an app’s checkout. Users perform scenarios to ensure functionality before launch.

40. Q: What is a business process reengineering?

 **A:** Business process reengineering redesigns workflows for major improvements, like automating manual tasks. It aims for efficiency, like reducing order processing time.

Machine Learning and Algorithms (40 Questions)

1. Q: What is supervised learning?

 **A:** Supervised learning trains models on labeled data, like predicting prices from house sizes. The model learns input-output mappings for predictions, like spam detection.

2. Q: What is unsupervised learning?

 **A:** Unsupervised learning finds patterns in unlabeled data, like clustering customers by behavior. It’s used for segmentation or anomaly detection without predefined outputs.

3. Q: What is overfitting in machine learning?

 **A:** Overfitting occurs when a model learns training data noise, performing poorly on new data. Prevent it with regularization, cross-validation, or more data.

4. Q: What is a decision tree?

 **A:** A decision tree splits data into branches based on features, like “age > 30” to predict purchases. It’s interpretable but can overfit without pruning.

5. Q: What is a confusion matrix?

 **A:** A confusion matrix shows classification performance, with true positives, true negatives, false positives, and false negatives. It evaluates model accuracy, like for disease detection.

6. Q: What is precision in classification?

 **A:** Precision is the ratio of true positives to predicted positives, like $TP / (TP + FP)$. High precision means fewer false positives, crucial for spam filters.

7. Q: What is recall in classification?

 **A:** Recall is the ratio of true positives to actual positives, like $TP / (TP + FN)$. High recall ensures most positives are caught, vital for medical diagnoses.

8. Q: What is the F1 score?

 **A:** The F1 score balances precision and recall, calculated as $2 * (\text{precision} * \text{recall}) / (\text{precision} + \text{recall})$. It’s used when both false positives and negatives matter.

9. Q: What is a random forest?

 **A:** A random forest combines multiple decision trees, averaging predictions to reduce overfitting. It’s robust for tasks like predicting customer churn.

10. Q: What is gradient descent?

 **A:** Gradient descent optimizes model parameters by minimizing a loss function, like adjusting weights in linear regression. It iteratively steps toward the minimum using gradients.

11. Q: What is a support vector machine (SVM)?

 **A:** SVM finds a hyperplane maximizing the margin between classes, like separating spam from non-spam emails. It handles high-dimensional data with kernels.

12. Q: What is k-means clustering?

 **A:** K-means clustering groups data into k clusters by minimizing distance to centroids. For example, it segments customers by purchase patterns.

13. Q: What is a neural network?

 **A:** A neural network mimics the brain, with layers of nodes processing inputs to outputs. It's used for complex tasks like image recognition.

14. Q: What is a loss function?

 **A:** A loss function measures model error, like mean squared error for regression. It guides optimization, minimizing the difference between predictions and actuals.

15. Q: What is cross-validation?

 **A:** Cross-validation splits data into folds, training on some and testing on others, like k-fold. It estimates model performance and reduces overfitting risks.

16. Q: What is feature scaling?

 **A:** Feature scaling standardizes or normalizes features, like scaling ages to 0-1. It ensures equal importance in algorithms like SVM or neural networks.

17. Q: What is a hyperparameter?

 **A:** A hyperparameter is a model setting, like the number of trees in a random forest. Tune them using grid search to optimize performance.

18. Q: What is logistic regression?

 **A:** Logistic regression predicts probabilities for binary outcomes, like customer churn (0 or 1). It uses a sigmoid function to model class probabilities.

19. Q: What is a ROC curve?

 **A:** A ROC curve plots true positive rate vs. false positive rate at various thresholds. The area under the curve (AUC) measures classification performance.

20. Q: What is dimensionality reduction?

 **A:** Dimensionality reduction simplifies data by reducing features, like PCA combining variables. It improves speed and reduces noise while preserving key information.

21. Q: What is principal component analysis (PCA)?

 **A:** PCA reduces dimensions by creating new variables (components) capturing maximum variance. It's used for visualization or speeding up machine learning models.

22. Q: What is a bias-variance tradeoff?

 **A:** The bias-variance tradeoff balances model complexity. High bias (simple models) underfits, while high variance (complex models) overfits. Optimal models balance both.

23. Q: What is a Naive Bayes classifier?

 **A:** Naive Bayes predicts classes using probability, assuming feature independence, like spam detection. It's fast and effective for text classification tasks.

24. Q: What is regularization?

 **A:** Regularization adds penalties to model complexity, like L1 (Lasso) or L2 (Ridge), to prevent overfitting. It shrinks coefficients in regression models.

25. Q: What is a learning rate in gradient descent?

 **A:** The learning rate controls step size in gradient descent, like 0.01. Too high causes instability; too low slows convergence. Tune for optimal training.

26. Q: What is a KNN algorithm?

 **A:** K-Nearest Neighbors (KNN) classifies data based on the k closest points, like predicting a customer's class by nearby data. It's simple but slow for large datasets.

27. Q: What is a feature in machine learning?

 **A:** A feature is an input variable, like "age" or "income" in a dataset. Good features improve model accuracy, like predicting house prices.

28. Q: What is a training set?

 **A:** A training set is data used to train a model, like labeled house prices. It's typically 70-80% of the dataset, with the rest for testing.

29. Q: What is a test set?

 **A:** A test set evaluates a trained model's performance, like predicting unseen house prices. It ensures the model generalizes to new data.

30. Q: What is a validation set?

 **A:** A validation set tunes model hyperparameters, like testing different k in KNN. It's separate from training and test sets to avoid bias.

31. Q: What is a decision boundary?

 **A:** A decision boundary separates classes in a model, like a line in SVM. It defines regions where the model predicts one class over another.

32. Q: What is ensemble learning?

 **A:** Ensemble learning combines models, like random forests or boosting, to improve accuracy. It reduces errors by averaging or weighting predictions.

33. Q: What is boosting in machine learning?

 **A:** Boosting builds models sequentially, focusing on errors, like in AdaBoost.

Each model corrects predecessors, improving accuracy for tasks like classification.

34. Q: What is bagging?

 A: Bagging (Bootstrap Aggregating) trains models on random data subsets, like random forests. It averages predictions to reduce variance and overfitting.

35. Q: What is a perceptron?

 A: A perceptron is a simple neural network unit, taking weighted inputs and applying an activation function. It's the building block for deep learning models.

36. Q: What is a kernel in SVM?

 A: A kernel transforms data into higher dimensions, like RBF kernel in SVM, to separate non-linear classes. It enables complex pattern recognition.

37. Q: What is a cost function?

 A: A cost function quantifies model error, like mean squared error. It's minimized during training to optimize parameters, like weights in regression.

38. Q: What is a gradient boosting machine?

 A: Gradient boosting builds trees sequentially, minimizing errors with gradients, like in XGBoost. It's powerful for structured data tasks, like predicting sales.

39. Q: What is a clustering algorithm?

 A: A clustering algorithm groups similar data, like k-means for customer segments. It's unsupervised, finding patterns without predefined labels.

40. Q: What is an activation function?

 A: An activation function decides a neuron's output, like sigmoid or ReLU in neural networks. It introduces non-linearity for complex pattern learning.

Data Engineering (40 Questions)

1. Q: What is ETL in data engineering?

 A: ETL (Extract, Transform, Load) extracts data from sources, transforms it (e.g., cleaning), and loads it into a destination, like a warehouse. It prepares data for analysis.

2. Q: What's the difference between batch and stream processing?

 A: Batch processing handles large data chunks at intervals, like nightly reports. Stream processing analyzes data in real-time, like live user clicks, for immediate insights.

3. Q: What is a data warehouse?

 **A:** A data warehouse stores structured, historical data for analysis, like sales over years. It's optimized for queries, supporting BI tools like Tableau.

4. Q: What is a data lake?

 **A:** A data lake stores raw, unstructured, or structured data, like logs or images. It's flexible for advanced analytics but requires processing before use.

5. Q: What is Apache Spark?

 **A:** Apache Spark is a distributed computing framework for big data. It processes large datasets in memory, like aggregating sales, faster than Hadoop MapReduce.

6. Q: What is a data pipeline?

 **A:** A data pipeline automates data flow from source to destination, like extracting logs, transforming, and loading into a warehouse. It ensures reliable data delivery.

7. Q: What is Hadoop?

 **A:** Hadoop is a framework for distributed storage (HDFS) and processing (MapReduce) of big data. It handles large-scale tasks, like log analysis, across clusters.

8. Q: What is a schema-on-read?

 **A:** Schema-on-read applies structure when reading data, like in data lakes. It's flexible, allowing different analyses without predefined schemas, unlike databases.

9. Q: What is a schema-on-write?

 **A:** Schema-on-write defines structure before storing data, like in data warehouses. It ensures consistency but requires upfront design, unlike data lakes.

10. Q: What is Apache Kafka?

 **A:** Apache Kafka is a streaming platform for handling real-time data feeds, like user events. It publishes and subscribes to data streams, ensuring scalability.

11. Q: What is data partitioning?

 **A:** Data partitioning splits data into chunks, like by date, for faster queries. It's used in databases or warehouses to improve performance.

12. Q: What is a star schema?

 **A:** A star schema organizes data with a central fact table (e.g., sales) linked to dimension tables (e.g., time, product). It simplifies queries in warehouses.

13. Q: What is a snowflake schema?

 **A:** A snowflake schema normalizes dimension tables in a star schema, like splitting product into sub-tables. It reduces redundancy but increases query complexity.

14. Q: What is data ingestion?

 **A:** Data ingestion collects data from sources, like APIs or files, into a system. It can be batch (e.g., daily) or streaming (e.g., real-time logs).

15. Q: What is data lineage?

 **A:** Data lineage tracks data's journey from source to destination, like from logs to reports. It ensures transparency and helps debug pipeline issues.

16. Q: What is Apache Airflow?

 **A:** Apache Airflow schedules and monitors data pipelines, defining workflows as DAGs (Directed Acyclic Graphs). It automates tasks, like ETL jobs.

17. Q: What is a data catalog?

 **A:** A data catalog organizes metadata, like table descriptions or sources. It helps analysts find and understand data, like in tools like Alation.

18. Q: What is data governance?

 **A:** Data governance sets policies for data quality, security, and access. For example, it ensures only authorized users query sensitive customer data.

19. Q: What is a data mart?

 **A:** A data mart is a subset of a data warehouse, focused on a department, like marketing. It's optimized for specific analytics, like campaign performance.

20. Q: What is data deduplication?

 **A:** Data deduplication removes duplicate records, like repeated customer entries. It ensures data quality and reduces storage in pipelines or warehouses.

21. Q: What is a distributed database?

 **A:** A distributed database stores data across multiple nodes, like Cassandra. It improves scalability and fault tolerance for large-scale applications.

22. Q: What is data compression?

 **A:** Data compression reduces data size, like using Parquet format. It saves storage and speeds up queries in big data systems like Spark.

23. Q: What is a data transformation?

 **A:** Data transformation converts data, like cleaning, aggregating, or formatting, during ETL. For example, it standardizes dates for consistent analysis.

24. Q: What is a NoSQL database?

 **A:** A NoSQL database handles unstructured or semi-structured data, like MongoDB for JSON. It's flexible for big data, unlike rigid SQL databases.

25. Q: What is a relational database?

 **A:** A relational database stores data in tables with relationships, like MySQL. It uses SQL for queries and ensures consistency with keys.

26. Q: What is data sharding?

 **A:** Data sharding splits a database across servers, like by customer ID ranges. It improves scalability and performance for large datasets.

27. Q: What is a data quality framework?

 **A:** A data quality framework ensures data accuracy, completeness, and consistency. It includes checks, like validating email formats, in pipelines.

28. Q: What is Apache Hive?

 **A:** Apache Hive queries big data in Hadoop using SQL-like syntax. It's used for analytics, like summarizing sales stored in HDFS.

29. Q: What is a data stream?

 **A:** A data stream is continuous real-time data, like website clicks. It's processed with tools like Kafka or Spark Streaming for instant insights.

30. Q: What is a batch job?

 **A:** A batch job processes data in chunks at scheduled times, like nightly sales reports. It's efficient for large, non-urgent datasets.

31. Q: What is data replication?

 **A:** Data replication copies data across systems, like mirroring a database for redundancy. It ensures availability and supports disaster recovery.

32. Q: What is a data model?

 **A:** A data model defines data structure and relationships, like tables in a database. It guides storage and querying, like a star schema.

33. Q: What is a cloud data warehouse?

 **A:** A cloud data warehouse, like Snowflake or BigQuery, stores and analyzes data online. It's scalable, cost-effective, and integrates with BI tools.

34. Q: What is data orchestration?

 **A:** Data orchestration manages pipeline tasks, like scheduling ETL jobs with Airflow. It ensures smooth data flow across systems and processes.

35. Q: What is a data sink?

 **A:** A data sink is the final destination in a pipeline, like a database or file. Data is loaded here after extraction and transformation.

36. Q: What is a data source?

 **A:** A data source provides raw data, like a database, API, or file. It's the starting point for ingestion in a data pipeline.

37. Q: What is data versioning?

 **A:** Data versioning tracks dataset changes, like snapshots of sales data. It supports reproducibility and auditing in analytics projects.

38. Q: What is a data lakehouse?

 **A:** A data lakehouse combines data lake flexibility with warehouse structure, like Delta Lake. It supports analytics and machine learning on unified data.

39. Q: What is a data integration platform?

 **A:** A data integration platform, like Talend, connects and transforms data from multiple sources. It streamlines ETL for analytics and reporting.

40. Q: What is a data retention policy?

 **A:** A data retention policy defines how long data is kept, like 7 years for financial records. It ensures compliance and manages storage costs.

Reporting and Documentation (40 Questions)

1. Q: What is a KPI in reporting?

 **A:** A Key Performance Indicator (KPI) measures performance, like sales revenue. It's tracked in dashboards to monitor progress toward business goals.

2. Q: Why is documentation important in data projects?

 **A:** Documentation records processes, data sources, and decisions, ensuring clarity. It helps teams troubleshoot, onboard, and reproduce work, reducing errors.

3. Q: What is a dashboard in reporting?

 **A:** A dashboard visualizes KPIs and metrics, like sales and profit, in one view. Tools like Tableau or Power BI create interactive dashboards for insights.

4. Q: What is a report specification?

 **A:** A report specification outlines report requirements, like data sources,

metrics, and format. It ensures the report meets stakeholder needs, like monthly sales summaries.

5. Q: What is data storytelling?

 **A:** Data storytelling presents insights compellingly, using visuals and narratives. For example, a dashboard with sales trends and annotations guides decision-making.

6. Q: What is a data glossary?

 **A:** A data glossary defines terms, like “revenue: total sales income.” It ensures consistent understanding across teams, reducing miscommunication.

7. Q: What is an executive summary in a report?

 **A:** An executive summary highlights key findings, like “sales grew 10%.” It’s a concise overview for decision-makers, placed at a report’s start.

8. Q: What is a metadata repository?

 **A:** A metadata repository stores data about data, like table schemas or update frequency. It aids discovery and management in reporting systems.

9. Q: What is a drill-down report?

 **A:** A drill-down report lets users explore details, like clicking sales by region to see cities. It’s interactive, built in tools like Power BI.

10. Q: What is a data quality report?

 **A:** A data quality report assesses metrics like accuracy or completeness, like missing customer IDs. It identifies issues for pipeline improvements.

11. Q: What is a scheduled report?

 **A:** A scheduled report runs automatically, like a weekly sales email. Tools like Power BI or Tableau Server deliver it to stakeholders regularly.

12. Q: What is a canned report?

 **A:** A canned report is pre-built with fixed parameters, like monthly revenue. It’s quick to access but less flexible than ad-hoc reports.

13. Q: What is an ad-hoc report?

 **A:** An ad-hoc report is created on-demand for specific questions, like sales for a new product. It’s tailored using tools like SQL or Tableau.

14. Q: What is a report template?

 **A:** A report template standardizes format and metrics, like a sales dashboard layout. It saves time and ensures consistency across reports.

15. Q: What is data lineage in reporting?

 **A:** Data lineage tracks data's path, like from source database to report. It ensures trust and helps debug errors in reporting pipelines.

16. Q: What is a data audit report?

 **A:** A data audit report reviews data usage and compliance, like access logs. It ensures adherence to policies, like GDPR, for security.

17. Q: What is a trend report?

 **A:** A trend report tracks metrics over time, like monthly sales growth. It uses charts to highlight patterns for strategic planning.

18. Q: What is a variance report?

 **A:** A variance report compares actual vs. expected metrics, like budget vs. spending. It identifies deviations for corrective actions.

19. Q: What is a summary report?

 **A:** A summary report condenses data into key insights, like total sales by region. It's concise, designed for quick decision-making.

20. Q: What is a detail report?

 **A:** A detail report lists granular data, like individual transactions. It's used for deep analysis, unlike summary reports with aggregated metrics.

21. Q: What is a data visualization best practice?

 **A:** Use clear visuals, avoid clutter, and choose appropriate charts, like bars for comparisons. Ensure colors are accessible and labels are readable.

22. Q: What is a report validation process?

 **A:** Report validation checks accuracy, like verifying sales totals against source data. It ensures reports are reliable before sharing with stakeholders.

23. Q: What is a performance report?

 **A:** A performance report tracks KPIs, like website traffic or sales targets. It evaluates progress and guides business decisions with data.

24. Q: What is a compliance report?

 **A:** A compliance report ensures adherence to regulations, like data privacy laws. It documents practices, like user data handling, for audits.

25. Q: What is a dashboard refresh rate?

 **A:** A dashboard refresh rate determines how often data updates, like daily or real-time. It's set in tools like Power BI for timely insights.

26. Q: What is a data annotation in reports?

 **A:** Data annotation adds notes to visuals, like explaining a sales spike. It clarifies insights for stakeholders in dashboards or reports.

27. Q: What is a report distribution list?

 **A:** A report distribution list specifies recipients, like managers getting weekly sales reports. It's managed in tools like Tableau Server for delivery.

28. Q: What is a snapshot report?

 **A:** A snapshot report captures data at a moment, like sales on a specific day. It's used for historical reference or compliance.

29. Q: What is a benchmark report?

 **A:** A benchmark report compares performance to standards, like sales vs. industry averages. It identifies areas for improvement or competitiveness.

30. Q: What is a report archive?

 **A:** A report archive stores historical reports, like past quarterly sales. It supports audits, compliance, or trend analysis over time.

31. Q: What is a data refresh in reporting?

 **A:** A data refresh updates report data, like pulling new sales from a database. Schedule it in tools like Power BI for current insights.

32. Q: What is a report filter?

 **A:** A report filter limits data, like showing sales for 2023. It's added in tools like Tableau to focus on relevant information.

33. Q: What is a report parameter?

 **A:** A report parameter lets users input values, like selecting a year for sales. It's set in tools like Power BI for dynamic reports.

34. Q: What is a cross-tab report?

 **A:** A cross-tab report summarizes data in a matrix, like sales by region and product. It's built in Excel or BI tools for comparisons.

35. Q: What is a report lifecycle?

 **A:** A report lifecycle covers creation, distribution, use, and archival. For example, a sales report is designed, shared monthly, and stored after use.

36. Q: What is a data export in reporting?

 **A:** A data export saves report data, like a dashboard to PDF or CSV. It's done in tools like Tableau for offline sharing or analysis.

37. Q: What is a report governance policy?

 **A:** A report governance policy sets rules for creation, access, and updates, like restricting sensitive data. It ensures consistency and security.

38. Q: What is a dynamic report?

 **A:** A dynamic report updates with user inputs or new data, like a dashboard with filters. It's built in BI tools for interactivity.

39. Q: What is a report audit trail?

 **A:** A report audit trail logs changes or access, like who viewed a sales dashboard. It ensures accountability and compliance in reporting.

40. Q: What is a report testing process?

 **A:** Report testing verifies functionality and accuracy, like checking dashboard calculations. It involves stakeholders to ensure reports meet requirements before deployment.