# Key concepts and Matching

# Key Concepts and Definitions for True/False and Matching Questions

---

# Table of Contents

---

# 1. Data Science Fundamentals

## Key Concepts:

- **Definition of Data Science**: An interdisciplinary field that uses scientific methods, processes, algorithms, and systems to extract knowledge and insights from structured and unstructured data.
- **Role of a Data Scientist**:
  - Collecting and cleaning data.

- Analyzing and interpreting data.
- Building predictive models.
- Communicating insights to stakeholders.
- **Data Science Lifecycle**:
    1. Problem Definition
    2. Data Collection
    3. Data Preparation
    4. Exploratory Data Analysis
    5. Modeling
    6. Evaluation
    7. Deployment
    8. Monitoring and Maintenance

## Possible True/False Statements:

- "Data Science is solely about building machine learning models." **(False)**
- "Exploratory Data Analysis is performed before data cleaning." **(False)**

## Matching Columns:

| Column A | Column B |
|---|---|
| Data Collection | Gathering data from various sources |
| Data Preparation | Cleaning and transforming data |
| Modeling | Building algorithms to predict outcomes |
| Deployment | Implementing models into production environments |

---

# 2. Data Collection and Preprocessing

## Key Concepts:

- **Data Collection Methods**:
    - Surveys
    - Web scraping
    - APIs
    - Sensors

- **Data Cleaning Techniques**:
    - Handling missing values (deletion, imputation)
    - Removing duplicates
    - Correcting inconsistencies
    - Outlier detection
- **Data Integration**: Combining data from multiple sources into a coherent dataset.
- **Data Transformation**:
    - Normalization
    - Encoding categorical variables
    - Aggregation

## Possible True/False Statements:

- "Data integration involves combining data from multiple sources." **(True)**
- "Normalization scales numerical data to have a mean of zero and standard deviation of one." **(True)**
- "One-hot encoding is used for encoding numerical variables." **(False)**

## Matching Columns:

| Column A | Column B |
| --- | --- |
| One-Hot Encoding | Encoding categorical variables into binary vectors |
| Imputation | Filling missing values with estimates |
| Standardization | Scaling data to have mean 0 and standard deviation 1 |
| Outlier Detection | Identifying data points significantly different |

# 3. Exploratory Data Analysis (EDA)

## Key Concepts:

- **Descriptive Statistics**:
    - Mean, Median, Mode
    - Variance, Standard Deviation
    - Skewness, Kurtosis
- **Visualization Techniques**:

- Histograms
- Scatter plots
- Box plots
- Heatmaps
- **Correlation Analysis**:
    - Pearson correlation coefficient
    - Spearman rank correlation

## Possible True/False Statements:

- "A histogram is used to display the distribution of a categorical variable." **(False)**
- "The Pearson correlation coefficient measures linear correlation between two variables." **(True)**
- "A box plot can help identify outliers in the data." **(True)**

## Matching Columns:

| Column A | Column B |
|---|---|
| Scatter Plot | Relationship between two continuous variables |
| Box Plot | Visual representation of the distribution's quartiles |
| Heatmap | Visualization of the correlation matrix |
| Skewness | Measure of the asymmetry of the probability distribution |

# 4. Statistical Foundations

## Key Concepts:

- **Measures of Central Tendency**:
    - Mean: Average value.
    - Median: Middle value in ordered data.
    - Mode: Most frequent value.
- **Measures of Dispersion**:
    - Range
    - Variance
    - Standard Deviation

- **Probability Distributions**:
    - Normal Distribution
    - Binomial Distribution
    - Poisson Distribution
- **Hypothesis Testing**:
    - Null hypothesis (H0)
    - Alternative hypothesis (H1)
    - p-value
    - Significance level (α)

## Possible True/False Statements:

- "Standard deviation is the square root of variance." **(True)**
- "In a normal distribution, the mean, median, and mode are equal." **(True)**
- "A p-value less than 0.05 typically leads to rejecting the null hypothesis at the 5% significance level." **(True)**

## Matching Columns:

| Column A | Column B |
| --- | --- |
| Null Hypothesis (H0) | Statement of no effect or no difference |
| Alternative Hypothesis | Statement that there is an effect or a difference |
| Type I Error | Rejecting a true null hypothesis (false positive) |
| Type II Error | Failing to reject a false null hypothesis (false negative) |

# 5. Data Visualization

## Key Concepts:

- **Principles of Effective Visualization**:
    - Clarity
    - Accuracy
    - Efficiency
    - Aesthetics
- **Visualization Tools and Libraries**:

- Matplotlib
- Seaborn
- Plotly
- Tableau
- **Chart Types and Their Uses**:
  - Line Chart: Trends over time
  - Bar Chart: Comparing categories
  - Pie Chart: Proportions of a whole
  - Heatmap: Correlation between variables

## Possible True/False Statements:

- "Seaborn is built on top of Matplotlib and provides a higher-level interface for statistical graphics." **(True)**
- "A pie chart is the best choice for showing trends over time." **(False)**
- "Heatmaps are useful for visualizing correlation matrices." **(True)**

## Matching Columns:

| Column A | Column B |
|---|---|
| Line Chart | Displaying data trends over intervals |
| Bar Chart | Comparing quantities across categories |
| Scatter Plot | Showing relationship between two variables |
| Histogram | Displaying the distribution of a dataset |

# 6. Machine Learning Basics

## Key Concepts:

- **Supervised Learning**: Learning from labeled data (input-output pairs).
- **Unsupervised Learning**: Finding patterns in unlabeled data.
- **Reinforcement Learning**: Learning optimal actions through rewards and penalties.
- **Overfitting**: When a model learns the training data too well, including noise, and performs poorly on new data.

- **Underfitting**: When a model is too simple and cannot capture the underlying trend of the data.

## Possible True/False Statements:

- "In supervised learning, the model learns from unlabeled data." **(False)**
- "Overfitting occurs when a model performs well on training data but poorly on unseen data." **(True)**
- "K-Means clustering is a supervised learning algorithm." **(False)**

## Matching Columns:

| Column A | Column B |
|---|---|
| Overfitting | Model fits training data too closely |
| Underfitting | Model is too simple to capture data patterns |
| Supervised Learning | Uses labeled data for training |
| Unsupervised Learning | Finds patterns without labeled responses |

# 7. Data Wrangling and Feature Engineering

## Key Concepts:

- **Data Wrangling**: Process of cleaning and unifying complex data sets for easy access and analysis.
- **Feature Engineering**: Creating new input features from existing ones to improve model performance.
- **Feature Selection**: Selecting the most relevant features to use in model construction.
- **Encoding Categorical Variables**:
  - One-Hot Encoding
  - Label Encoding

## Possible True/False Statements:

- "Feature engineering involves creating new features from existing data." **(True)**
- "Label encoding converts categorical variables into binary vectors." **(False)**
- "Feature selection can help reduce overfitting." **(True)**

## Matching Columns:

| Column A | Column B |
|---|---|
| One-Hot Encoding | Converts categories into binary columns |
| Feature Selection | Choosing relevant features for the model |
| Data Wrangling | Cleaning and unifying complex data sets |
| Principal Component Analysis | Dimensionality reduction technique |

---

# 8. Big Data Technologies

## Key Concepts:

- **Big Data Characteristics (5 V's)**:
    - Volume
    - Velocity
    - Variety
    - Veracity
    - Value
- **Hadoop Ecosystem Components**:
    - HDFS (Hadoop Distributed File System)
    - MapReduce
    - YARN (Yet Another Resource Negotiator)
- **Apache Spark**:
    - In-memory data processing
    - Supports batch and real-time analytics
- **NoSQL Databases**:
    - MongoDB
    - Cassandra
    - HBase

## Possible True/False Statements:

- "HDFS stands for Hadoop Distributed File System." **(True)**
- "Apache Spark is only capable of batch processing, not real-time analytics." **(False)**
- "NoSQL databases are designed to handle structured data only." **(False)**

## Matching Columns:

| Column A | Column B |
| --- | --- |
| HDFS | Distributed file system in Hadoop |
| MapReduce | Programming model for processing large data sets |
| Apache Spark | In-memory data processing framework |
| NoSQL Database | Non-relational database designed for large data volumes |

---

# 9. Ethics and Data Privacy

## Key Concepts:

- **Ethical Principles in Data Science**:
  - Privacy
  - Fairness
  - Transparency
  - Accountability
  - Security
- **Bias in Data and Models**:
  - Selection Bias
  - Confirmation Bias
  - Algorithmic Bias
- **Privacy Laws and Regulations**:
  - GDPR (General Data Protection Regulation)
  - CCPA (California Consumer Privacy Act)
  - HIPAA (Health Insurance Portability and Accountability Act)

## Possible True/False Statements:

- "Algorithmic bias can occur when training data is not representative of the population." **(True)**
- "GDPR is a privacy regulation enforced in the European Union." **(True)**
- "Anonymizing data completely removes the risk of re-identification." **(False)**

## Matching Columns:

| Column A | Column B |
| --- | --- |
| GDPR | European Union data protection regulation |
| Algorithmic Bias | Systematic errors in a computer system |
| Transparency | Openness in methods and decision-making processes |
| Informed Consent | Obtaining permission with full disclosure |

# 10. Data Mining Techniques

## Key Concepts:

- **Classification**: Predicting categorical labels.
- **Regression**: Predicting continuous values.
- **Clustering**: Grouping similar data points without pre-defined labels.
- **Association Rule Mining**: Discovering interesting relations between variables in large databases.
- **Anomaly Detection**: Identifying unusual data points.

## Possible True/False Statements:

- "K-Means is a clustering algorithm used for classification tasks." **(False)**
- "Association rule mining is commonly used in market basket analysis." **(True)**
- "Anomaly detection is used to find outliers in the data." **(True)**

## Matching Columns:

| Column A | Column B |
| --- | --- |
| Classification | Predicting discrete labels |
| Regression | Predicting continuous numerical values |
| Clustering | Grouping data without labels |
| Association Rule Mining | Finding relationships between variables |

# 11. Supervised Learning Algorithms

# Key Concepts:

- **Naïve Bayes Classifier**:
  - Based on Bayes' Theorem with an assumption of feature independence.
- **Support Vector Machine (SVM)**:
  - Finds the optimal hyperplane that separates classes.
- **Decision Trees**:
  - Tree-like model of decisions and their possible consequences.
- **Random Forest**:
  - Ensemble method using multiple decision trees.
- **Ensemble Learning**:
  - Combining predictions from multiple models.

## Possible True/False Statements:

- "Naïve Bayes assumes that all features are dependent on each other." **(False)**
- "Random Forest reduces overfitting compared to individual decision trees." **(True)**
- "SVM can only perform linear classification." **(False)**

## Matching Columns:

| Column A | Column B |
| --- | --- |
| Naïve Bayes | Probabilistic classifier with independence assumption |
| Support Vector Machine | Classifier that finds optimal separating hyperplane |
| Random Forest | Ensemble of decision trees |
| Ensemble Learning | Combining multiple models to improve performance |

---

# 12. Unsupervised Learning and Deep Learning

## Key Concepts:

- **Unsupervised Learning**:
  - Works with unlabeled data to find hidden patterns.
- **Deep Learning**:
  - Neural networks with multiple layers.

- **Convolutional Neural Networks (CNN)**:
  - Specialized for processing grid-like data (e.g., images).
- **Recurrent Neural Networks (RNN)**:
  - Designed for sequential data.
- **Long Short-Term Memory Networks (LSTM)**:
  - Type of RNN that can learn long-term dependencies.

## Possible True/False Statements:

- "Autoencoders are used for unsupervised learning tasks." **(True)**
- "CNNs are primarily used for natural language processing tasks." **(False)**
- "LSTMs are a type of RNN that can handle long-term dependencies." **(True)**

## Matching Columns:

| Column A | Column B |
| --- | --- |
| CNN | Neural network for image processing |
| RNN | Neural network for sequential data |
| LSTM | RNN variant handling long-term dependencies |
| Autoencoder | Neural network used for unsupervised learning |

# 13. Natural Language Processing (NLP)

## Key Concepts:

- **Tokenization**: Splitting text into words or sentences.
- **Sentiment Analysis**: Determining the emotional tone behind a body of text.
- **Lexicons in NLP**:
  - VADER (Valence Aware Dictionary for Sentiment Reasoning)
  - SentiWordNet
  - AFINN
- **Bag of Words**: Representing text as the frequency of words.
- **Part-of-Speech Tagging**: Assigning grammatical categories to words.

## Possible True/False Statements:

- "Tokenization involves combining multiple words into a single token." **(False)**
- "VADER is a lexicon and rule-based sentiment analysis tool." **(True)**
- "Bag of Words model considers the order of words in a sentence." **(False)**

## Matching Columns:

| Column A | Column B |
|---|---|
| Tokenization | Splitting text into smaller units |
| Sentiment Analysis | Determining emotional tone of text |
| Bag of Words | Text representation based on word frequency |
| Part-of-Speech Tagging | Assigning grammatical categories to words |

---

# 14. Explainable AI (XAI) and Large Language Models (LLMs)

## Key Concepts:

- **Explainable AI (XAI)**: Techniques that make the output of AI models understandable to humans.
- **LIME (Local Interpretable Model-Agnostic Explanations)**:
  - Explains individual predictions.
- **SHAP (SHapley Additive exPlanations)**:
  - Uses game theory for feature attribution.
- **Large Language Models (LLMs)**:
  - AI models trained on large text datasets (e.g., GPT, BERT).

## Possible True/False Statements:

- "LIME is specific to neural network models only." **(False)**
- "SHAP values can be used to interpret the contribution of each feature to the prediction." **(True)**
- "Large Language Models are only used for text generation." **(False)**

## Matching Columns:

| Column A | Column B |
| --- | --- |
| XAI | Making AI decisions understandable |
| LIME | Local explanations for model predictions |
| SHAP | Feature attribution using Shapley values |
| LLM | Large models trained on vast text data |

# 15. Ethics in Data Science

## Key Concepts:

- **Bias Types**:
    - Selection Bias
    - Measurement Bias
    - Algorithmic Bias
- **Privacy and Consent**:
    - Data ownership
    - Informed consent
- **Ethical Principles**:
    - Fairness
    - Transparency
    - Accountability
    - Security
- **Case Studies**:
    - Cambridge Analytica Scandal
    - Amazon's Biased Hiring Algorithm
    - Predictive Policing Issues

## Possible True/False Statements:

- "Algorithmic bias cannot be mitigated once the model is deployed." **(False)**
- "Informed consent requires that individuals are aware of how their data will be used." **(True)**
- "The Cambridge Analytica scandal involved the misuse of personal data for political advertising." **(True)**

## Matching Columns:

| Column A | Column B |
| --- | --- |
| Algorithmic Bias | Systematic errors due to biased data or assumptions |
| Informed Consent | Individuals agree to data use with full understanding |
| Fairness | Ensuring equitable treatment in models |
| Predictive Policing | Using data to forecast criminal activity |

---

# Additional Concepts for True/False and Matching

## Data Types:

- **Structured Data**: Data that is organized in a fixed format (e.g., tables).
- **Unstructured Data**: Data without a predefined data model (e.g., text, images).

## Machine Learning Terms:

- **Hyperparameters**: Configuration settings used to tune how models learn.
- **Cross-Validation**: Technique for assessing how a model will generalize to an independent dataset.

## Possible True/False Statements:

- "Cross-validation helps in reducing overfitting." **(True)**
- "Hyperparameters are learned during the training of the model." **(False)**

## Matching Columns:

| Column A | Column B |
| --- | --- |
| Hyperparameters | Settings configured before training |
| Overfitting | Model fits training data too closely |
| Cross-Validation | Technique for model validation |
| Structured Data | Data in a fixed format |

---

---

# Tips for True/False Questions:

- Read the statement carefully; look out for absolutes like "always" or "never."
- Consider the definitions and key characteristics of concepts.
- Think about exceptions to general rules.

# Tips for Matching Questions:

- Understand each term and its definition.
- Eliminate options that are clearly incorrect to narrow down choices.
- Be aware of similar-sounding terms and avoid confusion.

---

By familiarizing yourself with these key concepts and their correct associations, you'll be well-prepared to tackle true/false and matching questions on your exam. Good luck with your studies!