**Probability Distributions**: Probability distributions describe the probability of different outcomes of a random variable. Common probability distributions used in data science include the normal distribution, binomial distribution, and Poisson distribution. Understanding probability distributions is important for modeling and analyzing data.

**Statistical Measures**: Statistical measures are used to summarize and analyze data. These can include measures of central tendency such as mean, median, and mode, measures of dispersion such as variance and standard deviation, and measures of relationship such as correlation and covariance. These measures help to describe the characteristics of data and identify patterns or trends.

**Data Visualization for Descriptive Statistics**: Data visualization techniques, such as harchards, histogram, scatter plots, and box plots, can be upply evisually represent descriptive statistics and gain insights from data. Visualizations can provide a clear and intuitive way to understand the distribution, variability, and relationships within data.

## MACHINE LEARNING ALGORITHMS

## **Overview of Popular Machine Learning Algorithms**

Overview of popular machine learning algorithms, categorized into three main types: supervised learning, unsupervised learning, and reinforcement learning.

#### **Supervised Learning**

Supervised learning algorithms are trained on labeled data, where the algorithm learns to make predictions or decisions based on input-output pairs. The most common types of supervised learning algorithms are: a) Linear Regression: A statistical algorithm that models the relationship between a dependent variable and one or more it dependent variables as a

b) Decision Trees: Tree-Gased algorithmit that can be used for both classificate hand regression fusces. They make decisions based on feature values to predict the target variable.

c) Random Forests: An ensemble algorithm that combines multiple decision trees to make more accurate predictions.

Support Vector Machines (SVM): A powerful algorithm for d) classification and regression tasks that finds a hyperplane that best separates data into different classes.

e) K-Nearest Neighbors (KNN): A simple algorithm that classifies new data points based on the majority vote of their k-nearest neighbors in the training data.

a) **Q-Learning**: A model-free algorithm that learns to make optimal decisions by updating its action-value function based on the rewards it receives from the environment.

b) **Deep Q-Networks (DQNs)**: Q-learning algorithms that use deep neural networks to approximate the action-value function for handling high-dimensional state spaces.

c) Policy Gradient: An algorithm that directly learns a policy, or a mapping from states to actions, by optimizing the policy parameters using gradient ascent.

d) Actor-Critic: An algorithm that combines elements of both value-based and policy-based methods, where an agent has a poicy (actor) and a value function (critic) that are updated included. from age 23 of 70

23

experimentation and debugging of deep learning models. PyTorch is known for its flexibility and ease of use, making it popular among researchers and practitioners in the field of deep learning. PyTorch also supports automatic differentiation, which makes it easy to compute gradients for optimizing model parameters. PyTorch has gained significant popularity in recent years due to its strong support for deep learning research and its adoption by the research community.

Preview from Notesale.co.uk Page 28 of 70

**Apache Hive**: Apache Hive is an open-source data warehouse system that provides SQL-like query capabilities on top of large datasets stored in Hadoop HDFS. Hive translates SQL-like queries into MapReduce or Spark jobs, making it easy to analyze big data using familiar SQL syntax. Hive also supports custom user-defined functions (UDFs) and supports data partitioning, bucketing, and indexing for optimized query performance. Hive is widely used for ad-hoc queries, data analysis, and data warehousing in big data environments.

Preview from Notesale.co.uk Page 30 of 70

## DATA MINING AND TEXT MINING

#### **Techniques for Mining Structured and Unstructured Data**

Mining structured and unstructured data involves extracting valuable insights and patterns from data in order to gain meaningful insights and make informed decisions. Here are some common techniques used for mining structured and unstructured data:

Association Rules: Association rules mining is a technique used to discover interesting relationships or patterns in large datasets. It involves identifying co-occurring items or events in a dataset and uncovering patterns such as "If A occurs, then B is likely to occur as well." Association rules are commonly used in market basket analyse. Where the goal is to identify products that are often purchasily together. The Apriori algorithm and the FP-Growth algorithm are commonly used algorithms for association rules mining a **Q** 

**Clustering**: Clustering is a technique used to group similar objects together based on their similarities or dissimilarities in a dataset. Clustering algorithms aim to identify natural groupings or clusters in data without any predefined labels or categories. Clustering is commonly used in customer segmentation, image segmentation, and anomaly detection. Popular clustering algorithms include k-means, hierarchical clustering, and DBSCAN.

**Text Mining**: Text mining is a technique used to extract valuable information and insights from unstructured text data. Text data can include various types of information, such as customer reviews, social media posts,

# DATA INTEGRATION AND ETL (EXTRACT, TRANSFORM, LOAD)

Data integration and ETL (Extract, Transform, and Load) are critical processes in the field of data analytics and data management. Here's an overview of data integration and ETL:

**Data Integration**: Data integration is the process of combining data from different sources and making it available for analysis or other data processing tasks. It involves gathering data from multiple sources, such as databases, data warehouses, APIs, flat files, and external sources, and consolidating it into a unified format. Data integration can involve different types of data, including structured, semi-structured and unstructured data, and it may require data cleansing, validation, and transformation to ensure data quality and consistency.

**ETLATERAL, TransformQCad**: ETL is a common data integration process that involves extracting data from different sources, transforming it into a standardized format, and loading it into a target system, such as a data warehouse or a data lake. ETL is typically used in batch processing scenarios, where data is extracted from source systems periodically, transformed to meet the target system's requirements, and loaded into the target system for analysis or other data processing tasks.

The ETL process typically consists of the following steps:

**Extraction**: In this step, data is extracted from different sources, such as databases, APIs, flat files, and external sources. Data extraction may

**Data documentation**: Maintain comprehensive documentation of data integration processes, data mappings, data validation rules, and data lineage. Document data sources, data transformations, and data flow for easy reference and troubleshooting. This documentation serves as a reference for data integration developers, data analysts, and other stakeholders.



CCPA provides consumers with rights such as the right to know what personal information is collected, the right to opt-out of the sale of personal information, the right to delete personal information, and the right to nondiscrimination for exercising their privacy rights. CCPA also imposes requirements on businesses, such as providing clear privacy notices, maintaining a "Do Not Sell My Personal Information" link on their websites, and implementing reasonable security measures to protect data. Non-compliance with CCPA can result in fines and penalties, as well as potential lawsuits from consumers.

Apart from GDPR and CCPA, there are other data privacy regulations and frameworks in different jurisdictions, such as the Personal Information Protection and Electronic Document Act (PIPEDA) in Canada, the Personal Data Protection Act (PDPA) in Origapore, and the Privacy Act in Australia, anong others. These regulations are aimed at protecting the privacy rights of individuals and promoting responsible data handling practices by organizations. visualizations. This allows for real-time updates and interactive behavior based on changes in the data or user interactions.

**Interactivity**: D3.js provides extensive support for interactivity, allowing users to add interactive elements such as tooltips, hover effects, and zooming to their visualizations. Users can also add event listeners to respond to user interactions and update visualizations accordingly.

Customization: D3.js provides a high level of customization, allowing users to create unique visualizations with their own styles, colors, and layouts. Users have full control over the visual elements and can create custom animations, transitions, and effects to enhance the interactivity and user experience.

63

- 2. Work on a capstone project that involves acquiring, cleaning, exploring, analyzing, and visualizing data to derive meaningful insights.
- 3. Present and communicate the results of the capstone project to demonstrate proficiency in data science and big data analytics skills.

#### **Review and Exam Preparation**

- 1. Review the key concepts, techniques, and tools covered in the course.
- 2. Prepare for exams and assessments by reviewing course materials, notes, and assignments.
- 3. Practice solving data science problems and applying the larned skills through sample exercises and practice exercise.

# Optional Topics (Based on Course Curriculum)

Depending on the specific course curriculum, there may be additional optional topics that can be covered, such as natural language processing, time series analysis, data visualization with D3.js, deep learning, or advanced topics in big data analytics.