Assignment

Solve 10 questions.
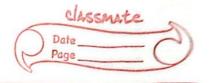
1.

4. Discuss: Statistics vs Data Mining vs Data Analytics vs. Data science.

→

Statistics

Statistics is the study of collecting, analyzing, interpreting, presenting and organizing data. It forms the foundation for many analytical methodologies and techniques. Traditional statistics focuses on hypothesis, testing, probability, regression and variance analysis to infer patterns and relationships within data. It's primary aim is to understand & explain the underlying phenomena represented by the data through well-defined mathematical principles & models.

Data mining

It involves exploring and analyzing large datasets to discover patterns, trends and relationships that might not to immediately apparent. It employs algorithms and techniques from machine learning statistics, and database systems to uncover hidden information. Data mining is heavily used in areas like market research, fraud detection and customer relationship management to generate insights that can inform decision-making processes.

Data Analytics

It is the process of examining datasets to draw conclusions about the information they contain. It encompasses a range of methods from descriptive analytics (summarizing past data) to predictive
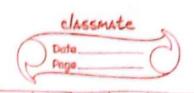
analytics (forecasting future events) and prescriptive analytics (recommending actions). Data analytics focuses on transforming raw data into actionable insights through the use of various tools & techniques, including statistical analysis, machine learning and data visualization

## Data Science

Data Science is an interdisciplinary field that combines elements of statistics, computer science and domain specific knowledge to extract meaningful insights and knowledge from data. It involves the entire data lifecycle, including data collection, cleaning, preprocessing, analysis and visualization. Data scientist use advanced techniques like machine-learning, DL and AI to build predictive model and automate decision making processes. Data science aims to solve complex problems and provide strategic insights that drive innovation & efficiency across various industries.

In summary, while statistics provide the theoretical foundation for understanding data, data mining focuses on discovering patterns within large datasets. Data analytics applies these techniques to generate actionable insights and data science integrates those approaches with advanced computational tools to solve complex problems and drive decision making.

2.

5. What is Supervised Learning? Give concrete examples of Regression and Classification.

→

Supervised Learning is a machine learning technique where a model is trained on a labeled dataset, learning to map input features to output labels. This enables model to make predictions on new, unseen data.

The two types of supervised learning are :-

| Regression | Classification |
|---|---|
| - Predicts a continuous numerical output based on input features. | Predicts a categorical output (class or group) based on input features |
| - In chemical engineering, used in predicting the yield of chemical reaction based on temperature, pressure and reactant concentrations. | In chemical engineering, classifying different types of based on their chemical composition. |
| - Usually used for predicting house price based on square footage, number of bedroom & location | Classifying emails as spam or not spam |

Supervised learning is valuable in chemical engineering for developing predictive models for processes like reaction yields, product quality or equipment fault detection.

**Q. Define Unsupervised learning. Give examples of Clustering and Anomaly Detection**

→

Unsupervised Learning is a machine learning technique where a model is trained on an unlabeled dataset, meaning the data points lack explicit output labels. The model aims to discover hidden patterns, structures or relationships within the data. Common techniques in unsupervised learning include clustering and anomaly detection.

| Clustering | Anomaly Detection |
|---|---|
| - Involves grouping similar data points together based on their inherent characteristics. This model identifies clusters of data points that share common features or patterns. | Involves identifying unusual or rare data points that deviate significantly from the majority of the data these anomalies may present outliers, errors, or events of interest. |
| Examples: | Examples: |
| - Customer Segmentation | - Fraud Detection |
| - Document Clustering | - Network Security |
| - Process Optimization in on chemical plant | - Anomaly Detection in Equipment Monitoring |

Unsupervised learning focuses on finding hidden patterns in unlabeled data. Clustering groups similar data points into clusters, with examples including customer segmentation and document clustering. Anomaly detection identifies unusual data points, with application such as fraud detection and network security.

4.

7. What is Reinforcement Learning? Provide example of Markov decision process. Q-learning & Monte-Carlo methods.

→

Reinforcement Learning (RL) is a type of machine learning where an agent learns to make decisions by interacting with an environment. The agent takes actions in the environment, and the environment provides feedback in the form of rewards or penalties. The agent's is to learn a policy that maximize its cumulative award over time.
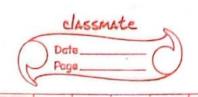
Components of RL:
- Agent - learner & decision-maker.
- Environment - The external system the agent interacts with.
- State: A representation of the current state in the environment.
- Action: A choice made by the agent that affects the environment.
- Reward: Feedback from the environment indicating the goodness or badness of the agent's action.

# Examples In Reinforcement Learning

## 1. Markov Decision Process (MDP)

A mathematical framework for modeling decision-making in situations where outcomes are partly random and partly under the control of a decision maker. In an MDP, the agent's decision-making process is divided into discrete time steps. At each time step, the agent is in a state, takes an action, receives a reward. The agent's goal is to find a policy (a mapping from states to actions) that maximizes the expected cumulative reward over time.

## 2. Q-learning

A model-free reinforcement learning algorithm that learns the value of taking actions in different states. The agent maintains a Q-table, which stores the expected future reward for each state-action pair. The agent updates the Q-table based on the rewards it receives, gradually learning the optimal policy.

## 3. Monte-Carlo methods.

A class of reinforcement learning algorithms that learn from complete episodes of interactions with the environment. The agent interacts with the environment until it reaches a terminal state, then updates its policy based on the cumulative reward it received during the episode.

RL can be applied to optimize control strategies for chemical processes. For example, an RL agent can learn to control the temp" & pressure of a reactor to yield or minimize energy consumption. The agent can also learn to schedule batch processes or optimize the operation of a distillation column.

5.

8. Discuss the importance of NumPy, SciPy, Matplotlib and Pandas in Data Analysis.

→

## NumPy

Numerical Python (NumPy) provides the foundation for numerical computations in Python. It introduces powerful array objects, enabling efficient storage and manipulation of large datasets. In chemical engineering, this is crucial for handling experimental data, simulation results, and model parameters.

It offers a wide range of mathematical functions optimized for array operations. These functions are essential for implementing and solving mathematical models in chemical engineering, such as mass and energy balances, reaction kinetics and transport phenomena.

## SciPy (Scientific Python)

It builds upon NumPy and provides additional functionality for te scientific and technical computing. It includes modules for optimization, linear algebra, integration, interpolation and signal processing and more. It is important in solving complex

engineering problems, analyzing experimental data, and optimizing process parameters.

Scipy's <u>odeint</u> function is particulary useful in ChE for solving ODEs that arise in dynamic models of chemical reactors, separation processes and other time dependent systems.

## Matplotlib

It is a versatile plotting library that allows for the creation of high-quality static, animated and interactive indivisualization. In ChE, effective visualization is essential for understanding complex data patterns, communicating results and gaining insights into process behavior.
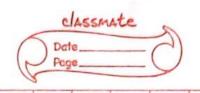
It offers extensive customization options enabling engineers to create plots tailored to their specific needs. This flexibility is valuable for presenting results in technical reports, publications & presentations.

## Pandas

Pandas provides data structures and data analysis tools designed to work with structural data (e.g tables, time series). It simplifies tasks like data cleaning, filtering, transforming & aggregation, which are common in chemical engineering data analysis.

Panda integrates went with other libraries like Numpy and Matplotlib, making it powerful tool for analyzing experimental data, exploring relationships between variables and preparing data for modeling & simulation

In summary, NumPy, SciPy, Matplotlib and Pandas form a robust ecosystem for data analysis in chemical engineering. They provide the essential tools for numerical computation, scientific computing, data visualization and data manipulation, enabling engineerings to effectively analyze data, develop and validate models and optimise processes.

6.

9. List some of the Machine learning Core libraries in Python.

→

- Scikit - learn

A comprehensive library for various machine learning algorithms including classification, regression, clustering and dimensionality reduction.

- TensorFlow

A powerful open-source library developed by Google for building & training deep neural networks.

- keras

A high level neural networks API that runs on top of TensorFlow, Theano or CNTK, providing a user-friendly interface for building & training deep learning models.

- PyTorch

An open-source machine learning library known for its dynamic computational graph & flexibility in building complex models.

These libraries provide a wide range of tools &

functionalities for implementing & evaluating machine learning models in Python, making them essential for data analysis & modeling tasks in chemical engineering