

Assignment -4

Machine learning.

1. What is clustering in machine learning?

Clustering is an **unsupervised learning technique** used **to group data points into clusters such that points in the same cluster are more similar** to each other than to those in other clusters.

It is used for **pattern recognition, exploratory data analysis, and data compression.**

2. Difference between supervised and unsupervised clustering?

Clustering is inherently an **unsupervised learning** task. However, the comparison between supervised and unsupervised techniques can be explained:

Supervised clustering involves **labelled data** where the **goal is to form clusters that match the provided labels.** This is rarely used and overlaps with classification.

Unsupervised clustering is the standard approach and works with **unlabelled data to identify inherent structures or patterns.**

3. Key applications of clustering algorithms?

Market segmentation

Image segmentation

Document categorization

Anomaly detection (e.g., fraud detection)

Social network analysis

Recommender systems

4. Describe the K-means clustering algorithm?

- **Initialization:** Choose k random centroids:
 - (a) random datapoint.
 - (b) centroid of all datapoints.
 - (c) any randomly new datapoints.
- **Assignment:** Assign each data point to the nearest centroid.
- **Update:** Recalculate centroids as the mean of all points assigned to them.
- **Repeat:** Continue the assignment and update steps until centroids stabilize or a maximum number of iterations is reached.

5. Advantages and disadvantages of K-means clustering?

Advantages:

- Simple and easy to implement.
- Computationally efficient.
- Works well on spherical clusters.

Disadvantages:

- Requires pre-specifying k.
- Sensitive to initialization and outliers.
- Poor performance on non-spherical or overlapping clusters.

6. How does hierarchical clustering work?

Hierarchical clustering builds a hierarchy of clusters using either:

- **Agglomerative approach** (bottom-up): Each data point starts as its own cluster, and clusters are merged iteratively based on similarity.
- **Divisive approach** (top-down): All data points start in one cluster, and splits are made iteratively.

7. Linkage criteria in hierarchical clustering?

- **Single Linkage**: Minimum distance between points in two clusters.
- **Complete Linkage**: Maximum distance between points in two clusters.
- **Average Linkage**: Average distance between all points in two clusters.
- **Centroid Linkage**: Distance between cluster centroids.

8. Explain the concept of DBSCAN clustering:

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) groups data points based on density. It identifies core points (dense regions) the datapoints have at least minimum number of datapoints 4, in its radius expands clusters around these points, and treats points in sparse regions as noise, border datapoint, outlier.

9. Parameters in DBSCAN clustering?

- **Epsilon (ϵ)**: The maximum distance between two points for them to be considered neighbors.

- **MinPts:** Minimum number of points required to form a dense region (core point).

10. Evaluating clustering algorithms?

Clustering evaluation focuses on how well the algorithm captures the structure of the data:

- **Internal Measures:** Evaluate within-cluster cohesion and between-cluster separation (e.g., silhouette score).
- **External Measures:** Compare results against ground truth (e.g., purity, adjusted Rand index).

11. What is the silhouette score, and how is it calculated?

The silhouette score measures how similar a point is to its cluster compared to other clusters:

$$S_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$

Where:

- a_i = average distance to points in the same cluster.
- b_i = average distance to points in the nearest cluster. Values range from -1 to 1, with higher values indicating better clustering.

12. Challenges of clustering high-dimensional data?

- **Curse of dimensionality:** Distance measures become less meaningful in high dimensions.
- **Sparsity:** Data points are more spread out, making it hard to define clusters.
- **Computational cost:** High-dimensional clustering requires more resources.

13. Explain the concept of density-based clustering

Density-based clustering identifies clusters as dense regions (CORE POINTS) of data points separated by sparser regions (BORDER POINTS). It is effective for irregularly shaped clusters and outlier detection.

14. How does Gaussian Mixture Model (GMM) clustering differ from K-means?

- **K-means:** Hard clustering where each point belongs to a single cluster.

- **GMM:** Soft clustering where each point has a probability of belonging to multiple clusters.
- GMM assumes data is generated from a mixture of Gaussian distributions.

15. Limitations of traditional clustering algorithms

- Difficulty handling non-spherical or overlapping clusters (e.g., K-means).
- Sensitivity to outliers and noise (e.g., K-means, hierarchical clustering).
- Scalability issues with large datasets (e.g., hierarchical clustering).
- Requires specifying hyperparameters like k or ϵ (e.g., K-means, DBSCAN).

16. Applications of Spectral Clustering

Spectral clustering is a graph-based clustering method that uses eigenvalues of a similarity matrix to identify clusters. Applications include:

- **Image segmentation:** Partitioning an image into distinct regions.
- **Social network analysis:** Detecting communities or groups in social networks.
- **Recommendation systems:** Grouping similar users or items.
- **Biological data:** Identifying gene groups or protein interactions.
- **Anomaly detection:** Detecting outliers in high-dimensional data.

17. Concept of Affinity Propagation

Affinity propagation is a clustering algorithm that identifies "exemplars" (data points that represent clusters). Unlike K-means, it does not require the number of clusters to be specified upfront:

1. **Input similarity matrix:** Defines pairwise similarities between data points.
2. **Message exchange:** Messages are iteratively exchanged between points to identify exemplars.
3. **Cluster formation:** Points are grouped around their exemplars.

Advantages:

- Automatically determines the number of clusters.
- Handles non-convex clusters well.

18.. Managing Categorical Variables in Clustering

Clustering categorical data involves using distance or similarity measures suited for non-numeric data:

- **Encoding techniques:**
 - One-hot encoding
 - Label encoding
- **Distance measures:**
 - Hamming distance: For binary variables.
 - Jaccard similarity: For sets or binary attributes.
 - Gower's distance: A combination for mixed data types.
- **Specialized algorithms:**
 - K-modes: Extends K-means for categorical data.
 - CLARA or PAM: Adapted for categorical and mixed data types.

19. Elbow Method for Determining the Optimal Number of Clusters

The elbow method identifies the ideal number of clusters (k) by plotting the **inertia** (within-cluster sum of squares) against different k values:

1. **Plot the graph:** Number of clusters (k) on the x-axis and inertia on the y-axis.
2. **Find the 'elbow point':** The point where adding more clusters results in a minimal reduction in inertia, resembling an elbow.

This point suggests the best trade-off between variance reduction and model complexity.

#to select K(centroid)>> elbow method

wcss = []

```
for k in range(1, 10):  
    kmeans = KMeans(n_clusters = k)  
    kmeans.fit(X_train)  
    wcss.append(kmeans.inertia_)  
  
wcss
```

```
#plot an elbow curve  
plt.plot(range(1, 10), wcss)  
plt.xticks(range(1, 10))  
plt.xlabel("No of clusters")  
plt.ylabel("wcss")  
plt.show()
```

20. Emerging Trends in Clustering Research

1. **Deep clustering:** Combining deep learning with clustering for feature extraction and representation learning.
2. **Scalable clustering:** Designing algorithms that handle massive datasets efficiently.
3. **Multi-view clustering:** Integrating data from multiple perspectives or modalities.
4. **Clustering with uncertainty:** Algorithms that consider probabilistic models and soft assignments.
5. **Explainable clustering:** Developing methods to make clustering results interpretable.
6. **Graph neural networks (GNNs):** Using GNNs for graph-based clustering.
7. **Federated clustering:** Clustering distributed data without centralizing it, maintaining privacy and scalability.

21. What is anomaly detection, and why is it important?

Anomaly detection involves **identifying unusual data points**, **patterns**, or **events** that **differ significantly** from **the majority of the data**. It is essential because anomalies often indicate critical situations such as:

- Fraudulent transactions.
- System failures.
- Security breaches.
- Medical conditions.

22. Types of anomalies encountered in anomaly detection

- **Point anomalies:** **Single data points** that are **significantly different** (e.g., a large transaction in a small-budget account).
- **Contextual anomalies:** **Data points** that are **unusual** in a **specific context** (e.g., temperature spikes in winter).
- **Collective anomalies:** **Groups of data points** that are **anomalous together** (e.g., sudden network traffic spikes).

23. Difference between supervised and unsupervised anomaly detection techniques

- **Supervised anomaly detection:**
 - Requires labelled data (normal and anomalous instances).
 - Examples: Classification models like logistic regression or neural networks.
- **Unsupervised anomaly detection:**
 - Works without labelled data.
 - Detects anomalies based on deviations from normal patterns.
 - Examples: Isolation Forest, DBSCAN, Local outlier factor detection.

24. Describe the Isolation Forest algorithm for anomaly detection

Isolation Forest isolates anomalies by recursively partitioning data:

- Randomly selects a feature and a split value.

- Anomalies require fewer splits to isolate because they are far from the majority.
- Outputs anomaly scores based on the average path length in a decision tree.

25. How does One-Class SVM work in anomaly detection?

One-Class SVM learns a decision boundary around normal data points in high-dimensional space:

- Maps data to a higher-dimensional space using a kernel function.
- Separates normal points from anomalies by maximizing the margin around the normal data.

26. Challenges of anomaly detection in high-dimensional data

- **Curse of dimensionality:** Distance metrics become less meaningful as dimensions increase.
- **Sparse anomalies:** Anomalies may be obscured in high-dimensional spaces.
- **Scalability:** Computation becomes expensive for large datasets.

27. Explain the concept of novelty detection?

Novelty detection identifies previously unseen data that differs from the training data. Unlike anomaly detection, which may focus on extreme cases, novelty detection assumes *new data points* are valid but different from known patterns. Applications include detecting new sensor failures or unknown customer behavior.

28. Real-world applications of anomaly detection?

- **Finance:** Fraudulent transaction detection.
- **Healthcare:** Identifying abnormal test results or symptoms.
- **Cybersecurity:** Intrusion detection in networks.
- **Manufacturing:** Predictive maintenance and fault detection.
- **Retail:** Detecting unusual shopping behaviours.
- **Environment:** Monitoring rare weather events.

29. Describe the Local Outlier Factor (LOF) algorithm?

The **Local Outlier Factor (LOF)** algorithm identifies anomalies by comparing the local density of a point to that of its neighbors:

- Computes the **local reachability density (LRD)** of each point, which represents how close a point is to its neighbors.
- An outlier has a significantly lower LRD than its neighbors.
- Outputs an LOF score, where values much greater than 1 indicate anomalies.

30. How do you evaluate the performance of an anomaly detection model?

- **Precision, Recall, F1-score:** Measures effectiveness in identifying anomalies.
- **Area Under the Receiver Operating Characteristic Curve (AUROC):** Evaluates model performance across thresholds.
- **Confusion Matrix:** Analyses false positives and false negatives.
- **Visual Inspection:** Plots clusters and anomalies when feasible.
- **Cost Analysis:** Assesses the cost of misclassification in the application context.

31. Role of feature engineering in anomaly detection

- **Feature scaling:** Normalizes features to ensure distance-based methods work effectively.
- **Feature selection:** Removes irrelevant or noisy features to improve performance.
- **Domain-specific features:** Incorporates knowledge from the application domain (e.g., time of day for transaction anomalies).
- **Dimensionality reduction:** Uses techniques like PCA to handle high-dimensional data.

32. Limitations of traditional anomaly detection methods

- Assumes specific data distributions (e.g., Gaussian).
- Sensitive to noise and outliers in the training set.
- Struggles with high-dimensional or dynamic datasets.
- Often requires manual tuning of hyperparameters.

33. Explain the concept of ensemble methods in anomaly detection?

Ensemble methods combine multiple anomaly detection models to improve robustness and accuracy:

- **Bagging:** Combines results from multiple models (e.g., Isolation Forest).
- **Boosting:** Sequentially improves detection by focusing on misclassified points.
- **Hybrid ensembles:** Mixes models based on different principles (e.g., distance-based, and density-based).

34. How does autoencoder-based anomaly detection work?

Autoencoders are neural networks trained to compress and reconstruct data:

- Normal data is reconstructed with minimal error.
- Anomalies have larger reconstruction errors as they deviate from learned patterns.
- Reconstruction error serves as the anomaly score.

35. Approaches for handling imbalanced data in anomaly detection

- **Oversampling:** Increase the representation of anomalies (e.g., SMOTE).
- **Undersampling:** Reduce the majority class size.
- **Cost-sensitive learning:** Penalizes misclassifying anomalies more heavily.
- **Synthetic data generation:** Uses GANs or similar methods to generate synthetic anomalies.

36. Describe the concept of semi-supervised anomaly detection

Semi-supervised anomaly detection models are trained on labelled normal data, with the goal of identifying deviations:

- Assumes anomalies are absent or rare in the training data.
- Uses approaches like one-class SVM or reconstruction-based techniques.

37. Trade-offs between false positives and false negatives in anomaly detection

- **False Positives:** Lead to unnecessary actions, wasting resources.
- **False Negatives:** Allow critical anomalies to go undetected.
- Applications determine the trade-off:
 - Fraud detection prioritizes minimizing false negatives.

- Medical applications often prioritize minimizing false positives for patient safety.

38. How do you interpret the results of an anomaly detection model?

- **Anomaly scores:** Understand how anomalies are ranked or scored.
- **Threshold analysis:** Choose a cutoff point to separate anomalies from normal data.
- **Feature contributions:** Analyse which features contributed to detecting anomalies.
- **Contextual understanding:** Validate anomalies against domain-specific knowledge.

39. Open research challenges in anomaly detection

- Handling dynamic and evolving data streams.
- Improving scalability for large datasets.
- Addressing high-dimensional, sparse data.
- Developing interpretable and explainable models.
- Reducing reliance on labelled data.

40. Explain the concept of contextual anomaly detection

Contextual anomaly detection identifies data points that are anomalous in a specific context but normal otherwise:

- Incorporates **contextual attributes** (e.g., time, location) and **behavioural attributes** (e.g., temperature value).
- Example: A feverish temperature is normal in summer but anomalous in winter.

41. What is time series analysis, and what are its key components?

Time series analysis studies data points collected sequentially over time to identify trends, patterns, and relationships.

Key components:

- **Trend:** Long-term direction of the data (e.g., upward sales).
- **Seasonality:** Repeating patterns (e.g., daily, or monthly cycles).

- **Noise:** Random variations in the data.
- **Stationarity:** Property where statistical characteristics (mean, variance) are constant over time.

42. Difference between univariate and multivariate time series analysis

- **Univariate time series:** Focuses on a single variable observed over time. Example: Monthly temperature readings.
- **Multivariate time series:** analyses multiple interdependent variables over time. Example: Temperature, humidity, and wind speed recorded simultaneously.

43. Process of time series decomposition?

Time series decomposition separates a time series into distinct components:

1. **Trend:** Long-term movement in the data.
2. **Seasonality:** Regular repeating patterns.
3. **Residuals:** Irregular fluctuations or noise. Decomposition can be additive ($y = T + S + R$) or multiplicative ($y = T \times S \times R$).

44. Main components of a time series decomposition

- **Trend:** Directional movement over time.
- **Seasonality:** Cyclical patterns within a fixed period.
- **Residual:** Random or unexplained variation.

45. Concept of stationarity in time series data

Stationarity refers to a time series whose statistical properties (mean, variance, autocorrelation) remain constant over time. Non-stationary data often have trends or seasonality, making them unsuitable for many modeling techniques.

46. Testing for stationarity in a time series

- **Visual inspection:** Look for consistent patterns in plots.
- **Statistical tests:**
 - **Augmented Dickey-Fuller (ADF) test:** Tests for unit roots.

- **Kwiatkowski-Phillips-Schmidt-Shin (KPSS) test:** Tests for trend stationarity.
- **Rolling statistics:** Compare mean and variance over time windows.

47. Discuss the ARIMA model

ARIMA (**A**utoregressive **I**ntegrated **M**oving **A**verage) models a time series by combining:

- **Autoregression (AR):** Relates current values to past values.
- **Integration (I):** Applies differencing to make the series stationary.
- **Moving Average (MA):** Models residual errors as a linear combination of past errors.

48. Parameters of the ARIMA model

- **p:** Number of lag observations in the autoregressive model.
- **d:** Degree of differencing applied to make the series stationary.
- **q:** Size of the moving average window.

49. Describe the SARIMA model

SARIMA (**S**easonal **A**utoregressive **I**ntegrated **M**oving **A**verage) extends ARIMA by incorporating seasonal components:

- **Seasonal AR, MA, and differencing terms** account for repeating patterns within a fixed period.
- **Parameters:** $(p, d, q) \times (P, D, Q, s)$ where P, D, Q are seasonal ARIMA terms and s is the season length.

50. Choosing the appropriate lag order in an ARIMA model

- Use **Autocorrelation Function (ACF)** and **Partial Autocorrelation Function (PACF)** plots:
 - **ACF** helps identify q (moving average terms).

- **PACF** helps identify p (autoregressive terms).
- Use criteria like **Akaike Information Criterion (AIC)** or **Bayesian Information Criterion (BIC)**.

51. Explain the concept of differencing in time series analysis

Differencing removes trends or seasonality to make a series stationary:

- **First differencing:** Subtracts each value from the previous one.
- **Seasonal differencing:** Subtracts values from the same point in the previous season.

52. What is the Box-Jenkins methodology?

The Box-Jenkins methodology provides a structured approach to time series modeling:

1. **Model identification:** Analyze ACF and PACF plots to select p , d , q .
2. **Parameter estimation:** Use statistical methods to estimate model parameters.
3. **Model diagnostics:** Validate the model by checking residuals for randomness.

53. Role of ACF and PACF plots in identifying ARIMA parameters

- **ACF plot:** Shows correlation of a time series with its lagged values, helping identify q .
- **PACF plot:** Shows partial correlation (after removing intermediate correlations), helping identify p .

54. Handling missing values in time series data

- **Interpolation:** Fill gaps using linear, spline, or polynomial interpolation.
- **Forward/Backward fill:** Propagate known values forward or backward.
- **Model-based imputation:** Use predictive models to estimate missing values.

- **Discarding incomplete data:** If gaps are small and won't significantly affect the analysis.

55. Concept of exponential smoothing

Exponential smoothing predicts future values by weighting past observations exponentially:

- **Simple exponential smoothing:** For data with no trend or seasonality.
- **Double exponential smoothing:** Adds trend components.
- **Triple exponential smoothing (Holt-Winters):** Adds trend and seasonality.

56. What is the Holt-Winters method, and when is it used?

The Holt-Winters method is an extension of exponential smoothing for seasonal data:

- **Additive:** Suitable for constant seasonal variations.
- **Multiplicative:** Suitable for proportional seasonal variations.
- Used in forecasting data with trend and seasonality, such as sales or temperature.

57. Challenges of Forecasting Long-Term Trends

- **Uncertainty:** Unpredictable events disrupt trends.
- **Model Drift:** Changing data patterns over time.
- **Noise:** Seasonal and random fluctuations.
- **External Factors:** Economic or social shifts.
- **Data Limitations:** Insufficient historical data.

58. Seasonality in Time Series

Recurring patterns at regular intervals due to cyclical factors (e.g., holiday sales, weather changes).

- **Additive Seasonality:** Constant variations.

- **Multiplicative Seasonality:** Variations proportional to data level.
- **Detection:** Decomposition, plots, or statistical tests.

59. Evaluating Forecast Models

- **Error Metrics:** MAE, RMSE, MAPE.
- **Residual Analysis:** Uncorrelated, zero-mean residuals.
- **Validation:** Walk-forward testing.
- **Bias:** Detect over-/under-forecasting.
- **Visuals:** Compare actual vs. predicted data.

60. Advanced Forecasting Techniques

- **ML/DL Models:** Random Forest, LSTM, Transformers.
- **Hybrid Models:** Combine ARIMA with ML.
- **Ensemble Methods:** Blend model predictions.
- **State-Space Models:** Trends via Kalman filters.
- **Exogenous Inputs:** Use external factors (e.g., weather).