

Lecture 1: Introduction to Machine Learning

Analytical Review, Short Notes, Analytical Questions, and 25+ Reasoned MCQs

Bulleted Short Notes

- **Machine Learning (ML):** Algorithms that improve performance automatically through experience/data.
- **Types of ML:**
 - **Supervised Learning:** Uses labeled data for tasks like classification/regression.
 - **Unsupervised Learning:** Finds structure/patterns in unlabeled data (e.g., clustering, dimensionality reduction).
 - **Reinforcement Learning:** Agent learns by interacting with environment, receiving rewards/penalties for actions.
- **ML Workflow:**
 - i. Data collection
 - ii. Data preprocessing/cleaning
 - iii. Feature engineering
 - iv. Model selection
 - v. Training
 - vi. Validation/testing
 - vii. Deployment
- **Overfitting vs Underfitting:**
 - *Overfitting:* Model memorizes training data, poor generalization to new data.
 - *Underfitting:* Model too simple, fails to capture trends in data.
- **Validation:**
 - Splitting data allows estimation of generalization performance, reduces risk of overfitting.
- **Feature Engineering:**
 - Transforming/creating features to improve learning and results.

- **Data Preprocessing:**

- Handling missing values, normalization/scaling, encoding categorical variables.

Analytical Short Questions (with Solutions)

Q1. What is the essential difference between supervised and unsupervised learning in terms of data structure and outcome?

A: Supervised learning uses labeled data to predict outputs; unsupervised learning uses unlabeled data to discover structure or patterns.

Q2. Provide a concrete real-world example for each of the three main ML types and justify your choice.

A:

- *Supervised:* Diagnosing disease from patient records (labels: disease/no disease).
- *Unsupervised:* Market basket analysis to find shopping patterns (no labels).
- *Reinforcement:* Training a robotic vacuum to clean efficiently by rewarding dirt collection.

Q3. How do overfitting and underfitting present themselves in both training and validation errors?

A:

- *Overfitting:* Low training error, high validation error.
- *Underfitting:* High error on both training and validation set.

Q4. Why do we split data into training and validation sets, and what problem does this solve?

A: To evaluate model generalization on unseen data and prevent overoptimistic performance estimates (data leakage).

Q5. What is the purpose of feature engineering, and give an example where it improves ML model performance?

A: Feature engineering creates or transforms features to highlight relationships; e.g., creating "age group" from "age" improves classification in demographic-based prediction.

Q6. If a model performs well on training data but poorly on validation data, what does this indicate, and what remedies could you apply?

A: Indicates overfitting; remedies include regularization, more data, simpler model, or data augmentation.

Q7. Why is normalization/scaling important for certain ML algorithms?

A: Ensures features contribute equally to distance-based or gradient-based algorithms, speeds up convergence, avoids dominance of large-scale features.

Q8. In what scenario is reinforcement learning a better fit than supervised learning?

A: When explicit labels are unavailable but feedback can be given as rewards or penalties (e.g., autonomous driving).

Q9. How can you detect and address underfitting in a model?

A: If both training and validation accuracy are low, increase model complexity, improve features, or reduce regularization.

Q10. What can happen if you evaluate a model only on training data?

A: You may significantly overestimate model performance due to memorization, not generalization.

Analytical MCQs (with Answers & Reasoning)

1. **Which of the following is NOT a fundamental type of machine learning?**

- A. Supervised Learning
- B. Unsupervised Learning
- C. Reinforcement Learning
- D. Relational Learning

Answer: D

Relational learning is not a primary ML type.

2. **In supervised learning, the training data must include:**

- A. Only input features
- B. Only output labels
- C. Both input features and output labels
- D. Only clusters

Answer: C

3. **Which is an example of unsupervised learning?**

- A. Spam detection
- B. Image clustering
- C. House price prediction
- D. Handwritten digit classification

Answer: B

4. **What is the main goal in reinforcement learning?**

- A. Minimize classification error
- B. Maximize cumulative reward
- C. Minimize clustering loss
- D. Maximize explained variance

Answer: B

5. **Which best indicates overfitting?**

- A. High training and validation accuracy
- B. Low training and validation accuracy
- C. High training, low validation accuracy
- D. Low training, high validation accuracy

Answer: C

6. **Validation set is used to:**

- A. Measure model's memorization
- B. Measure generalization on unseen data

- C. Compute feature importance
- D. Encrypt data

Answer: B

7. **If a model is underfitting, it is likely:**

- A. Too complex
- B. Too simple
- C. Well regularized
- D. Using too much data

Answer: B

8. **What does cross-validation help with?**

- A. Faster training
- B. Model validation and generalization
- C. Clustering
- D. Data augmentation

Answer: B

9. **Which scenario uses reinforcement learning?**

- A. Email classification
- B. Customer segmentation
- C. Robot learning to walk
- D. House price prediction

Answer: C

10. **A model that fits training data well but fails on new data is:**

- A. Underfitting
- B. Overfitting
- C. Regularized
- D. Well-calibrated

Answer: B

11. **Splitting data into training and validation sets helps to:**

- A. Prevent data leakage
- B. Reduce computational cost
- C. Tune hyperparameters
- D. Both A and C

Answer: D

12. **Feature engineering is important because:**

- A. Always increases complexity
- B. Can improve accuracy

- C. Prevents overfitting automatically
- D. Makes data collection easier

Answer: B

13. First step in ML pipeline:

- A. Model selection
- B. Data collection
- C. Evaluation
- D. Feature analysis

Answer: B

14. Primary risk of evaluating only on training data:

- A. Underestimating generalization error
- B. Overestimating computational cost
- C. Ignoring regularization
- D. Increasing bias

Answer: A

15. When is a model underfitting?

- A. High training, low validation accuracy
- B. Both errors high
- C. Both errors low
- D. High validation, low training accuracy

Answer: B

16. Benefit of data preprocessing:

- A. Increases bias
- B. Reduces noise/errors
- C. Makes training unnecessary
- D. Eliminates need for validation

Answer: B

17. Typical supervised learning example:

- A. Market basket analysis
- B. Fraud detection
- C. Topic modeling
- D. Dimensionality reduction

Answer: B

18. Why normalize features?

- A. Increases overfitting
- B. Ensures comparable scale

- C. Removes all outliers
- D. Creates more features

Answer: B

19. **TRUE about unsupervised learning:**

- A. Needs output labels
- B. For clustering
- C. Maximizes reward
- D. Needs reinforcement

Answer: B

20. **Data augmentation is most useful when:**

- A. Data is plentiful
- B. Overfitting is not a risk
- C. Training data is limited
- D. Only output labels are missing

Answer: C

21. **Overfitting leads to:**

- A. Reduced training accuracy
- B. Poor generalization
- C. Best for large datasets
- D. Never a concern

Answer: B

22. **Role of validation set:**

- A. Adjust model weights
- B. Evaluate during training
- C. Increase training set size
- D. Replace test set

Answer: B

23. **Reinforcement learning is preferred when:**

- A. Labeled data is abundant
- B. Only clustering is needed
- C. Actions lead to delayed rewards
- D. Feature selection is critical

Answer: C

24. **Purpose of feature engineering:**

- A. Increase overfitting
- B. Find better representations

- C. Remove need for training
- D. Make data collection easier

Answer: B

25. **Which is NOT a reason for validation?**

- A. To reduce overfitting
- B. Tune hyperparameters
- C. Data encryption
- D. Assess generalization

Answer: C

26. **If both training and validation errors are low:**

- A. Model may generalize well
- B. Model is underfitting
- C. Model is overfitting
- D. Model needs regularization

Answer: A

Lecture 2: Decision Trees

Analytical Review, Short Notes, Analytical Questions, and 25+ Reasoned MCQs

Bulleted Short Notes

- **Decision Trees:**
 - Tree-structured models for classification/regression.
 - Internal nodes: feature-based splits; leaves: output predictions.
- **Splitting Criteria:**
 - *Gini Impurity*: Measures frequency of misclassification.
 - *Entropy*: Measures information content/impurity.
 - *Information Gain*: Reduction in impurity after a split.
- **Tree Growth:**
 - Recursive partitioning of data until stopping criteria (max depth, min samples per leaf).
 - Axis-aligned decision boundaries.
- **Pruning:**
 - Removes branches to prevent overfitting.
 - Can be pre-pruning (early stopping) or post-pruning (after full growth).
- **Advantages:**
 - Interpretable, handles non-linear relationships, works with mixed data types.
- **Disadvantages:**
 - Prone to overfitting, sensitive to small data changes, greedy splitting may miss global optimum.

Analytical Short Questions (with Solutions)

Q1. How do Gini impurity and Entropy differ as splitting criteria in decision trees?

A:

- *Gini impurity* focuses on probability of incorrect classification (sum of squared class probabilities, subtracted from 1).
 - *Entropy* measures disorder—higher entropy means more class mix; calculated as $(-\sum p_i \log_2 p_i)$.
 - Both aim to maximize purity after a split but may rank splits differently due to mathematical form.
-

Q2. What is Information Gain and how is it used in building decision trees?

A:

Information Gain is the decrease in impurity (measured by entropy or Gini) after splitting a node.

$$[\text{Information Gain}] = \text{Parent Impurity} - \sum_k \frac{n_k}{n} \text{Child}_k \text{ Impurity}$$
 The split that maximizes Information Gain is chosen at each step.

Q3. Why are decision tree boundaries axis-aligned, and what implication does this have?

A:

Each split threshold operates on a single feature, resulting in splits that are perpendicular (axis-aligned) to feature axes.

Implication: Trees can't efficiently represent diagonal or curved boundaries unless very deep.

Q4. What is overfitting in decision trees, why does it occur, and how can it be mitigated?

A:

Overfitting occurs when the tree becomes too complex, memorizing training data and capturing noise.

Mitigation: Limit tree depth, set min samples per leaf, prune branches, use ensemble methods.

Q5. Explain the difference between pre-pruning and post-pruning.

A:

- *Pre-pruning*: Stop growing the tree early based on criteria (max depth, min samples).
 - *Post-pruning*: Grow the full tree, then remove branches that do not improve validation set performance.
-

Q6. Why are decision trees considered interpretable compared to other models?

A:

You can follow the path from root to leaf for any prediction, seeing which features and thresholds led to a decision.

Q7. What is the effect of increasing tree depth on bias and variance?

A:

Deeper trees decrease bias (fit data better) but increase variance (fit noise), risking overfitting.

Q8. How does greedy splitting in decision trees potentially lead to suboptimal solutions?

A:

At each node the split is locally optimal, but this may not lead to the globally optimal tree structure.

Q9. Describe a practical situation where a decision tree might be preferable over a linear model.

A:

When relationships between features and output are highly non-linear or involve categorical features with complex interactions.

Q10. What is the role of minimum samples per leaf in controlling a decision tree's complexity?

A:

Prevents splits that would create leaves with very few samples, which helps avoid overfitting from noise or outliers.

Analytical MCQs (with Answers & Reasoning)

1. **What is the main criterion used by decision trees to select a split?**

- A. Minimum sample size
- B. Maximum information gain
- C. Maximum tree depth
- D. Random feature selection

Answer: B

2. **Gini impurity measures:**

- A. The probability of correct classification
- B. The frequency of misclassification
- C. The number of classes
- D. The log of class probabilities

Answer: B

3. **Which splitting criterion is based on information theory?**

- A. Gini impurity
- B. Entropy
- C. Variance
- D. Pruning

Answer: B

4. **When does a decision tree stop splitting a node?**

- A. When information gain is maximized
- B. When all samples are of the same class
- C. When splitting increases impurity
- D. When tree is fully grown

Answer: B

5. **Why are decision tree boundaries axis-aligned?**

- A. Splits are made on individual features
- B. Splits are random
- C. Trees use PCA

D. All splits are diagonal

Answer: A

6. **Which is NOT an advantage of decision trees?**

- A. Interpretability
- B. Handling non-linear data
- C. Robustness to noise
- D. Mixed data type handling

Answer: C

7. **What is the effect of increasing tree depth?**

- A. Increases bias
- B. Decreases variance
- C. Increases variance
- D. Reduces overfitting

Answer: C

8. **Pre-pruning involves:**

- A. Growing tree fully, then trimming
- B. Stopping tree growth early
- C. Merging nodes after full growth
- D. Ignoring minimum samples

Answer: B

9. **What does post-pruning attempt to accomplish?**

- A. Increase overfitting
- B. Improve validation performance
- C. Make deeper trees
- D. Add more features

Answer: B

10. **Which is a weakness of decision trees?**

- A. Cannot handle categorical data
- B. Poor interpretability
- C. Sensitive to small data changes
- D. Always underfits

Answer: C

11. **The primary cause of overfitting in decision trees is:**

- A. Too few leaves
- B. Too deep trees
- C. Small validation set

D. Under-pruning

Answer: B

12. **Why are decision trees considered interpretable?**

- A. Use of linear decision surfaces
- B. Easily visualized paths from root to leaf
- C. Always shallow
- D. Use of regularization

Answer: B

13. **Which parameter can help control overfitting in decision trees?**

- A. Learning rate
- B. Minimum samples per leaf
- C. L1 regularization
- D. Feature scaling

Answer: B

14. **Which of the following is true about greedy splitting in decision trees?**

- A. Always achieves global optimum
- B. Only optimizes locally at each node
- C. Applies to ensemble methods only
- D. Reduces variance

Answer: B

15. **Which of the following would NOT be a suitable criterion for a decision tree split?**

- A. Information gain
- B. Gini impurity
- C. Mean squared error (for regression)
- D. Random guessing

Answer: D

16. **What is the purpose of pruning a decision tree?**

- A. To increase complexity
- B. To reduce overfitting
- C. To add more leaves
- D. To maximize variance

Answer: B

17. **If a tree is too shallow, it is likely to:**

- A. Overfit
- B. Underfit

- C. Have high variance
- D. Be random

Answer: B

18. **Which tree property is most responsible for its sensitivity to small data changes?**

- A. Greedy splitting
- B. Axis-aligned splits
- C. Pre-pruning
- D. Handling of mixed data types

Answer: A

19. **A tree that perfectly classifies the training data may:**

- A. Generalize well
- B. Be overfitting
- C. Be underfitting
- D. Have high bias

Answer: B

20. **Which is a practical use case for decision trees?**

- A. Image compression
- B. Non-linear classification with categorical variables
- C. Clustering
- D. Linear regression only

Answer: B

21. **If every split in a tree is made on the same feature, what does this indicate?**

- A. Feature is very informative
- B. Tree is random
- C. Tree is overfitting
- D. Data is underfitting

Answer: A

22. **Which describes a leaf node in a decision tree?**

- A. Decision point
- B. Contains predicted output
- C. Root of the tree
- D. Always binary

Answer: B

23. **Which is NOT a typical stopping criterion for tree growth?**

- A. Maximum depth

- B. Minimum information gain
- C. Minimum samples per split
- D. Maximum entropy

Answer: D

24. What is the main effect of post-pruning on a grown decision tree?

- A. Increases depth
- B. Reduces size and overfitting
- C. Adds more branches
- D. Increases number of leaves

Answer: B

25. Which is a disadvantage of decision trees compared to linear models?

- A. Cannot handle categorical data
- B. Susceptible to small data perturbations
- C. Cannot model non-linearities
- D. Require feature scaling

Answer: B

26. What happens if the minimum samples per split is set too high?

- A. Tree becomes too deep
- B. Tree becomes very shallow and may underfit
- C. Tree overfits
- D. Tree ignores categorical features

Answer: B

Lecture 2: Performance Measures

Analytical Review, Short Notes, Analytical Questions, and 25+ Reasoned MCQs

Bulleted Short Notes

- **Confusion Matrix:**
 - Summarizes classification performance: True Positive (TP), False Positive (FP), True Negative (TN), False Negative (FN).
 - **Key Metrics:**
 - **Accuracy:** $(TP + TN) / (TP + TN + FP + FN)$
 - **Precision:** $TP / (TP + FP)$
 - **Recall (Sensitivity):** $TP / (TP + FN)$
 - **Specificity:** $TN / (TN + FP)$
 - **F1-Score:** $2 \times (\text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall})$
 - **Balanced Accuracy:** $(\text{Sensitivity} + \text{Specificity}) / 2$
 - **ROC Curve & AUC:**
 - ROC: Plots TPR (Recall) vs. FPR (1 - Specificity); AUC is the area under ROC.
 - **Precision-Recall Curve:**
 - Useful for imbalanced classes.
 - **Class Imbalance:**
 - Accuracy can be misleading; focus on precision, recall, and F1.
 - **Multiclass Extension:**
 - Macro/micro averaging for metrics.
-

Analytical Short Questions (with Solutions)

Q1. Explain why accuracy is a misleading metric in imbalanced datasets and provide an example.

A:

If 95% of data is class A and 5% is class B, predicting all A yields 95% accuracy—good by accuracy but fails to detect minority class. Precision, recall, and F1-score reveal poor performance on class B.

Q2. Derive the formula for F1-score using precision and recall, and interpret its significance.

A:

$$F1 = 2 \times (\text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall}).$$

It balances precision and recall; high only when both are high.

Q3. What does the ROC curve show, and why is AUC a robust metric?

A:

ROC plots TPR vs. FPR at various thresholds.

AUC measures likelihood that a random positive is ranked above a random negative—threshold-independent.

Q4. How would you interpret a classifier with high precision but low recall?

A:

When it predicts positive, it's usually correct (few FPs), but it misses many actual positives (many FNs).

Q5. In multiclass classification, what is the difference between macro-averaged and micro-averaged F1-score?

A:

Macro: average F1 for each class, treating all classes equally.

Micro: aggregate contributions before calculating F1, favoring frequent classes.

Q6. Why is the confusion matrix essential beyond reporting a single metric?

A:

It reveals detailed error patterns: which classes are confused, where the model is weakest, and supports computing all other metrics.

Q7. If a model has high recall but low precision, what does that mean?

A:

It finds most positives (few FNs), but includes many false alarms (many FPs).

Q8. Describe a scenario where specificity is more important than sensitivity.

A:

In screening for a rare disease where false positives lead to unnecessary, costly procedures.

Q9. How would you use the precision-recall curve to select a threshold for a classifier?

A:

Choose the threshold balancing recall and precision for the desired tradeoff, especially important under class imbalance.

Q10. What would it mean if $AUC = 0.5$ for a classifier?

A:

Classifier is no better than random guessing.

Analytical MCQs (with Answers & Reasoning)

1. **What does the confusion matrix NOT show?**

- A. True Positives
- B. Recall
- C. False Negatives
- D. Feature importance

Answer: D

2. **Which metric is most affected by class imbalance?**

- A. Accuracy
- B. Recall
- C. F1-Score

D. Specificity

Answer: A

3. **Precision is defined as:**

A. $TP / (TP + FN)$

B. $TP / (TP + FP)$

C. $TN / (TN + FP)$

D. $(TP + TN) / \text{Total}$

Answer: B

4. **Recall is also known as:**

A. Specificity

B. Sensitivity

C. Precision

D. F1-score

Answer: B

5. **A classifier with high recall but low precision:**

A. Misses many actual positives

B. Has many false positives

C. Is always accurate

D. Has balanced performance

Answer: B

6. **F1-score is high only when:**

A. Precision is low, recall is high

B. Both precision and recall are high

C. Accuracy is high

D. Specificity is low

Answer: B

7. **Balanced accuracy is used to:**

A. Address class imbalance

B. Replace F1-score

C. Compute ROC

D. Ignore false negatives

Answer: A

8. **If a model predicts all negatives in a highly imbalanced binary dataset, accuracy will be:**

A. Low

B. High

- C. Zero
- D. Undefined

Answer: B

9. **AUC is best described as:**

- A. Likelihood a random positive is ranked above a negative
- B. The value of recall at optimum precision
- C. The mean squared error
- D. The threshold for maximum accuracy

Answer: A

10. **Which metric is most helpful when false negatives are very costly?**

- A. Precision
- B. Recall
- C. Specificity
- D. Accuracy

Answer: B

11. **The ROC curve is plotted with:**

- A. Precision vs. Recall
- B. TPR vs. FPR
- C. Accuracy vs. Threshold
- D. Specificity vs. Recall

Answer: B

12. **For a perfect classifier, the AUC is:**

- A. 0
- B. 0.5
- C. 1
- D. -1

Answer: C

13. **Which of the following is FALSE about F1-score?**

- A. It's harmonic mean of precision and recall
- B. It's sensitive to class imbalance
- C. It can be higher than accuracy
- D. It's always between 0 and 1

Answer: C

14. **Which metric is NOT derived directly from the confusion matrix?**

- A. Accuracy
- B. ROC AUC

C. Precision

D. Recall

Answer: B

15. **Macro-averaged F1-score:**

A. Weighs classes by frequency

B. Treats all classes equally

C. Ignores minority classes

D. Is threshold dependent

Answer: B

16. **If specificity is high and sensitivity is low, the model:**

A. Finds most positives

B. Has few false positives

C. Misses many positives

D. Both B and C

Answer: D

17. **Which scenario favors optimizing for precision?**

A. Spam detection (to avoid flagging legitimate mail)

B. Disease screening (to catch all cases)

C. Fraud detection (to minimize real loss)

D. Image recognition

Answer: A

18. **Micro-averaged F1-score:**

A. Aggregates all TP, FP, FN before computing

B. Averages F1 across classes

C. Ignores rare classes

D. Is higher than macro-F1 by definition

Answer: A

19. **Which is NOT a component of the confusion matrix?**

A. TP

B. FN

C. AUC

D. FP

Answer: C

20. **Precision-recall curve is most useful when:**

A. Classes are balanced

B. Data is high-dimensional

- C. One class is rare
- D. ROC AUC is 1

Answer: C

21. **If AUC = 0.5, the classifier is:**

- A. Perfect
- B. Random
- C. Useless
- D. Both B and C

Answer: D

22. **A high specificity means:**

- A. Few false positives
- B. Few false negatives
- C. Low recall
- D. High F1

Answer: A

23. **Which metric combines both precision and recall?**

- A. Accuracy
- B. F1-score
- C. ROC AUC
- D. Specificity

Answer: B

24. **For a rare disease, which metric is most likely to be misleading?**

- A. F1-score
- B. Precision
- C. Accuracy
- D. Recall

Answer: C

25. **If a model's confusion matrix shows many FNs and few FPs, which metric is likely low?**

- A. Precision
- B. Recall
- C. Specificity
- D. Accuracy

Answer: B

26. **Which value in a confusion matrix increases when the model is too conservative in positive predictions?**

A. TP

B. FN

C. FP

D. TN

Answer: B

Lecture 3: Classification

Analytical Review, Short Notes, Analytical Questions, and 25+ Reasoned MCQs

Bulleted Short Notes

- **Classification:**
 - Assigns input samples to discrete categories (labels/classes).
 - Types: Binary, Multiclass, Multilabel.
 - **Linear Classifiers:**
 - E.g., Logistic Regression, Perceptron.
 - Decision boundary: hyperplane in feature space.
 - **Logistic Regression:**
 - Predicts probability using sigmoid activation.
 - Loss: Cross-entropy (log-loss).
 - Output: Probability between 0 and 1.
 - **Perceptron:**
 - Simple linear classifier; updates weights if misclassified.
 - Works only if data is linearly separable.
 - **Decision Boundary:**
 - Linear for logistic regression and perceptron; non-linear for trees/neural networks.
 - **Evaluation:**
 - Accuracy, confusion matrix, ROC, precision, recall, F1.
 - **Non-linear Classification:**
 - Achieved via feature engineering, kernel methods, or deep models.
 - **Probabilistic Interpretation:**
 - Logistic regression outputs can be interpreted as class probabilities.
-

Analytical Short Questions (with Solutions)

Q1. Write the logistic regression decision function and explain its probabilistic interpretation.

A:

$$P(y=1|x) = \frac{\sigma(w^T x + b)}{1 + e^{-(w^T x + b)}}$$

The output gives the probability that input (x) belongs to class 1.

Q2. Compare and contrast hinge loss and cross-entropy loss for classification.

A:

- *Hinge loss* (used in SVM): penalizes points within the margin or misclassified, focuses on the decision boundary.
 - *Cross-entropy loss* (used in logistic regression): penalizes wrong probabilistic predictions, encourages correct probabilities, differentiable everywhere.
-

Q3. Why can't a single-layer perceptron solve the XOR problem?

A:

XOR is not linearly separable; a single linear boundary cannot split the classes as required for XOR.

Q4. How does a non-linear decision boundary arise in classification?

A:

By using feature engineering (polynomial features), kernel methods, or multi-layer (deep) models, boundaries can become non-linear.

Q5. What is the geometric interpretation of the perceptron's decision boundary?

A:

It is a hyperplane ($w^T x + b = 0$) that divides the feature space into two halves—one for each class.

Q6. How is the cross-entropy loss computed for a binary classifier?

A:

$$[L = -[y \log(p) + (1-y)\log(1-p)]]$$

where (p) is the predicted probability for class 1.

Q7. What evaluation metrics are most informative for imbalanced class classification?

A:

Precision, recall, and F1-score are better than accuracy.

Q8. What is the role of the sigmoid function in logistic regression?

A:

It converts the raw linear output into a probability between 0 and 1.

Q9. How does kernel trick enable non-linear classification?

A:

It implicitly maps data into higher-dimensional space where linear separation is possible, without explicit computation.

Q10. Why is probabilistic output important in many classification applications?

A:

Allows for uncertainty estimation, threshold tuning, and integration with probabilistic decision systems.

Analytical MCQs (with Answers & Reasoning)

1. **Which algorithm is a linear classifier?**

- A. Decision Tree
- B. Logistic Regression
- C. KNN

D. RBF SVM

Answer: B

2. **The decision boundary of a perceptron is:**

- A. Non-linear
- B. Linear
- C. Circular
- D. Not well-defined

Answer: B

3. **Which loss is used in logistic regression?**

- A. Hinge
- B. Squared error
- C. Cross-entropy
- D. Zero-one

Answer: C

4. **The perceptron learning rule updates weights when:**

- A. The prediction is correct
- B. The prediction is wrong
- C. The loss is zero
- D. The input is missing

Answer: B

5. **Why can't a single-layer perceptron solve XOR?**

- A. It is non-differentiable
- B. Data is not linearly separable
- C. Weights cannot be negative
- D. Input must be continuous

Answer: B

6. **What is the output of the sigmoid function?**

- A. Any real number
- B. Only 1 or 0
- C. Value between 0 and 1
- D. Integer

Answer: C

7. **What does cross-entropy loss penalize?**

- A. Incorrect class probabilities
- B. High accuracy
- C. Large weights

D. Proper calibration

Answer: A

8. **A non-linear decision boundary can be formed by:**

- A. Using more features
- B. Kernel methods
- C. Deep neural networks
- D. All of the above

Answer: D

9. **The main benefit of probabilistic output in classification is:**

- A. Simpler models
- B. Uncertainty estimation
- C. Always higher accuracy
- D. Less data needed

Answer: B

10. **Which metric is LEAST informative for highly imbalanced classes?**

- A. Accuracy
- B. Recall
- C. Precision
- D. F1-score

Answer: A

11. **The kernel trick enables:**

- A. Faster training
- B. Non-linear classification with linear algorithms
- C. Better feature scaling
- D. More regularization

Answer: B

12. **Which of the following is NOT a type of classification?**

- A. Binary
- B. Multiclass
- C. Regression
- D. Multilabel

Answer: C

13. **What is the shape of the decision boundary for logistic regression in 2D?**

- A. Line
- B. Circle
- C. Hyperbola

D. Spline

Answer: A

14. Which of these is most affected by threshold choice in classification?

A. Accuracy

B. Precision

C. Recall

D. All of the above

Answer: D

15. If a classifier outputs 0.7 for class 1, this means:

A. The sample is class 1

B. The sample is class 0

C. There is a 70% chance it is class 1

D. The loss is 0.7

Answer: C

16. Which is a limitation of logistic regression?

A. Only works with binary classes

B. Only models linear boundaries

C. Cannot use regularization

D. Cannot output probabilities

Answer: B

17. The perceptron algorithm stops updating when:

A. All samples are correctly classified

B. Loss increases

C. No weights are left

D. Data is non-linear

Answer: A

18. Hinge loss is most closely associated with:

A. Perceptron

B. Decision Tree

C. SVM

D. Logistic Regression

Answer: C

19. Which statement about multiclass classification is TRUE?

A. Only one class per sample

B. Samples can belong to multiple classes

C. Uses regression loss

D. Never uses probabilistic outputs

Answer: A

20. **What is a major benefit of the confusion matrix?**

- A. Shows only accuracy
- B. Shows detailed error patterns
- C. Increases loss
- D. Only for regression

Answer: B

21. **Which is a probabilistic classifier?**

- A. KNN
- B. Logistic Regression
- C. Perceptron
- D. SVM with linear kernel

Answer: B

22. **Why is regularization important in classification models?**

- A. Forces higher weights
- B. Prevents overfitting
- C. Ensures linear boundaries
- D. Reduces accuracy

Answer: B

23. **Which method can help create non-linear boundaries in logistic regression?**

- A. Add polynomial features
- B. Use raw features only
- C. Reduce data size
- D. Use only categorical features

Answer: A

24. **Which metric would be most important for a medical diagnosis classifier?**

- A. Recall
- B. Accuracy
- C. Specificity
- D. All are equally important

Answer: A

25. **What does the perceptron algorithm do if a sample is misclassified?**

- A. Nothing
- B. Updates weights to reduce error
- C. Changes activation function

D. Removes the sample

Answer: B

26. **In multilabel classification, each sample:**

A. Has exactly one label

B. Has no label

C. Can have multiple labels

D. Must be numeric

Answer: C

Lecture 4: Supervised Learning

Analytical Review, Short Notes, Analytical Questions, and 25+ Reasoned MCQs

Bulleted Short Notes

- **Supervised Learning:**

- Learning from labeled data to predict outputs for new inputs.
- Two main problems: Classification (predict discrete labels), Regression (predict continuous values).

- **Typical Workflow:**

- i. Collect labeled data
- ii. Preprocess and split into train/test (and often validation)
- iii. Model selection (e.g., logistic regression, decision tree, SVM)
- iv. Model training
- v. Model evaluation (accuracy, RMSE, etc.)
- vi. Hyperparameter tuning (cross-validation)

- **Bias-Variance Tradeoff:**

- High bias: Underfitting
- High variance: Overfitting
- Goal: Find a balance for best generalization

- **Feature Scaling:**

- Important for algorithms using distances or gradients (e.g., SVM, kNN, linear regression)

- **Cross-Validation:**

- Robust method for model evaluation and hyperparameter selection.

- **Loss Functions:**

- Classification: Cross-entropy, hinge loss
- Regression: Mean squared error, mean absolute error

- **Regularization:**

- Techniques (L1/L2) to prevent overfitting by penalizing large weights.
-

Analytical Short Questions (with Solutions)

Q1. How does supervised learning differ from unsupervised learning in terms of data and objectives?

A:

Supervised learning uses labeled data to learn input-output mappings, aiming to predict outputs for new inputs. Unsupervised learning finds structure in unlabeled data, such as clusters or associations, without explicit targets.

Q2. What does the bias-variance tradeoff mean, and how does it impact model selection?

A:

Bias is error from overly simple models (underfitting); variance is error from overly complex models (overfitting). The tradeoff is finding model complexity that minimizes total error and generalizes best.

Q3. Why is cross-validation superior to a single train/test split for model assessment?

A:

Cross-validation (e.g., k-fold) averages performance over multiple splits, providing a more reliable estimate of generalization and reducing dependence on one random split.

Q4. When is feature scaling critical and why?

A:

Critical for algorithms sensitive to feature magnitudes (e.g., kNN, SVM, linear/logistic regression) to ensure all features contribute equally and improve convergence.

Q5. In regression, compare mean squared error (MSE) and mean absolute error (MAE) as loss functions.

A:

MSE penalizes large errors more (sensitive to outliers), while MAE treats all errors linearly (more robust to outliers).

Q6. What is the role of regularization in supervised learning? Name two types.

A:

Regularization penalizes model complexity to reduce overfitting. L1 (lasso) encourages sparsity; L2 (ridge) shrinks weights smoothly.

Q7. Describe the process of hyperparameter tuning using cross-validation.

A:

Try different hyperparameters, evaluate each using cross-validation, select the set with the best average validation performance.

Q8. Why can a model perform well on training data but poorly on test data?

A:

Overfitting—model memorizes training data including noise, failing to generalize to unseen examples.

Q9. Give an example of a regression and a classification problem.

A:

Regression: Predicting house prices (continuous output).

Classification: Email spam detection (discrete labels: spam/not spam).

Q10. How do you detect if your model is underfitting?

A:

Both training and validation errors are high; model is too simple to capture underlying patterns.

Analytical MCQs (with Answers & Reasoning)

1. **Which of the following is a regression problem?**

- A. Predicting if an email is spam
- B. Predicting house price
- C. Assigning topic to a document
- D. Diagnosing a disease (yes/no)

Answer: B

2. **What is the main goal of supervised learning?**

- A. Cluster data
- B. Find associations
- C. Predict outputs for new inputs
- D. Find feature correlations

Answer: C

3. **For which algorithm is feature scaling most critical?**

- A. Decision Tree
- B. kNN
- C. Naive Bayes
- D. Random Forest

Answer: B

4. **Cross-validation helps to:**

- A. Increase bias
- B. Reduce overfitting
- C. Assess generalization reliably
- D. Only tune the test set

Answer: C

5. **High bias in a model leads to:**

- A. Overfitting
- B. Underfitting
- C. High variance
- D. Good generalization

Answer: B

6. **Which is NOT a typical loss function for regression?**

- A. Mean squared error
- B. Mean absolute error
- C. Cross-entropy
- D. Hinge loss

Answer: D

7. **L1 regularization encourages:**

- A. Large weights
- B. Sparse solutions
- C. Overfitting
- D. High bias

Answer: B

8. **If validation error is much higher than training error, the model is likely:**

- A. Underfitting
- B. Overfitting
- C. Well-generalized
- D. Linear

Answer: B

9. **Which method can best assess model robustness to random data splits?**

- A. Single train/test split
- B. Cross-validation
- C. Data augmentation
- D. Feature scaling

Answer: B

10. **Which is a discrete output problem?**

- A. Regression
- B. Classification
- C. Clustering
- D. PCA

Answer: B

11. **Mean absolute error is more robust than MSE to:**

- A. Outliers
- B. Large datasets
- C. Feature scaling
- D. Regularization

Answer: A

12. **Which regularization method penalizes the number of non-zero weights?**

- A. L1
- B. L2
- C. Both
- D. Neither

Answer: A

13. If your model underfits, you should:

- A. Increase model complexity
- B. Add more regularization
- C. Reduce number of features
- D. Lower learning rate

Answer: A

14. Which is NOT a supervised learning metric?

- A. RMSE
- B. F1-score
- C. Silhouette score
- D. Accuracy

Answer: C

15. What is the main purpose of the test set?

- A. Train the model
- B. Tune hyperparameters
- C. Assess final generalization
- D. Feature selection

Answer: C

16. Which step comes FIRST in the supervised learning workflow?

- A. Model evaluation
- B. Model training
- C. Data collection
- D. Hyperparameter tuning

Answer: C

17. Which technique is used to prevent overfitting?

- A. Increasing model size
- B. Regularization
- C. Ignoring validation set
- D. Lowering variance

Answer: B

18. If your model performs poorly on both train and test data, you should:

- A. Decrease complexity
- B. Increase complexity
- C. Add more regularization
- D. Use cross-validation

Answer: B

19. Which problem is best addressed by supervised learning?

- A. Market segmentation
- B. Predicting housing prices
- C. Finding latent topics
- D. Extracting principal components

Answer: B

20. Cross-entropy loss is usually used for:

- A. Regression
- B. Clustering
- C. Classification
- D. Dimensionality reduction

Answer: C

21. Which is a benefit of cross-validation?

- A. Uses less data
- B. Reduces test set size
- C. Gives more stable performance estimates
- D. Avoids feature scaling

Answer: C

22. High variance is indicated by:

- A. Similar train and test errors
- B. Large gap between train and test errors
- C. Low bias
- D. Both B and C

Answer: D

23. Which is NOT a type of supervised learning task?

- A. Regression
- B. Classification
- C. Clustering
- D. Time series forecasting

Answer: C

24. What is the effect of using too large a regularization parameter?

- A. Overfitting
- B. Underfitting
- C. No effect
- D. Increased variance

Answer: B

25. Which is the BEST strategy for hyperparameter tuning?

- A. Random guessing
- B. Using only test set
- C. Cross-validation
- D. Ignoring validation error

Answer: C

26. Feature scaling is NOT critical for:

- A. Linear regression
- B. SVM
- C. Decision trees
- D. kNN

Answer: C

Lecture 5: Shallow Neural Networks

Analytical Review, Short Notes, Analytical Questions, and 25+ Reasoned MCQs

Bulleted Short Notes

- **Shallow Neural Networks:**
 - Consist of input layer, one hidden layer, and output layer.
 - Each neuron computes a weighted sum of inputs plus bias, followed by an activation function.
 - **Perceptron:**
 - Single neuron; can only solve linearly separable problems.
 - **Activation Functions:**
 - Introduce non-linearity (e.g., sigmoid, tanh, ReLU).
 - **Forward Propagation:**
 - Inputs are processed through layers to produce predictions.
 - **Learning & Gradient Descent:**
 - Adjusts weights to minimize loss (typically via backpropagation).
 - **Linear Separability & XOR Problem:**
 - Shallow networks with non-linear activation can solve non-linearly separable problems; single-layer perceptron cannot.
 - **Limitations:**
 - Limited representation power compared to deep networks.
 - Prone to underfitting on complex data.
-

Analytical Short Questions (with Solutions)

Q1. Why can't a single-layer perceptron solve the XOR problem, and how does a shallow neural network overcome this?

A:

A single-layer perceptron can only form linear decision boundaries, but XOR is not linearly separable. A shallow network with a hidden layer and non-linear activation can represent the XOR function by combining multiple linear boundaries non-linearly.

Q2. What is the role of the activation function in a neural network?

A:

It introduces non-linearity, enabling the network to learn complex patterns and boundaries. Without it, the model would behave like a linear model regardless of layers.

Q3. Describe the process of forward propagation in a shallow neural network.

A:

Input features are multiplied by weights and summed (plus bias) at each neuron, passed through activation, then passed to the next layer, finally producing an output.

Q4. How does gradient descent update the weights in a neural network?

A:

By computing the gradient of the loss with respect to each weight, then adjusting weights in the opposite direction of the gradient to minimize loss.

Q5. What is meant by "linear separability," and why is it important in the context of neural networks?

A:

It means data can be separated by a straight line/hyperplane. Single-layer perceptrons require linear separability to classify correctly; networks with hidden layers do not.

Q6. What limitation do shallow neural networks have compared to deep neural networks?

A:

Shallow networks have limited capacity to represent complex functions and may underfit data with intricate patterns.

Q7. How does the choice of activation function affect learning in a shallow neural network?

A:

Non-linear activations (sigmoid, ReLU, tanh) allow modeling of complex relationships; linear activation reduces the network to a linear model.

Q8. What is the geometric interpretation of what a hidden layer does in a neural network?

A:

It transforms the original feature space, enabling the network to combine multiple linear boundaries into more complex, non-linear decision surfaces.

Q9. Explain the importance of bias terms in neural network neurons.

A:

Bias allows the activation function's threshold to be shifted, enabling the neuron to better fit the data.

Q10. Can increasing the number of hidden units in a shallow network always solve any classification problem? Explain.

A:

No; while more units increase capacity, shallow networks still cannot efficiently represent hierarchical or highly complex patterns compared to deep architectures.

Analytical MCQs (with Answers & Reasoning)

1. **A shallow neural network contains:**

- A. Only the input layer
- B. One hidden layer
- C. Multiple hidden layers
- D. No activation functions

Answer: B

2. **The main reason for using an activation function is to:**

- A. Normalize inputs
- B. Introduce non-linearity
- C. Increase linearity
- D. Reduce overfitting

Answer: B

3. **The perceptron can only solve problems that are:**

- A. Non-linear
- B. Linearly separable
- C. Hierarchical
- D. Multi-class

Answer: B

4. **Which activation function is NOT non-linear?**

- A. Sigmoid
- B. Tanh
- C. ReLU
- D. Identity

Answer: D

5. **Forward propagation refers to:**

- A. Computing gradients
- B. Passing inputs through the network to get outputs
- C. Updating weights
- D. Pruning neurons

Answer: B

6. **Gradient descent works by:**

- A. Increasing the loss
- B. Following the gradient uphill
- C. Minimizing the loss by adjusting weights
- D. Removing data points

Answer: C

7. **The XOR problem demonstrates the limitation of:**

- A. Deep networks
- B. Single-layer perceptrons
- C. Non-linear classifiers
- D. Multi-class models

Answer: B

8. **Bias terms in neural networks:**

- A. Have no effect
- B. Are always zero
- C. Allow shifting of activation thresholds
- D. Only matter in deep networks

Answer: C

9. **Which function is most likely to cause vanishing gradients in shallow networks?**

- A. ReLU
- B. Linear
- C. Sigmoid
- D. Identity

Answer: C

10. **Which is NOT a typical limitation of shallow neural networks?**

- A. Limited representation power
- B. Prone to underfitting on complex data
- C. Always overfit
- D. Cannot model hierarchical features

Answer: C

11. **A neuron computes:**

- A. Weighted sum + bias, then activation
- B. Only the sum of inputs
- C. Only activation of bias
- D. Product of inputs

Answer: A

12. **If you remove all activation functions from a shallow neural network, it becomes equivalent to:**

- A. A linear model
- B. A deep network
- C. kNN
- D. An SVM

Answer: A

13. **What is a key advantage of non-linear activation functions?**

- A. Make training faster
- B. Allow complex, non-linear decision boundaries
- C. Reduce the number of weights

D. Increase overfitting

Answer: B

14. Which best describes the learning process in neural networks?

- A. Random weight changes
- B. Adjusting weights to minimize loss
- C. Manual feature engineering
- D. Only using bias terms

Answer: B

15. Increasing the number of hidden units in a shallow network:

- A. Always solves any problem
- B. Can increase representational power, but has limits
- C. Makes the network linear
- D. Reduces training time

Answer: B

16. Which of the following enables a shallow network to solve non-linear problems?

- A. More data
- B. Non-linear activation and hidden layer
- C. Removing bias
- D. Using only output layer

Answer: B

17. Which scenario is NOT a limitation of shallow networks?

- A. Handling complex, hierarchical features
- B. Solving linearly separable problems
- C. Underfitting on complex data
- D. Representing deep structures

Answer: B

18. If a neuron uses the identity activation, the output is:

- A. Non-linear
- B. A threshold
- C. The input
- D. The weighted sum

Answer: D

19. A shallow neural network with enough units and non-linear activation can:

- A. Approximate any continuous function
- B. Only solve linear problems

- C. Never overfit
- D. Act as a clustering algorithm

Answer: A

20. **The process of changing weights to reduce prediction error is called:**

- A. Backpropagation
- B. Forward propagation
- C. Weight decay
- D. Normalization

Answer: A

21. **Vanishing gradients are more common in:**

- A. Deep networks
- B. Shallow networks
- C. Single-layer perceptrons
- D. Linear models

Answer: A

22. **The main function of a hidden layer is to:**

- A. Transform the feature space
- B. Store labels
- C. Normalize outputs
- D. Reduce number of weights

Answer: A

23. **Which function would likely lead to dead neurons in a shallow network?**

- A. ReLU
- B. Sigmoid
- C. Tanh
- D. Linear

Answer: A

24. **If all weights in a shallow neural network are initialized to zero:**

- A. Network can still learn
- B. All neurons learn the same features
- C. It speeds up training
- D. It prevents overfitting

Answer: B

25. **Which architecture has more representational power?**

- A. Shallow network
- B. Deep network

C. Single-layer perceptron

D. Linear regression

Answer: B

26. **Which is a key difference between shallow and deep networks?**

A. Use of bias

B. Number of layers

C. Use of activation functions

D. Learning rate

Answer: B

Lecture 6: Deep Neural Networks

Analytical Review, Short Notes, Analytical Questions, and 25+ Reasoned MCQs

Bulleted Short Notes

- **Deep Neural Networks (DNNs):**
 - Networks with multiple hidden layers (depth > 1).
 - Each layer extracts higher-level features.
 - **Architecture:**
 - Input \rightarrow Multiple hidden layers \rightarrow Output.
 - Each neuron: Weighted sum + bias, followed by activation.
 - **Forward and Backward Propagation:**
 - Forward: Compute outputs layer by layer.
 - Backward: Compute gradients of loss w.r.t. all weights (backpropagation).
 - **Non-linear Function Approximation:**
 - Deep networks can model complex, hierarchical patterns.
 - **Vanishing/Exploding Gradients:**
 - In deep networks, gradients can become very small (vanish) or very large (explode), making training difficult.
 - **Activations:**
 - ReLU helps mitigate vanishing gradient; sigmoid/tanh can cause it.
 - **Representational Power:**
 - More depth \rightarrow greater ability to approximate complex functions, but higher risk of overfitting.
 - **Regularization in DNNs:**
 - Dropout, L1/L2, batch normalization to prevent overfitting and stabilize training.
 - **Applications:**
 - Image, speech, text, and time series analysis.
-

Analytical Short Questions (with Solutions)

Q1. Why do deep neural networks have greater representational power than shallow networks?

A:

Each layer learns increasingly abstract features by composing non-linear transformations. More layers allow modeling of complex, hierarchical patterns unachievable by shallow networks.

Q2. What is the vanishing gradient problem and how does it affect deep network training?

A:

Gradients become very small as they are propagated back through many layers (especially with sigmoid/tanh), causing weights in earlier layers to update very little or not at all, thus stalling learning.

Q3. How does the ReLU activation function help address vanishing gradients?

A:

ReLU does not saturate in the positive domain, so its derivative is 1 for positive inputs, allowing gradients to flow backward without shrinking rapidly, making training deeper networks feasible.

Q4. Explain how backpropagation works in deep neural networks.

A:

Backpropagation applies the chain rule to compute gradients of the loss with respect to all weights, layer by layer from output to input, enabling efficient weight updates via gradient descent.

Q5. What is the tradeoff between depth and overfitting in deep neural networks?

A:

While depth increases capacity to model complex data, too many layers can memorize noise, resulting in overfitting. Regularization and large datasets are needed to counteract this.

Q6. Why is dropout used in training deep neural networks?

A:

Dropout randomly disables neurons during training, forcing the network to learn redundant, robust representations and reducing overfitting.

Q7. How does batch normalization stabilize and speed up deep network training?

A:

It normalizes inputs to each layer, reducing internal covariate shift, allowing higher learning rates, and mitigating vanishing/exploding gradients.

Q8. Give an example of a hierarchical feature learned in a deep network for image recognition.

A:

In early layers, the network may learn edges; in deeper layers, combinations of edges (shapes), and in final layers, object parts or entire objects.

Q9. What is the role of non-linearity in each layer of a deep network?

A:

Non-linear activation after each linear transformation allows the network to approximate complex, non-linear functions.

Q10. When might increasing depth hurt model performance?

A:

When data is insufficient, overfitting increases, or vanishing/exploding gradients prevent effective training. Also, depth beyond necessary adds computational cost without accuracy gain.

Analytical MCQs (with Answers & Reasoning)

1. **A deep neural network is defined as a network with:**

- A. No hidden layers
- B. One hidden layer
- C. Multiple hidden layers
- D. Only an input and output layer

Answer: C

2. **Which problem occurs when gradients become very small in deep networks?**

- A. Overfitting
- B. Vanishing gradients
- C. Exploding gradients
- D. Batch normalization

Answer: B

3. **The main purpose of backpropagation is to:**

- A. Compute outputs
- B. Update weights using gradient information
- C. Normalize data
- D. Reduce model depth

Answer: B

4. **Which activation helps mitigate the vanishing gradient problem?**

- A. Sigmoid
- B. Tanh
- C. ReLU
- D. Softmax

Answer: C

5. **Dropout is a technique to:**

- A. Increase depth
- B. Prevent overfitting
- C. Reduce training time
- D. Normalize gradients

Answer: B

6. **Deep neural networks are especially powerful because they:**

- A. Use only linear transformations
- B. Can model highly complex, hierarchical patterns
- C. Require less data than shallow networks

D. Always generalize better

Answer: B

7. **Which layer typically learns the most abstract features?**

A. Input

B. First hidden

C. Last hidden

D. Output

Answer: C

8. **Batch normalization improves training by:**

A. Increasing overfitting

B. Stabilizing activations

C. Removing hidden layers

D. Lowering learning rate

Answer: B

9. **Which scenario can cause exploding gradients?**

A. Too few layers

B. Large weights and deep networks

C. Using ReLU

D. Low learning rate

Answer: B

10. **Non-linear activation functions are necessary in deep networks to:**

A. Reduce computation

B. Approximate non-linear functions

C. Prevent overfitting

D. Limit memory use

Answer: B

11. **If a deep network overfits, you should try:**

A. Adding more layers

B. Removing dropout

C. Collecting more data

D. Increasing learning rate

Answer: C

12. **The main difference between shallow and deep networks is:**

A. Use of bias

B. Number of hidden layers

C. Input size

D. Output format

Answer: B

13. What does the chain rule enable in backpropagation?

- A. Forward computation
- B. Efficient gradient computation through layers
- C. Data normalization
- D. Weight initialization

Answer: B

14. Which is a disadvantage of very deep networks?

- A. Cannot model complex data
- B. More prone to vanishing/exploding gradients
- C. Always underfit
- D. Never require regularization

Answer: B

15. What is a hierarchical feature in DNNs?

- A. Feature learned from raw data only
- B. Feature composed of lower-level features
- C. Linear combination of inputs
- D. Random noise

Answer: B

16. When is ReLU activation less effective?

- A. When inputs are negative
- B. When inputs are positive
- C. In shallow networks
- D. With batch normalization

Answer: A

17. To reduce internal covariate shift, you should use:

- A. Dropout
- B. L2 regularization
- C. Batch normalization
- D. Sigmoid activation

Answer: C

18. Which technique is NOT specifically used to regularize deep networks?

- A. Dropout
- B. L1/L2 penalty
- C. Increasing learning rate

D. Early stopping

Answer: C

19. **The output of the last hidden layer in a DNN is often called:**

- A. Input
- B. Logit
- C. Feature representation
- D. Loss

Answer: C

20. **Which of the following is a sign that your deep network is underfitting?**

- A. Low train and test accuracy
- B. High train, low test accuracy
- C. High variance
- D. Exploding gradients

Answer: A

21. **What is the main computational challenge with deep networks?**

- A. Too few parameters
- B. Training time and vanishing/exploding gradients
- C. Lack of expressiveness
- D. No need for regularization

Answer: B

22. **Which of these is NOT a typical application of deep neural networks?**

- A. Image recognition
- B. Speech processing
- C. Sorting numbers
- D. Text analysis

Answer: C

23. **If you notice your DNN is not learning, what's a first thing to check?**

- A. Depth
- B. Data labels
- C. Learning rate and activation functions
- D. Output size

Answer: C

24. **Which of these is most directly related to the "depth" of a DNN?**

- A. Number of layers
- B. Number of features
- C. Number of neurons per layer

D. Size of output

Answer: A

25. **Which is a sign of exploding gradients?**

A. Model weights become extremely large or NaN

B. Loss decreases gradually

C. Training speed increases

D. Output is constant

Answer: A

26. **Which activation function is most likely to cause vanishing gradients?**

A. ReLU

B. Sigmoid

C. Linear

D. Softmax

Answer: B

Lecture 7: Backpropagation

Analytical Review, Short Notes, Analytical Questions, and 25+ Reasoned MCQs

Bulleted Short Notes

- **Backpropagation:**
 - Algorithm for efficiently computing gradients in neural networks.
 - Enables weight updates via gradient descent.
 - **Forward Pass:**
 - Compute outputs/predictions for given inputs.
 - **Backward Pass (Backpropagation):**
 - Apply chain rule to compute gradients of loss with respect to each weight.
 - Gradients are propagated backward layer by layer.
 - **Gradient Descent:**
 - Use computed gradients to update weights and minimize loss.
 - **Key Steps:**
 - i. Forward pass: compute activations and output.
 - ii. Compute loss.
 - iii. Backward pass: compute gradients layer by layer.
 - iv. Update weights with gradient descent.
 - **Efficiency:**
 - Backpropagation is much faster than naive differentiation; exploits chain rule and caching of intermediate results.
 - **Challenges:**
 - Vanishing/exploding gradients, especially in deep networks.
 - **Variants:**
 - Stochastic, mini-batch, and full-batch gradient descent.
-

Analytical Short Questions (with Solutions)

Q1. Why is backpropagation necessary for training deep neural networks?

A:

It efficiently computes the gradient of the loss function with respect to every parameter, enabling effective gradient-based optimization even in very deep architectures.

Q2. Explain how the chain rule is used in backpropagation.

A:

The chain rule allows calculation of gradients of composite functions by multiplying derivatives at each step, enabling backpropagation to compute derivatives of the loss with respect to each layer's parameters.

Q3. Describe the steps in a single iteration of the backpropagation algorithm.

A:

1. Forward pass: compute activations and predictions.
 2. Compute loss.
 3. Backward pass: calculate gradients from output to input layer.
 4. Update weights using the gradients.
-

Q4. What is gradient descent, and how does it relate to backpropagation?

A:

Gradient descent is an optimization algorithm that minimizes the loss by updating parameters in the direction of the negative gradient. Backpropagation computes these gradients in neural networks.

Q5. What problem do vanishing gradients cause in deep networks?

A:

Gradients become too small to update early layers, resulting in very slow or stalled learning for those layers.

Q6. How does mini-batch gradient descent differ from full-batch and stochastic methods?

A:

Mini-batch uses a subset of training data for each update (balance between noisy stochastic and stable full-batch), often improving convergence and computational efficiency.

Q7. How does caching intermediate results in the forward pass help during backpropagation?

A:

It avoids redundant calculations by reusing values needed for gradient computations, making the algorithm efficient.

Q8. Why might exploding gradients be problematic, and how are they commonly addressed?

A:

Exploding gradients cause unstable updates and can lead to NaN weights or divergence. Solutions include gradient clipping and careful weight initialization.

Q9. Can backpropagation be used for models with non-differentiable components? Explain.

A:

No; backpropagation requires differentiable operations in the computational graph. Non-differentiable components break the chain rule.

Q10. What is the primary mathematical advantage of backpropagation compared to naive differentiation?

A:

It leverages the chain rule and shared computation to compute all required gradients in time proportional to a single forward pass, rather than exponential time.

Analytical MCQs (with Answers & Reasoning)

1. **The main purpose of backpropagation is to:**

- A. Compute network outputs
- B. Compute gradients for all parameters efficiently
- C. Normalize inputs
- D. Randomly update weights

Answer: B

2. **Which mathematical principle is at the core of backpropagation?**

- A. Taylor expansion
- B. Chain rule of calculus
- C. Matrix inversion
- D. Eigen decomposition

Answer: B

3. **Which step comes first in backpropagation?**

- A. Backward pass
- B. Forward pass
- C. Weight update
- D. Gradient clipping

Answer: B

4. **Gradient descent uses backpropagation to:**

- A. Compute the loss
- B. Compute gradients and update weights
- C. Increase loss
- D. Normalize activations

Answer: B

5. **Vanishing gradients cause:**

- A. Rapid learning in all layers
- B. Early layers to stop learning
- C. Exploding weights
- D. Overfitting

Answer: B

6. **Which variant of gradient descent updates weights after seeing all training data?**

- A. Mini-batch
- B. Stochastic

- C. Full-batch
- D. Adam

Answer: C

7. **Caching intermediate values during the forward pass allows:**

- A. Faster gradient computation
- B. Slower training
- C. No effect on speed
- D. Overfitting

Answer: A

8. **Exploding gradients are usually addressed by:**

- A. Increasing learning rate
- B. Gradient clipping
- C. Removing hidden layers
- D. Using sigmoid everywhere

Answer: B

9. **Which of the following is necessary for backpropagation to work?**

- A. Non-differentiable operations
- B. Differentiable operations
- C. Random weights
- D. Large batch sizes

Answer: B

10. **Which of the following best describes the backward pass?**

- A. Computes network output
- B. Propagates gradients from output to input
- C. Applies dropout
- D. Normalizes weights

Answer: B

11. **Mini-batch gradient descent is preferred over full-batch because:**

- A. It is always slower
- B. Balances noise and efficiency
- C. Uses all data at once
- D. Never converges

Answer: B

12. **If gradients are NaN during training, a likely cause is:**

- A. Vanishing gradients
- B. Exploding gradients

- C. Small learning rate
- D. Batch normalization

Answer: B

13. **The chain rule in backpropagation allows:**

- A. Layer-by-layer computation of gradients
- B. Only output layer gradients
- C. Ignoring hidden layers
- D. Direct computation of weights

Answer: A

14. **Which is NOT a typical cause of vanishing gradients?**

- A. Many layers with sigmoid activation
- B. ReLU activations
- C. Small weights
- D. Deep architectures

Answer: B

15. **Gradient clipping is used to:**

- A. Prevent weights from growing too large
- B. Normalize inputs
- C. Reduce overfitting
- D. Increase learning rate

Answer: A

16. **What happens if you skip the backward pass in training?**

- A. Model still learns
- B. Weights are never updated
- C. Gradients are computed
- D. Loss is minimized

Answer: B

17. **Which loss is commonly minimized in neural network training?**

- A. Accuracy
- B. Cross-entropy
- C. Precision
- D. ROC AUC

Answer: B

18. **How does batch normalization indirectly help backpropagation?**

- A. Increases gradients
- B. Reduces internal covariate shift, stabilizing gradients

- C. Adds noise
- D. Removes bias

Answer: B

19. **Which is a benefit of mini-batch training?**

- A. Faster convergence and parallelization
- B. Always lower final loss
- C. Increased overfitting
- D. No effect on speed

Answer: A

20. **If your network is not learning and gradients are zero, likely cause is:**

- A. Exploding gradients
- B. Vanishing gradients
- C. High learning rate
- D. Dropout

Answer: B

21. **Which optimizer is based on gradient descent and uses backpropagation?**

- A. Adam
- B. Random Forest
- C. KNN
- D. PCA

Answer: A

22. **Backpropagation is most critical for:**

- A. Training neural networks
- B. Training SVMs
- C. Clustering
- D. Feature engineering

Answer: A

23. **Which of the following is NOT a challenge in backpropagation for deep networks?**

- A. Vanishing gradients
- B. Exploding gradients
- C. Efficient computation
- D. Memory usage

Answer: C

24. **Which step immediately follows computation of gradients in backpropagation?**

- A. Activation
- B. Weight update
- C. Forward pass
- D. Loss computation

Answer: B

25. **Backpropagation exploits which computational property for efficiency?**

- A. Randomization
- B. Reuse of intermediate results
- C. Only output gradients
- D. Linear activations

Answer: B

26. **Which activation is most likely to cause vanishing gradients?**

- A. Sigmoid
- B. ReLU
- C. Linear
- D. Softmax

Answer: A

Lecture 17: Support Vector Machines (SVM)

Analytical Review, Short Notes, Analytical Questions, and 25+ Reasoned MCQs

Bulleted Short Notes

- **Support Vector Machines (SVM):**

- Supervised learning models for classification and regression.
- Aim to find the optimal hyperplane that maximally separates classes in feature space.
- Support vectors: data points closest to the hyperplane, most influential in defining the decision boundary.
- Margin: distance between hyperplane and nearest data points (support vectors); SVM maximizes this margin.
- **Hard margin:** No misclassification allowed (only works if data is perfectly separable).
- **Soft margin:** Allows some misclassification (uses regularization parameter C to balance margin size and misclassification).

- **Kernels:**

- Allow SVMs to perform non-linear classification by mapping data into higher-dimensional space.
- Common kernels: linear, polynomial, radial basis function (RBF), sigmoid.

- **Advantages:**

- Effective in high-dimensional spaces.
- Can model complex (non-linear) boundaries with kernels.
- Only support vectors matter—efficient use of data.

- **Disadvantages:**

- Not ideal for very large datasets (slow training).
 - Choice of kernel and tuning parameters (C , γ) is crucial.
 - Less interpretable than simpler models.
-

Analytical Short Questions (with Solutions)

Q1. What is the main goal of an SVM in classification tasks?

A:

To find the hyperplane that maximizes the margin between different classes, thereby achieving the best possible separation.

Q2. Why are support vectors important in SVM?

A:

They are the data points closest to the decision boundary and entirely determine the position and orientation of the hyperplane.

Q3. How does the regularization parameter C affect the SVM decision boundary?

A:

A small C increases the margin but allows more misclassifications (regularization), while a large C tries to classify all points correctly but may lead to overfitting.

Q4. What is the role of the kernel in SVM?

A:

The kernel enables SVM to perform non-linear classification by implicitly mapping input features into higher-dimensional space.

Q5. Give an example when you would choose an RBF kernel over a linear kernel.

A:

When the data is not linearly separable, and you suspect that a non-linear boundary is needed for good classification.

Q6. What is the margin in SVM, and why does maximizing it help generalization?

A:

The margin is the distance from the hyperplane to the nearest data points (support vectors). Maximizing it reduces the model's sensitivity to noise and improves generalization.

Q7. Why might SVMs be less suitable for very large datasets?

A:

Because training complexity depends on the number of data points, and SVMs require solving a quadratic optimization problem, which is computationally expensive for large datasets.

Q8. How does SVM handle multi-class classification?

A:

Via strategies like one-vs-rest or one-vs-one, combining several binary SVM classifiers.

Q9. What is the effect of a very high gamma value in an RBF kernel?

A:

The decision boundary becomes very sensitive to individual data points, leading to overfitting.

Q10. How can you select the best kernel and parameters for an SVM?

A:

Use cross-validation to compare performance of different kernels and hyperparameter values (C , γ).

Analytical MCQs (with Answers & Reasoning)

1. **The primary goal of SVM is to:**

- A. Minimize classification error only
- B. Maximize margin between classes

- C. Reduce number of features
- D. Increase number of support vectors

Answer: B

2. **Support vectors in SVM are:**

- A. All training samples
- B. Points closest to the decision boundary
- C. Randomly chosen points
- D. Outliers

Answer: B

3. **Which parameter controls the trade-off between maximizing margin and minimizing classification error in SVM?**

- A. Kernel
- B. Gamma
- C. C
- D. Learning rate

Answer: C

4. **The RBF kernel is useful when:**

- A. Data is linearly separable
- B. Data requires a non-linear decision boundary
- C. Features are categorical
- D. Only regression is needed

Answer: B

5. **A very high value of C in SVM leads to:**

- A. Wider margin, more misclassification
- B. Narrower margin, less misclassification
- C. Always underfitting
- D. No effect

Answer: B

6. **Kernel trick in SVM allows:**

- A. Faster training
- B. Nonlinear classification without explicit mapping
- C. More interpretable models
- D. Always linear decision boundaries

Answer: B

7. **Which kernel is best for linearly separable data?**

- A. Linear

- B. Polynomial
- C. RBF
- D. Sigmoid

Answer: A

8. **SVMs are less suitable for:**

- A. High-dimensional data
- B. Large datasets
- C. Data with clear margins
- D. Binary classification

Answer: B

9. **SVM handles multi-class classification by:**

- A. Single SVM
- B. One-vs-rest or one-vs-one strategies
- C. Only with linear kernel
- D. Using clustering

Answer: B

10. **The margin in SVM refers to:**

- A. Total number of features
- B. Distance between hyperplane and nearest points
- C. Learning rate
- D. Number of kernels

Answer: B

11. **The output of an SVM for a new sample is:**

- A. Probability score
- B. Distance from the hyperplane (can be converted to class label)
- C. Always 0 or 1
- D. Regression coefficient

Answer: B

12. **Which of the following is NOT a kernel function?**

- A. RBF
- B. Linear
- C. Decision tree
- D. Polynomial

Answer: C

13. **A soft margin SVM allows:**

- A. No misclassification

- B. Some misclassification
- C. Only linear separation
- D. No support vectors

Answer: B

14. Which method is typically used to tune SVM hyperparameters?

- A. Grid search with cross-validation
- B. Random guessing
- C. Manual selection
- D. Use defaults only

Answer: A

15. Which is an advantage of SVM?

- A. Handles high-dimensional data well
- B. Needs large datasets for good performance
- C. Always interpretable
- D. Only works for regression

Answer: A

16. For non-linear SVM, the kernel function:

- A. Must be linear
- B. Maps data to higher dimensions
- C. Is only for regression
- D. Increases model interpretability

Answer: B

17. Which SVM parameter controls the flexibility of the RBF kernel?

- A. C
- B. gamma
- C. Margin
- D. Alpha

Answer: B

18. The “kernel trick” avoids:

- A. Explicit computation in high-dimensional space
- B. Parameter tuning
- C. Support vector calculation
- D. Overfitting

Answer: A

19. Which is a disadvantage of SVMs?

- A. Not effective in high-dimensional space

- B. Training is slow for large datasets
- C. Cannot use kernels
- D. Only works for two classes

Answer: B

20. **Which scenario is NOT ideal for SVMs?**

- A. Small dataset, high dimension
- B. Large dataset, low dimension
- C. Non-linear class boundaries
- D. Few features, clear margin

Answer: B

21. **Which is a typical use of SVMs?**

- A. Text classification
- B. Image recognition
- C. Bioinformatics
- D. All of the above

Answer: D

22. **Increasing gamma in RBF kernel SVM:**

- A. Makes boundary smoother
- B. Makes boundary more sensitive to data (risk of overfitting)
- C. Has no effect
- D. Always improves accuracy

Answer: B

23. **Which is NOT true about support vectors?**

- A. They define the decision boundary
- B. All data points are support vectors
- C. They are closest to the hyperplane
- D. Removing them changes the boundary

Answer: B

24. **SVMs can be used for:**

- A. Classification only
- B. Regression only
- C. Both classification and regression
- D. Clustering

Answer: C

25. **Which kernel is most commonly used for complex, non-linear problems?**

- A. RBF

- B. Linear
- C. Polynomial
- D. Sigmoid

Answer: A

26. **The main computational challenge with SVMs is:**

- A. Memory and training time scale poorly with number of samples
- B. Lack of kernel functions
- C. Only works for small feature spaces
- D. Always underfits

Answer: A

SVM Part-I (Support Vector Machines) – Critical Analysis, Analytical Questions, and MCQs

Slidewise Critical Analysis and Analytical Short Questions (with Answers)

Slide 1: Introduction to SVM

- **Critical Analysis:**
Introduces the concept and context of SVM as a supervised learning algorithm used for classification (and sometimes regression). SVM is known for its ability to find an optimal boundary/hyperplane between classes.
 - **Analytical Question:**
Q: Why is SVM considered a powerful classifier compared to simple linear classifiers?
A: Because SVM finds the optimal hyperplane that maximizes the margin, leading to better generalization on unseen data.
-

Slide 2: Linear Separability

- **Critical Analysis:**
Discusses the assumption that data classes can be separated by a straight line (2D) or hyperplane (higher dimensions). Not all datasets are linearly separable.
 - **Analytical Question:**
Q: What is linear separability, and why is it important in the context of SVM?
A: Linear separability means that a single straight line (or hyperplane) can separate all classes. SVM relies on this property for its basic formulation but can be extended for non-linear cases.
-

Slide 3: The SVM Classifier/Hyperplane

- **Critical Analysis:**
Defines the mathematical form of the hyperplane and the concept of classifying data points using the sign of a linear function.
 - **Analytical Question:**
Q: How does the sign of the decision function relate to class membership in SVM?
A: If $(w^T x + b > 0)$, the point is assigned to one class; if (< 0) , to the other class.
-

Slide 4: Margin and Support Vectors

- **Critical Analysis:**
Introduces the concept of the margin (distance between hyperplane and nearest points) and support vectors (points on the margin).
 - **Analytical Question:**
Q: Why are support vectors crucial for defining the SVM decision boundary?
A: Because they are the only points that affect the position and orientation of the optimal hyperplane.
-

Slide 5: Maximizing the Margin

- **Critical Analysis:**
Explains that maximizing the margin minimizes the model's generalization error and leads to a more robust classifier.
 - **Analytical Question:**
Q: What is the relationship between margin size and model generalization?
A: Larger margins generally lead to better generalization and less overfitting.
-

Slide 6: Optimization Problem Formulation

- **Critical Analysis:**
Presents the hard-margin SVM optimization: minimize $(|w|^2)$ subject to correct classification constraints.

- **Analytical Question:**

Q: Why is minimizing $(|w|^2)$ equivalent to maximizing the margin?

A: The margin is inversely proportional to $(|w|)$; minimizing $(|w|^2)$ increases the margin.

Slide 7: Lagrangian and Dual Formulation

- **Critical Analysis:**

SVM uses Lagrangian multipliers to move from the primal to the dual problem, which is easier to solve and enables the kernel trick.

- **Analytical Question:**

Q: What is the benefit of solving the dual formulation in SVM?

A: It simplifies computation (especially for high-dimensional data) and allows the use of kernels for non-linear separation.

Slide 8: Karush-Kuhn-Tucker (KKT) Conditions

- **Critical Analysis:**

Details the optimality conditions for the SVM solution; KKT conditions help identify support vectors.

- **Analytical Question:**

Q: How do KKT conditions help in identifying support vectors in SVM?

A: Support vectors correspond to data points where the Lagrange multipliers are non-zero (KKT complementary slackness).

Slide 9: Non-separable Case and Soft Margin

- **Critical Analysis:**

Introduces slack variables and the regularization parameter (C) to allow misclassifications (soft margin) for non-separable data.

- **Analytical Question:**

Q: How does the value of (C) affect the SVM decision boundary in soft-margin SVM?

A: Large (C) penalizes misclassification more, leading to smaller margin; small (C) allows more misclassifications and a wider margin.

Slide 10: Summary and Next Steps

- **Critical Analysis:**
Recaps the main points and typically hints at kernel methods for handling non-linear data (to be discussed in Part II).
- **Analytical Question:**
Q: Why are kernel methods a natural extension after understanding linear SVM?
A: Because real-world data is often not linearly separable, and kernels allow SVM to find optimal hyperplanes in transformed (higher-dimensional) spaces.

25+ Analytical MCQs (with Reasoning and Answers)

1. **SVM seeks a decision boundary that:**
A. Minimizes classification error on training data
B. Maximizes margin between classes
C. Minimizes number of support vectors
D. Maximizes data dimensionality
Answer: B. (Wider margin improves generalization.)
2. **Support vectors are:**
A. Points farthest from the hyperplane
B. Points closest to the hyperplane
C. All points in the dataset
D. Points inside the margin only
Answer: B. (Support vectors define the margin.)
3. **The margin in SVM is:**
A. The sum of distances to all points
B. The width between the closest points of opposite classes
C. The minimum distance from hyperplane to support vectors
D. The number of misclassified points
Answer: C.
4. **In the SVM optimization, minimizing ($|w|^2$) leads to:**
A. Narrower margin
B. Wider margin
C. More support vectors

D. Increased overfitting

Answer: B.

5. **The dual formulation benefits SVM by:**

- A. Reducing computation for small datasets
- B. Enabling use of kernels for non-linear data
- C. Increasing primal constraints
- D. Ignoring support vectors

Answer: B.

6. **If all data points satisfy their margin constraints, their Lagrange multipliers are:**

- A. Zero
- B. Non-zero
- C. Negative
- D. Infinite

Answer: A. (Only support vectors have non-zero multipliers.)

7. **The regularization parameter (C) in soft-margin SVM:**

- A. Controls margin width only
- B. Penalizes margin violations
- C. Is always set to zero
- D. Does not affect optimization

Answer: B.

8. **Larger values of (C) in SVM result in:**

- A. More misclassification
- B. Fewer misclassifications but narrower margin
- C. Increased margin width
- D. Ignoring support vectors

Answer: B.

9. **KKT conditions are used to:**

- A. Find dual variables
- B. Guarantee optimality of the SVM solution
- C. Count support vectors
- D. Standardize features

Answer: B.

10. **If data is not linearly separable, SVM can:**

- A. Only use hard margin
- B. Use soft margin or kernel trick

- C. Not be applied
- D. Require data removal

Answer: B.

11. **A hyperplane is defined by:**

- A. (w) and (b) such that $(w^T x + b = 0)$
- B. The distance between two random points
- C. The sum of all feature vectors
- D. The product of all feature vectors

Answer: A.

12. **Which points influence the SVM solution?**

- A. All points
- B. Only support vectors
- C. Only misclassified points
- D. Only points far from margin

Answer: B.

13. **The SVM decision function for a new point uses:**

- A. Only the bias term
- B. Weighted sum over support vectors
- C. Sum of all Lagrange multipliers
- D. The mean of all points

Answer: B.

14. **If a data point has a slack variable $(\xi_i > 0)$, it is:**

- A. Correctly classified with margin
- B. On the margin
- C. Violating the margin
- D. Not part of the solution

Answer: C.

15. **Why is maximizing the margin desirable in SVM?**

- A. It reduces computational cost
- B. It improves generalization to unseen data
- C. It increases overfitting
- D. It guarantees perfect accuracy

Answer: B.

16. **When is the SVM solution unique?**

- A. When data is degenerate
- B. When support vectors are unique and not collinear

- C. When all points are support vectors
- D. When $C = 0$

Answer: B.

17. The kernel trick allows SVMs to:

- A. Reduce dimensionality
- B. Operate in high-dimensional feature spaces without explicit mapping
- C. Ignore support vectors
- D. Maximize training error

Answer: B.

18. Which statement about SVMs is FALSE?

- A. They can be used for regression
- B. They are sensitive to feature scaling
- C. They always require kernel methods
- D. They produce sparse solutions

Answer: C.

19. In the dual SVM formulation, the solution is expressed in terms of:

- A. Feature weights only
- B. Lagrange multipliers and support vectors
- C. Only the bias
- D. Only the data means

Answer: B.

20. The sign of the SVM decision function for a point (x) determines:

- A. Number of support vectors
- B. Class assignment of (x)
- C. Value of C
- D. Margin width

Answer: B.

21. Which scenario requires soft-margin SVM?

- A. Perfectly separable data
- B. Data with overlap/noise
- C. Data with very high dimensions
- D. Data with only one class

Answer: B.

22. Support vectors have what property regarding the margin?

- A. They lie outside the margin
- B. They are at or within the margin

- C. They are always misclassified
- D. They are always at zero slack

Answer: B.

23. **The bias term, (b), in SVM:**

- A. Moves the hyperplane closer to the origin
- B. Shifts the hyperplane to fit data better
- C. Is always zero
- D. Has no effect

Answer: B.

24. **Which best describes the SVM loss function for hard margin?**

- A. Hinge loss
- B. Zero-one loss
- C. Exponential loss
- D. Quadratic loss

Answer: A.

25. **A point with Lagrange multiplier ($\alpha_i = 0$):**

- A. Is a support vector
- B. Is not a support vector
- C. Is always misclassified
- D. Must be at the origin

Answer: B.

26. **The main constraint in SVM optimization ensures:**

- A. Points are maximally distant from hyperplane
- B. Points are on correct side of margin
- C. Points are in same class
- D. All data is normalized

Answer: B

PCA (Principal Component Analysis) – Critical Analysis, Analytical Questions, and Analytical MCQs

Slidewise Critical Analysis and Analytical Short Questions (with Answers)

Slide 1: Introduction to PCA

- **Critical Analysis:**
PCA is a statistical method for dimensionality reduction, data visualization, and noise filtering, widely used in exploratory data analysis and preprocessing for machine learning.
 - **Analytical Question:**
Q: What are the main goals of performing PCA on a dataset?
A: To reduce dimensionality while retaining as much variance as possible, and to discover new uncorrelated features (principal components).
-

Slide 2: The Need for Dimensionality Reduction

- **Critical Analysis:**
High-dimensional data can suffer from the "curse of dimensionality," making analysis, visualization, and modeling challenging.
 - **Analytical Question:**
Q: Why might high-dimensional data lead to poor model performance?
A: High-dimensional data increases computational cost and risk of overfitting, and may contain redundant or irrelevant features.
-

Slide 3: Mathematical Formulation of PCA

- **Critical Analysis:**
PCA projects data onto orthogonal directions (principal components) that maximize variance, typically found by eigendecomposition of the covariance matrix.
 - **Analytical Question:**
Q: How are principal components mathematically determined?
A: As the eigenvectors of the data's covariance matrix, ordered by their corresponding eigenvalues.
-

Slide 4: Variance Maximization

- **Critical Analysis:**
The first principal component is the direction of maximum variance; each subsequent component is orthogonal to previous ones and captures the next highest variance.
 - **Analytical Question:**
Q: Why does PCA choose directions of maximum variance?
A: Because directions of maximum variance retain the most information in the data.
-

Slide 5: Computation Steps in PCA

- **Critical Analysis:**
The PCA process includes mean-centering the data, calculating the covariance matrix, performing eigendecomposition/SVD, and projecting data onto selected components.
 - **Analytical Question:**
Q: Why is it important to center the data before applying PCA?
A: Centering ensures that the first principal component aligns with the direction of maximum variance, not biased by the mean.
-

Slide 6: Choosing the Number of Components

- **Critical Analysis:**
The number of components is often chosen so that a specified proportion of variance (e.g., 95%) is retained.

- **Analytical Question:**
Q: How do you decide the optimal number of principal components to retain?
A: By analyzing the explained variance ratio and selecting enough components to reach a desired threshold (e.g., via scree plot).
-

Slide 7: Geometric Interpretation

- **Critical Analysis:**
PCA finds a new coordinate system where the axes (PCs) are the directions of greatest variance, and data can be projected onto this reduced space.
 - **Analytical Question:**
Q: What is the geometric meaning of projecting data onto the first principal component?
A: It means representing each data point by its coordinate along the direction of maximum variance.
-

Slide 8: Applications of PCA

- **Critical Analysis:**
PCA is applied in image compression, data visualization, noise reduction, and as a preprocessing step for other algorithms.
 - **Analytical Question:**
Q: Give two practical scenarios where PCA is beneficial.
A: Visualizing high-dimensional gene expression data; compressing images by reducing pixel dimensions.
-

Slide 9: Limitations of PCA

- **Critical Analysis:**
PCA assumes linear relationships and is sensitive to scaling, outliers, and may not capture non-linear structure.
- **Analytical Question:**
Q: Why might PCA fail to identify important features in some datasets?
A: Because it only captures linear correlations and can be misled by outliers or differing scales.

Slide 10: Summary and Extensions

- **Critical Analysis:**
Summarizes the strengths and weaknesses of PCA and may introduce extensions (kernel PCA, sparse PCA) for non-linear or structured data.
- **Analytical Question:**
Q: What is kernel PCA and when should it be used?
A: Kernel PCA is a non-linear extension that uses kernel functions to capture complex structure in data not captured by linear PCA.

25+ Analytical MCQs (with Reasoning and Answers)

1. **The main objective of PCA is to:**
A. Maximize the number of dimensions
B. Minimize data variance
C. Reduce the number of features while preserving variance
D. Find clusters in the data
Answer: C
2. **Principal components are:**
A. Original variables
B. Linear combinations of original variables
C. Random projections
D. Always equal in number to features
Answer: B
3. **The first principal component:**
A. Has the smallest variance
B. Is always the same as the first feature
C. Captures the direction of maximum variance
D. Is orthogonal to the last component only
Answer: C
4. **PCA relies on computing:**
A. The covariance matrix
B. The distance matrix
C. The correlation matrix only

D. The mean vector only

Answer: A

5. **If data is not centered before PCA:**

- A. Results may be biased by the mean
- B. Principal components remain unchanged
- C. Covariance matrix is not needed
- D. PCA cannot be performed

Answer: A

6. **Which of the following is NOT a typical PCA application?**

- A. Feature extraction
- B. Supervised classification
- C. Data visualization
- D. Noise reduction

Answer: B

7. **The eigenvalues in PCA represent:**

- A. The amount of variance captured by each component
- B. The direction of the PCs
- C. The mean of the data
- D. The number of features

Answer: A

8. **How is the optimal number of PCs usually determined?**

- A. By the number of features
- B. By cumulative explained variance
- C. By PCA always using two components
- D. By the number of observations

Answer: B

9. **PCA is sensitive to:**

- A. Scaling of features
- B. Outliers
- C. Both A and B
- D. The number of observations only

Answer: C

10. **The principal components are always:**

- A. Correlated
- B. Orthogonal (uncorrelated)
- C. Equal to original features

D. Chosen arbitrarily

Answer: B

11. If the first two PCs explain 85% of the variance, this means:

- A. The remaining PCs are not needed for most analyses
- B. Data can safely be projected onto two dimensions
- C. Both A and B
- D. All PCs must be retained

Answer: C

12. Which step is NOT part of PCA computation?

- A. Centering the data
- B. Calculating covariance matrix
- C. Performing SVD or eigendecomposition
- D. Normalizing eigenvalues to sum to one

Answer: D

13. Kernel PCA is used when:

- A. Data relationships are non-linear
- B. Data is always linearly separable
- C. Data is categorical
- D. Covariance matrix is not computable

Answer: A

14. Which property of PCA makes it useful for visualization?

- A. It reduces data to two or three dimensions while retaining variance
- B. It increases the number of features
- C. It always clusters data
- D. It removes all outliers

Answer: A

15. Which of the following is TRUE about the directions of principal components?

- A. Each PC is orthogonal to all others
- B. Each PC must be parallel
- C. Only the first PC is orthogonal
- D. PCs are chosen at random

Answer: A

16. A major limitation of PCA is:

- A. It cannot reduce data
- B. It cannot handle linear data

- C. It cannot capture non-linear structure
- D. It always increases variance

Answer: C

17. **PCA will perform poorly if:**

- A. Variables are measured on different scales and not standardized
- B. Variables are already uncorrelated
- C. Data has no noise
- D. There is only one feature

Answer: A

18. **Scree plot is used in PCA to:**

- A. Visualize the proportion of variance explained by each component
- B. Plot original data
- C. Show the mean of the data
- D. Identify missing data

Answer: A

19. **Which of the following best explains why PCA can be used for noise reduction?**

- A. Noise often appears in lower-variance components, which can be discarded
- B. Noise is always in the first component
- C. PCA amplifies noise
- D. All PCs represent noise

Answer: A

20. **The loading vector in PCA indicates:**

- A. How much each original variable contributes to a principal component
- B. The mean of each variable
- C. The number of samples
- D. The scaling applied

Answer: A

21. **If two features are highly correlated, PCA will:**

- A. Combine them into a single principal component
- B. Remove one feature
- C. Ignore the correlation
- D. Increase their variance

Answer: A

22. **PCA assumes the principal directions are those that:**

- A. Minimize variance

- B. Maximize variance
- C. Are parallel to the axes
- D. Are random

Answer: B

23. If the cumulative explained variance of the first k PCs is only 50%, this suggests:

- A. More PCs are needed to adequately represent the data
- B. The data is low dimensional
- C. PCA is not needed
- D. The first PC suffices

Answer: A

24. Sparse PCA is used when:

- A. Interpretability of PCs is desired
- B. PCs must load on all variables
- C. Data is categorical
- D. Only non-linear patterns are present

Answer: A

25. Which of the following is NOT a step in standard PCA?

- A. Standardize or center the data
- B. Compute the covariance matrix
- C. Project data onto principal components
- D. Cluster data using k-means

Answer: D

26. The transformation from the original space to the PCA space is:

- A. Non-invertible
- B. Linear
- C. Non-linear
- D. Unsupervised

Answer: B

PCA Part-II & Clustering – Critical Analysis, Analytical Questions, and Analytical MCQs

Slidewise Critical Analysis and Analytical Short Questions (with Answers)

Slide 1: Introduction to Clustering & PCA-II

- **Critical Analysis:**
Introduces the bridge between dimensionality reduction (PCA) and unsupervised learning (clustering). Highlights why clustering is meaningful after reducing dimensionality.
 - **Analytical Question:**
Q: Why is PCA commonly used before clustering in data analysis pipelines?
A: PCA reduces noise and redundancy, making clusters more distinct and improving clustering performance.
-

Slide 2: Clustering Overview

- **Critical Analysis:**
Explains the goal of clustering: grouping similar data points without labels. May introduce the concept of similarity/distance metrics.
 - **Analytical Question:**
Q: What is the main challenge in clustering high-dimensional data?
A: High dimensionality can obscure natural groupings due to the curse of dimensionality and increased noise.
-

Slide 3: Types of Clustering Algorithms

- **Critical Analysis:**
Lists main clustering approaches: partitioning (e.g., k-means), hierarchical, density-based, etc. Each has unique strengths and weaknesses.
 - **Analytical Question:**
Q: Which clustering method is best suited for finding arbitrarily shaped clusters?
A: Density-based methods (e.g., DBSCAN) are effective for arbitrarily shaped clusters.
-

Slide 4: K-Means Clustering

- **Critical Analysis:**
Describes k-means as a popular, efficient clustering method that partitions data into k clusters by minimizing within-cluster variance.
 - **Analytical Question:**
Q: Why is the choice of initial centroids important in k-means clustering?
A: Poor initialization can lead to suboptimal or inconsistent clustering results.
-

Slide 5: The Role of PCA in Clustering

- **Critical Analysis:**
Shows how PCA simplifies the feature space, improving clustering quality and computational efficiency.
 - **Analytical Question:**
Q: How does PCA help improve the performance of clustering algorithms?
A: By reducing dimensionality, PCA removes noise and redundant features, making clusters more separable.
-

Slide 6: PCA for Visualization of Clusters

- **Critical Analysis:**
Emphasizes using the first two or three principal components to project high-dimensional data for visualization.

- **Analytical Question:**
Q: Why is PCA useful for cluster visualization?
A: It enables meaningful 2D/3D projections where cluster structure can be visually assessed.
-

Slide 7: Clustering Evaluation Metrics

- **Critical Analysis:**
Introduces internal (e.g., Silhouette Score, Davies-Bouldin) and external (e.g., Purity, ARI) evaluation metrics for clustering results.
 - **Analytical Question:**
Q: What does a high silhouette score indicate about clustering quality?
A: That clusters are well-separated and data points are close to their own cluster center.
-

Slide 8: Limitations and Practical Tips

- **Critical Analysis:**
Discusses limitations, such as sensitivity to outliers, need for scaling, and the risk of losing important information during PCA.
 - **Analytical Question:**
Q: What is a potential drawback of applying PCA before clustering?
A: Important clustering information might be lost if too much variance is discarded.
-

Slide 9: Real-world Applications

- **Critical Analysis:**
Provides examples where PCA and clustering together enable meaningful insights (e.g., customer segmentation, gene expression analysis).
 - **Analytical Question:**
Q: Give an example where combining PCA and clustering is especially valuable.
A: In genetics, where high-dimensional gene expression data is reduced by PCA for clearer clustering of cell types.
-

Slide 10: Summary and Outlook

- **Critical Analysis:**
Summarizes the synergy between PCA and clustering, and hints at advanced topics (e.g., kernel PCA, spectral clustering).
 - **Analytical Question:**
Q: What is one way to address non-linear structures in clustering after PCA?
A: Use non-linear dimensionality reduction (e.g., kernel PCA) or clustering methods that can handle complex shapes (e.g., DBSCAN).
-

25+ Analytical MCQs (with Reasoning and Answers)

1. **The main purpose of applying PCA before clustering is to:**
A. Increase the number of features
B. Reduce noise and redundancy
C. Merge all clusters into one
D. Remove all data labels
Answer: B
2. **Which clustering method is most sensitive to initial conditions?**
A. Hierarchical
B. DBSCAN
C. K-means
D. Spectral
Answer: C
3. **A high silhouette score indicates:**
A. Poorly defined clusters
B. Overlapping clusters
C. Well-separated, coherent clusters
D. Random grouping
Answer: C
4. **What is a limitation of k-means clustering?**
A. It can only find circular (convex) clusters
B. It works for all distance metrics
C. It requires labeled data
D. It does not require initialization
Answer: A

5. **PCA is particularly useful for clustering when:**

- A. Data is low-dimensional
- B. Features are highly correlated
- C. Data has no noise
- D. All features are categorical

Answer: B

6. **Which metric is internal to clustering evaluation?**

- A. Adjusted Rand Index
- B. Silhouette Score
- C. Purity
- D. NMI

Answer: B

7. **DBSCAN is preferred over k-means when:**

- A. Clusters are of similar size
- B. Clusters have irregular shapes
- C. Data is one-dimensional
- D. Data is perfectly separated

Answer: B

8. **If PCA is used for visualization, how many components are typically selected?**

- A. 1
- B. 2 or 3
- C. All
- D. 10

Answer: B

9. **When could PCA harm clustering performance?**

- A. When important variance for clusters is discarded
- B. When all variance is retained
- C. When data is visualized
- D. When clusters are very distinct

Answer: A

10. **Hierarchical clustering differs from k-means by:**

- A. Not requiring pre-specified cluster number
- B. Being non-iterative
- C. Producing a dendrogram

D. All of the above

Answer: D

11. Which of the following is not a clustering algorithm?

A. k-means

B. DBSCAN

C. Agglomerative

D. PCA

Answer: D

12. The principal components used for clustering are chosen based on:

A. Maximum explained variance

B. Minimum eigenvalue

C. Random selection

D. Feature names

Answer: A

13. If clusters overlap after PCA, a possible solution is to:

A. Increase the number of components

B. Try a different clustering algorithm

C. Use non-linear dimensionality reduction

D. Any of the above

Answer: D

14. The curse of dimensionality refers to:

A. Increased noise and reduced cluster separability in high dimensions

B. Decreased computation cost

C. More interpretable clusters

D. Fewer outliers

Answer: A

15. Cluster purity is an example of:

A. Internal metric

B. External metric (requires ground truth)

C. Distance function

D. Initialization method

Answer: B

16. Spectral clustering is often used when:

A. Data is linearly separable

B. Clusters have complex structures

C. There is only one cluster

D. PCA is not performed

Answer: B

17. After PCA, clustering may be improved because:

A. Less noise and fewer irrelevant variables

B. Data is always perfectly separable

C. All clusters are spherical

D. More clusters are found

Answer: A

18. Which step is not typical in a PCA-clustering workflow?

A. Standardizing features

B. Applying PCA

C. Assigning clusters

D. Labeling clusters with known classes

Answer: D

19. The Davies-Bouldin Index is best interpreted as:

A. Lower values indicate better clustering

B. Higher values are better

C. Used only for supervised learning

D. A measure of explained variance

Answer: A

20. If your data contains outliers, which clustering method is more robust?

A. K-means

B. DBSCAN

C. Hierarchical

D. K-means++

Answer: B

21. Which of the following can be an artifact of reducing dimensions too aggressively with PCA before clustering?

A. Loss of cluster structure

B. Improved cluster purity

C. More clusters

D. Perfect separation

Answer: A

22. A dendrogram is most associated with:

A. Hierarchical clustering

B. K-means

C. DBSCAN

D. PCA

Answer: A

23. **Why should data be scaled before PCA and clustering?**

A. To ensure features contribute equally

B. To improve computational speed

C. To remove outliers

D. To increase variance

Answer: A

24. **If you want to identify the optimal number of clusters, you might use:**

A. Scree plot

B. Elbow method

C. PCA loadings

D. Covariance matrix

Answer: B

25. **Combining PCA and clustering is especially useful when:**

A. The dataset is high-dimensional and redundant

B. All features are binary

C. Only supervised learning is required

D. There is only one cluster

Answer: A

26. **Which clustering evaluation metric does NOT require ground truth labels?**

A. Silhouette Score

B. Adjusted Rand Index

C. Purity

D. NMI

Answer: A

Digital Image Forensics – Critical Analysis, Analytical Questions, and Analytical MCQs

Slidewise Critical Analysis and Analytical Short Questions (with Answers)

Slide 1: Introduction to Digital Image Forensics

- **Critical Analysis:**
Introduces the field of digital image forensics, which deals with detecting manipulation, verifying authenticity, and tracing the origin of digital images.
 - **Analytical Question:**
Q: Why is digital image forensics increasingly important in today's digital world?
A: The ease of image editing and dissemination makes it vital to verify image authenticity for legal, journalistic, and scientific purposes.
-

Slide 2: Types of Image Forgeries

- **Critical Analysis:**
Discusses common types: splicing, copy-move (cloning), retouching, and image synthesis (deepfakes).
 - **Analytical Question:**
Q: How does copy-move forgery differ from splicing?
A: Copy-move uses parts of the same image, while splicing combines parts from different images.
-

Slide 3: Image Manipulation Techniques

- **Critical Analysis:**
Explains common editing methods: cropping, scaling, color enhancement, compression artifacts, and deep learning-based alterations.
 - **Analytical Question:**
Q: Which manipulation technique is often hardest to detect and why?
A: Deep learning-based manipulations (deepfakes) can be very subtle and realistic, making detection challenging.
-

Slide 4: Forensic Analysis Pipeline

- **Critical Analysis:**
Outlines the typical workflow: image acquisition, pre-processing, feature extraction, detection/classification, and reporting.
 - **Analytical Question:**
Q: Why is feature extraction a crucial step in forensic analysis?
A: It helps identify subtle statistical or structural changes caused by manipulation.
-

Slide 5: Passive vs. Active Forensics

- **Critical Analysis:**
Explains passive (no prior info, e.g., statistical analysis) and active (embedded watermarks/fingerprints) approaches.
 - **Analytical Question:**
Q: What is a key limitation of active forensic methods?
A: They require images to be pre-embedded with watermarks or signatures, which is rarely the case for found images.
-

Slide 6: Copy-Move Forgery Detection

- **Critical Analysis:**
Discusses block-based and keypoint-based detection methods, their strengths/weaknesses, and challenges with post-processing.

- **Analytical Question:**
Q: Why are keypoint-based methods robust to simple geometric transformations?
A: They rely on invariant feature descriptors that are stable under rotation and scaling.
-

Slide 7: Splicing Detection

- **Critical Analysis:**
Covers detection techniques such as edge analysis, color inconsistencies, and machine learning.
 - **Analytical Question:**
Q: How can color inconsistencies indicate splicing?
A: Differing lighting or camera characteristics between spliced regions may produce detectable anomalies.
-

Slide 8: Deep Learning in Image Forensics

- **Critical Analysis:**
Examines the role of CNNs and transfer learning in detecting subtle, complex manipulations.
 - **Analytical Question:**
Q: What is a key advantage of deep learning over traditional feature-based methods?
A: Deep learning can automatically discover complex patterns and adapt to new types of manipulations.
-

Slide 9: Evaluation and Challenges

- **Critical Analysis:**
Discusses issues like dataset bias, generalization to unseen manipulations, and adversarial attacks.
 - **Analytical Question:**
Q: Why is generalization a major challenge for forensic detectors?
A: Detectors may overfit to known manipulations and fail to detect novel ones.
-

Slide 10: Legal and Ethical Considerations

- **Critical Analysis:**
Explores implications for law, privacy, and ethics, including potential for both good (authenticity) and harm (false positives).
 - **Analytical Question:**
Q: How can digital forensics tools themselves be misused?
A: By fabricating evidence or falsely accusing individuals based on flawed analysis.
-

25+ Analytical MCQs (with Reasoning and Answers)

1. **Digital image forensics primarily aims to:**
A. Edit images efficiently
B. Detect and analyze image manipulations
C. Increase image resolution
D. Compress images for storage
Answer: B
2. **Which type of forgery involves copying a region within the same image?**
A. Splicing
B. Copy-move
C. Retouching
D. Morphing
Answer: B
3. **Passive forensic methods rely on:**
A. Embedded watermarks
B. Statistical analysis of the image itself
C. Prior knowledge of the device
D. User passwords
Answer: B
4. **Which is a limitation of block-based copy-move detection?**
A. Sensitive to compression noise
B. Cannot handle geometric changes
C. Slow on large images
D. All of the above
Answer: D

5. **Keypoint-based forgery detection is robust because:**

- A. It ignores all transformations
- B. It uses features invariant to scaling and rotation
- C. It works only on color images
- D. It requires original images

Answer: B

6. **Splicing detection may analyze:**

- A. Color inconsistencies
- B. Edges and boundaries
- C. Texture transitions
- D. All of the above

Answer: D

7. **A main challenge for deep learning in image forensics is:**

- A. Overfitting to training manipulations
- B. Lack of computational power
- C. Too simple architecture
- D. Inability to process JPEG

Answer: A

8. **Active forensics typically fails when:**

- A. Watermarks are missing
- B. Images are grayscale
- C. Images are high-resolution
- D. There is no internet connection

Answer: A

9. **Feature extraction in forensics is used to:**

- A. Compress the image
- B. Identify manipulation traces
- C. Restore original images
- D. Blur sensitive regions

Answer: B

10. **Which is a legal risk in digital image forensics?**

- A. Privacy invasion
- B. False positives in detection
- C. Both A and B
- D. Higher storage requirements

Answer: C

11. A deepfake is best described as:

- A. A compressed image
- B. An image synthesized or altered by AI
- C. A watermark
- D. A raw sensor image

Answer: B

12. Which of the following is NOT a forgery detection technique?

- A. Lighting analysis
- B. JPEG artifact analysis
- C. Histogram equalization
- D. Machine learning classification

Answer: C

13. Image resampling detection aims to find:

- A. Changes in image dimensions
- B. Interpolation artifacts from resizing
- C. Color balance errors
- D. Missing metadata

Answer: B

14. A robust forensic method should:

- A. Work only on known manipulations
- B. Generalize to unseen manipulations
- C. Require original camera files
- D. Ignore image content

Answer: B

15. The main advantage of passive forensic techniques is:

- A. They work without prior preparation
- B. They need watermarks
- C. They are always faster
- D. They guarantee detection

Answer: A

16. Forensic tools are vulnerable to adversarial attacks, meaning:

- A. They can be tricked by specially crafted manipulations
- B. They never make mistakes
- C. They are immune to new editing techniques
- D. They always detect splicing

Answer: A

17. **A forensic method that uses camera sensor patterns is called:**

- A. Passive statistical
- B. Sensor pattern noise analysis
- C. Block matching
- D. Histogram equalization

Answer: B

18. **Which is a drawback of deep learning forensics?**

- A. Requires large labeled datasets
- B. Cannot process large images
- C. Is always explainable
- D. Ignores pixel values

Answer: A

19. **Tampering localization aims to:**

- A. Identify where manipulation occurred in the image
- B. Detect image format
- C. Restore original images
- D. Compress images

Answer: A

20. **Which is a sign of JPEG forgery?**

- A. Double compression artifacts
- B. Perfect color balance
- C. High sharpness
- D. No metadata

Answer: A

21. **Which method is best for detecting copy-move forgeries involving rotation?**

- A. Block matching
- B. Keypoint-based methods
- C. Color histogram
- D. Metadata analysis

Answer: B

22. **Splicing often introduces:**

- A. Lighting inconsistencies
- B. Sensor pattern inconsistencies
- C. Edge mismatches
- D. All of the above

Answer: D

23. Which is a practical challenge for real-world image forensics?

- A. Diverse manipulation techniques
- B. Large-scale datasets
- C. Lack of ground truth
- D. All of the above

Answer: D

24. A false positive in forensics implies:

- A. A real forgery missed
- B. An authentic image flagged as fake
- C. Perfect detection
- D. None of the above

Answer: B

25. One ethical concern with forensics tools is:

- A. Misuse for targeting innocent people
- B. High computational cost
- C. Lack of user interface
- D. Training time

Answer: A

26. Which is true about deepfakes?

- A. They are always easy to detect
- B. They can be generated by GANs
- C. They are always illegal
- D. They cannot be detected by forensics

Answer: B

Evaluation Measures 2 – Critical Analysis, Analytical Questions, and Analytical MCQs

Slidewise Critical Analysis and Analytical Short Questions (with Answers)

Slide 1: Introduction to Evaluation Measures (Part 2)

- **Critical Analysis:**
This slide sets up a deeper exploration of evaluating model performance, likely focusing on metrics beyond accuracy, such as those for imbalanced data or regression tasks.
 - **Analytical Question:**
Q: Why might accuracy alone be insufficient for model evaluation?
A: Accuracy does not account for class imbalance and may misrepresent model performance when classes are not evenly distributed.
-

Slide 2: Confusion Matrix Recap

- **Critical Analysis:**
Revisits the confusion matrix, emphasizing the importance of TP, TN, FP, and FN in calculating other metrics.
 - **Analytical Question:**
Q: Which elements of the confusion matrix are essential to compute recall?
A: True Positives (TP) and False Negatives (FN).
-

Slide 3: Precision, Recall, and F1-Score

- **Critical Analysis:**
Explains these metrics and their importance for tasks like information retrieval, medical diagnosis, etc.
 - **Analytical Question:**
Q: When is precision more important than recall?
A: When the cost of a false positive is high, as in spam detection.
-

Slide 4: ROC Curve and AUC

- **Critical Analysis:**
Introduces the ROC curve (TPR vs. FPR) and Area Under the Curve (AUC) as threshold-independent metrics.
 - **Analytical Question:**
Q: What does an AUC of 0.5 indicate about a classifier?
A: The classifier performs no better than random guessing.
-

Slide 5: Precision-Recall Curve

- **Critical Analysis:**
Highlights the usefulness of this curve in imbalanced datasets.
 - **Analytical Question:**
Q: Why can the precision-recall curve be more informative than ROC in imbalanced datasets?
A: Because it focuses on the positive class and is less affected by the abundance of true negatives.
-

Slide 6: Regression Evaluation Metrics

- **Critical Analysis:**
Covers metrics like Mean Absolute Error (MAE), Mean Squared Error (MSE), and R2 score.

- **Analytical Question:**
Q: What does the R^2 score represent in regression?
A: The proportion of variance in the dependent variable explained by the model.
-

Slide 7: Cohen's Kappa and Matthews Correlation Coefficient (MCC)

- **Critical Analysis:**
Presents advanced evaluation measures that account for chance agreement and provide balanced assessment, especially for imbalanced data.
 - **Analytical Question:**
Q: What is the advantage of MCC over accuracy in imbalanced datasets?
A: MCC provides a more balanced measure even when classes are of very different sizes.
-

Slide 8: Metric Selection Guidelines

- **Critical Analysis:**
Discusses how to choose the right metric based on the problem, data distribution, and business needs.
 - **Analytical Question:**
Q: What should guide the choice of evaluation metric?
A: The problem's context, data characteristics, and the cost of different types of errors.
-

Slide 9: Common Pitfalls in Model Evaluation

- **Critical Analysis:**
Warns against overreliance on a single metric, improper validation, or ignoring data leakage.
 - **Analytical Question:**
Q: What is a common pitfall when evaluating a model on the same data used for training?
A: It leads to overfitting and overestimation of model performance.
-

Slide 10: Summary and Best Practices

- **Critical Analysis:**
Summarizes key points and offers best practices: use multiple metrics, cross-validate, and always consider the business context.
 - **Analytical Question:**
Q: Why is it good practice to report multiple evaluation metrics?
A: It provides a more comprehensive view of model performance and highlights trade-offs.
-

25+ Analytical MCQs (with Reasoning and Answers)

1. **Which metric is most appropriate for imbalanced binary classification?**
A. Accuracy
B. Precision-Recall AUC
C. Mean Squared Error
D. R2 Score
Answer: B
2. **If a model has high precision but low recall, this suggests:**
A. Most positive predictions are correct, but many positives are missed
B. Model rarely makes false positives
C. Model misses many actual positives
D. Both A and C
Answer: D
3. **The ROC curve plots:**
A. Precision vs. Recall
B. True Positive Rate vs. False Positive Rate
C. Accuracy vs. Error
D. Recall vs. Specificity
Answer: B
4. **Which element is NOT in a confusion matrix?**
A. True Positive
B. False Positive
C. Mean Squared Error
D. False Negative
Answer: C

5. **A classifier with AUC = 1.0 means:**

- A. Random performance
- B. Perfect discrimination between classes
- C. No true positives
- D. All negatives predicted as positives

Answer: B

6. **The F1-score is the harmonic mean of:**

- A. Precision and recall
- B. Recall and specificity
- C. Accuracy and recall
- D. ROC and AUC

Answer: A

7. **Which metric is least affected by class imbalance?**

- A. Accuracy
- B. MCC
- C. Precision
- D. Recall

Answer: B

8. **Mean Absolute Error (MAE) is used for:**

- A. Classification
- B. Regression
- C. Clustering
- D. Feature selection

Answer: B

9. **An R2 score below 0 indicates:**

- A. Model is perfect
- B. Model is worse than predicting the mean
- C. Model is better than mean prediction
- D. Not possible

Answer: B

10. **Cohen's kappa accounts for:**

- A. Chance agreement
- B. Only true positives
- C. Only false negatives
- D. Regression bias

Answer: A

11. Which metric is most useful for medical diagnosis where missing a disease is costly?

- A. Precision
- B. Recall
- C. F1-score
- D. MAE

Answer: B

12. If precision = 1, recall = 0.5, what is F1-score?

- A. 0.75
- B. 1.0
- C. 0.67
- D. 0.5

Answer: C ($F1 = 2 \times 0.5 / (1 + 0.5) = 0.67$)

13. What does a high false positive rate indicate?

- A. Many negatives are incorrectly classified as positives
- B. Model rarely predicts positive
- C. All positives are correct
- D. None of the above

Answer: A

14. Precision-Recall curve is preferred over ROC when:

- A. Classes are balanced
- B. Positive class is rare
- C. Negative class is rare
- D. Both classes are rare

Answer: B

15. Which metric combines all confusion matrix elements in a balanced way?

- A. Accuracy
- B. MCC
- C. Recall
- D. F1-score

Answer: B

16. Cross-validation helps to:

- A. Avoid overfitting in evaluation
- B. Increase training data
- C. Reduce model complexity

D. Improve data imbalance

Answer: A

17. **AUC is threshold-independent because:**

A. It considers all possible thresholds

B. It fixes a single threshold

C. It is only for regression

D. It ignores true negatives

Answer: A

18. **Which metric is not threshold dependent?**

A. Accuracy

B. Precision

C. AUC

D. Recall

Answer: C

19. **If a model's F1-score is low despite high accuracy, this likely means:**

A. Severe class imbalance

B. Model is overfitting

C. Model is underfitting

D. F1 and accuracy are always similar

Answer: A

20. **Which metric penalizes large errors more heavily?**

A. MAE

B. MSE

C. R2

D. F1

Answer: B

21. **Which metric is best when both false positives and false negatives are costly?**

A. Precision

B. Recall

C. F1-score

D. Specificity

Answer: C

22. **Data leakage can result in:**

A. Overoptimistic evaluation metrics

B. Underfitting

- C. Lower accuracy
- D. Higher bias

Answer: A

23. **A perfect F1-score is:**

- A. 0
- B. 1
- C. Undefined
- D. -1

Answer: B

24. **Which metric can be used for multi-class classification?**

- A. Accuracy
- B. Macro-averaged F1-score
- C. MCC
- D. All of the above

Answer: D

25. **Precision is calculated as:**

- A. $TP / (TP + FP)$
- B. $TP / (TP + FN)$
- C. $TN / (TN + FP)$
- D. $TP / (TP + TN)$

Answer: A

26. **A model with high recall but low precision is likely to:**

- A. Classify most positives correctly, but also many false positives
- B. Miss many true positives
- C. Have no false positives
- D. Have the same F1 and accuracy

Answer: A

Evaluation Measures 3 – Critical Analysis, Analytical Questions, and Analytical MCQs

Slidewise Critical Analysis and Analytical Short Questions (with Answers)

Slide 1: Introduction to Advanced Evaluation Measures

- **Critical Analysis:**
This slide introduces further evaluation metrics, possibly for multi-class, multi-label, or ranking problems, expanding on earlier parts.
 - **Analytical Question:**
Q: Why are advanced evaluation metrics required beyond accuracy, precision, and recall?
A: Complex tasks (multi-class, multi-label, ranking) require more nuanced metrics to capture model performance characteristics.
-

Slide 2: Multi-Class Classification Metrics

- **Critical Analysis:**
Discusses accuracy, macro/micro-averaged precision, recall, and F1-score for multi-class settings.
 - **Analytical Question:**
Q: What is the difference between macro and micro-averaging for F1-score in multi-class classification?
A: Macro averages metrics for each class equally; micro aggregates contributions before averaging, weighting by class frequency.
-

Slide 3: Multi-Label Evaluation Metrics

- **Critical Analysis:**
Explains hamming loss, subset accuracy, and label-based/macro-micro averaged metrics for multi-label classification.
 - **Analytical Question:**
Q: When is hamming loss preferred over subset accuracy in multi-label tasks?
A: When partial correctness is important—hamming loss penalizes each label independently, while subset accuracy is all-or-nothing.
-

Slide 4: Ranking Metrics (e.g., MAP, NDCG, MRR)

- **Critical Analysis:**
Introduces ranking evaluation: Mean Average Precision (MAP), Normalized Discounted Cumulative Gain (NDCG), Mean Reciprocal Rank (MRR), for evaluating search/recommendation.
 - **Analytical Question:**
Q: Why is NDCG preferred for graded relevance in ranking tasks?
A: NDCG accounts for both the relevance and the position of results, rewarding highly relevant items at top ranks.
-

Slide 5: Calibration and Probability Metrics

- **Critical Analysis:**
Covers Brier score, calibration curves, and log-loss as measures of predicted probability quality.
 - **Analytical Question:**
Q: What does a low Brier score indicate about a classifier's probability estimates?
A: Predictions are well-calibrated and close to true probabilities.
-

Slide 6: Cost-Sensitive and Custom Metrics

- **Critical Analysis:**
Discusses the need for custom or cost-sensitive metrics in domains where different errors have different impacts.
 - **Analytical Question:**
Q: Why might you use a cost-sensitive metric in fraud detection?
A: Because false negatives (missed fraud) may cost much more than false positives.
-

Slide 7: Cross-Validation and Model Selection Criteria

- **Critical Analysis:**
Reviews model selection using cross-validation and metrics like AIC, BIC for balancing fit and complexity.
 - **Analytical Question:**
Q: What is the main advantage of using cross-validation for metric estimation?
A: It provides a more robust estimate by averaging over multiple data splits.
-

Slide 8: Pitfalls and Best Practices for Advanced Metrics

- **Critical Analysis:**
Warns about overfitting to metrics, metric selection bias, and misinterpreting scores.
 - **Analytical Question:**
Q: How can metric selection bias affect model evaluation?
A: Over-optimizing for a single metric may ignore other important aspects of performance.
-

Slide 9: Putting It All Together—Case Study/Examples

- **Critical Analysis:**
Gives a real-world example of matching metrics to business needs.
- **Analytical Question:**
Q: How should metric choice align with application requirements in real-world projects?

A: The chosen metric should reflect the business goal and risk associated with different errors.

Slide 10: Summary and Recommendations

- **Critical Analysis:**
Summarizes the importance of metric diversity, alignment with objectives, and continuous evaluation.
 - **Analytical Question:**
Q: Why is it important to periodically revisit your evaluation metrics as a project evolves?
A: Changing requirements or data may mean different metrics become more relevant or necessary.
-

25+ Analytical MCQs (with Reasoning and Answers)

1. **Macro-averaged F1-score is best used when:**
A. Classes are balanced
B. All classes are equally important
C. Majority class dominates
D. There is only one class
Answer: B
2. **Micro-averaged F1-score gives more weight to:**
A. Minority classes
B. Majority classes
C. Each class equally
D. Random predictions
Answer: B
3. **Hamming loss in multi-label classification measures:**
A. Fraction of labels incorrectly predicted
B. Fraction of samples with any wrong label
C. Fraction of correct labels
D. Probability estimates
Answer: A

4. **Subset accuracy in multi-label classification requires:**

- A. All labels to be correct for a sample
- B. At least one label to be correct
- C. Labels to be independent
- D. Majority labels to be correct

Answer: A

5. **Which ranking metric considers both the order and relevance of results?**

- A. Precision@k
- B. MAP
- C. NDCG
- D. Hamming loss

Answer: C

6. **A low Brier score indicates:**

- A. Poor calibration
- B. Good probability estimates
- C. High error
- D. Overfitting

Answer: B

7. **In cost-sensitive evaluation, a higher cost is assigned to:**

- A. All errors equally
- B. Specific types of errors
- C. Correct predictions
- D. Random predictions

Answer: B

8. **Which metric is most appropriate for evaluating search engines?**

- A. NDCG
- B. Hamming loss
- C. Macro-F1
- D. AUC

Answer: A

9. **Mean Reciprocal Rank (MRR) is most informative when:**

- A. Only the first relevant result matters
- B. All relevant results are needed
- C. There are no relevant results
- D. Scores are continuous

Answer: A

10. **If a classifier's predicted probabilities are poorly calibrated, which will likely be high?**

- A. Accuracy
- B. Brier score
- C. Macro-F1
- D. Subset accuracy

Answer: B

11. **AIC and BIC are used for:**

- A. Model complexity penalization
- B. Regression accuracy
- C. Probability calibration
- D. Ranking results

Answer: A

12. **A model that performs well on micro-F1 but poorly on macro-F1 suggests:**

- A. It favors frequent classes
- B. It handles all classes equally
- C. It is overfitting
- D. None of the above

Answer: A

13. **Which metric is most punishing for partial errors in multi-label tasks?**

- A. Subset accuracy
- B. Hamming loss
- C. MAP
- D. ROC-AUC

Answer: A

14. **If two models have similar accuracy, but one has better NDCG, this means:**

- A. The better NDCG model ranks relevant items higher
- B. Both models are identical
- C. Accuracy is always superior
- D. The first model is better

Answer: A

15. **Label ranking average precision is used in:**

- A. Multi-label ranking
- B. Binary classification
- C. Regression

D. Clustering

Answer: A

16. **Why is cross-validation important for metric estimation?**

- A. Reduces overfitting risk
- B. Provides more reliable estimates
- C. Accounts for data variability
- D. All of the above

Answer: D

17. **Which metric is threshold-independent?**

- A. NDCG
- B. ROC-AUC
- C. Macro-F1
- D. Subset accuracy

Answer: B

18. **A high log-loss value indicates:**

- A. Model's probability predictions are far from true labels
- B. Model is well-calibrated
- C. Model has perfect accuracy
- D. Model is underfitting

Answer: A

19. **Which metric is suitable for evaluating recommendation systems with graded relevance?**

- A. Precision
- B. MAP
- C. NDCG
- D. Hamming loss

Answer: C

20. **Overfitting to a metric can result in:**

- A. Poor generalization
- B. Good generalization
- C. Lower performance on other important metrics
- D. Both A and C

Answer: D

21. **Metric selection bias may occur if:**

- A. Only one metric is optimized
- B. Multiple metrics are reported

- C. Cross-validation is used
- D. None of the above

Answer: A

22. **Why is it important to match metric choice to business needs?**

- A. Some errors are costlier than others
- B. All metrics are equivalent
- C. Metrics never impact business
- D. Only accuracy matters

Answer: A

23. **Label-wise macro averaging in multi-label learning:**

- A. Treats each label equally
- B. Weighs labels by frequency
- C. Ignores rare labels
- D. Is always higher than micro

Answer: A

24. **A well-calibrated model's predicted probabilities:**

- A. Match observed frequencies
- B. Always sum to one
- C. Are random
- D. Overestimate confidence

Answer: A

25. **The Brier score is minimized when:**

- A. Predicted probabilities are close to true outcomes
- B. Model is overfit
- C. Predictions are always 0.5
- D. All predictions are wrong

Answer: A

26. **Which metric is least appropriate for multi-label classification?**

- A. Hamming loss
- B. Subset accuracy
- C. Standard (binary) accuracy
- D. Macro-F1

Answer: C

Ranking – Critical Analysis, Analytical Questions, and Analytical MCQs

Slidewise Critical Analysis and Analytical Short Questions (with Answers)

Slide 1: Introduction to Ranking Problems

- **Critical Analysis:**
Introduces ranking as a machine learning and information retrieval task, where the goal is to order items (documents, products, etc.) by relevance or preference.
 - **Analytical Question:**
Q: Why are ranking problems different from standard classification or regression?
A: Because the output is an ordered list, not a single label or value, and evaluation focuses on the quality of the ordering.
-

Slide 2: Applications of Ranking

- **Critical Analysis:**
Highlights applications such as search engines, recommendation systems, and ad placement, where the ordering of results directly impacts user experience.
 - **Analytical Question:**
Q: Give an example of a real-world system that relies on effective ranking.
A: Web search engines, which must rank web pages by relevance to a user query.
-

Slide 3: Types of Ranking Problems

- **Critical Analysis:**
Describes pointwise, pairwise, and listwise approaches to ranking. Each approach differs in how it models the ranking problem and what loss functions it uses.
 - **Analytical Question:**
Q: What is the main difference between pairwise and listwise ranking approaches?
A: Pairwise compares pairs of items to learn relative orderings, while listwise optimizes the ranking of entire lists.
-

Slide 4: Evaluation Metrics for Ranking

- **Critical Analysis:**
Introduces key metrics: Precision@K, Recall@K, Mean Average Precision (MAP), NDCG (Normalized Discounted Cumulative Gain), and Mean Reciprocal Rank (MRR).
 - **Analytical Question:**
Q: Why is NDCG particularly well-suited for ranking tasks with graded relevance?
A: Because it weights items according to their relevance and their rank position, rewarding highly relevant items at the top.
-

Slide 5: Precision@K and Recall@K

- **Critical Analysis:**
Focuses on measuring the quality of the top-K results, which is important when only the top results matter to users.
 - **Analytical Question:**
Q: When is Precision@K preferred over overall precision?
A: When the user only interacts with the top-K ranked items, as in search or recommendation tasks.
-

Slide 6: NDCG and MAP in Depth

- **Critical Analysis:**
Explains how NDCG normalizes DCG to allow for comparison across queries, and how MAP summarizes precision across recall levels.
 - **Analytical Question:**
Q: How does NDCG differ from MAP in evaluating ranked lists?
A: NDCG accounts for the position and graded relevance, while MAP averages precision across all relevant items.
-

Slide 7: Learning to Rank Algorithms

- **Critical Analysis:**
Discusses popular algorithms: RankNet, LambdaRank, LambdaMART, and SVM-Rank, and their use of neural networks or SVMs for ranking.
 - **Analytical Question:**
Q: What is the advantage of listwise learning-to-rank algorithms?
A: They optimize the quality of the entire ranked list rather than local pairwise comparisons.
-

Slide 8: Challenges in Ranking

- **Critical Analysis:**
Discusses issues such as data sparsity, position bias, and scalability in large datasets.
 - **Analytical Question:**
Q: What is position bias and how does it affect ranking evaluation?
A: Users are more likely to click items higher in a list, regardless of true relevance, potentially biasing evaluation metrics.
-

Slide 9: Practical Tips and Best Practices

- **Critical Analysis:**
Offers advice on feature engineering, handling queries with few or no relevant items, and combining multiple metrics for robust evaluation.

- **Analytical Question:**

Q: Why is it important to use multiple metrics when evaluating ranking systems?

A: Different metrics capture different aspects of ranking quality, such as relevance, order, and grade of items.

Slide 10: Summary and Outlook

- **Critical Analysis:**

Summarizes the field and points to ongoing research in learning to rank, unbiased evaluation, and fairness in ranking.

- **Analytical Question:**

Q: What is a current research direction in ranking systems?

A: Developing unbiased evaluation methods and ensuring fairness in ranked outputs.

25+ Analytical MCQs (with Reasoning and Answers)

1. **Which metric is most appropriate for evaluating the top-ranked results in a search engine?**

- A. Accuracy
- B. Precision@K
- C. Mean Absolute Error
- D. Brier Score

Answer: B

2. **NDCG is preferred over MAP when:**

- A. All relevant items have equal importance
- B. Graded relevance levels exist
- C. Only binary relevance is considered
- D. There is no need to normalize scores

Answer: B

3. **A pointwise ranking approach treats ranking as:**

- A. Classification or regression for each item individually
- B. Comparing item pairs
- C. Optimizing the whole list

D. None of the above

Answer: A

4. **Pairwise ranking approaches learn by:**

- A. Assigning absolute scores to items
- B. Comparing pairs of items to determine which should rank higher
- C. Sorting the entire list
- D. Ignoring labels

Answer: B

5. **Which is NOT a ranking metric?**

- A. Precision@K
- B. Mean Reciprocal Rank
- C. ROC-AUC
- D. NDCG

Answer: C

6. **In NDCG, the discount factor is used to:**

- A. Penalize relevant items at lower ranks
- B. Reward irrelevant items
- C. Normalize relevance scores
- D. Compute average precision

Answer: A

7. **Mean Reciprocal Rank (MRR) measures:**

- A. The reciprocal of the rank of the first relevant item
- B. The average number of relevant items
- C. The maximum possible rank
- D. Overall accuracy

Answer: A

8. **Why is ranking more challenging than classification?**

- A. Requires ordering, not just labeling
- B. Needs to consider pairwise/listwise relationships
- C. Evaluation depends on the order of items
- D. All of the above

Answer: D

9. **Listwise learning-to-rank approaches optimize:**

- A. The correctness of the entire ranked list
- B. Individual items' scores
- C. The number of features

D. Only pairwise preferences

Answer: A

10. Position bias in ranking evaluation refers to:

- A. Irrelevant items appearing at the top
- B. Users clicking higher-ranked items more, regardless of relevance
- C. Overfitting the training data
- D. Ignoring user feedback

Answer: B

11. MAP is most appropriate when:

- A. Relevance is graded
- B. Only the first relevant item matters
- C. Every relevant item should be retrieved, regardless of position
- D. There are few queries

Answer: C

12. What is a limitation of Precision@K?

- A. It ignores items beyond the Kth position
- B. It considers all items
- C. It is threshold-independent
- D. It measures recall

Answer: A

13. Which algorithm is NOT used for learning to rank?

- A. RankNet
- B. LambdaMART
- C. SVM-Rank
- D. K-means

Answer: D

14. In ranking, data sparsity refers to:

- A. Few items per query
- B. Many items but few relevance labels
- C. Non-numeric features
- D. Balanced classes

Answer: B

15. If a user only wants the top result, which metric is most relevant?

- A. Recall@K
- B. MRR
- C. MAP

D. NDCG

Answer: B

16. **A high NDCG@10 score indicates:**

- A. Relevant items are ranked highly in the top 10
- B. Irrelevant items are at the top
- C. Many items are missing
- D. Only one item is relevant

Answer: A

17. **Learning to rank algorithms can be trained using:**

- A. Click data
- B. Human relevance judgments
- C. Both A and B
- D. Only labeled pairs

Answer: C

18. **Why use multiple ranking metrics in evaluation?**

- A. To capture different user needs and behaviors
- B. Because all metrics are equivalent
- C. To increase complexity
- D. To avoid reporting low scores

Answer: A

19. **Which metric is most sensitive to the presence of highly relevant items at lower ranks?**

- A. Precision@1
- B. NDCG
- C. MAP
- D. Recall@K

Answer: B

20. **A challenge in learning to rank is:**

- A. Handling missing or noisy relevance labels
- B. Overfitting to training queries
- C. Scalability for large datasets
- D. All of the above

Answer: D

21. **In pairwise ranking, the loss function is based on:**

- A. Absolute difference between scores
- B. Whether the pair is correctly ordered

- C. The sum of all item scores
- D. The number of features

Answer: B

22. **Click data for ranking evaluation is:**

- A. Always unbiased
- B. Prone to position bias
- C. Not useful for training
- D. Equivalent to random labels

Answer: B

23. **Which method addresses the issue of graded relevance in ranking evaluation?**

- A. MAP
- B. NDCG
- C. MRR
- D. Precision@K

Answer: B

24. **Why might listwise approaches outperform pairwise in some settings?**

- A. They directly optimize the full order of results
- B. They require less data
- C. They use only the top result
- D. They ignore relevance labels

Answer: A

25. **A high MAP score means:**

- A. Most relevant items are retrieved early in the ranking
- B. Only the first relevant item is considered
- C. All items are highly ranked
- D. All irrelevant items are at the top

Answer: A

26. **Fairness in ranking evaluation concerns:**

- A. Ensuring all groups are fairly represented in the top results
- B. Only maximizing precision
- C. Ignoring user feedback
- D. Overfitting to one group

Answer: A

How Not to Get Fooled by Machine Learning – Critical Analysis, Analytical Questions, and Analytical MCQs

Slidewise Critical Analysis and Analytical Short Questions (with Answers)

Slide 1: Introduction – The Risks of ML Misinterpretation

- **Critical Analysis:**
Opens with the prevalence of overtrusting or misinterpreting ML results and the need for critical thinking in data science.
 - **Analytical Question:**
Q: Why is it important to approach ML results with skepticism?
A: Because models can be misleading due to data issues, spurious correlations, or overfitting, leading to incorrect conclusions.
-

Slide 2: Common Pitfalls in ML

- **Critical Analysis:**
Lists pitfalls such as overfitting, data leakage, confirmation bias, and improper validation.
 - **Analytical Question:**
Q: What is data leakage and why is it dangerous?
A: Data leakage occurs when information from outside the training dataset is used to create the model, leading to overly optimistic performance estimates.
-

Slide 3: The Illusion of High Accuracy

- **Critical Analysis:**
Shows how high accuracy can be misleading, especially with imbalanced datasets or irrelevant baselines.
 - **Analytical Question:**
Q: Why might a model report high accuracy yet be useless in practice?
A: In imbalanced data, predicting the majority class achieves high accuracy without learning relevant patterns.
-

Slide 4: Spurious Correlations and Causality

- **Critical Analysis:**
Warns against interpreting correlation as causation and highlights risks of spurious relationships in data.
 - **Analytical Question:**
Q: What is a spurious correlation, and how can it fool an ML practitioner?
A: It's a statistical relationship that appears by chance or due to confounding variables, leading to false causal interpretations.
-

Slide 5: Data Snooping and P-Hacking

- **Critical Analysis:**
Explains dangers of excessive hypothesis testing and cherry-picking results to achieve statistical significance.
 - **Analytical Question:**
Q: How does p-hacking distort model evaluation?
A: By selectively reporting only significant results, it increases the risk of false discoveries.
-

Slide 6: Improper Validation Techniques

- **Critical Analysis:**
Discusses the necessity of proper train-test splits, cross-validation, and the dangers of testing on training data.

- **Analytical Question:**
Q: Why is cross-validation preferred over a single train-test split?
A: It provides a more robust estimate of model generalization by averaging performance over multiple splits.
-

Slide 7: The Problem of Overfitting

- **Critical Analysis:**
Highlights how complex models can fit noise instead of signal and underlines the need for regularization and validation.
 - **Analytical Question:**
Q: How can you detect overfitting in a model?
A: If training performance is high but test/validation performance is much lower.
-

Slide 8: Confirmation Bias and Interpretability

- **Critical Analysis:**
Addresses the risk of seeing what you want to see in results and the need for model interpretability.
 - **Analytical Question:**
Q: How does confirmation bias manifest in ML analysis?
A: By unconsciously favoring results that support one's hypothesis and ignoring contradictory evidence.
-

Slide 9: Ethical and Societal Risks

- **Critical Analysis:**
Examines risks like bias amplification, discrimination, and unintended consequences when ML is blindly trusted.
 - **Analytical Question:**
Q: Give an example of societal harm caused by uncritical application of ML.
A: Biased hiring algorithms that perpetuate discrimination.
-

Slide 10: Best Practices for Robust ML

- **Critical Analysis:**
Summarizes strategies: rigorous validation, transparency, domain expertise, and skepticism.
 - **Analytical Question:**
Q: Name two best practices to avoid being fooled by ML.
A: Use cross-validation and always inspect for data leakage.
-

25+ Analytical MCQs (with Reasoning and Answers)

1. **Which is a common pitfall that leads to overestimating ML model performance?**
A. Using validation data for both training and testing
B. Cross-validation
C. Regularization
D. Feature selection
Answer: A
2. **A model that achieves high accuracy on training data but low accuracy on test data is likely:**
A. Overfit
B. Underfit
C. Well regularized
D. Robust
Answer: A
3. **Data leakage occurs when:**
A. Test data is used during model training
B. Data is split randomly
C. Cross-validation is performed
D. Features are standardized
Answer: A
4. **Spurious correlations can be detected by:**
A. Domain knowledge and careful hypothesis testing
B. Ignoring all relationships
C. Only using accuracy

D. Training deeper models

Answer: A

5. **Why is accuracy a misleading metric for imbalanced datasets?**

A. It ignores class distribution

B. It is always high

C. It is the only metric

D. It accounts for precision

Answer: A

6. **Which practice helps prevent overfitting?**

A. Using more features

B. Cross-validation

C. Training on test data

D. P-hacking

Answer: B

7. **P-hacking refers to:**

A. Tuning model parameters for best validation score

B. Selecting only statistically significant results from multiple tests

C. Using deep learning

D. Data normalization

Answer: B

8. **Confirmation bias in ML can be mitigated by:**

A. Pre-registering hypotheses and using blinded analysis

B. Ignoring contradictory evidence

C. Only reporting positive results

D. Training larger models

Answer: A

9. **The danger of data snooping is:**

A. Overestimating model generalizability

B. Underestimating model complexity

C. Reducing dimensionality

D. Using regularization

Answer: A

10. **Which is a sign of model overfitting?**

A. High variance between train and test performance

B. Identical train and test error

C. Low training accuracy

D. High bias

Answer: A

11. Why should you avoid making decisions based solely on high cross-validated accuracy?

A. Other performance aspects may be ignored

B. It guarantees causality

C. It prevents bias

D. It avoids feature leakage

Answer: A

12. Which scenario is most likely to result in model bias?

A. Using an unrepresentative training dataset

B. Using cross-validation

C. Oversampling minority class

D. Data normalization

Answer: A

13. Causality in ML models can be reliably inferred by:

A. Randomized controlled experiments

B. Observing high correlation

C. Achieving high test accuracy

D. Using deep learning

Answer: A

14. Ethical harm from ML can arise when:

A. Models amplify societal biases

B. Models are validated

C. Regularization is used

D. Models are explainable

Answer: A

15. Which is a best practice for preventing data leakage?

A. Splitting data before feature engineering

B. Performing feature selection on whole dataset

C. Using all data for cross-validation

D. Tuning on test set

Answer: A

16. Improper validation techniques can cause:

A. Overoptimistic performance metrics

B. Model interpretability

- C. Domain knowledge
- D. Robust testing

Answer: A

17. The illusion of high accuracy can be most dangerous when:

- A. Dataset is imbalanced
- B. Model is underfit
- C. Data is normalized
- D. Regularization is used

Answer: A

18. Why is domain expertise critical in ML projects?

- A. To detect spurious or non-causal relationships
- B. To automate feature engineering
- C. To maximize model complexity
- D. To avoid regularization

Answer: A

19. Which step can help ensure reproducible ML results?

- A. Random seed fixing and documentation
- B. Selective reporting
- C. Ignoring failed experiments
- D. Data snooping

Answer: A

20. A model that picks up on confounding variables is likely:

- A. Spurious and unreliable
- B. Robust and generalizable
- C. Highly interpretable
- D. Regularized

Answer: A

21. Overfitting is less likely when:

- A. The model is simple and validated
- B. The model is highly complex
- C. The training set is small
- D. The test data is used for training

Answer: A

22. Which is a sign of confirmation bias in ML?

- A. Only reporting results that confirm expectations
- B. Rigorous cross-validation

- C. Using hold-out validation
- D. Feature scaling

Answer: A

23. **What is the risk of not considering ethical impacts of ML?**

- A. Unintended discrimination or harm
- B. More accurate models
- C. Improved robustness
- D. Faster computation

Answer: A

24. **Which validation approach is most robust against random data splits?**

- A. K-fold cross-validation
- B. Single hold-out split
- C. Training on all data
- D. No validation

Answer: A

25. **A well-documented ML workflow helps prevent:**

- A. Reproducibility errors and hidden pitfalls
- B. Overfitting
- C. Underfitting
- D. Data normalization

Answer: A

26. **A key reason to maintain skepticism about ML outcomes is:**

- A. Models can pick up patterns that are not meaningful or causal
- B. Models are always unbiased
- C. Validation is unnecessary
- D. Cross-validation is always perfect

Answer: A