

Machine Learning

Q1 Define Bias and Variance.

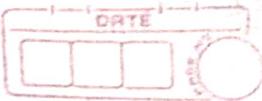
(Bias)

- * It is simply defined as inability of model because of that there is some difference ~~of error~~ occurring b/w the model's predicted value and actual value.
- * These difference b/w actual or expected values and the predicted values are known as error or bias error.
- * Bias (\hat{Y}) = $E(\hat{Y}) - Y$

Variance

- * Variance is the measure of spread in data from its mean position.
- * It is amount by which the performance of a predictive model changes when it is trained on different subsets of training data.
- * Variance = $E[(\hat{Y} - E(\hat{Y}))^2]$

Q2 What is structure of learning?



late

box

Q3 How to make predictions in linear regression

- * Making prediction is as simple as solving the equation for a specific input.
- * Let's make this concrete with an example.
Imagine we are predicting weight(Y) from height(X).
- * Our linear regression model representation for this problem would be
$$Y = B_0 + B_1(x)$$
$$\text{Weight} = B_0 + B_1(\text{Height})$$

Q4 What is No-Free lunch theorem?

→ This theorem states that there is no one model that works best for every problem. The assumption of a great model for one problem may not hold for another problem, so it is common in ML to try multiple models and find one that works best for a particular problem.

Q5 What do you mean by Stopping Criterion in CART?

→ The most frequent halting method is to utilize a min amount of training data allocated to every leaf node. If the count is smaller than the specified threshold, the split is rejected and also the node is considered the last leaf node.

Q6 What are the libraries used to do visualization.

* Matplotlib * plotly * Altair
* Seaborn * GGplot * Geoplotlib

Q7 Explain framework of inductive learning.

* Training example

→ A sample from X , including its o/p from target function.

* Target function

→ A mapping function f from X to $f(X)$

* Hypothesis

→ Approximation of f , candidate function.



- * Classifier
 - learning program builds a classifier that is used to classify
- * Learner
 - process that creates the classifier

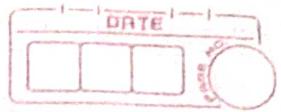
Q8 Define Training set & Test Set

Training Set

- * It is used to fit the model to learn the parameters of the model.
- * Larger in size as compared to test sets
- * Slower on larger datasets but job can be done in parallel using multiprocessing

Test Set

- * It is used to test whether model can generalize well on unseen data
- * Smaller in size as compared to training set
- * Faster than training sets



Q9 Define overfitting & underfitting.

Overfitting

- It occurs when our ML model tries to cover all the data points or more than the required data points present in given dataset. Because of this, the model starts caching noise and inaccurate values present in the dataset, and all these factors reduce the efficiency and accuracy of the model.
- The overfitted model has high variance & low bias.

Underfitting

- It occurs when our ML model is not able to capture the underlying trend of the data. To avoid the overfitting in the model, the fed of training data can be stopped at an early stage, due to which the model may not learn enough from the training data. As a result, it may fail to find the best fit of the dominant trend in the data.
- The underfitted model has low variance and high bias.
- It can neither model the training data nor generalize to new data.



Q10 What is aggregation?

- * Aggregation is the process to reduce the many measurements into a few values/statistics
- * You can do this aggregation in many different ways, the mean & 95% confidence interval is just one of many options to summarize the several measurements.

Q11 Give any 3 applications of ML?

- * Fraud detection
- * Speech recognition
- * Stock market analysis & forecasting

Q12 Define Gini ~~cost~~ cost

- * It is a proportion of impurity or inequality in statistical and monetary settings.
- * In ML, Gini cost is utilized as an impurity measure in decision tree algorithms for classification tasks.
- * It is calculated by subtracting the sum of the Squared probabilities of each class from one.

Q13 What do you mean by Gradient Descent?

- * When there are one or more inputs, you can use a process of optimizing the values of coefficients by iteratively minimizing the error of the model on your training data. This operation is called Gradient Descent and works by starting with zero values for each coefficient.
- * When using this method, you must select a learning rate (α) parameter that determines the size of the improvement step to take on each iteration of the procedure.
- * Gradient descent is often taught using a linear regression model because it is relatively straightforward to understand. In practise, it is useful when you have a very large dataset either in the no. of rows or the no. of columns that may not fit into memory.

Q14 Explain limit overfitting. How to handle the sequence?

- * Both overfitting and underfitting can lead to poor model performance. But by far the most common problem in applied machine learning is overfitting.
- * Overfitting is such a problem because the evaluation of ML on training data is different from evaluation we actually care the most about, namely how well the algorithm performs on unseen data.
- * There are 2 imp. techniques that you can use when evaluating ML algo to limit overfitting:
 - use a resampling technique to estimate model accuracy.
 - hold back a validation dataset

Q15 How Random forest using ML can influence the efficiency of algorithm?

→ High Accuracy

→ Random Forest models generally have high accuracy because they are an ensemble of multiple decision trees.

→ Reduced Overfitting

→ Random forests are less prone to overfitting compared to individual decision trees.

→ Handling Missing Values

→ Random Forests can handle missing values in the dataset w/o the need for imputation.

→ Robustness

→ Random Forest is robust to outliers & noise in data. It can handle noisy data and still provide accurate predictions, reducing the need for extensive data cleaning efforts and making the overall process more efficient.

Q16 How do you check the accuracy in ML model.

= We can calculate the accuracy of any model by dividing the correctly predicted problems by the total no. of predictions made.

$$\text{Accuracy} = \frac{\text{no. of correct predictions}}{\text{total no. of predictions}}$$

The above formula is very useful for calculating the accuracy of any model. It provides a simple understanding of a binary classification problem.

Q19 How Precision & Recall calculated in ML

Precision

- Precision measures the %age of predictions made by model that are correct.
- Precision = $\frac{TP}{TP + FP}$

→ TP = True Positive

→ FP = False Positive

TP

Recall

- Recall measures the %age of data points that were correctly identified by model.

$$\rightarrow \text{Recall} = \frac{TP}{TP + FN}$$

→ TP = True Positive

→ FN = False Negative

Q20 How CART model can be created from data?

- * Creating a CART model involves selecting i/p variable & split points on those variables until a suitable tree is constructed.
- * The selection of which i/p variable to use and the specific split or cut-point is chosen using a greedy algorithm to minimize a cost function.
- * Tree construction ends using a predefined stopping criterion, such as min. no. of training instances assigned to each node of the tree.
- * Tree can be constructed using following algorithm:-
 - Greedy algorithm
 - Stopping Criterion
 - Tree pruning

Q21 How errors are calculated in linear Regression

→ Linear regression most often uses mean square error (MSE) to calculate the error of the model. MSE is calculated by :-

1. measuring the distance of the observed y-values from the predicted y-values at each value of x
2. squaring each of these distances
3. calculating the mean of each of the squared distances.

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

n = no. of data points
 y_i = observed values

\hat{y}_i = predicted values

Q22 Explain Tree Pruning

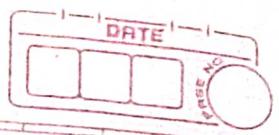
- * Tree pruning in ML refers to the process of reducing the size of a decision tree by removing certain branches of the tree that do not provide significant information or do not improve the model's predictive accuracy.
- * It is used to avoid overfitting, where the tree fits the training data too closely and performs poorly on unseen or new data.
- * 2 types of pruning techniques
 - pre-pruning
 - post-pruning
- * Pruning helps in creating simpler and more interpretable models that generalize well for unseen data, as it reduces the risk of overfitting and improves tree's efficiency.

Q23 Explain learning methods of ML

- * Batch learning
 - In many cases, we have end-to-end ML systems in which we need to train the model in one go by using whole available training data. Such kind of learning method is called batch learning.

* Online learning

- In this method, the training data is supplied in multiple incremental batches, called mini batches to algorithm.



* Instance based learning
→ this method builds the model by doing generalisation based on i/p data.

* Model based learning
→ In the method an iterative process takes place on ML models that are built based on various model parameters called hyperparameters and in which i/p data is used to extract features.

Q29 Linear Algorithm | Non-linear algorithm

Assume a linear relationship b/w input feature & O/p predictions.

Represent data using straight lines (in 2D) or hyperplanes (in multiple dimensions).

Computationally efficient & easy to interpret, suitable for simpler relationships in data.

Example] linear regression
linear SVM

Do not assume a linear relationship b/w i/p features & o/p predictions.

Can capture complex non-linear patterns & relationships in the data.

Computationally intensive and may require larger data.

Example] decision trees
neural networks



Q5 Explain Generalization in ML.

- * Generalization refers to how well the concepts learned by a ML model apply to specific examples not seen by the model when it was learning.
- * The goal of good ML model is to generalize well from the training data to any data from the problem domain. This allows us to make predictions in the future on data the model has never seen.
- * There is a terminology used in ML when we talk about ^{how} well a ML model learns & generalizes to new data, namely overfitting & underfitting.
- * Overfitting & underfitting are the 2 biggest causes for poor performance of ML algo.

Q6 Explain Linear Regression Model

- linear regression is a linear model e.g. a model that assumes the linear relationship b/w the input variables (x) and the single o/p variable (y).
- More specifically, that y can be calculated from a linear combination of i/p variables (x). When there is a Single i/p variable (x), the method is referred to as simple linear regression.
$$y = \beta_0 + \beta_1(x)$$
- When there are multiple i/p variables the method is referred as multiple linear regression.

* Different techniques can be used to prepare or train the linear regression equation for data, the most common of which is ordinary least squares.

Q38 How ML & AI are coordinated?

* AI's big picture

→ AI is the overarching concept where machine simulate human intelligence, including problem solving & decision-making

* ML's specialized role

→ ML is a subset of AI. It involves algorithms learning from data, improving tasks w/o being explicitly programmed.

* Data-powered learning

→ ML algo analyze data, identifying patterns and making predictions. This data-driven learning enhances AI's abilities.

* Continuous Improvement

→ ML enables AI Systems to learn from new data, ensuring they adapt and improve, making them more intelligent & effective over time.