

## **MACHINE LEARNING**

**In Q1 to Q11, only one option is correct, choose the correct option:**

1. Which of the following methods do we use to find the best fit line for data in Linear Regression?

- A) Least Square Error   B) Maximum Likelihood   C) Logarithmic Loss   D) Both A and B

**ANS: A) Least Square Error**

2. Which of the following statement is true about outliers in linear regression?

- A) Linear regression is sensitive to outliers      B) linear regression is not sensitive to outliers  
C) Can't say      D) none of these

**ANS: A) Linear regression is sensitive to outliers**

3. A line falls from left to right if a slope is \_\_\_\_\_?

- A) Positive      B) Negative      C) Zero      D) Undefined

**ANS: B) Negative**

4. Which of the following will have symmetric relation between dependent variable and independent variable?

- A) Regression      B) Correlation      C) Both of them      D) None of these

**ANS: C) Both of them**

5. Which of the following is the reason for over fitting condition?

- A) High bias and high variance      B) Low bias and low variance  
C) Low bias and high variance      D) none of these

**ANS: C) Low bias and high variance**

6. If output involves label then that model is called as:

- A) Descriptive model      B) Predictive modal  
C) Reinforcement learning      D) All of the above

**ANS: B) Predictive modal**

7. Lasso and Ridge regression techniques belong to \_\_\_\_\_?

- A) Cross validation      B) Removing outliers      C) SMOTE      D) Regularization

**ANS: D) Regularization**

8. To overcome with imbalance dataset which technique can be used?

- A) Cross validation      B) Regularization      C) Kernel      D) SMOTE

**ANS: A) Cross validation**

9. The AUC Receiver Operator Characteristic (AUCROC) curve is an evaluation metric for binary classification problems. It uses \_\_\_\_\_ to make graph?

- A) TPR and FPR
- B) Sensitivity and precision
- C) Sensitivity and Specificity
- D) Recall and precision

**ANS: C) Sensitivity and Specificity**

10. In AUC Receiver Operator Characteristic (AUCROC) curve for the better model area under the curve should be less.

- A) True
- B) False

**ANS: B) False**

11. Pick the feature extraction from below:

- A) Construction bag of words from an email
- B) Apply PCA to project high dimensional data
- C) Removing stop words
- D) All of the above

**ANS: D) All of the above**

**In Q12, more than one options are correct, choose all the correct options:**

12. Which of the following is true about Normal Equation used to compute the coefficient of the Linear Regression?

- A) We don't have to choose the learning rate.
- B) It becomes slow when number of features is very large.
- C) We need to iterate.
- D) It does not make use of dependent variable.

**ANS: A) We don't have to choose the learning rate.**

**B) It becomes slow when number of features is very large.**

**Q13 and Q15 are subjective answer type questions, Answer them briefly.**

**13. Explain the term regularization?**

**ANS:** Regularization is the process of making something regular or acceptable. This is precisely why we employ it in applied machine learning. Regularization is the process of regularising or shrinking the coefficients towards zero in the context of machine learning. To prevent overfitting, regularisation discourages learning a more complex or flexible model.

The basic idea is to penalize the complex models i.e. adding a complexity term that would give a bigger loss for complex models. To understand it, let's consider a simple relation for linear regression. Mathematically, it is stated as below:

$$Y \approx W_0 + W_1 X_1 + W_2 X_2 + \dots + W_P X_P$$

Where  $Y$  is the learned relation i.e. the value to be predicted.

$X_1, X_2, \dots, X_P$ , are the features deciding the value of  $Y$ .

$W_1, W_2, \dots, W_P$ , are the weights attached to the features  $X_1, X_2, \dots, X_P$  respectively.

$W_0$  represents the bias.

To fit a model that accurately predicts the value of  $Y$ , we need a loss function as well as optimised parameters, such as bias and weights.

The residual sum of squares loss function is commonly used in linear regression (RSS). It can be written as follows using the above-mentioned linear regression relationship:

$$RSS = \sum_{i=1}^m (Y_i - W_0 - \sum_{j=1}^n W_j X_{ji})^2$$

RSS is also known as the linear regression objective without regularisation.

The model will now learn using this loss function. It will adjust the weights based on our training data (coefficients). If our dataset is noisy, it will suffer from overfitting, and the estimated coefficients will not generalise to new data. This is where regularization comes into action. It regularizes these learned estimates towards zero by penalizing the magnitude of coefficients.

#### 14. Which particular algorithms are used for regularization?

**ANS:** The following are the algorithms are used for Regularization

- Ridge Regression
- LASSO (Least Absolute Shrinkage and Selection Operator) Regression
- Dropout

##### **Ridge Regression (L2 Regularization)**

Ridge regression is also called L2 norm or regularization.

When using this technique, we add the sum of weight's square to a loss function and thus create a new loss function which is denoted thus:

$$\text{Loss} = \sum_{j=1}^m \left( Y_i - W_0 - \sum_{i=1}^n W_i X_{ji} \right)^2 + \lambda \sum_{i=1}^n W_i^2$$

As seen above, the original loss function is modified by adding normalized weights. Here normalized weights are in the form of squares.

You may have noticed parameters  $\lambda$  along with normalized weights.  $\lambda$  is the parameter that needs to be tuned using a cross-validation dataset. When you use  $\lambda=0$ , it returns the residual sum of square as loss function which you chose initially. For a very high value of  $\lambda$ , loss will ignore core loss function and minimize weight's square and will end up taking the parameters' value as zero.

Now the parameters are learned using a modified loss function. To minimize the above function, parameters need to be as small as possible. Thus, L2 norm prevents weights from rising too high.

### Lasso Regression (L1 Regularization)

Also called lasso regression and denoted as below:

$$\text{Loss} = \sum_{j=1}^m \left( Y_i - W_o - \sum_{i=1}^n W_i X_{ji} \right)^2 + \lambda \sum_{i=1}^n |W_i|$$

This technique is different from ridge regression as it uses absolute weight values for normalization.  $\lambda$  is again a tuning parameter and behaves in the same as it does when using ridge regression.

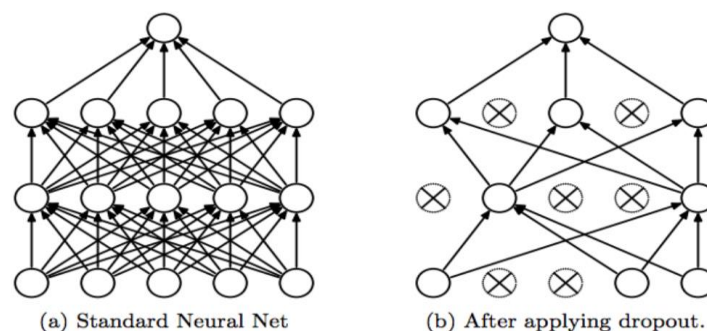
As loss function only considers absolute weights, optimization algorithms penalize higher weight values.

In ridge regression, loss function along with the optimization algorithm brings parameters near to zero but not actually zero, while lasso eliminates less important features and sets respective weight values to zero. Thus, lasso also performs feature selection along with regularization.

### Dropout

Dropout is a regularization technique used in neural networks. It prevents complex co-adaptations from other neurons.

In neural nets, fully connected layers are more prone to overfit on training data. Using dropout, you can drop connections with 1-p probability for each of the specified layers. Where p is called keep probability parameter and which needs to be tuned.



With dropout, you are left with a reduced network as dropped out neurons are left out during that training iteration.

Dropout decreases overfitting by avoiding training all the neurons on the complete training data in one go. It also improves training speed and learns more robust internal functions that generalize better on unseen data. However, it is important to note that Dropout takes more

epochs to train compared to training without Dropout (If you have 10000 observations in your training data, then using 10000 examples for training is considered as 1 epoch).

Along with Dropout, neural networks can be regularized also using L1 and L2 norms. Apart from that, if you are working on an image dataset, image augmentation can also be used as a regularization method.

For real-world applications, it is a must that a model performs well on unseen data. The techniques we discussed can help you make your model learn rather than just memorize.

### **15. Explain the term error present in linear regression equation?**

**ANS:** An error term is a residual variable produced by a statistical or mathematical model, which is created when the model does not fully represent the actual relationship between the independent variables and the dependent variables. As a result of this incomplete relationship, the error term is the amount at which the equation may differ during empirical analysis.

The error term is also known as the residual, disturbance, or remainder term, and is variously represented in models by the letters  $e$ ,  $\epsilon$ , or  $u$ .

- An error term appears in a statistical model, like a regression model, to indicate the uncertainty in the model.
- The error term is a residual variable that accounts for a lack of perfect goodness of fit.
- Heteroskedastic refers to a condition in which the variance of the residual term, or error term, in a regression model varies widely.

### **Error Term Use in a Formula**

An error term essentially means that the model is not completely accurate and results in differing results during real-world applications. For example, assume there is a multiple linear regression function that takes the following form:

$$Y = \alpha X + \beta \rho + \epsilon$$

Where:

$\alpha, \beta$  = Constant parameters

$X, \rho$  = Independent variables

$\epsilon$  = Error term

Linear regression is a form of analysis that relates to current trends experienced by a particular security or index by providing a relationship between a dependent and independent variables, such as the price of a security and the passage of time, resulting in a trend line that can be used as a predictive model.

A linear regression exhibits less delay than that experienced with a moving average, as the line is fit to the data points instead of based on the averages within the data. This allows the line to change more quickly and dramatically than a line based on numerical averaging of the available data points.