# MACHINE LEARNING

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- Option(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
3. A line falls from left to right if a slope is _____?
A) Positive B) Negative
C) Zero D) Undefined

**Ans- Option(A) Linear regression is sensitive to outliers**

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- Option(B) Correlation**

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- Option(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- Option(A) Descriptive model**

7. Lasso and Ridge regression techniques belong to _____?
A) Cross validation B) Removing outliers
C) SMOTE D) Regularization

**Ans- Option(D) Regularization**

8. To overcome with imbalance dataset which technique can be used?
A) Cross validation B) Regularization
C) Kernel D) SMOTE

**Ans- Option(D) SMOTE**

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- Option(A) TPR and FPR**

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

**Ans- Option(B) False**

11. Pick the feature extraction from below:
A) Construction bag of words from a email
B) Apply PCA to project high dimensional data
C) Removing stop words
D) Forward selection

**Ans- Option(B) Apply PCA to project high dimensional data**

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- Option(A,B,C,D)**

MACHINE LEARNING  WORKSHEET ASSIGNMENT

**13. Explain the term regularization?**

Regularization refers to techniques that are used to calibrate machine learning models in order to minimize the adjusted loss function and prevent overfitting or underfitting.

Using Regularization, we can fit our machine learning model appropriately on a given test set and hence reduce the errors in it.
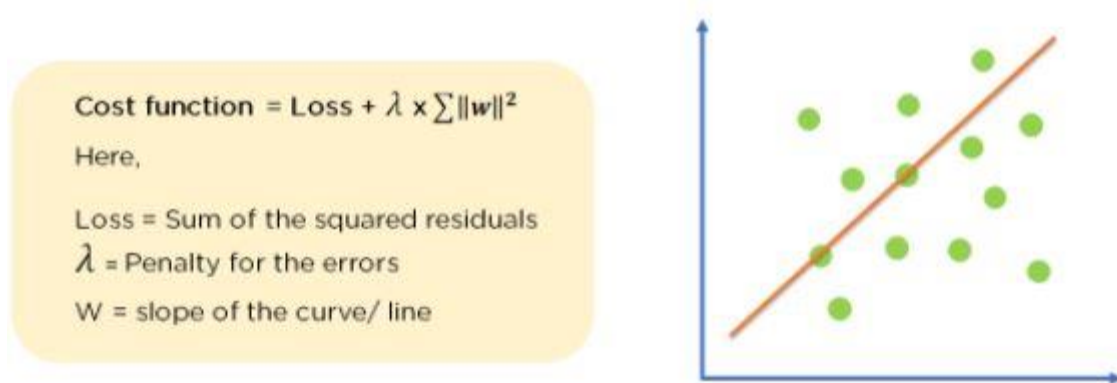
# Regularization Techniques

There are two main types of regularization techniques: Ridge Regularization and Lasso Regularization.

# Ridge Regularization :

Also known as Ridge Regression, it modifies the over-fitted or under fitted models by adding the penalty equivalent to the sum of the squares of the magnitude of coefficients.

This means that the mathematical function representing our machine learning model is minimized and coefficients are calculated. The magnitude of coefficients is squared and added. Ridge Regression performs regularization by shrinking the coefficients present. The function depicted below shows the cost function of ridge regression :



Cost function = Loss + $\lambda \times \sum \|w\|^2$

Here,

Loss = Sum of the squared residuals

$\lambda$ = Penalty for the errors

W = slope of the curve/ line

Cost Function of Ridge Regression

In the cost function, the penalty term is represented by Lambda $\lambda$. By changing the values of the penalty function, we are controlling the penalty term. The higher the penalty, it reduces

the magnitude of coefficients. It shrinks the parameters. Therefore, it is used to prevent multicollinearity, and it reduces the model complexity by coefficient shrinkage.

Comparing the two models, with all data points, we can see that the Ridge regression line fits the model more accurately than the linear regression line.

# Lasso Regression

It modifies the over-fitted or under-fitted models by adding the penalty equivalent to the sum of the absolute values of coefficients.

Lasso regression also performs coefficient minimization, but instead of squaring the magnitudes of the coefficients, it takes the true values of coefficients. This means that the coefficient sum can also be 0, because of the presence of negative coefficients.

We can control the coefficient values by controlling the penalty terms, just like we did in Ridge Regression.

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

Ans-

**Regularization**

The term 'regularization' refers to a set of techniques that regularizes learning from particular features for traditional algorithms or neurons in the case of neural network algorithms.It normalizes and moderates weights attached to a feature or a neuron so that algorithms do not rely on just a few features or neurons to predict the result. This technique helps to avoid the problem of overfitting.To understand regularization, let's consider a simple case of linear regression. Mathematically, linear regression is stated as below

$$= w_0 + w_1x_1 + w_2x_2 + \ldots + w_nx_n$$

where y is the value to be predicted;

$x_1, x_2, \ldots, x_n$ are features that decides the value of y;

$w_0$ is the bias;

$w_1$, $w_2$, ….., $w_n$ are the weights attached to $x_1$, $x_2$, …., $x_n$ relatively.

Now to build a model that accurately predicts the *y* value, we need to optimize above mentioned bias and weights.

To do so, we need to use a loss function and find optimized parameters using gradient descent algorithms and its variants.

To know more about building a machine learning application and the process, check out below blog:

How to Develop Machine Learning Applications for Business

The loss function called 'the residual sum of square' is mostly used for linear regression. Here's what it looks like :

$$\text{RSS} = \sum_{j=1}^{m}\left(Yi - Wo - \sum_{i=1}^{n} Wi\, Xji\right)^2$$

Next, we will learn bias (or intercept) and weights (also identified as parameters and coefficients) using the optimization algorithm (gradient descent) and data. If your dataset does have noise in it, it will face overfitting problem and learned parameters will not generalize well on unseen data.

To avoid this, you will need to regularize or normalize your weights for better learning.

There are three main regularization techniques, namely:

1. Ridge Regression (L2 Norm)

2. Lasso (L1 Norm)

3. Dropout

Ridge and Lasso can be used for any algorithms involving weight parameters, including neural nets. Dropout is primarily used in any kind of neural networks e.g. ANN, DNN, CNN or RNN to moderate the learning. Let's take a closer look at each of the techniques.

**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:

$$\text{Loss} = \sum_{j=1}^{m}\left(Yi - Wo - \sum_{i=1}^{n} Wi\,Xji\right)^2 + \lambda \sum_{i=1}^{n} Wi^2$$

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:

$$\text{Loss} = \sum_{j=1}^{m}\left(Yi - Wo - \sum_{i=1}^{n} Wi\,Xji\right)^2 + \lambda \sum_{i=1}^{n} |Wi|$$

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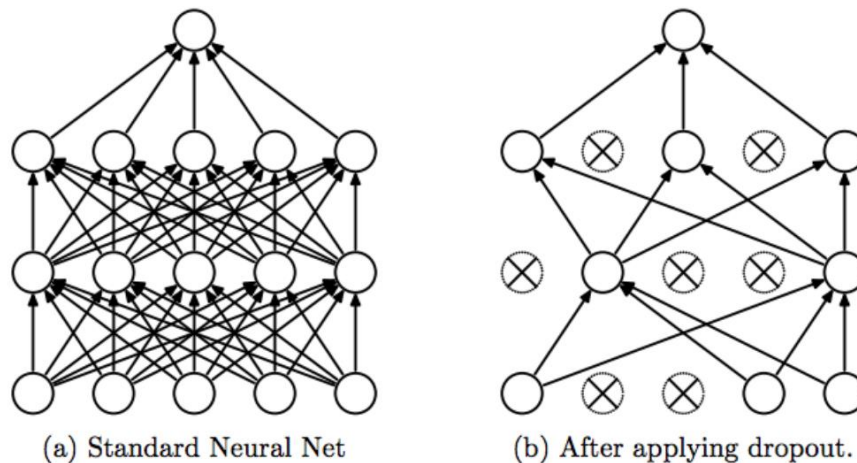
As loss function only considers absolute weights, optimization algorithms penalize higher weight values.

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.



(a) Standard Neural Net          (b) After applying dropout.

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..

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

**Ans-** The standard error of the regression (S), also known as the standard error of the estimate, represents the average distance that the observed values fall from the regression line. Conveniently, it tells you how wrong the regression model is on average using the units of the response variable. Smaller values are better because it indicates that the observations are closer to the fitted line.