

ANSWERS to Worksheet Set1

PYTHON ASSIGNMENT

Question Number	Answer		
1	C		
2	B		
3	C		
4	A		
5	D		
6	C		
7	A		
8	C		
9	A	C	
10	A	B	

STATISTICS ASSIGNMENT

QNo.	ANSWER
1	a
2	a
3	b
4	d
5	c
6	b
7	b
8	a
9	c
10	<p>Normal Distribution is also called as Gaussian Distribution. It is a probability distribution that is symmetric about the mean. It shows that data near the mean are more frequent in occurrence than the data far from the mean.</p> <p>In graphical form the normal distribution will appear as a bell curve.</p> <p>In a normal distribution mean is 0 and standard deviation is 1.</p>
11	<p>The real-world data often has a lot of missing values. The cause of missing values can be data corruption or failure to record data.</p> <p>There are some imputation techniques to remove missing data in a dataset-</p> <ol style="list-style-type: none">1. Deleting rows with missing values2. Imputation using (mean/median) values3. Imputation using (most frequent) or (constant/zero) values.4. Imputation using K-NN method5. Imputation using Multivariate Imputation by Chained Equation (MICE)6. Imputation using Deep learning(DataWig)

12	<p>A/B testing also called split testing, originated from the randomized control trials in Statistics, is one of the most popular ways for Businesses to test new UX features, new versions of a product, or an algorithm to decide whether your business should launch that new product/feature or not.</p> <p>The idea behind A/B testing is that you show the varied version of the product to a sample of customers (the experimental group) and the existing version of the product to another sample of customers (the control group). Then, the difference in product performance in experimental/treatment versus the control group is tracked, to identify the effect of this new version(s) of the product on the performance of the product.</p> <p>So, the goal is then to track the metric during the test period and find out whether there is a <i>difference</i> in the performance of the product and what type of difference is it.</p>
13	<p>No, mean imputation of missing data is not an acceptable practice. It is a simple solution to remove missing data but there are some disadvantages of using it –</p> <ol style="list-style-type: none"> 1. Mean imputation does not preserve the relationship among variables. 2. Mean imputation leads to an underestimate of standard errors.
14	<p>Linear Regression is a basic and commonly used type of predictive analysis. The main idea of regression is to examine two things –</p> <ol style="list-style-type: none"> 1. does a set of predictor variables do a good job In predicting an outcome(dependent) variable. 2. which variables in particular are significant predictors of the outcome variable and in what way do they impact the outcome variable. <p>Simple formula for regression is $y = c + bx$,</p>

	<p>Where:</p> <p>y = estimated dependent variable score</p> <p>c = constant</p> <p>b = regression coefficient and</p> <p>x = score on the independent variable</p>
15	<p>There are main two types of Statistics –</p> <ul style="list-style-type: none"> A. Descriptive Statistics B. Inferential Statistics <p><u>Descriptive Statistics</u></p> <p>The branch of statistics that focuses on collecting, summarizing and presenting a set of data.</p> <p>Ex. Average age of student in a class, height, weight of students in a particular section etc..</p> <p><u>Inferential Statistics</u></p> <p>The branch of statistics that analyzes sample data to draw conclusions about a population.</p> <p>Ex. A sample survey of people in a city who voted in a particular election, and it can be considered as a population.</p>

MACHINE LEARNING ASSIGNMENT

QNo.	ANSWER		
1	A		
2	A		
3	C		
4	B		
5	C		
6	B		
7	D		
8	D		
9	A		
10	B		
11	B		
12	A	B	
13	<p><u>REGULARIZATION-</u></p> <p>Regularization is one of the most important concepts in Machine Learning. It is a technique to prevent the model from overfitting by adding an extra information to it.</p> <p>Sometimes the machine learning model performs well with the training data but does not perform well with the test data. It means the model is not able to predict the output when deals with unseen data by introducing noise in the output, and hence the model is called overfitted. This problem can be deal with the help of a regularization technique.</p> <p>This technique can be used in such a way that it will allow to maintain all variables or features in the model by reducing the magnitude of the variables. Hence, it maintains accuracy as well as a generalization of the model.</p> <p>The linear regression equation is –</p> $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n + b$		

	<p>In the above equation, Y represents the value to be predicted</p> <p>X_1, X_2, \dots, X_n are the features for Y.</p> <p>$\beta_0, \beta_1, \dots, \beta_n$ are the weights or magnitude attached to the features, respectively. Here represents the bias of the model, and b represents the intercept.</p>
14	<p>There are two main techniques or algorithms used for regularization-</p> <p>1. <u>Ridge Regression</u></p> <p>It is also called as L2 regularization. It is used to reduce the complexity of the model.</p> <p>Ridge regression is one of the types of linear regression in which a small amount of bias is introduced so that we can get better long-term predictions.</p> <p>A general linear or polynomial regression will fail if there is high collinearity between the independent variables, so to solve such problems, Ridge regression can be used.</p> <p>2. <u>Lasso Regression</u></p> <p>This type of regression is also called L1 regularization.</p> <p>Lasso regression is another regularization technique to reduce the complexity of the model. It stands for <i>Least Absolute and Selection Operator</i>.</p>

	<p>It is similar to the Ridge Regression except that the penalty term contains only the absolute weights instead of a square of weights.</p> <p>Lasso regression can help us to reduce the overfitting in the model as well as the feature selection.</p> <p><u>Key Difference between Ridge Regression and Lasso Regression</u></p> <ul style="list-style-type: none"> ○ Ridge regression is mostly used to reduce the overfitting in the model, and it includes all the features present in the model. It reduces the complexity of the model by shrinking the coefficients. ○ Lasso regression helps to reduce the overfitting in the model as well as feature selection.
15	<p>An error term represents the margin of error within a statistical model; it refers to the sum of the deviations within the regression line, which provides an explanation for the difference between the theoretical value of the model and the actual observed results.</p> <p>The regression line is used as a point of analysis when attempting to determine the correlation between one independent variable and one dependent variable.</p> <p>Error Term Use in a Formula</p> <p>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:</p> $Y = \alpha X + \beta p + \epsilon$

	<p>where:</p> <p>α, β=Constant parameters</p> <p>X, ρ=Independent variables</p> <p>ϵ=Error term</p>
--	--