

m c @

Q1 a

Q16 = c

Q2 b -

Q17 = a

Q3 -

Q18 = a

Q4 A

Q19 = a

Q5 c, ~~a~~

Q20 = b

Q6 d

Q21 = b

Q7 b

Q22 =

Q8 a

Q23 = a

Q9 b

Q24 = d

Q10 b

Q25 = d

Q11 a

Q26 = b

Q12 c

Q27 = a

Q13 a, b, c

Q28

Q14 a, b

Q15 a

S.B @

Q1 For missing values we can use multiple steps like we can fill the missing values by using mean, median, mode (depends on the context), and also we can also use the random method. i.e. for ~~avoiding~~ maintaining the randomness of the data we can use the ~~or~~ fill the missing value by using random filling method i.e.

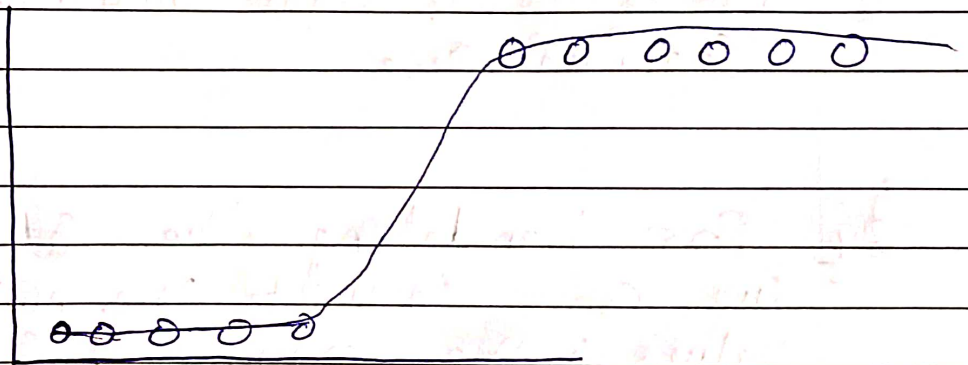
$np.random.randint(\text{mean} - \text{stdv}, \text{mean} + \text{stdv})$

where stdv is the standard deviation

and if the numbers of the missing value is very low as compared to the shape of the dataset then we can simply drop that dataset.

And if the number of the missing value is very high then we can simply drop that column.

Q1/ for making churn prediction model in telecommunications company we will use logistic regression model because our target ~~cat~~ column will be in classification i.e. in "Yes" or "No" because we will have only 2 values in our target column. Whether the client will churn i.e. "Yes" or it will remain on our platform i.e. "No".



logistic regression is a classification algorithm which helps to predict a binary outcome.

Q3 For segmenting the Customers into different groups based on the purchasing behaviour we can use Group by method so that we can filter different customers based on the purchasing behaviour like we can group by our customers on the basis of salary, age, need etc. so, the benefit we will get that we can target the right people with so that our company or organization can grow, it also help us for getting the right audience from the crowded one.

Q4 For handling the Categorical Variables we can convert the ~~new~~ Categorical values into the numerical one like if the dataset contains a column named Sex i.e male or female then we can convert this column into numerical by using 1 for male and 0 for female. ~~we~~ we can also do manually ~~by~~ with the help of rename function or we can use the LABEL ENCODER from the SKlearn library. the syntax will be

```
from sklearn.preprocessing import Label Encoder
```

LE = LabelEncoder()

$$df["Colname"] = LE.fit_transform(df["Colname"])$$

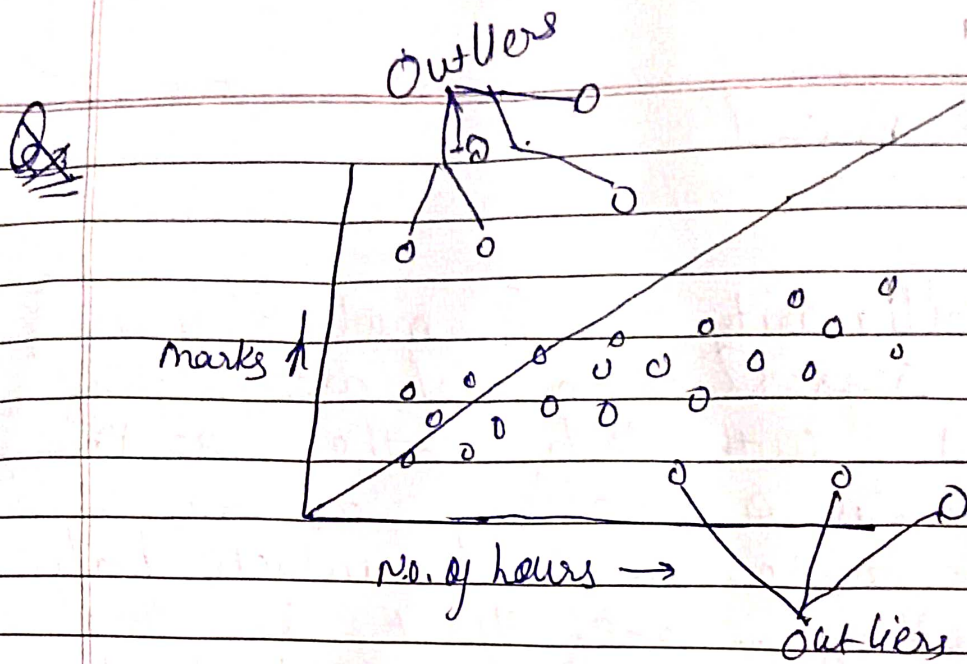
Concept Based

Q1 multicollinearity = In multicollinearity two regression model are highly correlated with each other so this will create a problem as we will not be able to find which feature is making more effect on the target. So it will lead to reduce the accuracy of our model so, we can use correlation matrix for handling multicollinearity.

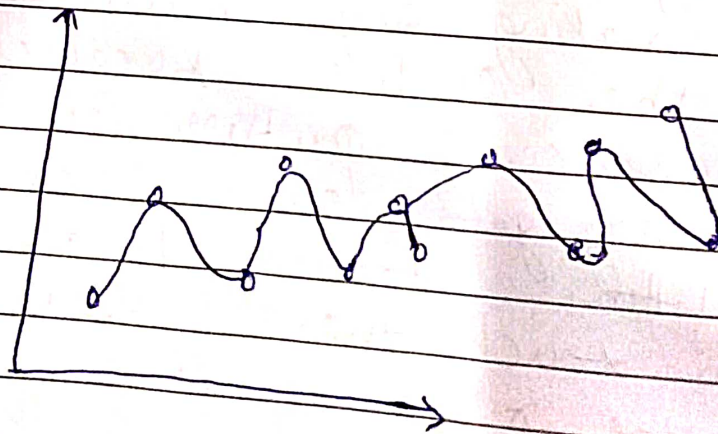
Q2 Outlier play a major role in making a ~~no~~ effective prediction model because before making any model we should handle the outlier. because due to the presence of outlier our model's accuracy gets down because in ~~first~~ fitting training phase our model tries to touch these outliers so our ~~model~~ the regression line gets deviated from its original path that leads to a bad model.

We can detect outliers by using the following methods:

- ① Scatterplot
- ② box plot
- ③ Z-score method.



Q3 Overfitting a phenomenon in which our regression line tries to touch each and every datapoint present in the dataset, it occurs by increasing the degree of the model.



Q4 For dealing with the non-linear relationship between predictors and target variable we can use polynomial regression as it gives when we have (n) graph non-linear relationship in our dataset then polynomial regression can be fitted



~~syntax~~

~~from sklearn.linear_model import PolynomialRegression~~

~~poly = PolynomialFeatures(degree=n)~~

~~poly.fit(X_train, y_train)~~

if captures the non-linear pattern with the help of degree we provide in the polynomial feature.