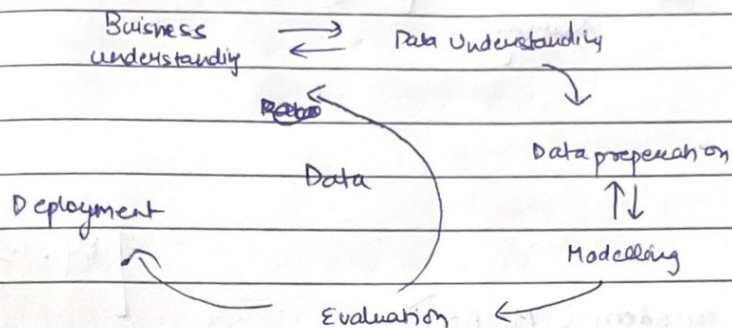


CRISP DM Framework

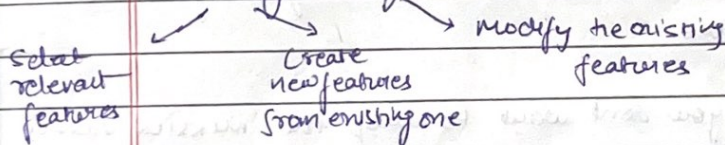
↳ cross industry standard for data mining.

DATE: / /
PAGE: /



Data preparation

- ↳ Missing value treatment
- ↳ Duplicates values (Remove)
- ↳ Outliers (Boxplot → Used to detect outlier, Detection & treatment)
- ↳ Feature engineering



① f_1, f_2, y (we'll check the correlation, the most correlated data will be used as relevant feature)

② Distance | time | speed

Data Encoding

* Null / Missing value treatment

① Missing completely at random (MCAR)

Independent of the data that we have or the data that is missing

Students Marks

A	39
B	39
C	-
D	30
E	31
F	-

→ No dependence of C on the people who appeared for exams & who didn't appeared.

② Missing at random (MAR)

↳ Some dependency on observed data.

DATE: / /
PAGE:

Name	Age	En: Girls not giving contact no. in the open survey.
V	21	
A	25	
S		

[`df.isnull().sum()`]

③ Missing not at random (MNAR)

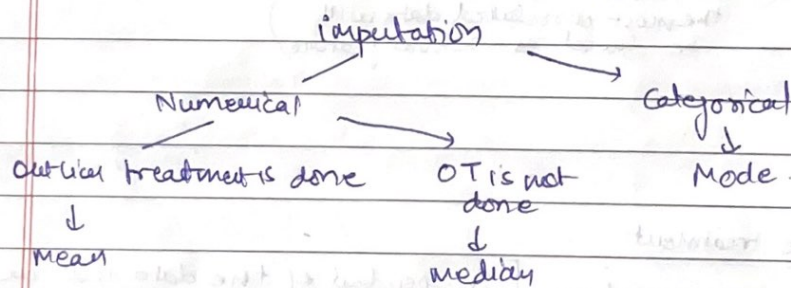
↳ Some dependency on observed data

* How to treat null values

↳ Never do missing value treatment without connecting to business team.

* if missing value $< 1\%$ and the data is large, you can drop the missing value.

* if data is small and you don't want to drop the missing values (bet small or large) you will do imputation.



Student wants

50 Here you cannot use mean & median
60 because if student appears for the exam
70
80
90

* impute missing value with a constant (Replace with zero 0) or the number that is not possible for that particular column)

* Create a new column with a flag if the column is missing and input some const. value in the missing column.

Name	Marks	Name Absent
A	28	0
B	-	1
C	30	0
D	-	0

Good Write