

Dealing with categorical data

We need a technique that can represent the categorical data in numerical form considering,

- ① Mathematical nature of the data is not disturbed.
- ② The domain nature and pattern of the data remains unchanged.

OHE (One Hot Encoding) → Dummy Variables.

city
BOM
DEL
CHN
BOM
CHN
CHN

- ① Extract the unique values of the column
['BOM', 'CHN', 'DEL']

- ② Sort the above list in ascending order.
['BOM', 'CHN', 'DEL']

- ③ Create DUMMY VARIABLES based on above list and show which variable is active for the given record using binary '1'.

city	BOM	CHN	DEL
BOM	1	0	0
DEL	0	0	1
CHN	0	1	0
BOM	1	0	0
CHN	0	1	0
CHN	0	1	0

Hypothesis Testing

Pre-requisite

① Significance Level [SL] / α values ^(alpha)

→ % error value the project can tolerate

→ usually SL values is determined in two ways:

① Assumed values: 0.05, 0.01, 0.1
5% 1% 10%

② Identify the nearby values of error using cross-validation method.
(practical SL)

② Confidence Level (CL): (how much minimum accuracy expected from the model report...).

$$CL = 1 - SL.$$

determines the minimum evaluation metric threshold to be achieved by the trained data product.

(trained model, analysis's report, strategy report...)

③ p-value: (probability value)

calculated SL value derived from statistical testing formulae/tools

Samples



Statistics

(Approximation)

[There can be errors or inaccuracies in the result or prediction]

- How much % accuracy must be maintained by the process?
- How much % errors in prediction or in analysis report, the project can tolerate?

eg: $SL = 0.2$

out of 10 predictions done by the data product, the business shall allow 2 incorrect predictions

Statistical Testing

test to identify if the given statistical feature is present in my dataset / column.

Test for Normalization

check if the given column has normal distribution

Test for Correlation

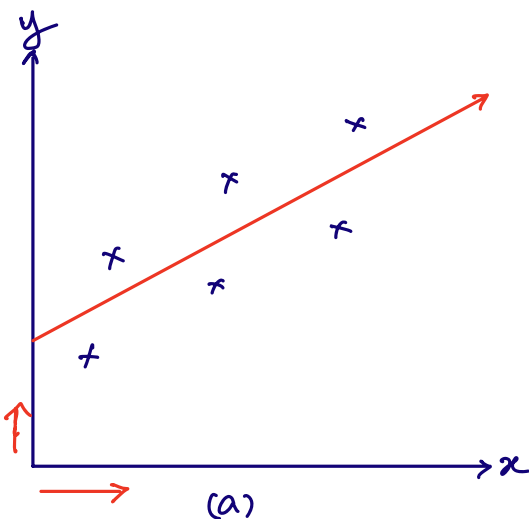
check if correlation exists between two variables / columns.

Test for Feature elimination

removing irrelevant / less significant independent variables.

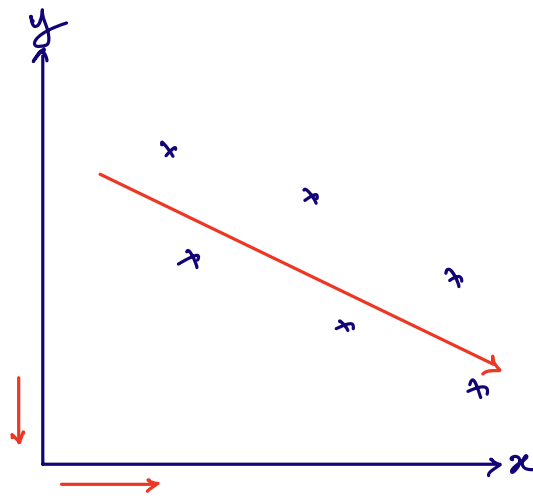
Correlation

Correlation is all about identifying understanding whether the given two variables / columns in the dataset has a linear relationship.



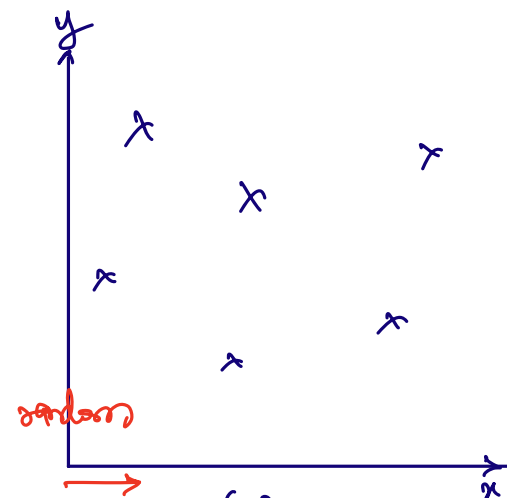
(a) Positive Correlation

with increase in x , y increases.



(b) Negative Correlation

with the increase in x , y decreases.



(c)

No correlation

with increase in x , y is random

Goal:

determine whether the evidence provided by the sample data is STRONG enough to support the said assumption.

(assumption)
Hypothesis
Testing

Goal: Test if col1 and col2 have linear relationship.

☹️ $H_0 \rightarrow$ col1 and col2 has NO LINEAR RELATIONSHIP

😊 $H_a \rightarrow$ col1 and col2 has LINEAR RELATIONSHIP.

who win?
via statistical testing

Null Hypothesis (H_0)



Alternate Hypothesis (H_1) (H_a)



Goal: Check if the given liquid is an ACID or BASE

Note: when it comes to formulating the hypothesis, always remember:

① Hypothesis must be formed as a BINARY OUTCOME.

② Only ONE STATISTICAL feature must be tested at a time.

Hypothesis 1 \rightarrow Check if the given liquid is ACID.

☹️ $H_0 \rightarrow$ given liquid is NOT an ACID.

😊 $H_a \rightarrow$ given liquid is an ACID.

Hypothesis 2 \rightarrow Check if the given liquid is BASE.

☹️ $H_0 \rightarrow$ given liquid is NOT a BASE.

😊 $H_a \rightarrow$ given liquid is a BASE.

statistical Testing

