

(3)

Develop algorithm field of study machine learning.  $\rightarrow$  learnable knowledge

1. Data mining storage - HDD, flash memory

2. abstraction

} different steps in  
ML

$\hookrightarrow$  Data cleaning

$\hookrightarrow$  Feature extraction  $\rightarrow$  80% training data, 20% test data

$\hookrightarrow$  standardization

3. Generalization  $\rightarrow$  future prediction

4. Evaluation

1. Data is very important which is stored in anywhere.

2. Abstraction give the meaning to the data.

Data cleaning  $\rightarrow$  remove the missing value.

Feature extraction  $\rightarrow$  eg: bomb boom gender not empty

3. Generalization  $\rightarrow$  transform the data, training the data

into a model. Write models.

4. Evaluation  $\rightarrow$  test the unseen data, the model is performing well or not.

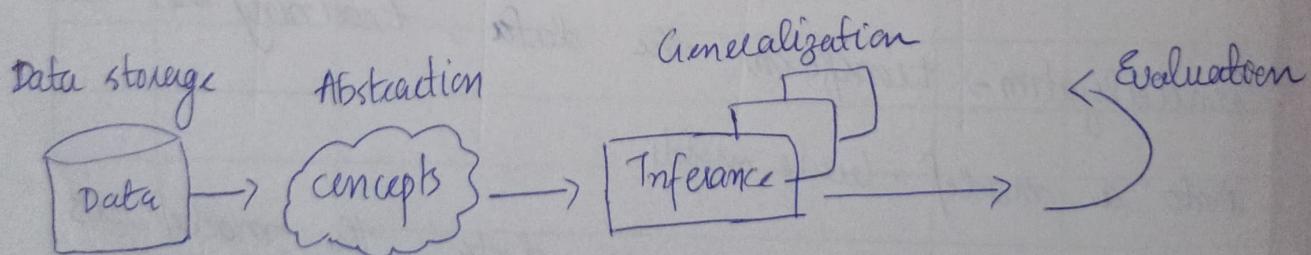
1. Overfitting
2. Underfitting

1. The model perform well on the training data but the model perform poor on the test data.
- Model fails to analyze the underline pattern of data.
  - Model ~~does not~~ fit the data

	sqft	rooms		
eg: city	1500	3	2	4,000 0
city	1480	3	2	1,000 0 X

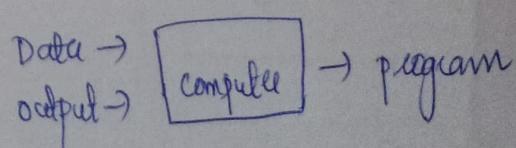
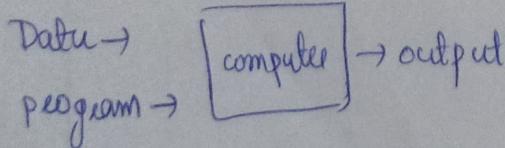
2. The model perform both poorly on training and test data

Good fit - Model perform both well on training & test data



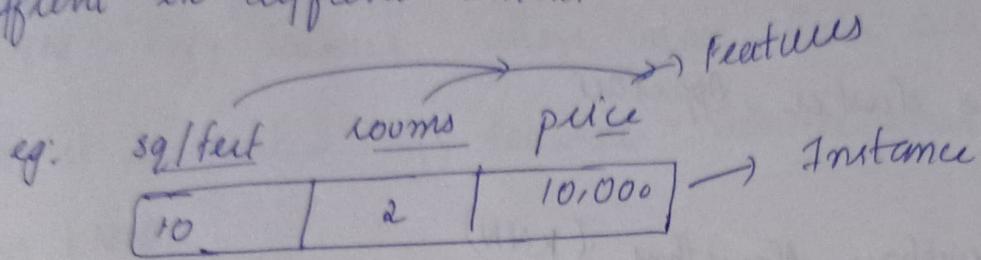
Traditional

Machine learning



## Input Data

dataset - is a collection of unit of observation - which is different in different context



Features can be different types :

① Numeric eg: height, weight

② categorical / Nominal

eg: male/female - gender

③ ordinal → education qualification - SSLC, 12<sup>th</sup>, UG, PG

variant of categorical eg: size of shirt

{ different type of data }

There are a type ML Model:

Predictive & Descriptive :

→ supervised → classification

Predictive : It is used for future predictions and the data consists of features and target variables (class label)

Descriptive : No feature is superior over other, which does not have a target variable.  
'unsupervised', clustering

predictive

↳ SUM, kNN, Decision tree → classification different algo

Descriptive

↳ k-means cluster, Apriori

## k-Nearest Neighbour Algorithm (kNN)

- supervised classification algorithm
- imp concept is to find the similarity of neighbor data

eg: Tomato : Sweetness = 6  
                                    Gumminess = 4

Ingredient	Sweetness	Gumminess	Food type	Distance to the food
Grape	8	5	Fruit	$\sqrt{(6-8)^2 + (4-5)^2} = \sqrt{5}$
Greenbeam	3	7	vegetable	$\sqrt{(6-3)^2 + (4-7)^2} = \sqrt{18}$
Nuts	3	6	protein	$\sqrt{(6-3)^2 + (4-6)^2} = \sqrt{18}$
Orange	7	3	Fruit	$\sqrt{(6-7)^2 + (4-3)^2} = \sqrt{2}$

→ Identify distance measurement

→ k value should not too much low & too much high

→ k=1, choose needed value  $\sqrt{2}$ , so we predict the tomato is a fruit.

→ k=3, majority class has priority

→ k value should be optimal, we should select the no. of folds

$$kNN : - \text{distance}(p_1, q_1) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

evaluation

→ kNN - lazy learner, when new data occur it only performs,  
so it is called lazy learner. Generalization stage does not  
exist in kNN.

It only stores the data, when a data occurs it calculates  
and analyse.

$$\text{dest}(p, q) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_n - q_n)^2}$$

$k = \sqrt{4} = 2$ , so we consider  $\sqrt{2}$  and  $\sqrt{5}$  so it's predicted that  
food type is fruit.

• kNN algorithm is a simple & effective algorithm that makes no  
assumption about the undefined data distribution. The letter  
 $k$  is a variable term implied attorney number of neighbors  
should be used.

In this example locating tomatoes nearest neighbors  
requires a nearest function, that measures the similarity b/w  
2 instances. Traditionally kNN algorithm uses 'euclidean  
distance' for measuring the similarity

- Simple and effective classification algorithm
- New data arrives at testing phase so that learning phase is bad.  
learning phase >> testing phase.
- No generalization and abstraction
- No model creation
- When missing values are there additional processing is needed
- Make no assumptions

✓  $k$  = No. of neighbours to classify to unlabeled data

✓ Euclidean Euclidean distance  $\rightarrow$  (most used for measurement)

Manhattan & Minkowski  $\rightarrow$  (less used)

✓ lazy learning & instant learning & rule learning  $\rightarrow$  ENN

Q. Based on a survey conducted in our institution students are classified based on the 2 attributes: academics excellence and other activities. Given the following data identify the classification of a student with  $x=5$  and  $y=7$  using kNN algorithm. (choose  $k=3$ )

$x$ (Academic Excellence)	$y$ (Other Activities)	Z (Classification)
8	6	outstanding
5	6	good
7	3	good
6	9	outstanding

Ans: distance measurement:

$$\sqrt{(5-8)^2 + (7-6)^2} = \sqrt{9+1} = \sqrt{10} \quad \text{outstanding}$$

$$\sqrt{(5-5)^2 + (7-6)^2} = \sqrt{0+1} = \sqrt{1} \quad \text{good}$$

$$\sqrt{(5-7)^2 + (7-3)^2} = \sqrt{4+16} = \sqrt{20} \quad \text{good}$$

$$\sqrt{(5-6)^2 + (7-9)^2} = \sqrt{1+4} = \sqrt{5} \quad \text{outstanding.}$$

$k=3$  so that the nearest neighbor  $\sqrt{10}, \sqrt{1}$  and  $\sqrt{5}$  outstanding good outstanding ie, outstanding

## Probilistic Reasoning

### Naive Bayes classification

eg: probability that today might be a rainy day (70% chance)  
10 days takes with the same and similar condition of today.  
then 7 days are rainy and concluded that  $7/10$  i.e., 70%

application:

eg: weather forecasting, Text classification (spam /not spam)

↳ probability that the incoming msg is spam or ham

↳ Anomaly detection in a network

↳ diagnosis of a disease, based on symptoms

• deals with probability of likelihood of an event.

(↳ Based on many diff trends [likelihood of an event])

eg: coin tossing  $\rightarrow$  tail / event  $\rightarrow$  head

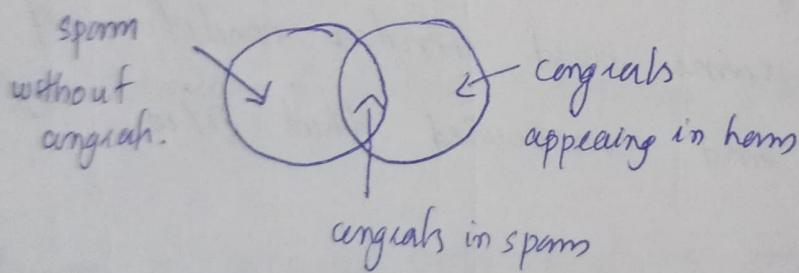
trial: Incoming msg , event  $\rightarrow$  spam

- 'event' and 'trial' main dependency of likelihood
- All the features have importance  
even if the feature poses small role even though less impact they also considered
- Mutually exclusive & exhaustive event  
i.e. Both events does not come together  $\rightarrow$  Either spam/harm

example:

~~eg~~ spam mail - 20% harm mail  $\rightarrow$  80%

$$i.e., 100 - 80 \rightarrow 20\%$$



- probability that msg is spam and msg contain congrats  
(dependent event)

"dependent event" are (probabilistic learning) considered as the basis for predictive modeling.

- Independent event.

- probability (spam  $\cap$  congratulation)

Independent event



$$P(\text{spam}) \times P(\text{congratulation})$$

$$.20 \times .05$$

$$= 0.01 \rightarrow 1\%$$

spam mig provided congratulation occur  $\rightarrow$  dependent went.

Buyer,

$$\checkmark P(A|B) = \frac{P(A \cap B)}{P(B)}$$

conditional probability

A mail which is spam and which includes congrats msg.

one type

$\hookrightarrow$  posterior probability

: A new instance occurs how the evaluation occurs.

$$\checkmark P(A \cap B) = P(A|B) \times P(B)$$

$$P(A \cap \bar{B}) = P(\bar{B} \cap A)$$

$$P(B|A) = \frac{P(B \cap A)}{P(A)}$$

$$P(\bar{B} \cap A) = P(B|A) \times P(A).$$

i.e.,

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

prior probability  $\rightarrow$  eg: In spam msg there is already congrats.  
msg.

$$P(\text{spam} | \text{congrats}) = \frac{P(\text{congrats} | \text{spam}) P(\text{spam})}{P(\text{congrats})}$$

posterior probability

likelihood

prior probability

marginal likelihood.

### Naive Bayes

problem

supervised

chills	Running Nose	Headache	Fever	Has Flu
Y	N	mild	Y	N
Y	Y	No	N	Y
Y	N	strong	Y	Y
N	Y	mild	Y	Y
N	N	No	N	N
N	Y	strong	Y	Y
N	Y	strong	N	N
Y	Y	mild	Y	Y

chills	Fever	Headache	Running Nose	Has Flu
Y	Y	mild	No	?

→ posterior probability

solution:

$$\text{Has Flu Yes} = 5/8 = 0.625 \leftarrow \text{prior probability}$$

$$\text{Has Flu No} = 3/8 = 0.375$$

$$P(\text{Flu} | \text{chills} = Y) = \frac{3}{5} = 0.6$$

$$P(\text{Fever} = Y | \text{Flu} = Y) = \frac{4}{5} = 0.8$$

$$P(\text{Headache} = \text{mild} | \text{Flu} = Y) = \frac{2}{5} = 0.4$$

$$P(\text{Running Nose} = \text{No} | \text{Flu} = Y) = 1/5 = 0.2$$

$$\therefore (0.6 \times 0.8 \times 0.4 \times 0.2) \times 0.625$$

$$\begin{matrix} \text{condition} \\ \text{probability} \end{matrix} = \underline{\underline{0.024}}$$

$$P(\text{Flu} = N | \text{chills} = Y) = 1/3 = 0.33$$

$$P(\text{Fever} = Y | \text{Flu} = N) = 1/3 = 0.33$$

$$P(\text{Headache} = \text{mild} | \text{Flu} = N) = 1/3 = 0.33$$

$$P(\text{Running Nose} = \text{No} | \text{Flu} = N) = 2/3 = 0.66$$

$$\therefore (0.33 \times 0.33 \times 0.33 \times 0.66) \times 0.375$$

$$= \underline{\underline{0.02371}} \times 0.375 = 0.009$$

$$\begin{matrix} 0.33 & \times \\ 0.33 & \times \\ \hline 0.99 \end{matrix}$$

$$0.990$$

$$\begin{matrix} 0.990 & \times \\ \hline 0.9989 \end{matrix}$$

$$0.02371$$

$$0.375$$

so here  $0.024 > 0.009$  i.e., has Flu = Y

### problem

- Estimate conditional probabilities of each attributes {color, legs, height, smelly} for the species classes {M, H} using the data given in the table
- Using these probabilities estimate the probability value for the new instance {color=green, legs=2, height = tall and smelly = No?}

No	Color	Legs	Height	Smelly	Species
1	white	3	short	yes	M
2	green	2	tall	No	M
3	green	3	short	yes	M
4	white	3	short	yes	M
5	green	2	short	No	H
6	white	2	tall	No	H
7	white	2	tall	No	H
8	white	2	short	yes	H
9					

### solution

color	legs	height	smelly	species
green	2	tall	No	?

Has species M =  $4/8 = .5$

Has species H =  $4/8 = .5$

$$\begin{array}{l}
 p(M=4 | \text{color} = \text{Green}) = 2/4 = .5 \\
 p(M=4 | \text{legs} = 2) = 1/4 = .25 \\
 p(M=4 | \text{Height} = \text{Tall}) = 1/4 = .25 \\
 p(M=4 | \text{smelly} = \text{No}) = 1/4 = .25
 \end{array}
 \quad
 \begin{array}{l}
 p(H=4 | \text{color} = \text{Green}) = 1/4 = .25 \\
 p(H=4 | \text{legs} = 2) = 4/4 = 1 \\
 p(H=4 | \text{Height} = \text{Tall}) = 2/4 = .5 \\
 p(H=4 | \text{smelly} = \text{No}) = 3/4 = .75
 \end{array}$$

$$\therefore (.5 \times .25 \times .25 \times .25) \times .5 \\
 = 0.0039025 \times .5 \\
 = 0.00390$$

$$\therefore (0.25 \times 1 \times .5 \times .75) \times .5 \\
 = 0.09375 \times 0.5 \\
 = 0.046875$$

i.e.,  $0.046875 > 0.00390$   $\therefore \underline{\text{species}} = H$

Name Bayes Algorithm gives all the features that important.

### Laplace Estimation / Laplace Smoothing

Fruit	Yellow	Sweet	Hung	Total
Mango	550	950	0	650
Banana	400	300	350	400
Others	50	850	400	150
Total	800	850	400	1200

Laplace estimator also known as Laplace smoothing is a technique used in probability of statistics to handle the problem of zero

probability in categorical data. It involves adding a small value (often 1) to each count of observed icons which ensures that no probability is ever zero. Even if some events have not been observed in examples, Laplace estimator ensures that all possible events have a non-zero probability making the model more robust especially when dealing with more dataset.

$$P(\text{mango}) = \frac{650}{1200} = 0.542$$

$$P(\text{Banana}) = \frac{400}{1200} = 0.333$$

$$P(\text{Other}) = \frac{150}{1200} = 0.125$$

$P(\text{Facet}/\text{yellow-sweet, long})$

$$P(\text{yellow/mango}) = \frac{3050}{650} =$$

$$P(\text{Sweet/mango}) = \frac{450}{650}$$

$$P(\text{long/mango}) = \frac{650}{650} =$$

when we multiple it become zero so avoid this,

$$\frac{O}{650} = \frac{O+1}{650+k} \quad k=2$$



Binary Category.

$$P(\text{yellow/mango}) = \frac{350}{650} = \frac{300+1}{650+2} = \frac{351}{652}$$

$$= 0.539$$

$$\therefore P(\text{sweet/mango}) = \frac{451}{652} = \underline{\underline{0.692}} \quad 0.692$$

$$\therefore P(\text{long/mango}) = \frac{1}{652} = \underline{\underline{0.002}}$$

$$P(\text{yellow/banana}) = \frac{400}{400} = \frac{401}{402} = \underline{\underline{0.998}}$$

$$P(\text{sweet/banana}) = \frac{301}{402} = \underline{\underline{0.749}}$$

$$P(\text{long/banana}) = \frac{851}{402} = 0.873$$

$$\rightarrow P(\text{mango/yellow, sweet, long})$$

$$= P(\text{mango}) = P(\text{yellow/mango}) \times$$

$$P(\text{sweet/mango}) \times$$

$$P(\text{long/mango})$$

$$= 0.542 \times 0.00050 \times 0.531 \times 0.002 = 0.000$$

$$0.542 \times 0.002 \times 0.692 \times 0.539 = \underline{\underline{0.000050}}$$

$$P(\text{Banana}) = 0.333 \times 0.998 \times 0.749 \times 0.873$$

=

$$P(\text{other}) =$$

$$P(\text{yellow/other}) = \frac{51}{152} = 0.335$$

$$P(\text{sweet/other}) = \frac{851}{152} = 0.598$$

$$P(\text{long/other}) = \frac{401}{152} = 0.657$$

$$\text{plotter}) = 0.125 \times 0.335 \times 5.598 \times 2.657$$

KN  
normal  
size M

(2)

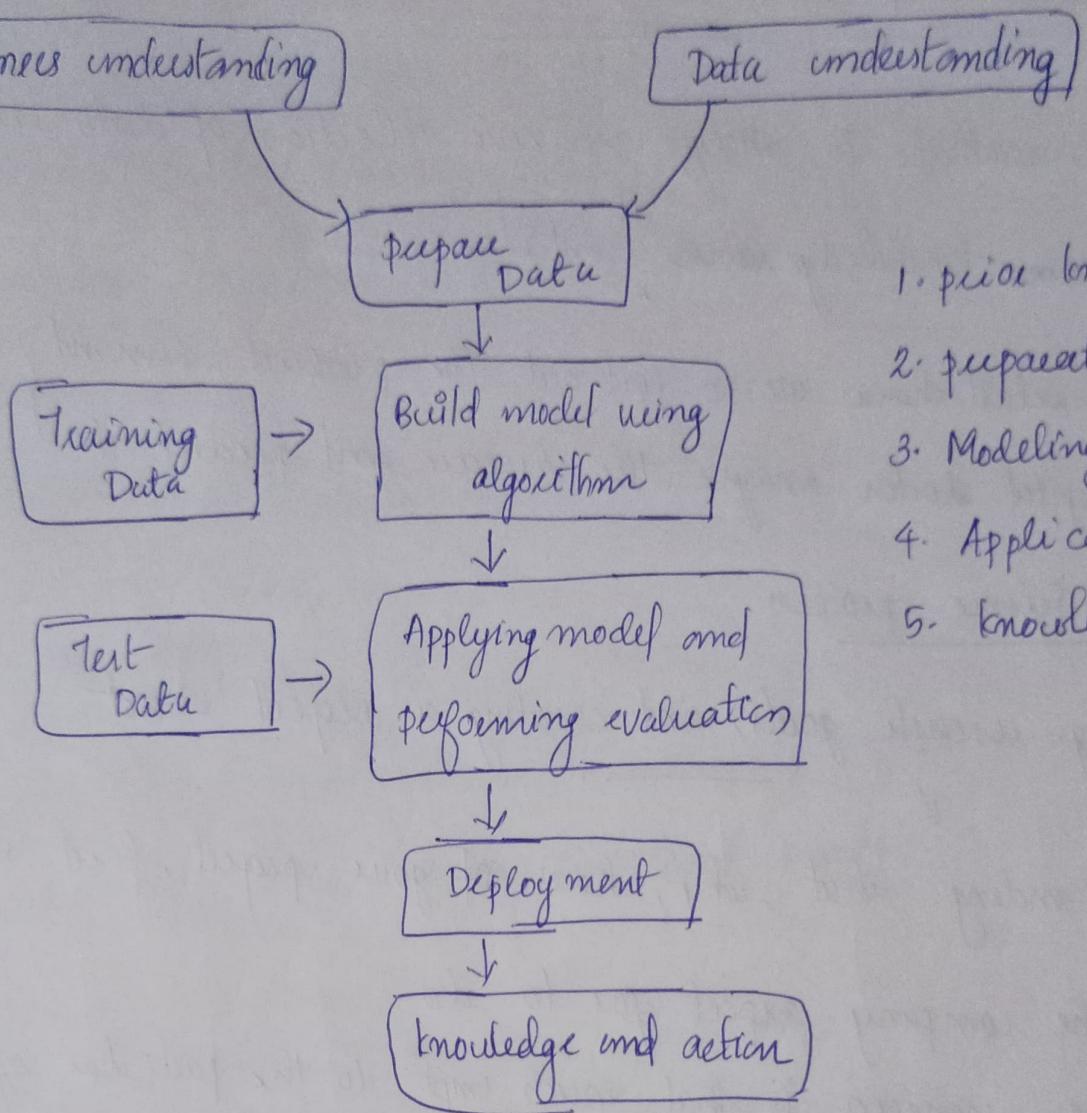
## Data Science Process

It is a collection of technique used to extract values from data.

- We can find patterns from data
- " " " connection within data } through data science
- " " " Relations within data

The method is based on evidence/empirical knowledge/ past historical knowledge/ observations)  $\Rightarrow$  science

Data science is commonly referred as knowledge discovery, data mining, predictive analytics, machine learning



1. prior knowledge
2. preparation
3. Modeling
4. Application
5. knowledge.

### 1. prior knowledge

- The inf which is already known about subject  
eg: credit interest rate for a person
  - The inf about data.
  - prior knowledge about the subject area and the data.
- 10,0000 → 5000 → defaulters  
↓  
income, credit history

## Data Science and Data Science Process

It involves methods to analyse massive collection of data and extract some knowledge <sup>it</sup> which it contains.

eg: In a retail chain we to find out the product demand? data  
In hospital doctor analyse the disease and predict if patient

### \* Data Science process

1. Defining research goals, and creating a project chart



• understanding what, why, how of your project, first step.

→ what the company expect you to do

→ why the scenario is that much imp to the particular project

→ how it is going to be / Is it a biggest strategic picture or a lone wolf (single individual idea)

eg: Decent defaulter



By following the ~~first~~ three steps we come up with a well defined research goal, well defined deliverables and action plan.

→ we can achieve.

creating a project chart:

→ which helps mainly the clients.

① A clear research goal

② The project mission and context

It is a schedule or document which helps the client to give more clarity.

- ③ How you are going to perform your analysis
- ④ what resources you expect to use
- ⑤ proof that is an achievable project or a proof of concepts
- ⑥ Deliverables and measures of success
- ⑦ Timeline

First step is time consuming.

## 2. Retrieving Data

- Start with data stored within the company. If it is not enough they collect the data from the third party.

- Database - (storage)
- Data marts - (subset of data warehouse)
- Data warehouse - (to store and analysing the data)
- Data lakes - (natural data / raw data storage)

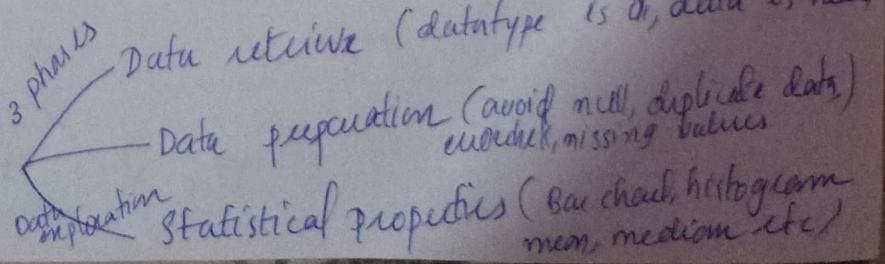
If the data is not enough

(Data is used for companies)

- Don't be afraid to shop around

Data is considered as precious, there is regulation & policies to retrieve the data.

- Data as an asset
- Do quality checks
- Investigate Data



### 3. Data Preparation - Cleaning, Integrating & Transforming data.

#### Cleaning

- we need to sanitize (clean) the data.
  - Data should be sanitized, so that model performs well on this data.
  - Interpretation error (eg: february 30 → that is logically incorrect, : human have 360 age)
  - Inconsistency (eg: take 2 column which represent female, in one column it write female other column write F)
- Data should be represented in a consistent way and should be interpreted correctly. Then occurs 2 types of errors
- a. Interpretation error:

In this case the data is represented in such a way that it is not logically correct. eg: age > 350, february 30  
(this kind of data should not happen)

- b. Inconsistency error:

The data is not represented in a consistent manner eg: representing 'female' in one column and 'F' in another column.

## Error Description

- mistakes during data entry
- Redundant whitespace
- Impossible values
- Outliers

↙

## Possible solutions

- Manual overrules
- use string function (strip(), lstrip(), rstrip())
- Manual overrules
- Validation and if erroneous, treat occurring values

(An observation that seems to be distinct from other observations)

- Different units of measurement  
(height in cm & m)

eg: height of student

5-6 F

J  
entry occurs 1 of  
↙

first we validate  
if it is wrong  
then we convert

## Integrating Data

- Data joining  
→ Data appending

client	item	Month
--------	------	-------

client	Region
--------	--------

client	item	Month	Region
--------	------	-------	--------

} Data joining  
common attribute  
(client)

Data appending - fields are same, additional storage

client	item	month
--------	------	-------

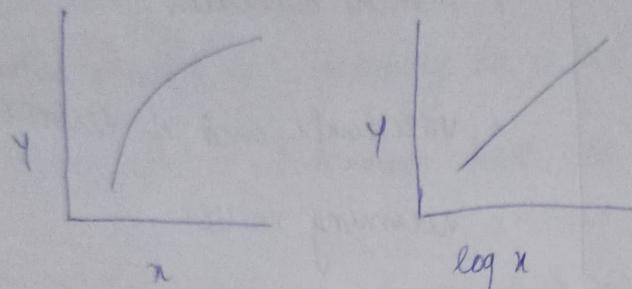
client	item	month
--------	------	-------

client	item	month
--------	------	-------

Data storage  
is problem  
↓  
overcoming  
virtual  
storage

## Transforming Data

- Reducing the no. of variables is one way of transformation
- Identifying the relationship b/w variables.



- Turning variables to dummies. Dummy variables can take a variables True or False.

eg:

customer	year	Gender	sales		customer	year	sales	Male	Femal
1	2015	F	10	⇒	1	2015	10	0	1
2	2015	M	8		1	2016	11	0	1
1	2016	F	11		2	2015	8	1	0

## 4. Exploratory Data Analysis

- We take a deep dive into data and use graphical techniques to gain understanding about your data.
- linking and brushing:

• linking refers to coordination b/w multiple visual representations

• Brushing refers to highlighting or focusing on specific data point in one visualization.

## 5. Building the model

- ① first we need to select the modeling technique and selecting the variables (feature)

eg: we just import the library we don't need to build from the scratch.

### ② Execution of the model

### ③ Diagnosing and comparing the models

- other factors:  
→ whether the model is to be used in the production environment  
→ whether the model requires maintenance

## 2- Execution of the model

### → Model fit : (undesigned patterns should be undetected)

A good model fit means that the model accurately captures the underlying patterns

eg:  $R^2 = 85\%$  (best fit → which means model understood that house price and location is depend on house price)

house price → houseprice | house location

$R^2 = 0$ , didn't analyze the model

$R^2 = 1$

R-squared indicates the proportion of variances in the dependent

model should be based on uncertainty eg: weather forecasting suitable nowadays

variable that is predictable from the independent variable.

e.g.: If you are predicting house prices not all houses are ~~are~~ in the same, there is a variation b/w prices based on factors like size, location etc.

$R^2=1$ , indicates a better fit

$R^2=0$ , indicates a model doesn't explained much of the variation

→ Model Comparison: It involves evaluating multiple models to determine which model is best fit for data or best at making the predictions.

## 6. Presenting Findings and Building applications on top of them

Presenting results to the stakeholders and individualizing your analysis process for repetitive reuse and integration with other tools.