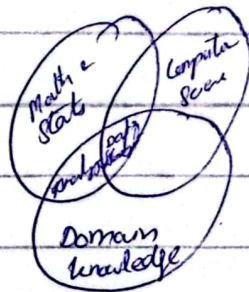


## Lecture 15: Data Science

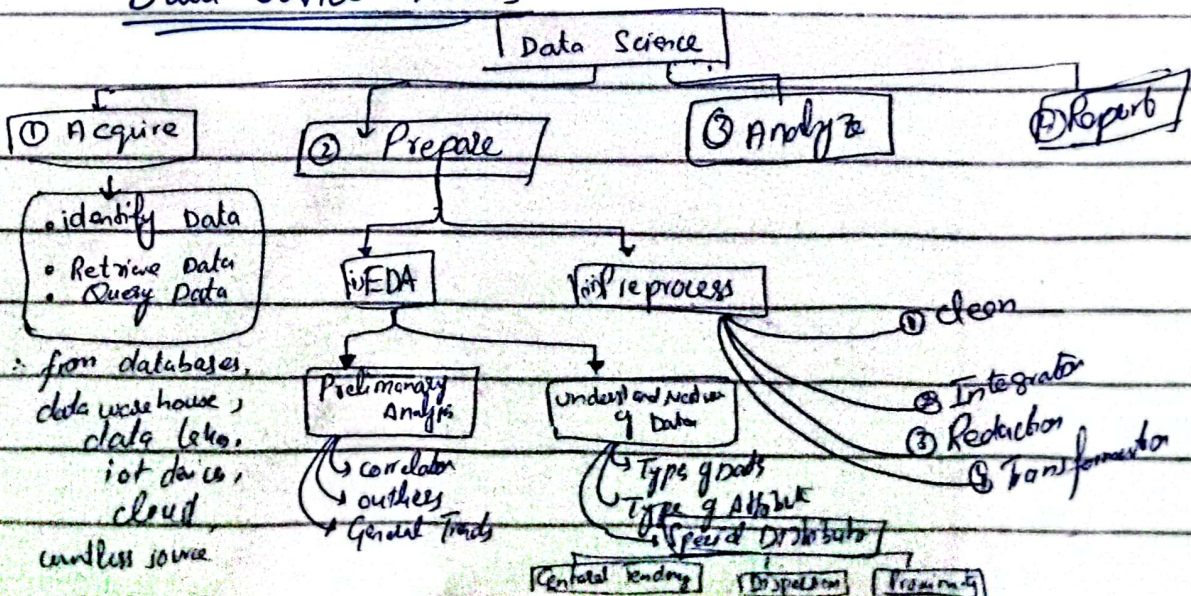
- multi-disciplinary data (statistics, computer science)
- domain knowledge where data is coming
- Extracting insights & knowledge from Data
- Today, data processing is backbone of all.  
(Like social media)  
↳ user preferences
- Data is like "Gold Mine"

### Data Science

" is the field of extracting knowledge and insights from data through a blend of statistics, programming, machine learning & domain expertise.



### Data Science Process





∴ Type of Attribute ??

• what type of data attributes have ??

• Statistical Distribution

## Structured Data

↳ organized & follow defined format

Attribute types usually stored in databases / spreadsheets

→ Nominal data (Pi) / categorical (no order)

→ Binary → (with two categories only)

→ Ordinal (order) smoker / non-smoker

like T-shirt size (S, M, L)

grades (A, B, C)

→ Numerical (Quantitative and measurable)

↳ can perform operation

With absolute zero

like experience

No-absolute zero

like temperature

discrete

absolute like  
2, 5, zero,

continuous

within range of values

## Statistical distributions

① Central Tendency (مركز)

”بیانی و بیانگر فوجک“

Three important measures

→ mean (average)

→ median (middle value of sorted data)

→ Mode (most frequent value)



## ② Measure of Dispersion (تفرق)

(How data is spread from middle value)

important measure

(1) standard deviation & variance (if we use mean)

(2) inter quartile Range (IQR)  
(if we use median)

## ③ Measure of proximity (قریب قریب)

(on 2 attributes usually)

(How value are close to each other)

important measures

① dot product

② Cosine Similarity

③ Correlation

Interest

→ Relationship Strength (strong/weak)

→ Relationship direction (same/opposite)

Data Distribution

Symmetric

Asymmetric (skewed)

Positively skewed

Negatively skewed

## ② Preprocess Data

(i) Pre-processing

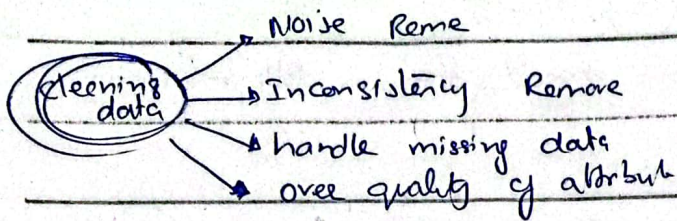
① Clean

(To have quality data)

(accurate, consistent, no missing value, operation etc)

(manage)





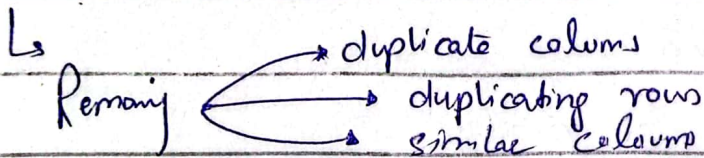
## ② Data integration

(Technique to merging data that are come from multiple sources)

## ③ Data Reduction

(i.e.  $\downarrow$ )

(duplicate columns, columns with similar information)



## Dimensionality Reduction

- To reduce similar data
- drop not relevant data

Technique  $\Rightarrow$  ① PCA (Principal Component Analysis)  
 $\downarrow$   
based on eigen-decomposition,  
is singular value decomposition

## ② P-SNE

✓ Low quality data removal

## ④ Data Transformation