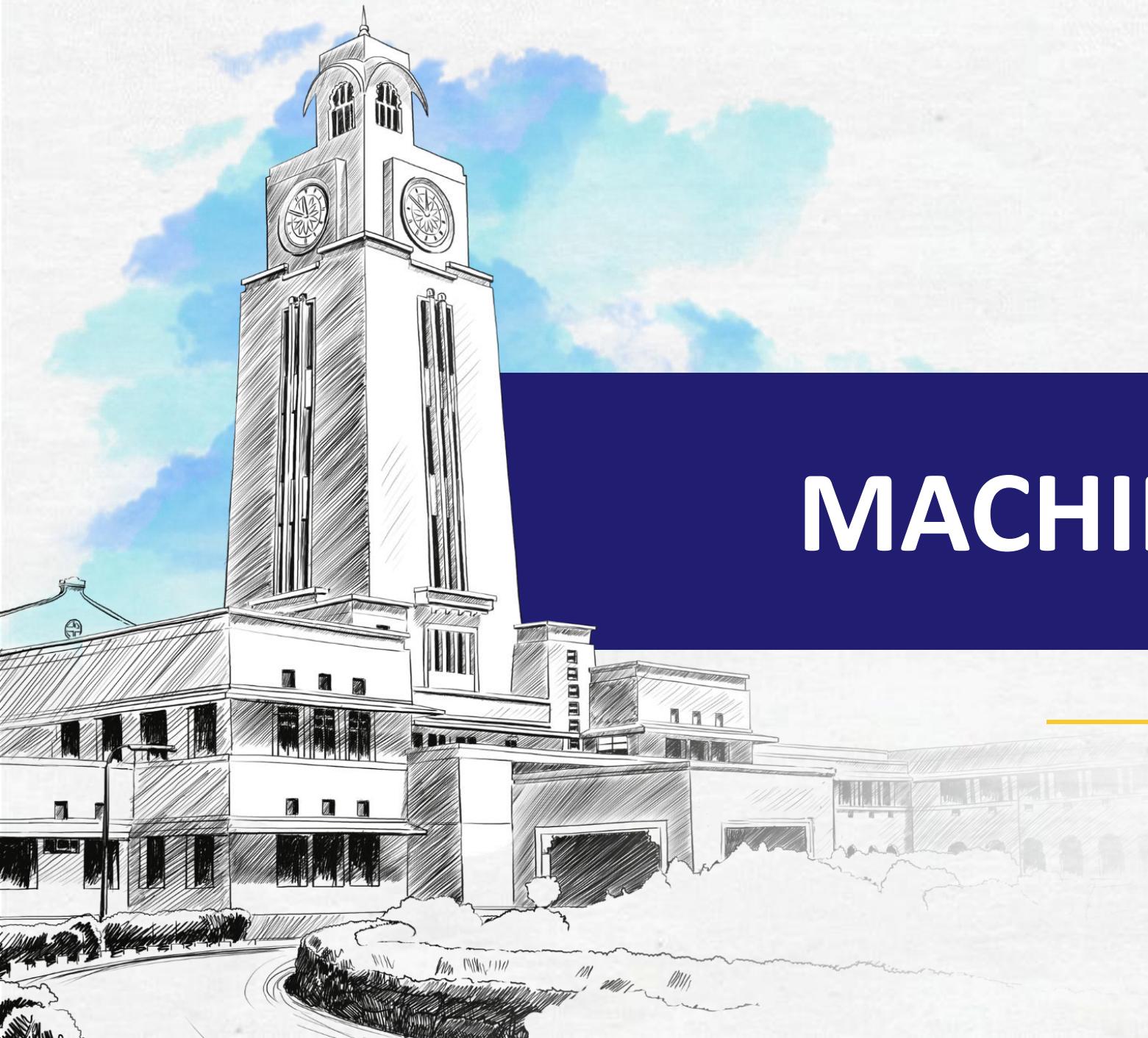




BITS Pilani
Pilani | Dubai | Goa | Hyderabad

MACHINE LEARNING



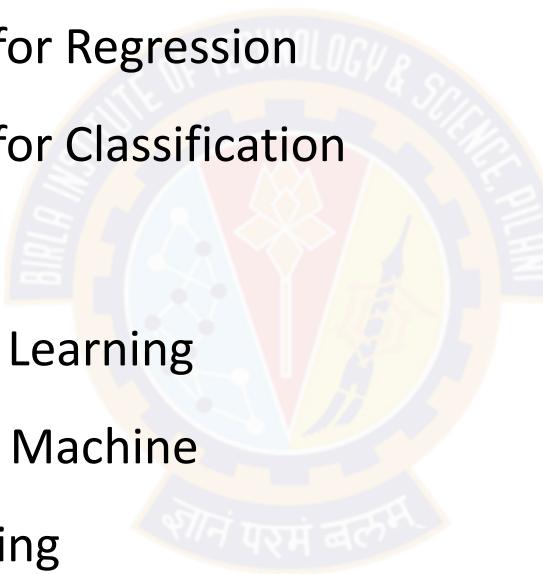


Session 2

(3rd August, 2025)

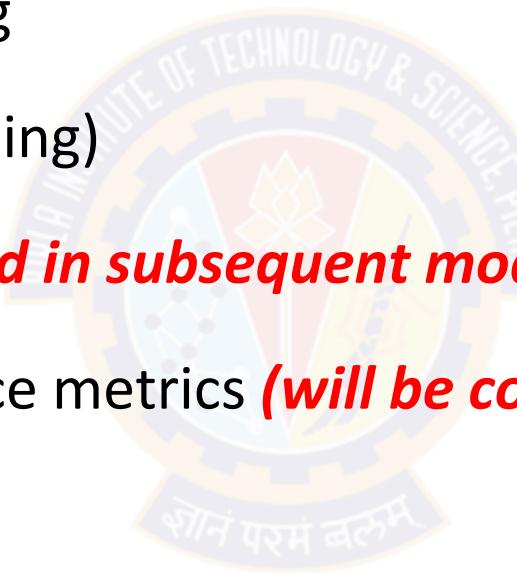
Course Plan

- M1 Introduction
- M2 Machine learning Workflow
- M3 Linear Models for Regression
- M4 Linear Models for Classification
- M5 Decision Tree
- M6 Instance Based Learning
- M7 Support Vector Machine
- M8 Bayesian Learning
- M9 Ensemble Learning
- M10 Unsupervised Learning
- M11 Machine Learning Model Evaluation/Comparison



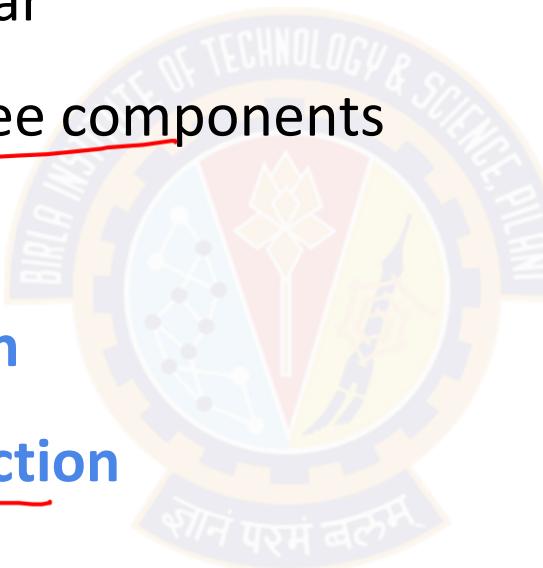
Agenda

- Role of Data
- Data Preprocessing / wrangling
- Data skewness removal (sampling)
- Model Training (*will be covered in subsequent modules*)
- Model Testing and performance metrics (*will be covered in subsequent modules*)

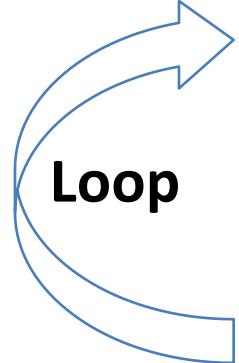


ML in a Nutshell

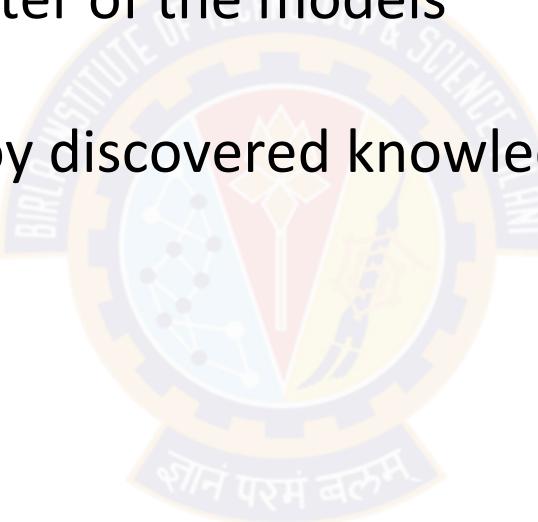
- Tens of thousands of machine learning algorithms
 - Hundreds new every year
- Every ML algorithm has three components
 - ✓ ○ **Data Representation**
 - ✓ ○ **Parameter Optimization**
 - ✓ ○ **Model Evaluation, Selection**



ML in Practice



- Understand domain, prior knowledge, and goals
- Data integration, selection, cleaning, pre-processing, etc.
- Learn optimal parameter of the models
- Interpret results
- Consolidate and deploy discovered knowledge



Definition of Data

- Collection of **data objects** and their **attributes**
- An **attribute** is a property or characteristic of an object
 - Examples: eye color of a person, temperature, etc.
 - aka variable, field, characteristic, dimension, or feature
- A collection of attributes describe an **object**
 - Object is also known as record, point, case, sample, entity, or instance

Attributes

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No ✓
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Objects

Types of Attributes

- There are different types of attributes

- Nominal

2025aimlb1001@liv

- Examples: ID numbers, zip codes

- Ordinal

- Examples: rankings (e.g., taste of potato chips on a scale from 1-10), grades, height {tall, medium, short}

- Interval

No
True zero value

- Examples: calendar dates, temperatures in Celsius or Fahrenheit.

- Ratio

0°K → no energy

- Examples: temperature in Kelvin, length, counts, elapsed time (e.g., time to run a race)

$$\text{Age} = 3 \text{ (son)}$$

$$\overline{10 \text{ cm}} = 2(5 \text{ cm})$$

A B
No absence of temp.
 $0^\circ\text{C} \rightarrow$ Frozen state

Properties of Attribute Values

- The type of an attribute depends on which of the following properties/operations it possesses:

Difference Between Ratio and Interval

➤ Is it physically meaningful to say that a temperature of 10° is twice that of 5° on

- the Celsius scale?

Not meaningful

- the Fahrenheit scale?

NMF

- the Kelvin scale?

Meaningful \rightarrow true zero exists

➤ Consider measuring the height above average

- If Bill's height is three inches above average and Bob's height is six inches above average, then

would we say that Bob is twice as tall as Bill?

- Is this situation analogous to that of temperature?

Attribute Type	Description	Examples	Operations
Categorical Qualitative	Nominal	Nominal attribute values only distinguish. ($=$, \neq)	zip codes, employee ID numbers, eye color, sex: { <i>male</i> , <i>female</i> }
	Ordinal	Ordinal attribute values also order objects. ($<$, $>$)	hardness of minerals, { <i>good</i> , <i>better</i> , <i>best</i> }, grades, street numbers
Numeric Quantitative	Interval	For interval attributes, differences between values are meaningful. (+, -)	calendar dates, temperature in Celsius or Fahrenheit
	Ratio	For ratio variables, both differences and ratios are meaningful. (*, /)	temperature in Kelvin, monetary quantities, counts, age, mass, length, current

This categorization of attributes is due to S. S. Stevens

Attribute Type	Transformation	Comments
Nominal	Any permutation of values	If all employee ID numbers were reassigned, would it make any difference?
	An order preserving change of values, i.e., $new_value = f(old_value)$ where f is a monotonic function	An attribute encompassing the notion of good, better best can be represented equally well by the values {1, 2, 3} or by { 0.5, 1, 10}.
Interval	$new_value = a * old_value + b$ where a and b are constants	Thus, the Fahrenheit and Celsius temperature scales differ in terms of where their zero value is and the size of a unit (degree).
	$new_value = a * old_value$	Length can be measured in meters or feet.

This categorization of attributes is due to S. S. Stevens

Discrete and Continuous Attributes

▪ Discrete Attribute

Counting

- Has only a finite or countably infinite set of values
- Examples: zip codes, counts, or the set of words in a collection of documents
- Often represented as integer variables.
- Note: **binary attributes** are a special case of discrete attributes

▪ Continuous Attribute

measuring

- Has real numbers as attribute values
- Examples: temperature, height, or weight.
- Practically, real values can only be measured and represented using a finite number of digits.
- Continuous attributes are typically represented as floating-point variables.

$101.1^{\circ}F$

1.57 cm

Important Characteristics of Data

- Dimensionality (number of attributes)
 - High dimensional data brings a number of challenges

- Sparsity

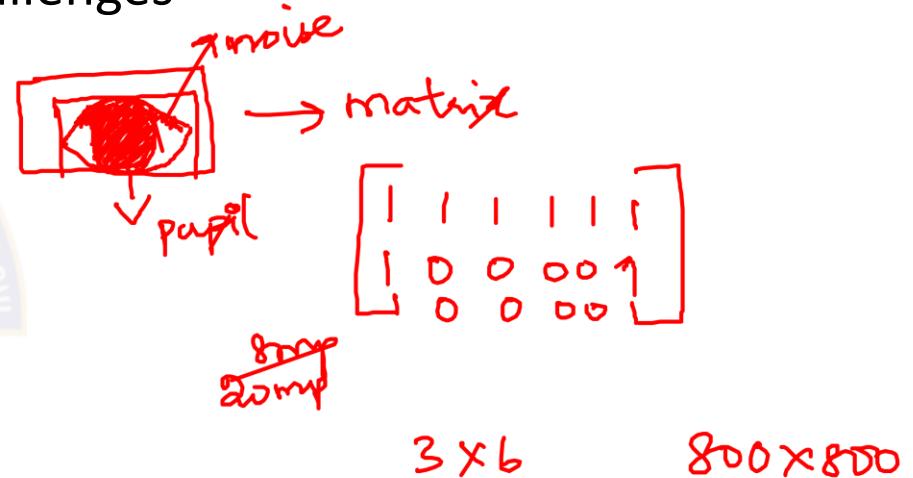
- Only presence counts

- Resolution

- Patterns depend on the scale

- Size

- Type of analysis may depend on size of data



1024×1024

Data Types

- Relational/Object
- Transactional Data
- Document Data
- Web & Social Network Data
- Spatial Data
- Time Series
- Sequence Data



The diagram illustrates various data types mapped onto a sample dataset. The data types are:

- Discrete Numeric (points to ServiceRating)
- Asymmetric Binary (points to IsPriorityCustomer)
- Ordinal (points to CardType)
- Continuous Numeric (points to CreditScore)
- Symmetric Binary (points to isMultipleAccountHolder)

	ServiceRating	IsPriority Customer	CardType	Credit Score	isMultipleAccount Holder
Jack	5	Yes	Platinum	7.5	Yes
Jill	2	Yes	Gold	8.2	No
John	9	No D	Gold	7	Yes

Data Types

- Relational/Object
- Transactional Data
- Document Data
- Web & Social Network
- Spatial Data
- Time Series
- Sequence Data



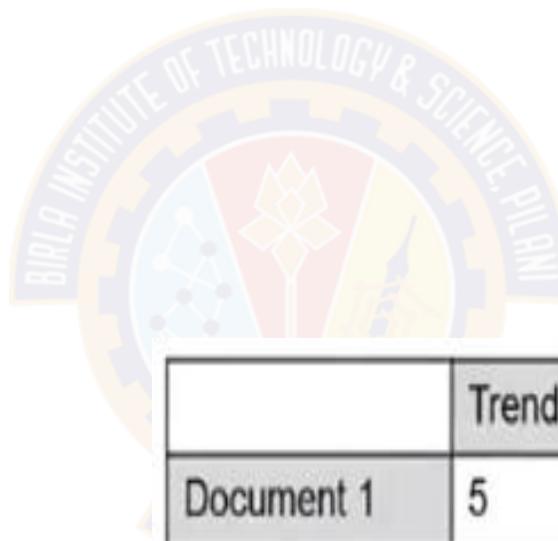
	Purchase 1	Purchase 2
Jack	Paper, Pen, Medicine	Milk, Bread, Egg, Milk			
Jill	Rice, Medicine, Vegetable, Milk	Rice, Egg, Vegetable, Milk	??		
John	Bread, Jam, Butter , Jam	Milk, Bread, Pasta, Medicine			

Transactional data

	Items Bought
Transaction 1	Paper, Pen, Medicine
Transaction 2	Rice, Medicine, Vegetable, Milk
Transaction 3	Milk, Bread, Egg, Milk
Transaction 4	Bread, Jam, Butter , Jam

Data Types

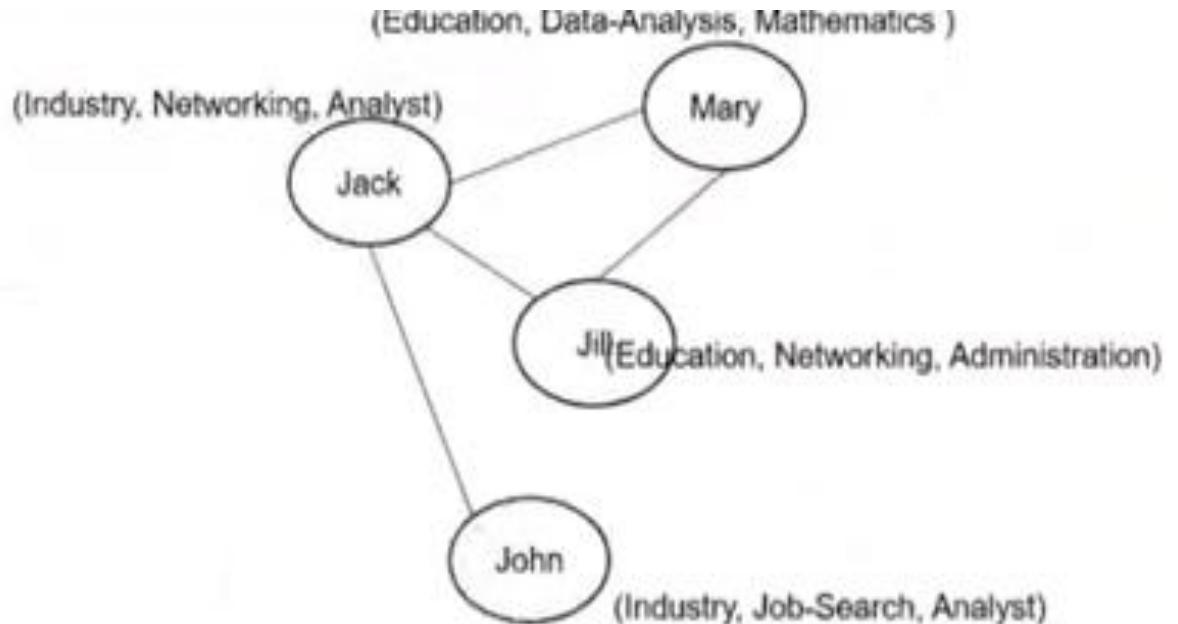
- Relational/Object
- Transactional Data
- Document Data
- Web & Social Network Data
- Spatial Data
- Time Series
- Sequence Data



	Trend	Data	Story	Mining	Cloth
Document 1	5	10	4	8	0
Document 2	5	5	8	0	7
Document 3	2	8	2	4	0

Data Types

- Relational/Object
- Transactional Data
- Document Data
- Web & Social Network Data
- Spatial Data
- Time Series



Categorical Nominal

	Work-Field	Purpose of Connect	Domain of work	No.of Connections	Link to parent	...
John	Industry	Job-Search	Analyst	1	Jack	
Mary	Education	Data-Analysis	Mathematics	2	Jack, Jill	

Data Types

- Relational/Object
- Transactional Data
- Document Data
- Web & Social Network Data
- Spatial Data
- Time Series ✓
- Sequence Data



User	Call Type	Call Duration	Time Stamp	Tower Cell ID	Latitude	Longitude
9341959679262440000	Voice	10	2019-11-20 14:15:01	123456	12.97	77.58
9341959679262440000	Text	0	2019-11-19 11:10:09	123456	12.73	77.82
9221959659362440000	Voice	10	2019-11-20 14:15:01	324576	19.07	72.87

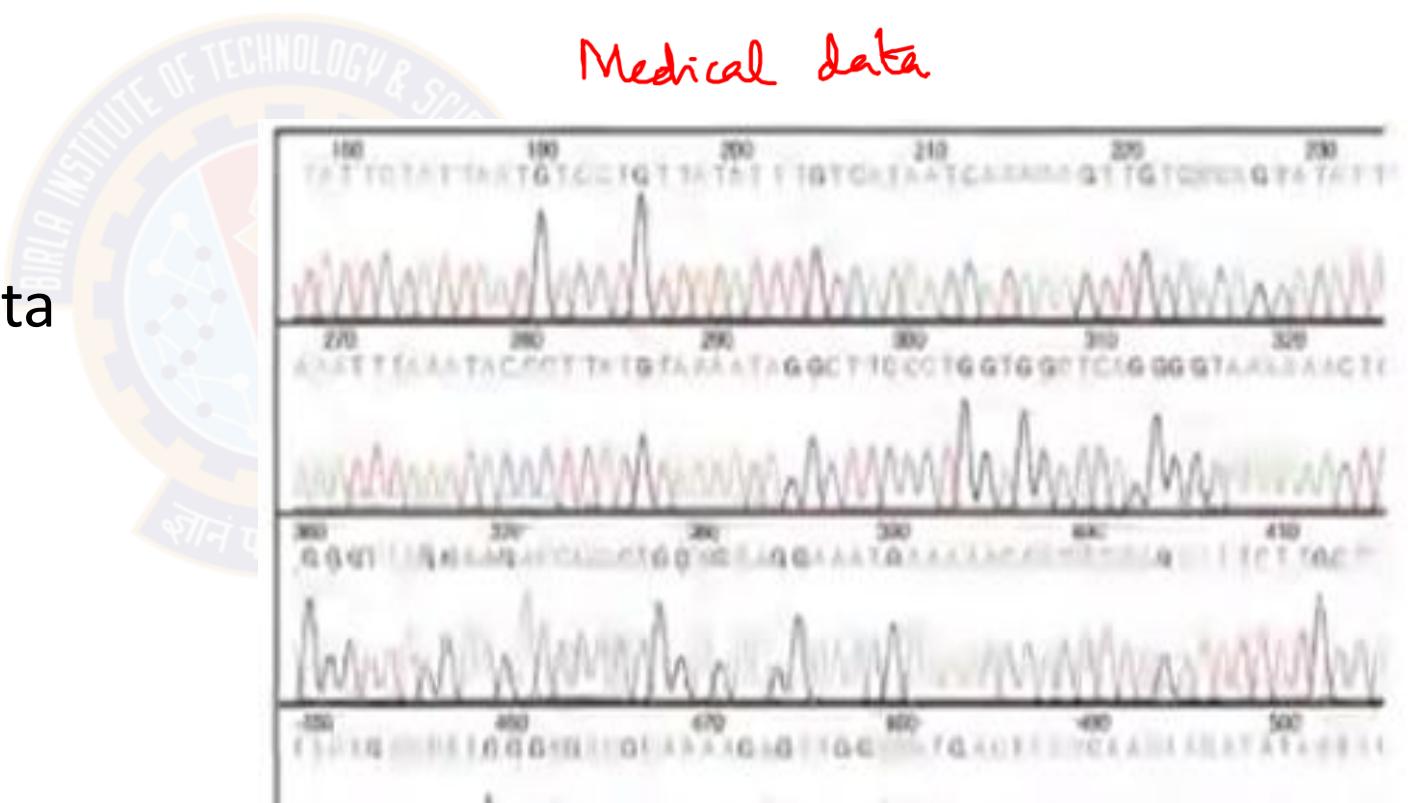
Data Types

- Relational/Object
- Transactional Data
- Document Data
- Web & Social Network Data
- Spatial Data
- Time Series
- Sequence Data



Data Types

- Relational/Object
- Transactional Data
- Document Data
- Web & Social Network Data
- Spatial Data
- Time Series
- Sequence Data



Case Study – 1

Identify the Data Types and attribute types

A bank wishes to analyze its customer base for targeted marketing and needs to segment the customers based on its account information with its branch. Post analysis it might be interested to target potential customers of high income level possessing Titanium card types.

Object

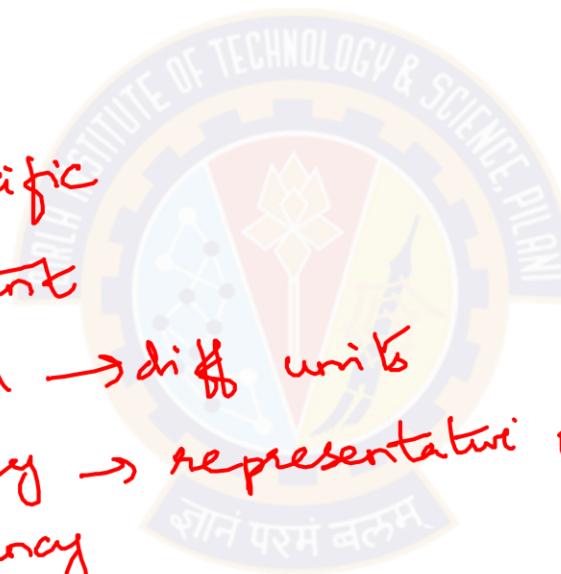
Name	Gender	Service Rating	Is Priority Customer?	Card Type	Credit Score	Is Multiple Account Holder	Income Level	Region
Jack	Male	5	Yes	Platinum	7.5	Yes	Upper	BGLR
Jill	Female	2	Yes	Gold	8.2	No	Middle	DELHI
John	Male	9	No	Gold	7	Yes	Lower	BGLR
Mary	Male	6	No	Gold	6.0	No	Lower	BGLR

Silver, Gold, platinum
hierarchy vs these → ordinal

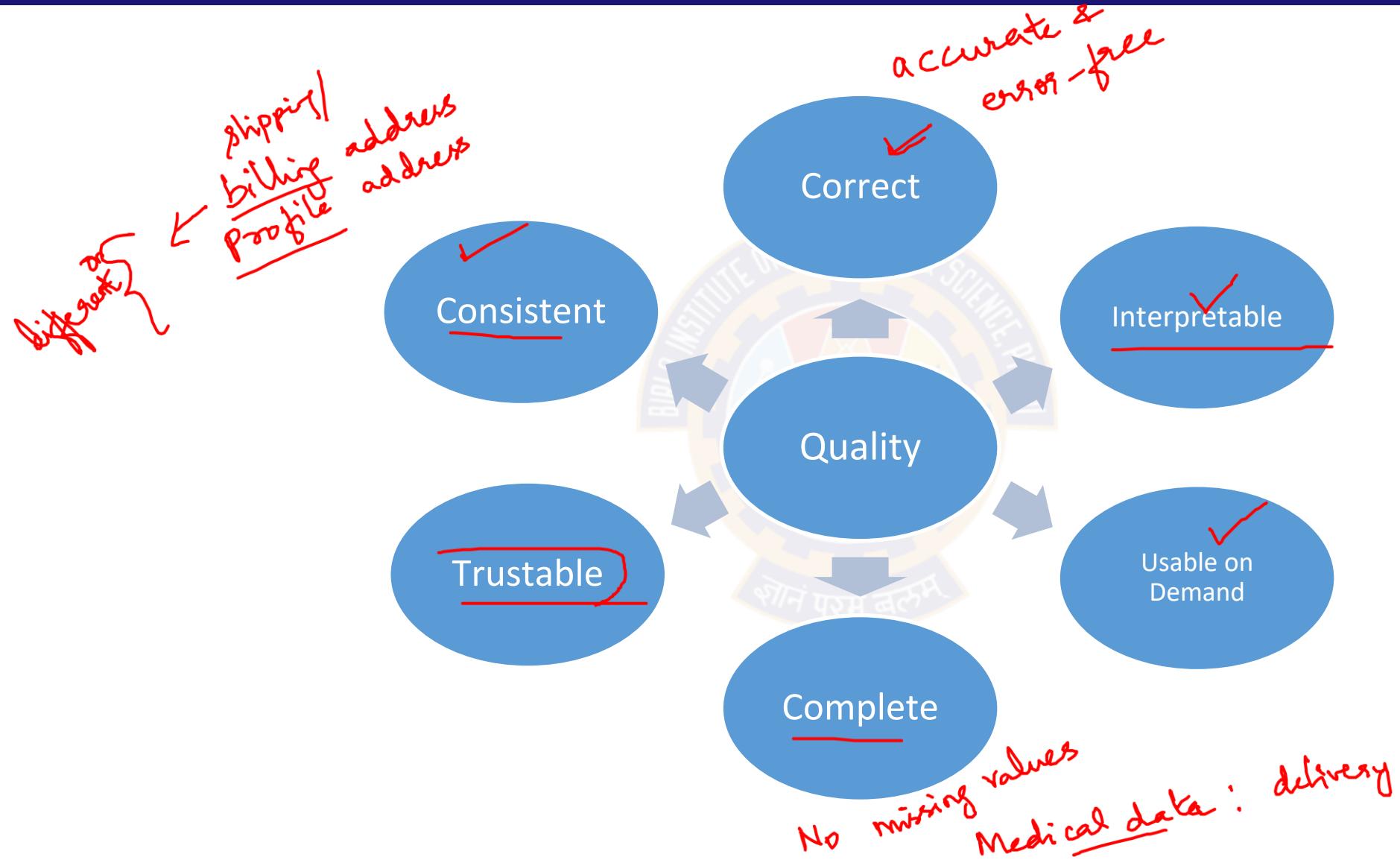
Issues

- 1) missing data
- 2) Noise / outliers
- 3) format
- 4) region specific
- 5) inconsistent
- 6) standard → diff units
- 7) relevancy → representativity of sample
- 8) redundancy
- 9)

Quality



Data Quality

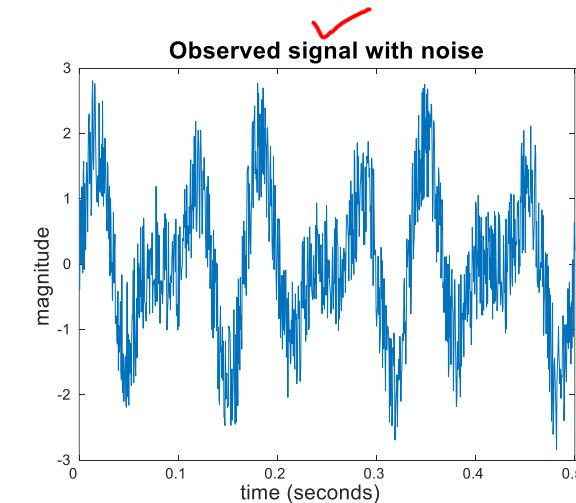
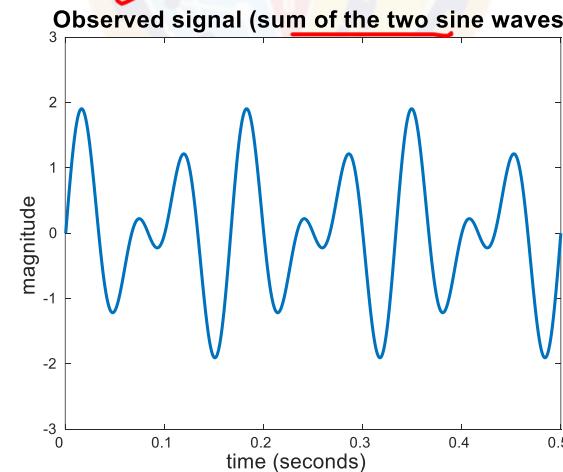
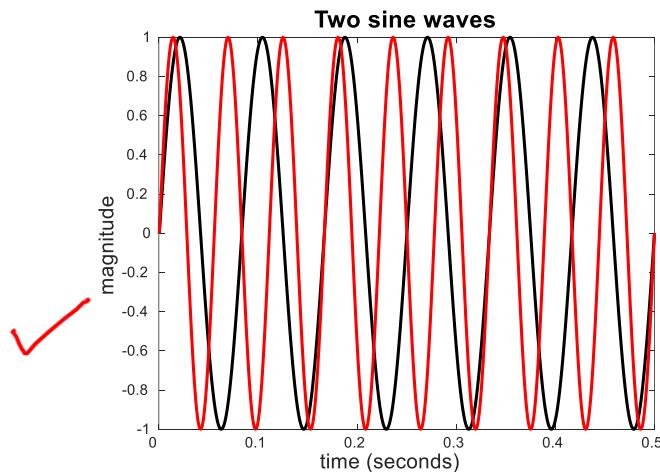


Data Quality

- Poor data quality negatively affects many data processing efforts
- ML example: a classification model for detecting people who are loan risks is built using poor data
 - Some credit-worthy candidates are denied loans
 - More loans are given to individuals that default
- What kinds of data quality problems?
- How can we detect problems with the data?
- What can we do about these problems?
- Examples of data quality problems:
 - Noise and outliers
 - Wrong data ✓
 - Fake data ✓ *Voting* → *multiple region*
 - Missing values
 - Duplicate data → *redundancy*

Noise

- For objects, noise is an extraneous object
- For attributes, noise refers to modification of original values
 - Examples: distortion of a person's voice when talking on a poor phone and "snow" on television screen
 - The figures below show two sine waves of the same magnitude and different frequencies, the waves combined, and the two sine waves with random noise
 - The magnitude and shape of the original signal is distorted

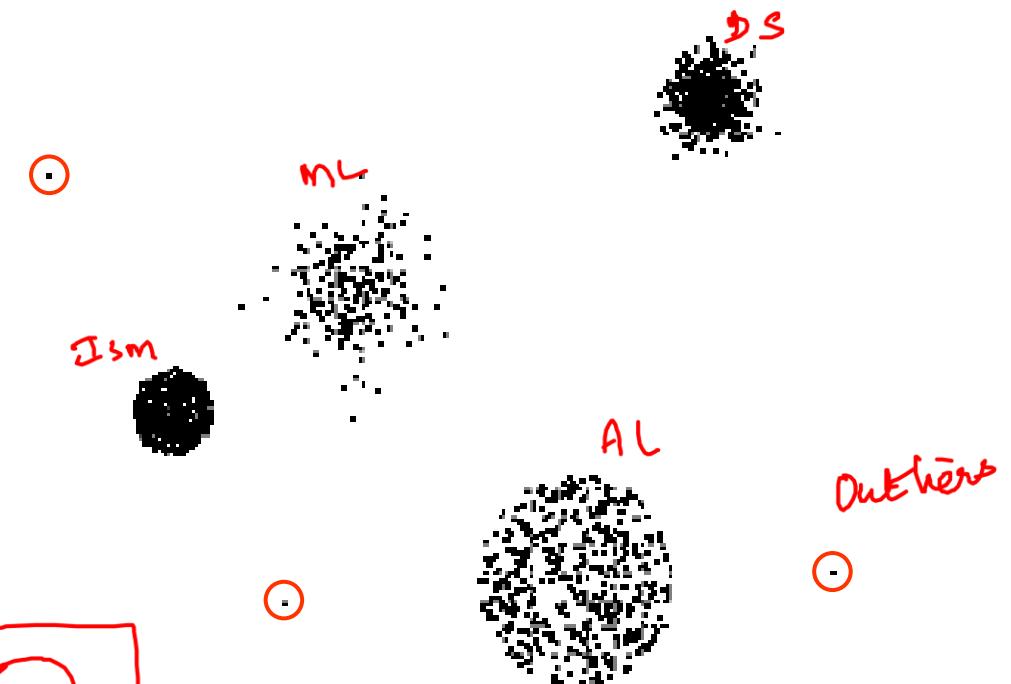
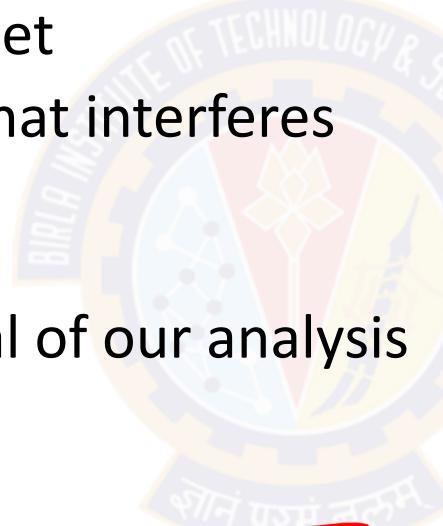
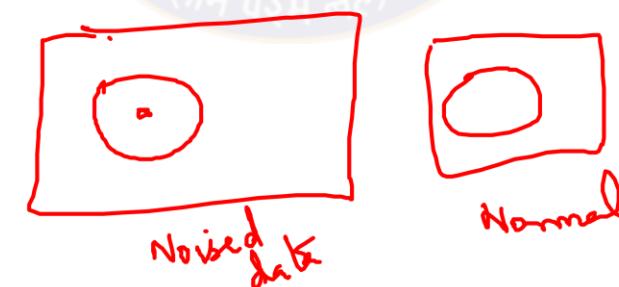


Outliers

- **Outliers** are data objects with characteristics that are considerably different than most of the other data objects in the data set

- **Case 1:** Outliers are noise that interferes with data analysis

- **Case 2:** Outliers are the goal of our analysis
 - Credit card fraud
 - Intrusion detection



Missing Values

➤ Reasons for missing values

- Information is not collected
 - (e.g., people decline to give their age and weight)
- Attributes may not be applicable to all cases
 - (e.g., annual income is not applicable to children)

➤ Handling missing values

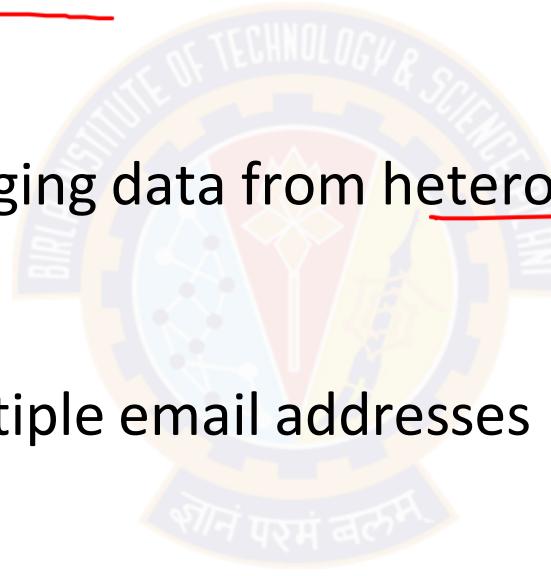
- Eliminate data objects or variables
- Estimate missing values
 - Example: time series of temperature
 - Example: census results
- Ignore the missing value during analysis

Test 1 Test 2 Test 3 Mark 5
4 5 — Avg / Best / mode

Past data 3 days
Thu 31°C
Fri —
Sat 30°C
Sun $\frac{31+30}{2}$

Duplicate Data

- Data set may include data objects that are duplicates, or almost duplicates of one another
 - Major issue when merging data from heterogeneous sources
- Examples:
 - Same person with multiple email addresses
- Data cleaning
 - Process of dealing with duplicate data issues

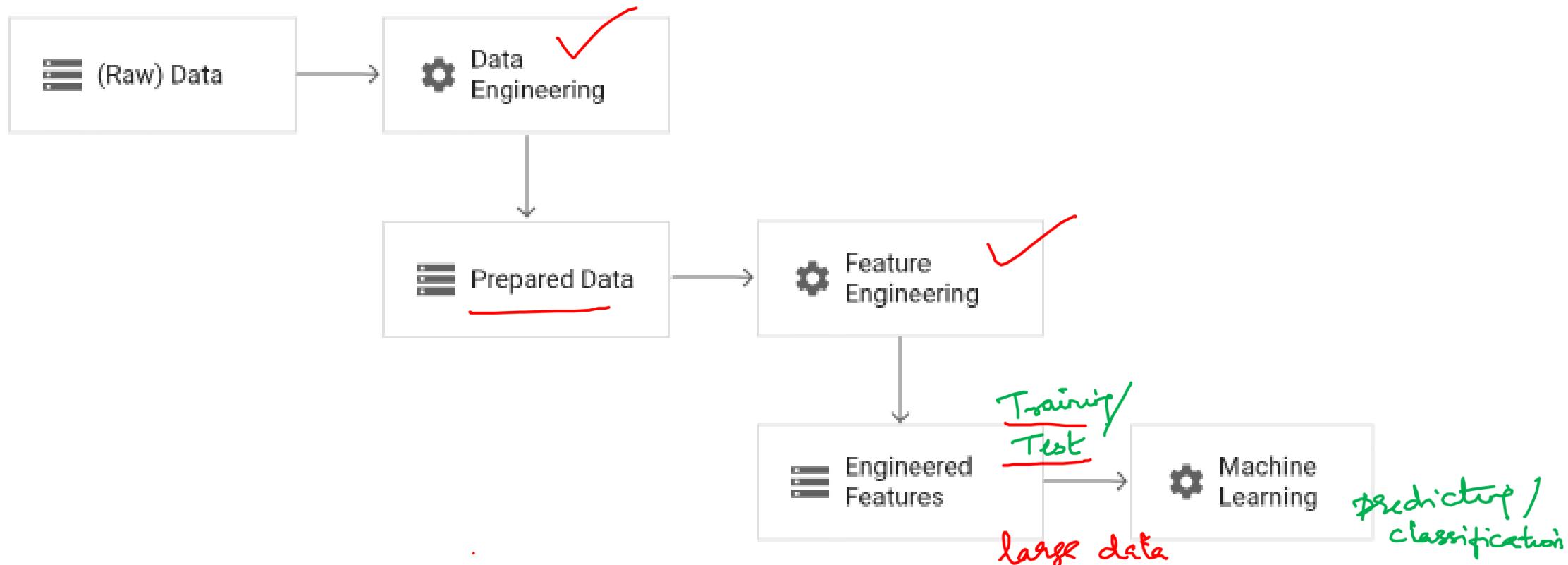


Super market
↓ post
discounts

unique features → discard redundant

Preprocessing

- Preprocessing the data for ML involves both data engineering and feature engineering
- Data engineering : process of converting raw data into prepared data.
- Feature engineering : tunes the prepared data to create the features that are expected by the ML model



① Raw data:

	Name	<u>Gender</u>	Age	Address	Purchase Amt	Date
	John	0 Male	28	"123 Jane St"	105.8	05-8-25
	ARK	1 Female		"453 Thomas St"		7 th Aug 25

missing values diff. format

② Data Engineering: Cleaning, transforming & format the raw data

Name	Gender	<u>Age</u>	Address	Purchase Amt	Date
John	Male	28	123 Jane St	105.8	05-8-25
ARK	Female	30	453 Thomas St	85	07-8-25

Analysis → regular customer

↓ take it bill

3) Feature Engineering \rightarrow creating new features / transforming existing ones to improve the model performance

Eg:
① Encoding $\rightarrow [0 \rightarrow \text{Male}, 1 \rightarrow \text{Female}]$

② Create new ones \rightarrow Date \rightarrow Purchase list for the Month

③ Binning numerical values [Eg: Age: 10-20, 20-30]

④ Scaling : normalization or standardisation

Eg:

Gender (encoded)	Age	Days of Week
0	20-30	Tuesday
1	30-40	Thursday

Stage	Description	Example
Raw Data	Initial data collection	Raw CSV with missing and messy values
Data Engineering	Clean, fill, transform	Fix missing values, format dates, deduplicate
Prepared Data	Structured data	Clean, tabular format ready for features
Feature Engineering	Add/modify features	Encode gender, extract day of week, scale
Engineered Features	Final model input	Numerical/categorical model-ready variables
Machine Learning	Model training/testing	Predict behavior based on engineered features

Case study

•BITS WILP is in collaboration with multiple IT companies interested to upskill and level skill their employee through inducting them in tailored Mtech AIML program. Over a year of successful completion , the student are yet to complete another one semester and enroll in Dissertation to complete the program with certification.

similar academic background irrespective of time of enrollment, seems to score more or less in same range in every semester. Accounting department requires to complete few academic year closure documentation for which , they would have to bill the collaborative organization based on the prospective no. of students who might be eligible for project semester. As of current semester the students have completed their exams but the process is pending for grading. As Data analyst help accounts team to get necessary information with the given available data across all the collaborative program.

Challenge 1 : Insufficient Training Data.

Idea : Trade-off algorithm vs Data readiness

AttributesOfInterest	✓
Name	
Gender	
Age	
DataOfBirth	
Organisation	
JobTitle	
NatureOfJob	
EntranceScore	
EligibilityScore	
PreviousDegree	
WILPBatch	
Section	
ISM	✓
MFML	✓
ACI	
ML	
NLP	
.....	

Data Pre-processing

- Data Aggregation
- Data cleansing
- Instances selection and partitioning
- Feature tuning



Aggregation

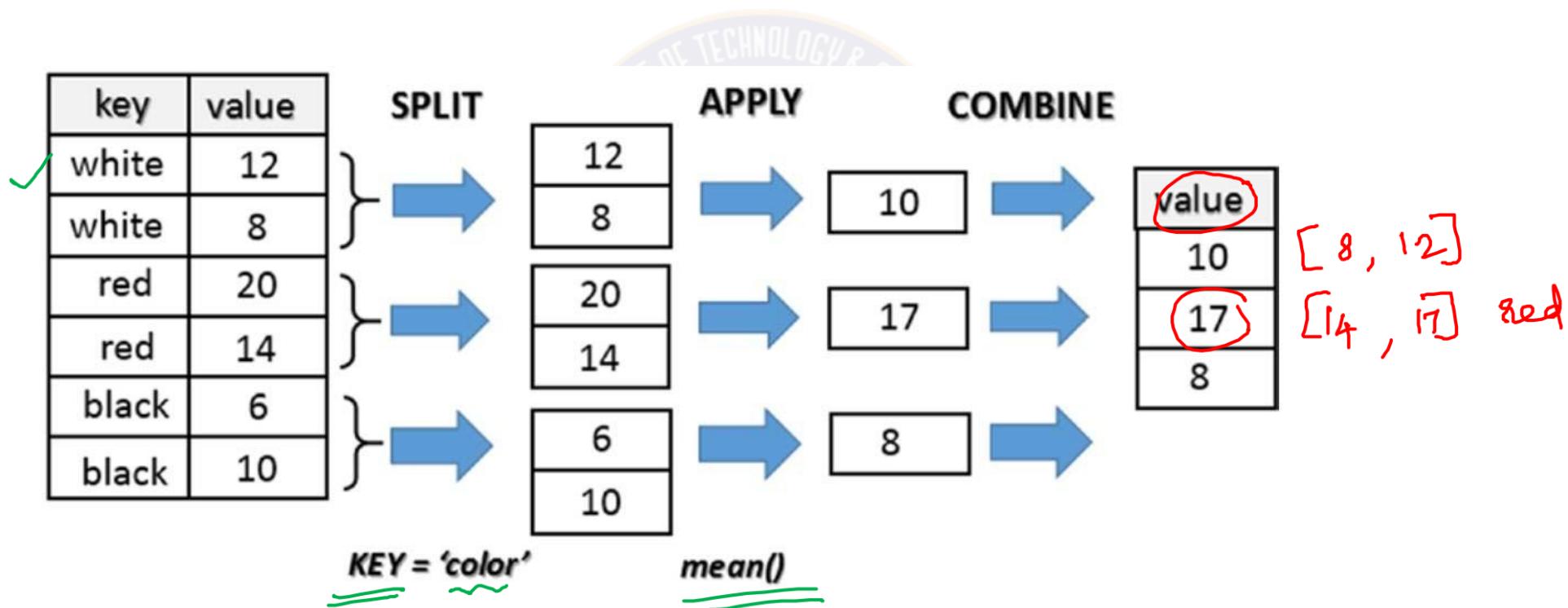
- Combining two or more attributes (or objects) into a single attribute (or object)
- Purpose
 - Data reduction -reduce the number of attributes or objects
 - Change of scale
 - Cities aggregated into regions, states, countries, etc. → *store things*
 - Days aggregated into weeks, months, or years
 - More “stable” data - aggregated data tends to have less variability

Table 2.4. Data set containing information about customer purchases.

Transaction ID	Item	Store Location	Date	Price	...
:	:	:	:	:	
101123	Watch	Chicago	09/06/04	\$25.99	...
101123	Battery	Chicago	09/06/04	\$5.99	...
101124	Shoes	Minneapolis	09/06/04	\$75.00	...
:	:	:	:	:	

Data Aggregation

Python Group By Example



Data cleansing

- Removing or correcting records of corrupted or invalid values from raw data
 - NOISY: containing noise, errors, or outliers .
 - e.g., Salary="−10" (an error)
 - INCONSISTENT: containing discrepancies in codes or names, e.g.,
 - Age="42", Birthday="03/07/2010"
 - Was rating "1, 2, 3", now rating "A, B, C"
 - discrepancy between duplicate records
 - INTENTIONAL (e.g., disguised missing data)
 - Jan. 1 as everyone's birthday
- Removing records that are missing a large number of columns
- Duplicate data

Data cleansing



A mistake or a millionaire?

Missing values

Inconsistent duplicate entries



Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	10000K	Yes
6	No	NULL	60K	No
7	Yes	Divorced	220K	NULL
8	No	Single	85K	Yes
9	No	Married	90K	No
	No	Single	90K	No

Data cleansing

Imputing Missing values

	DATE	air_mv	air_mv_zero	air_mv_previous	air_mv_mean	air_expand
1	JAN49	112	112	112	112	112
2	FEB49	118	118	118	118	118
3	MAR49	132	132	132	132	132
4	APR49	129	129	129	129	129
5	MAY49		0	129	284.54385965	128.29783049
6	JUN49	135	135	135	135	135
7	JUL49		0	135	284.54385965	144.73734152
8	AUG49	148	148	148	148	148
9	SEP49	136	136	136	136	136
10	OCT49	119	119	119	119	119
11	NOV49		0	119	284.54385965	116.19900978
12	DEC49	118	118	118	118	118
13	JAN50	115	115	115	115	115
14	FEB50	126	126	126	126	126
15	MAR50	141	141	141	141	141

Data cleansing

Handling outliers (univariate)

- IQR ✓
 - Outliers are usually, a value higher/lower than $1.5 \times \text{IQR}$
- Z-score method (3 sigma)

Data → normally distributed
non-normal



$$Z = \frac{x - \mu}{\sigma} \quad \text{statistical technique}$$

If data points with Z-Score > 3 or < -3 are considered as Outliers

Data cleansing

Handling outliers (univariate) using IQR

❖ Interquartile Range (IQR):

- IQR = Q3 - Q1 (where Q1 is the 25th percentile and Q3 is the 75th percentile)

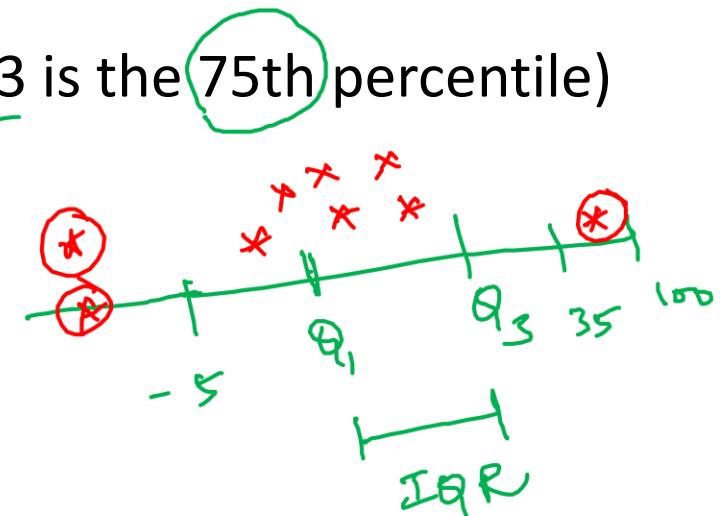
the value below
which 25% of
the data falls

❖ Outlier Detection:

- Lower Bound:** Q1 - 1.5 * IQR
- Upper Bound:** Q3 + 1.5 * IQR

❖ Example:

- If Q1 = 10 and Q3 = 20, then IQR = 10 = $Q_3 - Q_1$
- Lower Bound = $10 - 1.5 * 10 = -5$ ✓
- Upper Bound = $20 + 1.5 * 10 = 35$ ✓
- Data points < -5 or > 35 are outliers



Exercise

➤ Find the outlier in the following data using Inter-Quartile Range.

❖ Data = 10, 12, 11, 15, 11, 14, 13, 17, 12, 22, 14, 11 ✓

$n = 12$

1. Sort : 10, 11, 11, 11, 12, 12, 13, 14, 14, 15, 17 22

2. Median: $(12+13)/2=12.5=Q_2$

3. $Q_1=11$ (25^{th} percentile)

4. $Q_3=14.5$ (75^{th} percentile)

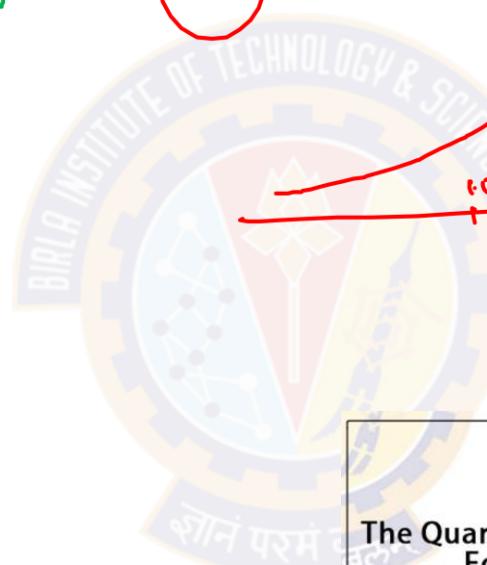
5. $IQR=Q_3-Q_1=3.5$ ✓

6. ~~Min~~ ^{LRL} = $Q_1 - 1.5IQR = 5.75$

7. ~~Max~~ ^{URL} = $Q_3 + 1.5IQR = 19.75$

• Outlier = 22

$$Q_1 = \frac{1}{4} (12+1)^{\text{th}} \text{ term} = \frac{13}{4}^{\text{th}} \text{ term}$$



Quartile Formula

The Quartile Formula = $\frac{1}{4} (n + 1)^{\text{th}}$ term
For Q₁

The Quartile Formula
For Q₃ = $\frac{3}{4} (n + 1)^{\text{th}}$ term

The Quartile Formula
For Q₂ = Q₃ – Q₁ (Equivalent to Median)

Data cleansing

Handling outliers (univariate) using 3 sigma ✓

➤ **3 Sigma Rule:** Based on the properties of a normal distribution

- **Mean (μ) and Standard Deviation (σ)

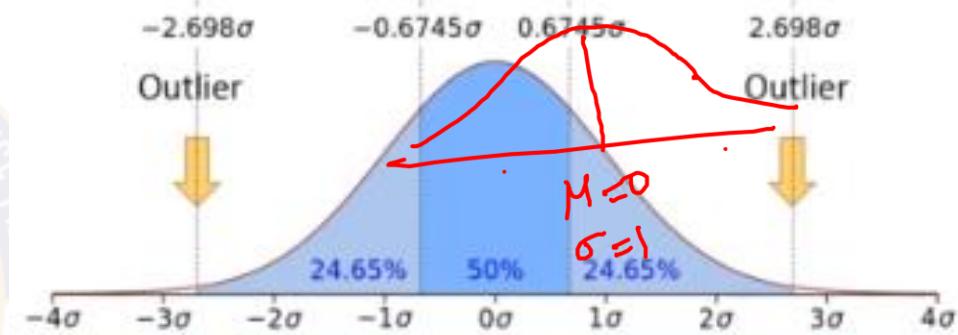
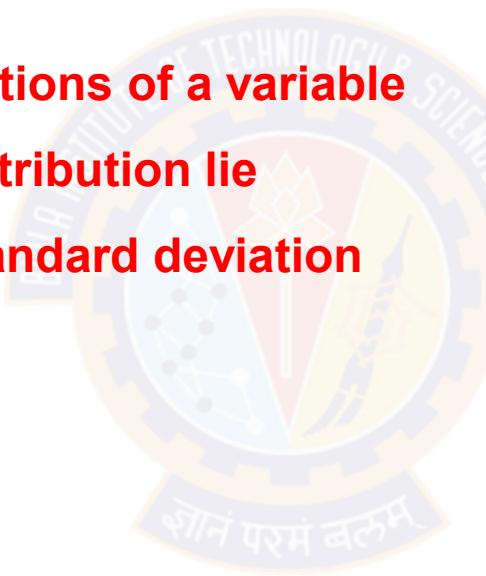
➤ **99% of the observations of a variable following a normal distribution lie within mean $\pm 3 \times$ standard deviation**

➤ **Outlier Detection:**

- Lower Bound: $\mu - 3\sigma$
- Upper Bound: $\mu + 3\sigma$ ✓

➤ **Example Calculation:**

- If $\mu = 50$ and $\sigma = 5$, then:
 - Lower Bound = $50 - 3 * 5 = 35$ ✓
 - Upper Bound = $50 + 3 * 5 = 65$ ✓
- Data points < 35 or > 65 are outliers



68% data lies within $\pm 1\sigma$
95% " " " $\pm 2\sigma$
99.7% " " " $\pm 3\sigma$

'6σ' standard

Instances selection and partitioning

training, evaluation (validation), test sets

Challenge 2 : Non-representative Training Data .

Idea : Training Data be representative of the new cases we want to generalize

- Small sample size leads to sampling noise. Increase sampling size.
- If sampling process is flawed, even large sample size can lead to sampling bias

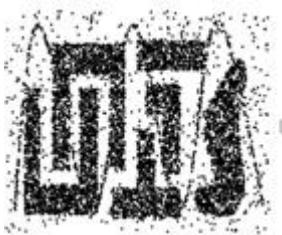
pattern / group

under represented
over represented

The key principle for effective sampling is the following:

- Using a sample will work almost as well as using the entire data set, if the sample is representative
- A sample is representative if it has approximately the same properties (of interest) as the original set of data

X not good representative



8000 points



2000 Points



500 Points

Example:

Problem: You're building a model to predict loan defaults, but your training data is from just one state or income group.

- **Small Sample:** May lead to noise → unreliable model
- **Biased Sample:** Even if large, won't generalize to other states or income groups

Good Practice:

- Ensure geographic, demographic, and economic diversity
- Align training data distribution with target population

Instances selection and partitioning

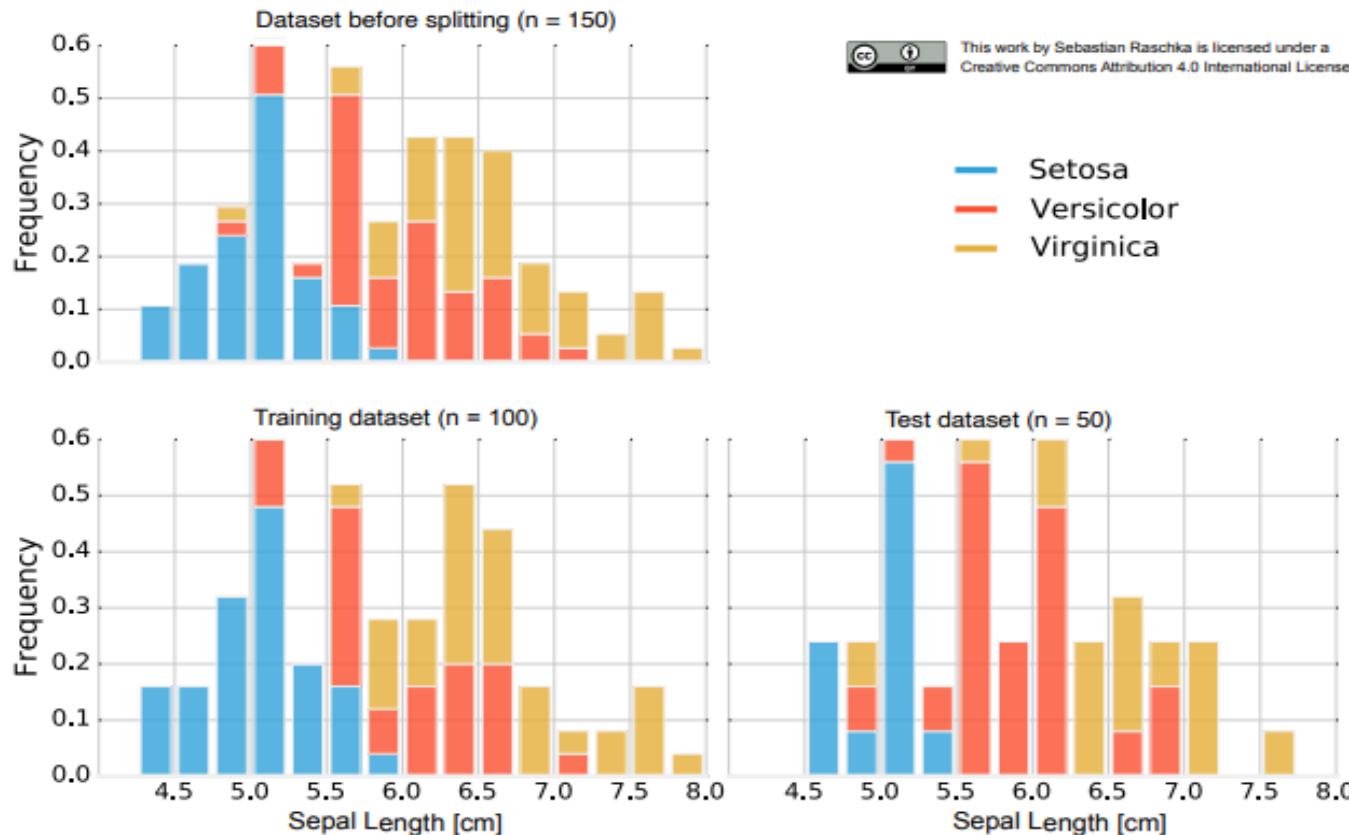
training, evaluation (validation), test sets

- Sampling is the main technique employed for data reduction.
 - It is often used for both the preliminary investigation of the data and the final data analysis.
- Statisticians often sample because obtaining the entire set of data of interest is too expensive or time consuming.
- Sampling is typically used in data mining because processing the entire set of data of interest is too expensive or time consuming.

Instances selection and partitioning

Sampling

Issues with Subsampling (Independence Violation)



IRIS Dataset of Flowers

50 Setosa,
50 Versicolor,
50 Virginica

- Random subsampling can assign 2/3 (100) to training set and 1/3 (50) to the test set
- Training set → 38 x Setosa, 28 x Versicolor, 34 x Virginica
- Test set → 12 x Setosa, 22 x Versicolor, 16 x Virginica

Instances selection and partitioning

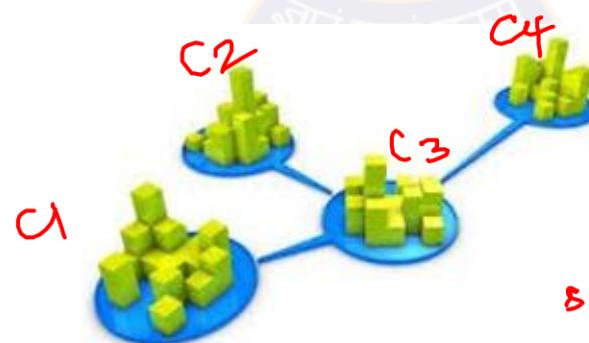
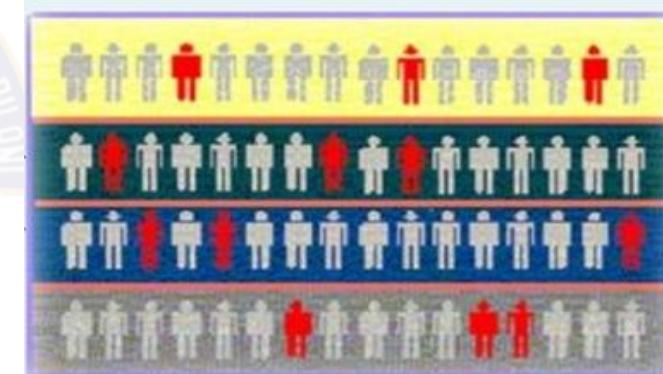
Sampling - Frequently Used

Simple Random Type



heterogeneous

Stratified Sampling Type



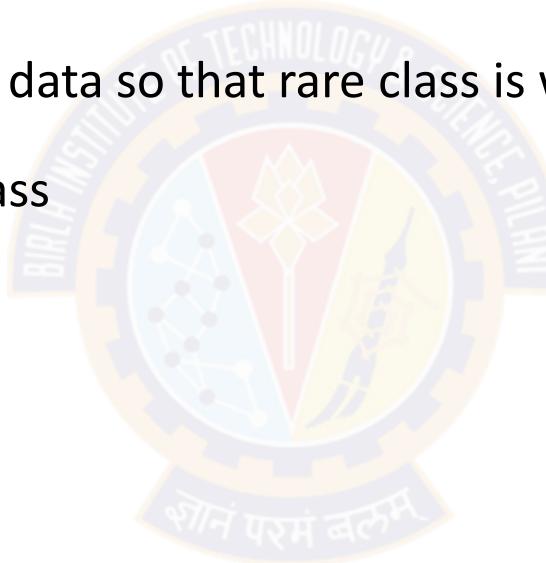
Clustered Sampling Type

$$\text{sample} = \{C_1, C_2, C_3, C_4\}$$

Instances selection and partitioning

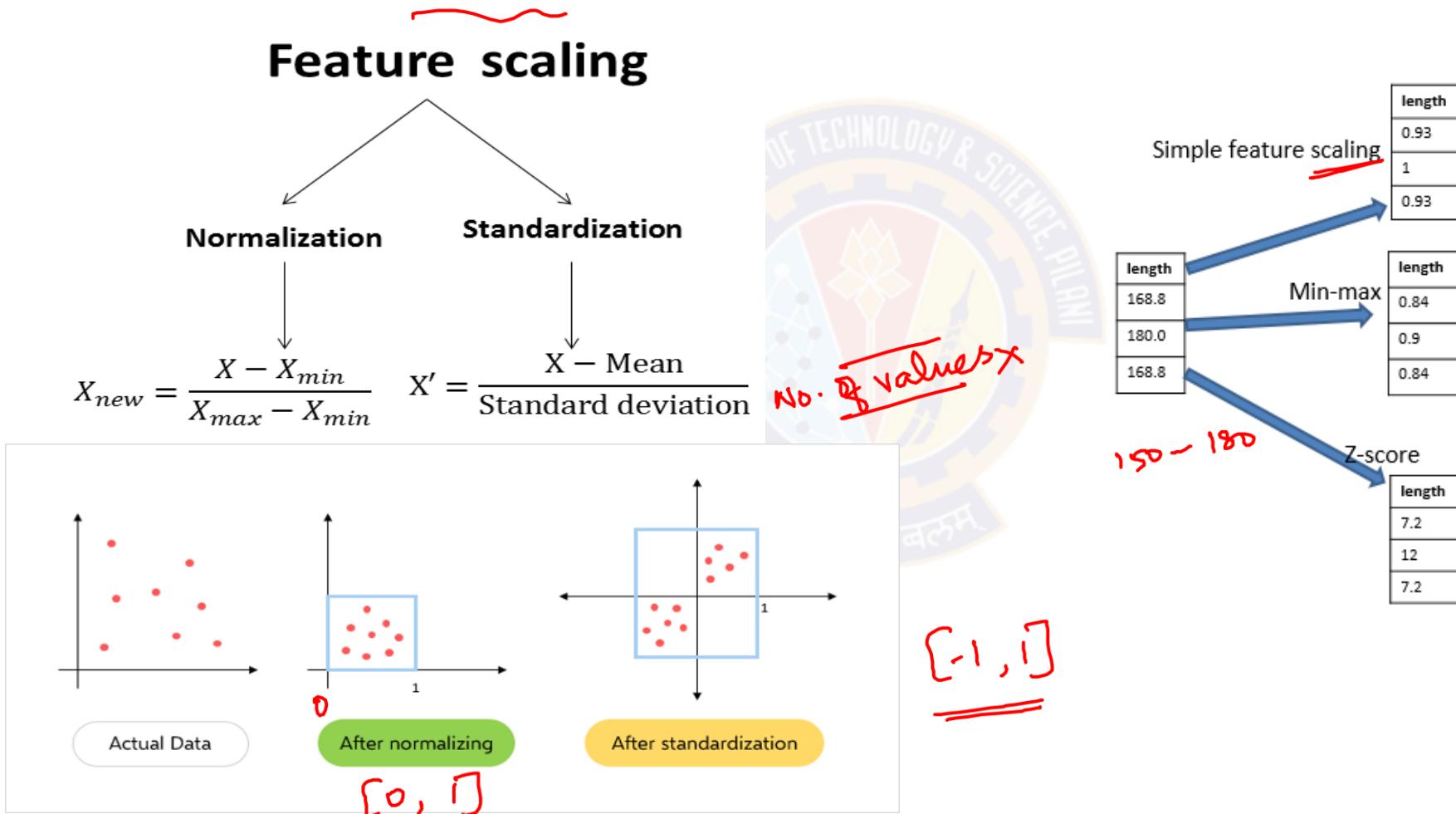
Sampling - Imbalanced Training Set

- **Scenario :** Building Classifiers with Imbalanced Training Set
- Modify the distribution of training data so that rare class is well-represented in training set
 - Under sample the majority class
 - Over sample the rare class



Feature tuning Feature Scaling

To map the continuous values from one range to target range to easily compare and fit in apt distribution to enable statistical processing



Note: Scaling the target values is generally not required

Feature tuning

Feature Scaling - Normalization Vs Standardization

- Normalization
 - when approximate upper and lower bounds on data is known
 - When data is approximately uniformly distributed across that range. E.g age. Not to be used on skewed attribute e.g. income
 - when the algorithms do not make assumptions about the data distribution e.g. (KNN, NN)
 - scales in a range of [0,1] or [-1,1]
- Standardization
 - used when algorithms make assumptions about the data distribution (Gaussian distribution)
 - not bounded by range
 - less affected by outliers

Z-score
ज्ञान परमं बलम्

Note:

- Fit the scalers to the training data only
- Use them to transform the training set and the test set

Feature tuning

Feature Scaling - Normalization Vs Standardization

- Min-max normalization: to $[new_min_A, new_max_A]$

$$v' = \frac{v - min_A}{max_A - min_A} (new_max_A - new_min_A) + new_min_A$$

- Ex. Let income range \$12,000 to \$98,000 normalized to [0.0, 1.0]. Then \$73,000 is mapped to

$$\frac{73,600 - 12,000}{98,000 - 12,000} (1.0 - 0) + 0 = 0.716$$

- Z-score normalization/Standardization (μ : mean, σ : standard deviation):

$$v' = \frac{v - \mu_A}{\sigma_A}$$

- Ex. Let $\mu = 54,000$, $\sigma = 16,000$. Then $\frac{73,600 - 54,000}{16,000} = 1.225$

- **Normalization by decimal scaling**

$$v' = \frac{v}{10^j} \quad \text{Where } j \text{ is the smallest integer such that } \text{Max}(|v'|) < 1$$

Feature Engineering

- Feature engineering needed for coming up with a good set of features - Irrelevant Features
- Feature extraction [10 attributes]
 - Dimensionality reduction → PCA
- Feature selection
 - more useful features to train on among existing features.
- Feature Construction
 - Combine existing features to produce a more useful one
- Feature Transformation



Eg: loan defaulter's list
Obj: missing payments ✓
?? statement of accts for last 6 months

Case study

Input:

WILP student details enrolled in MTech AIML program.

Analysis:

Predict the GPA of the AIML students in Semester3 to estimate
the no. of students who might enroll in dissertation

Observation:

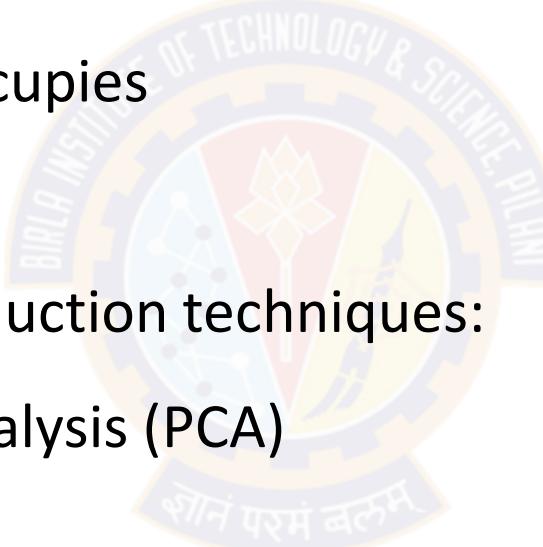
Students with similar educational background tend to perform
same in the exams

AttributesOfInterest
Name
Gender
Age
DataOfBirth
Organisation
JobTitle
NatureOfJob
EntranceScore
EligibilityScore
PreviousDegree
WILPBatch
Section
ISM
MFML
ACI
ML
NLP
.....

Feature Engineering - Extraction

Curse of Dimensionality

- Reducing the number of features by creating lower-dimension
- When dimensionality increases, data becomes increasingly sparse in the space that it occupies
- Solution : Dimensionality Reduction techniques:
e.g Principal Components Analysis (PCA)



AttributesOfInterest
Name
Gender
Age
DataOfBirth
Organisation
JobTitle
NatureOfJob
EntranceScore
EligibilityScore
PreviousDegree
WILPBatch
Section
ISM
MFML
ACI
ML
NLP
.....

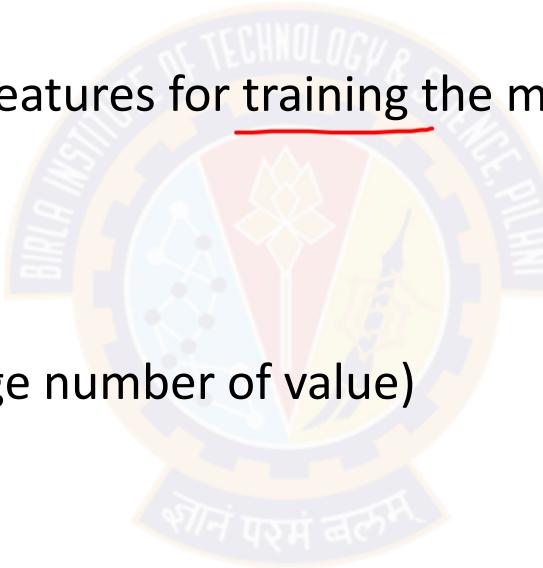
Feature Engineering - Selection

- Selecting a subset of the input features for training the model

Handle Redundant features

Remove Irrelevant feature

dropping features (missing a large number of value)



```
dataframe= dataframe.drop(['COLNAME-1','COLNAME-2'],axis=1)
```

AttributesOfInterest
Name
Gender
Age
DateOfBirth
Organisation
JobTitle
NatureOfJob
EntranceScore
EligibilityScore
PreviousDegree ✓
WILPBatch
Section
ISM
MFML
ACI
ML ✓
NLP ✓
.....

Feature Engineering - construction

➤ Creating new features by using techniques

- Polynomial expansion (by using univariate mathematical functions)
- Feature crossing (to capture feature interactions)
- Features can also be constructed by using business logic from the domain of the ML use case.

$x \rightarrow \text{sq.ft of house}$
Maintenance fee $\rightarrow \text{Rs } 3/\text{sq.ft}$

$(3x) = \text{Balance Amt}$
 $\text{EMI} \times \text{No. months} \rightarrow \text{paid}$

Marking schemes \rightarrow Answer Key
1 $\rightarrow 2M$
2 $\rightarrow 2M$

Feature Engineering - Transform

AttributesOfInterest
PreviousDegree ✓
SEM-1-Total .
SEM-2-Total .
SEM-3-Total .
CGPA .
<u>isEligibleForDissertation</u>

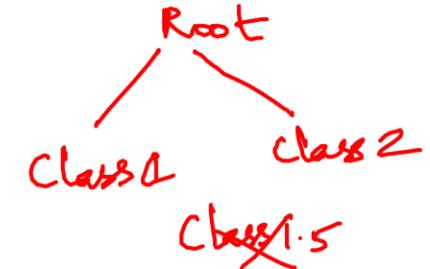
AttributesOfInterest
PreviousDegree
SEM-1-GPA ✓
SEM-2-GPA ✓
SEM-3-GPA ✓
CGPA
<u>isEligibleForDissertation</u>

AttributesOfInterest
PreviousDegree
S1-isComplete✓
S2-isComplete
S3-isComplete
CGPA
<u>isEligibleForDissertation</u>

Feature Engineering - Transform

Encoding Numerical Features

- **Discretization** : Convert continuous attribute into a discrete attribute
 - Naive Bayes, decision trees and their ensembles including Random forest, Minimum distance classifiers or KNN prefer discrete features.
 - Also known as binning' or 'bucketing'
 - To handle outliers.
 - To improve the value spread i.e., spread of data
- Discretization involves converting the raw values of a numeric attribute (e.g., age) into
 - interval labels (e.g., 0–10, 11–20, etc.) OR
 - conceptual labels (e.g., youth, adult, senior)

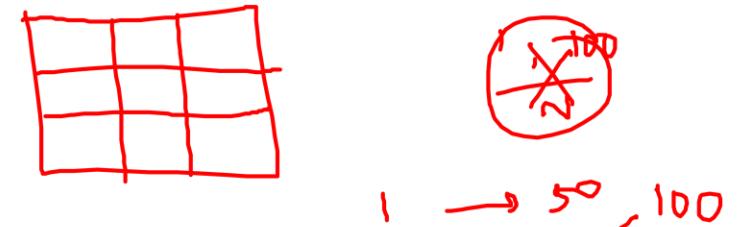


Feature Engineering - Transform

Encoding Numerical Features

Simple Discretization: Binning

- Equal-width (distance) partitioning
 - Divides the range into N intervals of equal size: uniform grid
 - if A and B are the lowest and highest values of the attribute, the width of intervals will be: $W = (B - A)/N$.
 - The most straightforward, but outliers may dominate presentation
 - Skewed data is not handled well
- Equal-depth (frequency) partitioning
 - Divides the range into N intervals, each containing approximately same number of samples
 - Good data scaling



10, 10, 10

Feature Engineering - Transform

Encoding Categorical Features

- Binarization maps a categorical attribute into one or more binary variables - One Hot/Dummy Encoding

The diagram illustrates the process of binarizing a categorical feature. On the left, there is a table with two columns: 'Car' and 'Fuel'. The 'Fuel' column contains categorical values: Gas, Diesel, Gas, Gas, and gas. An arrow points from this table to a second table on the right. The second table has five columns: 'Car', 'Fuel', an empty column, 'Gas', and 'Diesel'. The 'Fuel' column from the first table is mapped to the second table's 'Fuel' column. The other four columns represent binary encoding for 'Gas' and 'Diesel'. For row A, 'Gas' is 1 and 'Diesel' is 0. For row B, 'Gas' is 0 and 'Diesel' is 1. For rows C and D, 'Gas' is 1 and 'Diesel' is 0.

Car	Fuel		Gas	Diesel
A	Gas	1	0
B	Diesel	0	1
C	Gas	1	0
D	gas	1	0

- Categorical features to a numeric representation - Label Encoding

The diagram illustrates the process of label encoding a categorical feature. On the left, there is a table with a single column labeled 'Fuel'. The values in this column are: Fuel, Gas, Diesel, Gas, and gas. An arrow points from this table to a second table on the right. The second table also has a single column labeled 'Fuel'. The values in this column are: 1, 2, 1, 2, and 3 respectively. This shows that each categorical value is assigned a unique integer label.

Fuel
1
2
1
2
3

Problem Type 3

Pre-Processing

- ~~To~~
- A marketing domain has launched their APP products tailored for different categories of student's population in a city to get feedback. The focus group has given following feedback. How do you propose to ready the data for analysis?

Features liked	Features to improve	Do you have similar app?	How much do you pay for existing app per month?	Rate game (1-10)	Rate social media connect (1-5)	Shopping facility (1-10)	How would you recommend this to a friend	Education
Graphics	Usability	No	\$35	8	5	4	High	School
Interactive	More features		Rs.500	8	3	6	Low	College
Graphics		Yes	Rs.250	7	5	6	Medium	School
Cheap	Creativity	Yes	\$20	7	5	8	High	College

Question : Identify the basic preprocessing & data cleaning required for this case



Challenges of Machine Learning

- Training Data
 - Insufficient
 - Non representative
- Model Selection
 - Overfitting
 - Underfitting
- Validation and Testing



IMPORTANT NOTE TO THE STUDENTS:

More on this slide will be discussed by faculty only in later modules 4,5,6.... On appropriate sections



Few Terminologies

(To interpret the jargons in the prescribed text book)

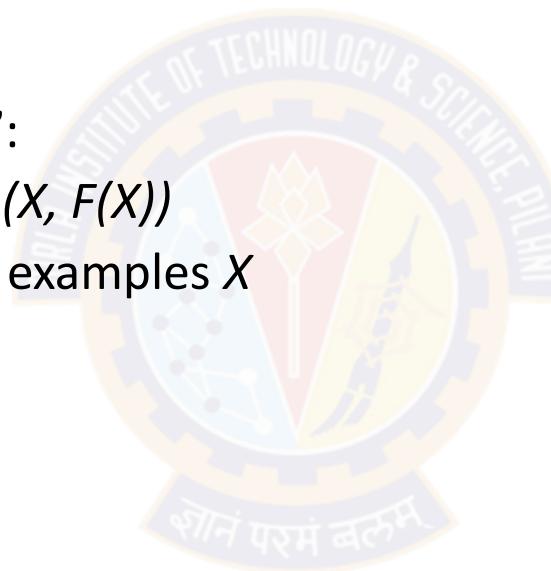
Terminologies

- **Training example.** An example of the form $\langle \mathbf{x}, f(\mathbf{x}) \rangle$.
- **Target function (target concept).** The true function f .
- **Hypothesis.** A proposed function h believed to be similar to f .
- **Concept.** A boolean function. Examples for which $f(\mathbf{x}) = 1$ are called **positive examples** or **positive instances** of the concept. Examples for which $f(\mathbf{x}) = 0$ are called **negative examples** or **negative instances**.
- **Classifier.** A discrete-valued function. The possible values $f(\mathbf{x}) \in \{1, \dots, K\}$ are called the **classes** or **class labels**.
- **Hypothesis Space.** The space of all hypotheses that can, in principle, be output by a learning algorithm.
- **Version Space.** The space of all hypotheses in the hypothesis space that have not yet been ruled out by a training example.

Amount taken	Period	Credit Score	Defaulter
40 lakhs	5 years	1000	No
10 Lakhs	5 months	550	YES
80 Lakhs	3 years	950	No
20 Lakhs	4 years	1500	No

Inductive Learning Hypothesis

- Any hypothesis found to approximate the target function well over a sufficiently large set of training examples will also approximate the target function well over other unobserved examples
- Inductive learning or “Prediction”:
 - Given examples of a function $(X, F(X))$
 - Predict function $F(X)$ for new examples X
- Classification
 $F(X) = \text{Discrete}$
- Regression
 $F(X) = \text{Continuous}$
- Probability estimation
 $F(X) = \text{Probability}(X):$



Inductive Learning Hypothesis

- Target Concept

- Discrete : $f(x) \in \{\text{Yes, No, Maybe}\}$
- Continuous : $f(x) \in [20-100]$
- Probability Estimation : $f(x) \in [0-1]$

Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport?
Sunny	Warm	Normal	Strong	Warm	Same	Yes
Sunny	Warm	High	Strong	Warm	Same	Yes
Rainy	Cold	High	Strong	Warm	Change	No
Sunny	Warm	High	Strong	Cool	Change	Yes

Inductive Learning Hypothesis

- Target Concept

- Discrete : $f(x) \in \{\text{Yes, No, Maybe}\}$
- Continuous : $f(x) \in [20-100]$
- Probability Estimation : $f(x) \in [0-1]$

Sky	AirTemp	Altitude	Wind	Water	Forecast	Humidity
Sunny	Warm	Normal	Strong	Warm	Same	60
Sunny	Warm	High	Strong	Warm	Same	75
Rainy	Cold	High	Strong	Warm	Change	70
Sunny	Warm	High	Strong	Cool	Change	45

Inductive Learning Hypothesis

- **Target Concept**

- **Discrete** : $f(x) \in \{\text{Yes, No, Maybe}\}$
- **Continuous** : $f(x) \in [20-100]$
- **Probability Estimation** : $f(x) \in [0-1]$

Sky	AirTemp	Humidity	Wind	Water	Forecast	$P(\text{EnjoySport} = \text{Yes})$
Sunny	Warm	Normal	Strong	Warm	Same	0.95
Sunny	Warm	High	Strong	Warm	Same	0.7
Rainy	Cold	High	Strong	Warm	Change	0.5
Sunny	Warm	High	Strong	Cool	Change	0.6

Hypothesis

Example	Sky	AirTemp	Humidity	Wind	Water	Forecast	EnjoySport
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes

- One possible hypothesis $(?, \text{Cold}, \text{High}, ?, ?, ?)$
- The most general hypothesis—that every day is a positive example— $(?, ?, ?, ?, ?, ?)$
- The most specific possible hypothesis—that no day is a positive example $(\phi, \phi, \phi, \phi, \phi, \phi)$



Thank you