

# ACDS Lecture Series

Lecture - 3

CSIR

## Machine Learning

**G. N. Sastry and Team**

ADVANCED COMPUTATION AND DATA SCIENCES (ACDS) DIVISION

CSIR-North East Institute of Science and Technology, Jorhat, Assam, India



Acquiring new or existing knowledge. Observe a phenomena and recognize a pattern.



Unconscious  
Incompetence



Conscious  
Incompetence



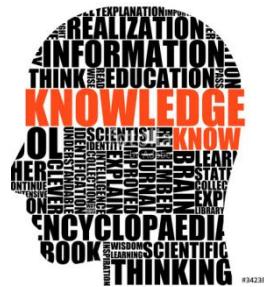
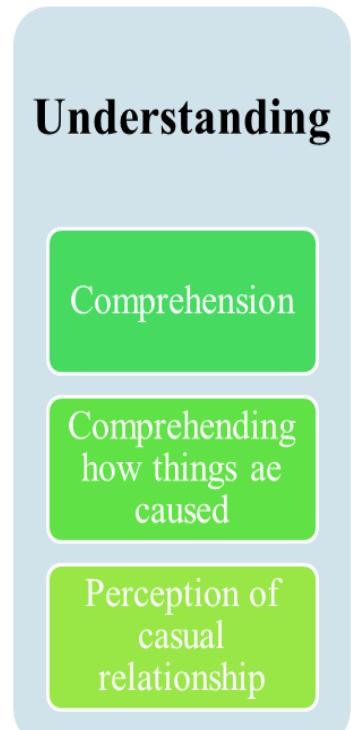
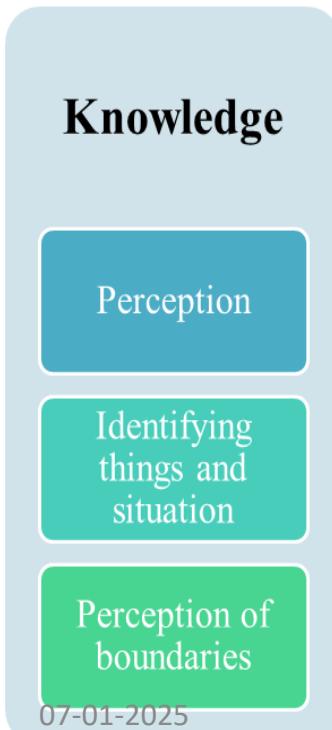
Conscious  
Competence



Unconscious  
Competence

## UNDERSTANDING

Ability to perform flexibility with what is taught and applying it beyond the knowledge acquired

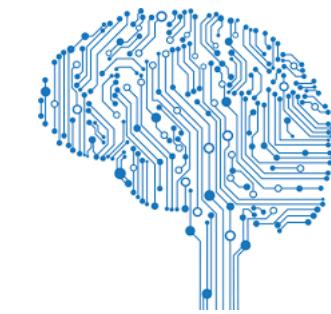


Understanding gained from learning and experiences

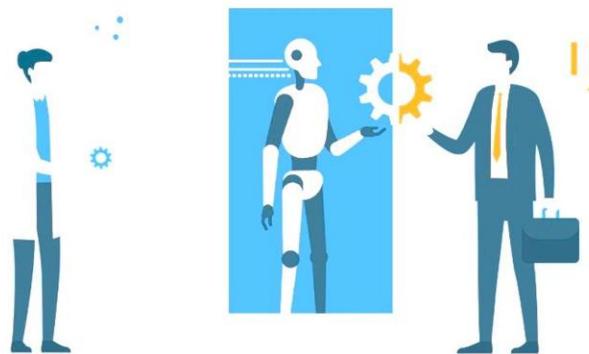
- What is going on?
- What do you notice?
- Where does this fit?
- What are you curious about it?
- Why do you think so?

# Can a Machine Mimic the Process of Human Learning?

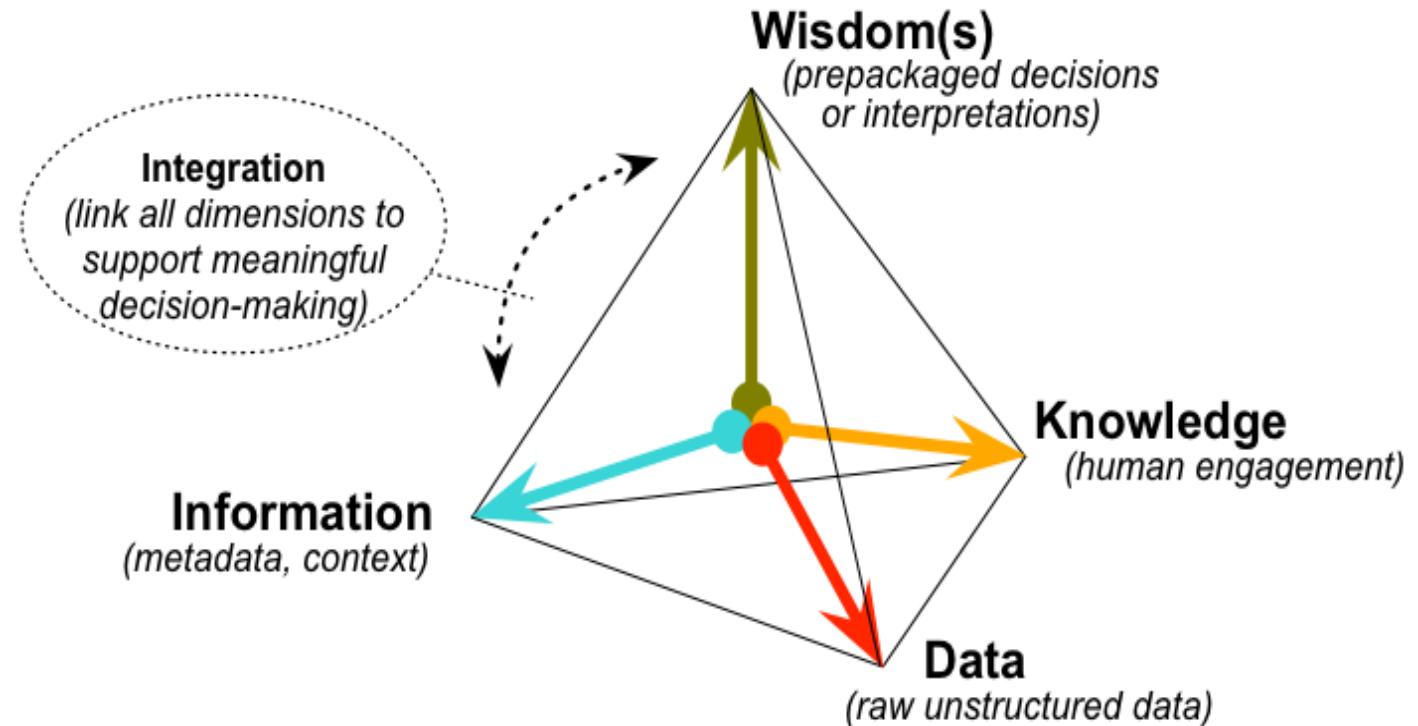
ACDS, CSIR-NEIST



## MACHINE LEARNING



Whether it's Human Learning or Machine Learning — both involve observations and executing about a thing or a process or phenomenon.

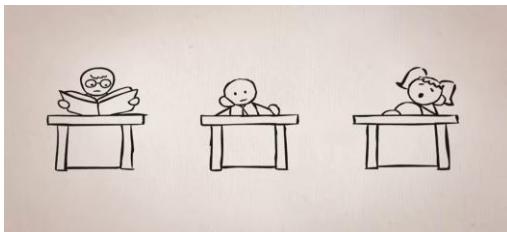
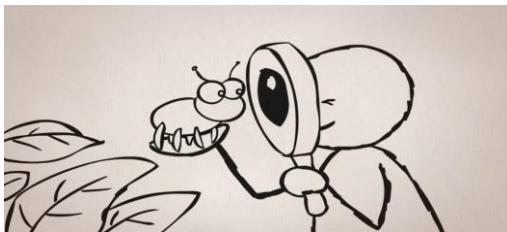


Nevertheless, it is humans who build, operate, and understand the tools that turn data into information, and information into knowledge, and knowledge into wisdom.

# Depiction of Learning for Human and Machine

ACDSD, CSIR-NEIST

Human  
Learning



Cognition/Model

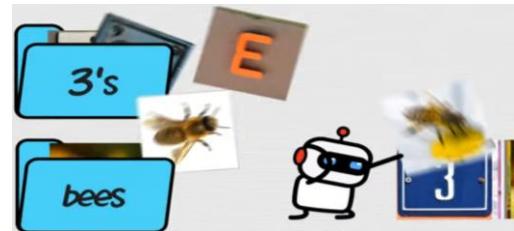
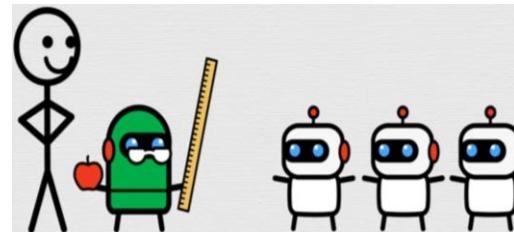
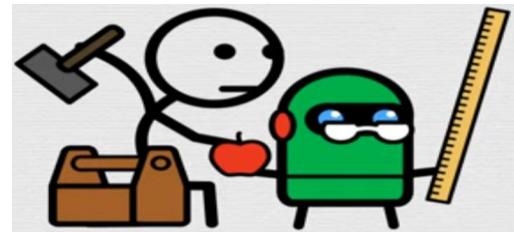
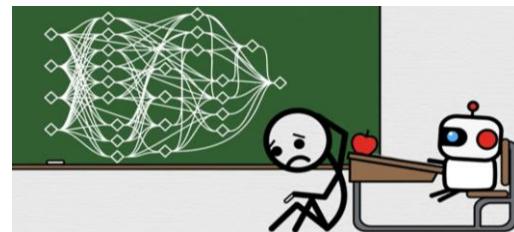
Data Feeding

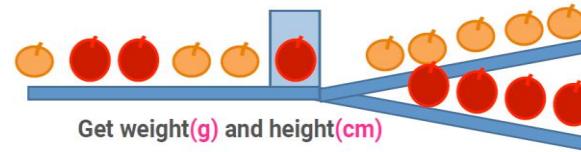
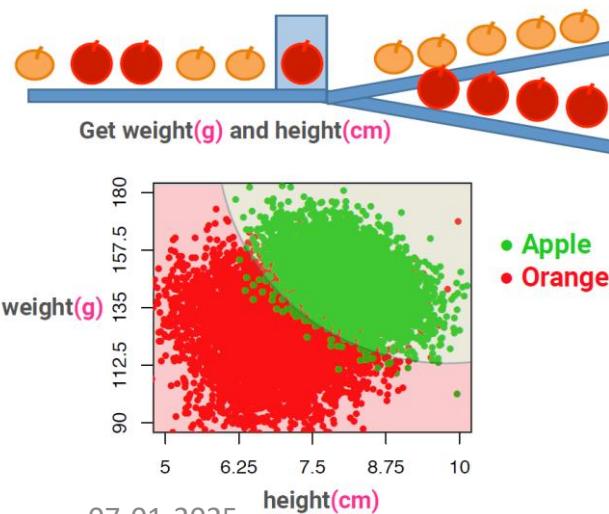
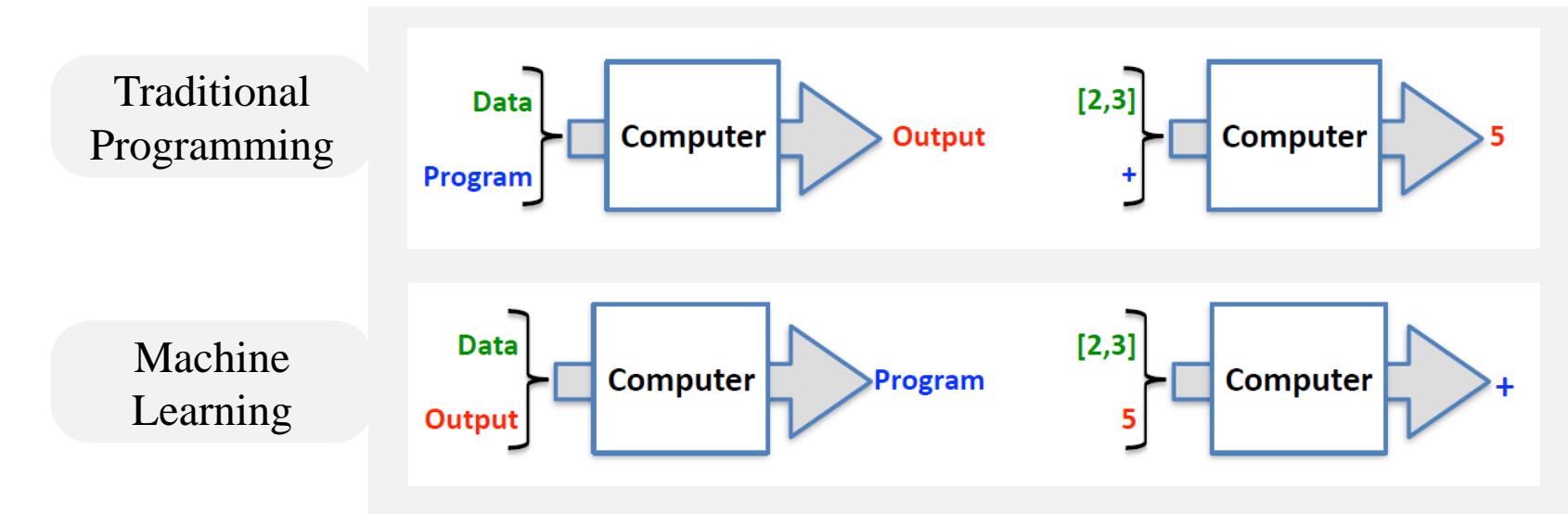
Improvisations

Training

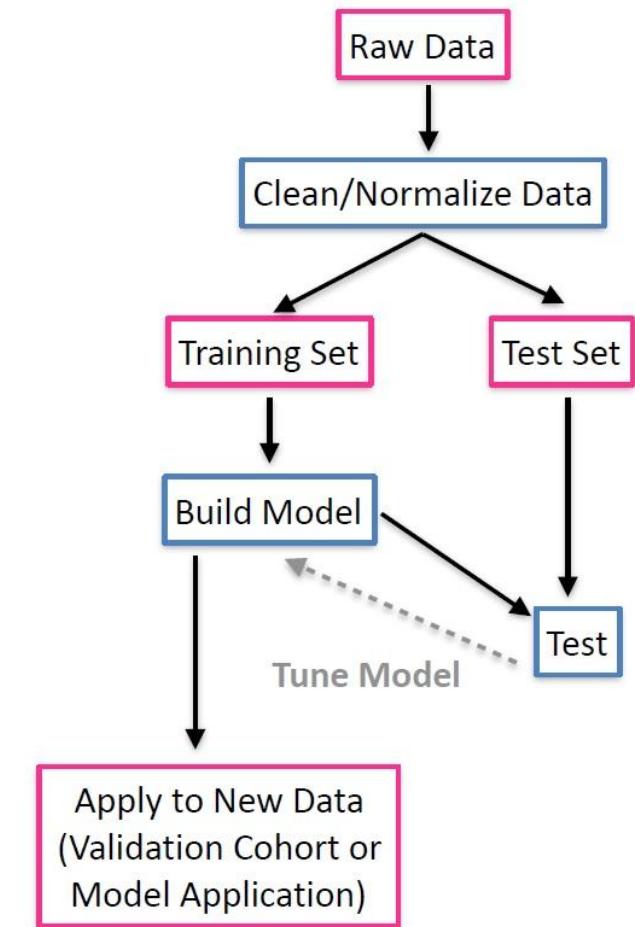
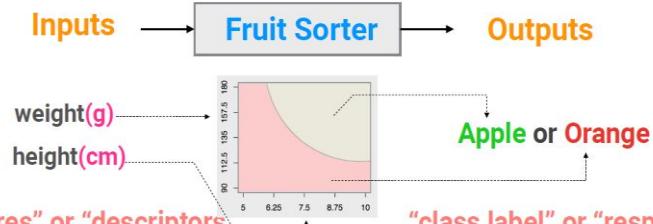
Expertise

Machine  
Learning





Now we got a computer program for this classification problem.



## Unsupervised Learning

### Clustering

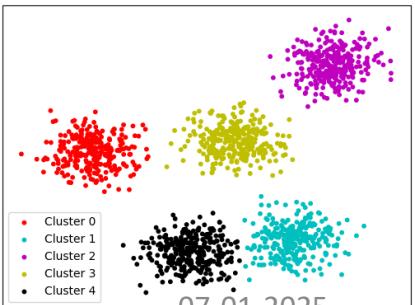
$$r_{nk} \in \{0, 1\}$$

$$r_{nk} = \begin{cases} 1 & \text{if datapoint } x_n \text{ assigned to cluster } k \\ 0 & \text{Otherwise} \end{cases}$$

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

$$r_{nk} = \begin{cases} 1 & \text{if } k = \operatorname{argmin}_j \|x_n - \mu_k\|^2 \\ 0 & \text{Otherwise} \end{cases}$$

$$\frac{\partial}{\partial J} = 2 \sum_{n=1}^N r_{nk} (x_n - \mu_k) = 0 \quad \mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$$



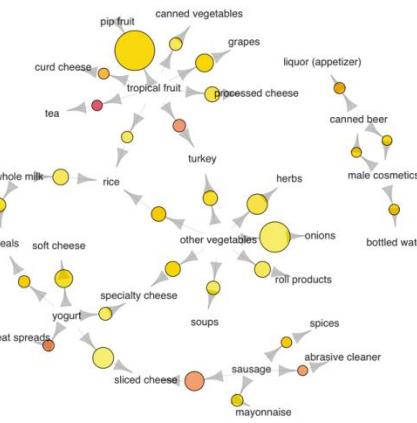
### Association Rules

Rule:  $x \Rightarrow y$

$$\text{Support } (S) = \frac{\text{freq}(x, y)}{N}$$

$$\text{Confidence } (C) = \frac{\text{freq}(x, y)}{\text{freq}(x)}$$

$$\text{Lift } (L) = \frac{\text{Support}}{\text{Support}(x) \times \text{Support}(y)}$$



## Supervised Learning

### Classification

#### [Linear classification]

- Logistic / Softmax regression
- Linear discriminant analysis
- Naive Bayes classifiers
- Perceptron
- Linear Support Vector Machines (SVM)

#### [Nonlinear classification]

- k-nearest neighbor classifiers
- Decision trees (Classification trees)
- Polynomial classifiers / Factorization machines
- Tree ensemble classifiers
  - Random Forest classifiers
  - Extra Trees classifiers
  - Gradient Boosted Decision Trees (GBDT)
- Kernel method classifiers
  - Support Vector Machines (SVM)
  - Gaussian process classifiers
  - Neural network (Deep learning) classifiers
    - Multi-layer perceptrons (MLP)
    - Convolutional networks (CNN)
      - VGG (OxfordNet)
      - Inception (GoogLeNet)
      - ResNet / ResNeXt
      - DenseNet
    - Recurrent networks (RNN)

### Regression

#### [Linear regression]

- Least squares regression
- Principal component regression
- Partial Least Squares (PLS) regression
- Penalized linear regression
  - LASSO regression (L1-penalized)
  - Ridge regression (L2-penalized)
  - ElasticNet regression (L1 & L2-penalized)

#### [Nonlinear regression]

- k-nearest neighbor regressors
- Decision trees (Regression trees, Model trees)
- Polynomial regressors / Factorization machines
- Tree ensemble regressors
  - Random Forest regressors
  - Extra Trees regressors
  - Gradient Boosted Regression Trees (GBRT)
- Kernel method regressors
  - Support Vector Regression (SVR)
  - Kernel Ridge Regression
- Gaussian process regressors
- Neural network (Deep learning) regressors
  - Multi-layer perceptrons (MLP)
  - Convolutional networks (CNN)
    - VGG (OxfordNet)
    - Inception (GoogLeNet)
    - ResNet / ResNeXt
    - DenseNet
  - Recurrent networks (RNN)

## Unsupervised Learning

### Clustering

- k-means
- Hierarchical clustering
- Gaussian mixtures
- Spectral methods
- DBSCAN

### Decomposition

- Principal component analysis (PCA)
- Independent component analysis (ICA)
- Canonical correlation analysis (CCA)
- Nonnegative matrix factorization (NMF)
- Latent Dirichlet allocation (LDA)

### Manifold learning

- Multidimensional scaling (MDS)
- Self-organizing maps (SOM)
- Isomap
- Locally linear embedding (LLE)
- Spectral embedding (Laplacian eigenmaps)
- t-distributed Stochastic Neighbor Embedding (t-SNE)
- Autoencoders

### Density estimation

## Others

### Semi supervised learning

### Ranking

### Transfer learning

### K-shot learning

### Domain adaptation

### Multitask learning

### Reinforcement learning

### Imitation learning

### Active learning

### Model-based optimization

### Time series/Sequence models

### Probabilistic inference

### (Bayesian, Generative, Graphical)

### Causal inference

### Online/Incremental learning

### Anomaly/Outlier detection

### Ensemble learning

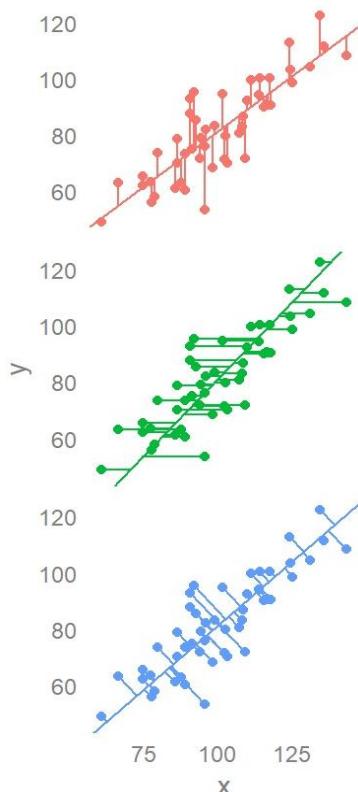
### Relational/Network learning

### Representation learning

### Structured prediction

### Meta Learning

:



## Regression

$$y = mx + c$$

$$m_{new} = m_{current} - k \frac{\delta}{\delta m} E(m, c)$$

$$c_{new} = c_{current} - k \frac{\delta}{\delta c} E(m, c)$$

## Supervised Learning

## Classification

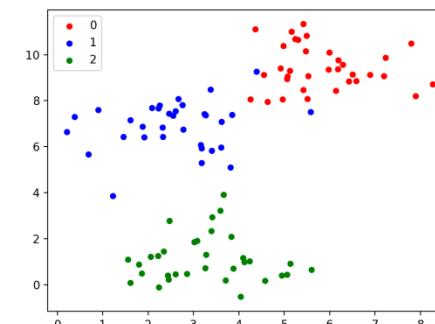
$$\text{class} = \begin{cases} C = 1 & \text{if } P(C = 1|x_1, x_2) > P(C = 0|x_1, x_2) \\ C = 0 & \text{Otherwise} \end{cases}$$

$$P(C|x) = \frac{P(C)p(x|C)}{p(x)}$$

$$P(C_i|x) = \frac{p(x|C_i)P(C_i)}{\sum_{k=1}^K p(x|C_k)P(C_k)}$$

*Bayes' Rule*

*Bayes Classifier*



## Unsupervised Learning

### Clustering

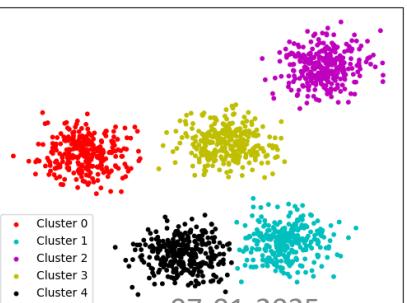
$$r_{nk} \in \{0, 1\}$$

$$r_{nk} = \begin{cases} 1 & \text{if datapoint } x_n \text{ assigned to cluster } k \\ 0 & \text{Otherwise} \end{cases}$$

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

$$r_{nk} = \begin{cases} 1 & \text{if } k = \operatorname{argmin}_j \|x_n - \mu_k\|^2 \\ 0 & \text{Otherwise} \end{cases}$$

$$\frac{\partial}{\partial J} = 2 \sum_{n=1}^N r_{nk} (x_n - \mu_k) = 0 \quad \mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}}$$



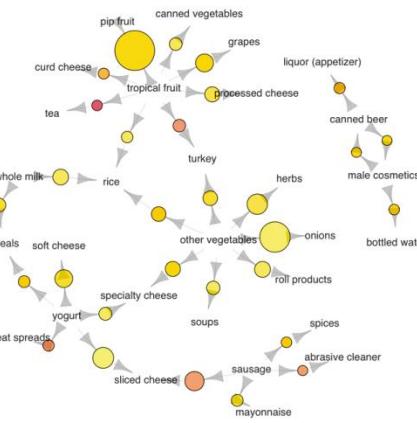
### Association Rules

Rule:  $x \Rightarrow y$

$$\text{Support } (S) = \frac{\text{freq}(x, y)}{N}$$

$$\text{Confidence } (C) = \frac{\text{freq}(x, y)}{\text{freq}(x)}$$

$$\text{Lift } (L) = \frac{\text{Support}}{\text{Support}(x) \times \text{Support}(y)}$$



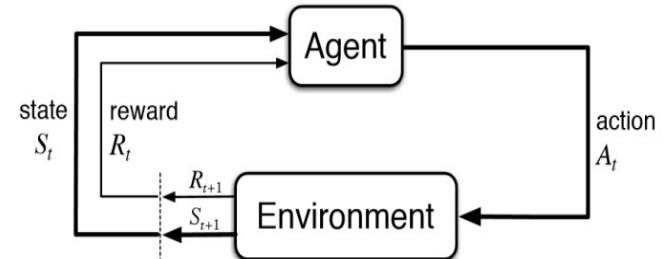
Reward function,  $r(s, a) \rightarrow R$

Markov dynamics,  $P(s_t^*|s_{t-1}, a_{t-1}, s_{t-2}, a_{t-2} \dots) = P(s_t|s_{t-i}, a_{t-i}) = P(S)$

Markov Decision Process (MDP):  $\langle S, A, R, T, \gamma \rangle$

$S = \text{set of states}$ ,  $A = \text{actions}$ ,  $R = r(s, a)$ ,  $T = p(s'|s, a)$ ,  $\gamma = (0, 1)$

Policy,  $\pi: s \rightarrow a$

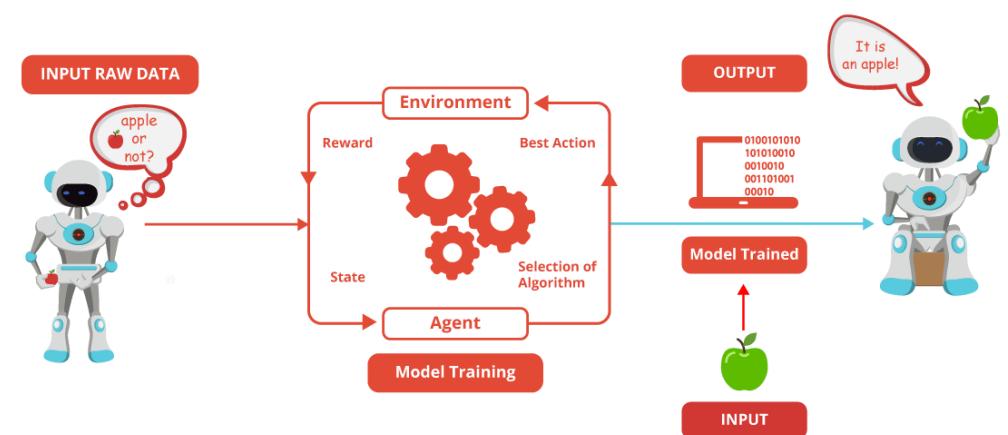


*Model based: directly estimate and use R & T*

*Value based: use value function*

$$V^\pi = E \left[ \sum_{i=1}^T r_{t+1} \right]$$

*Policy based:*  $\arg \max_{\pi \in \Pi} V^\pi$



Human disease = Genetics + Environment + Life Style

## Cancer

Men:- 1:2 risk of developing

1:4 risk of dying

Women:- 1:3 risk of developing

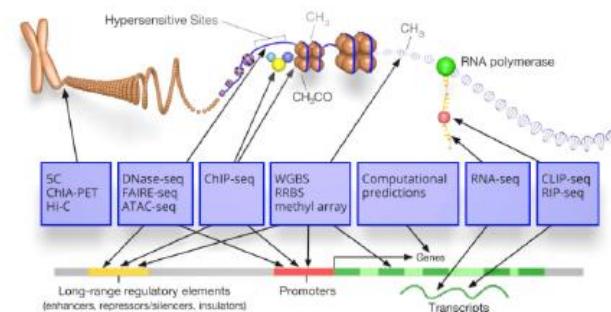
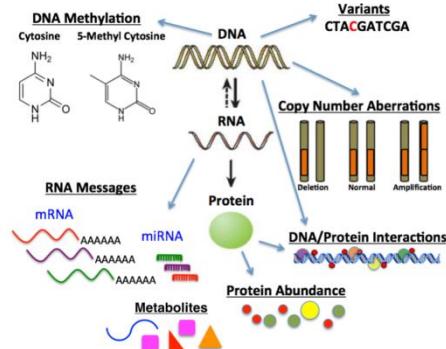
1:5 risk of dying

## Complex Problems

1. Which patient have high risk?
2. Early biomarkers?
3. Who can be long term or short term survivors?
4. What chemotherapeutic might benefit a patient?

*Goal of Genomics is to identify genetics/genomic variation related with disease to improve patient care*

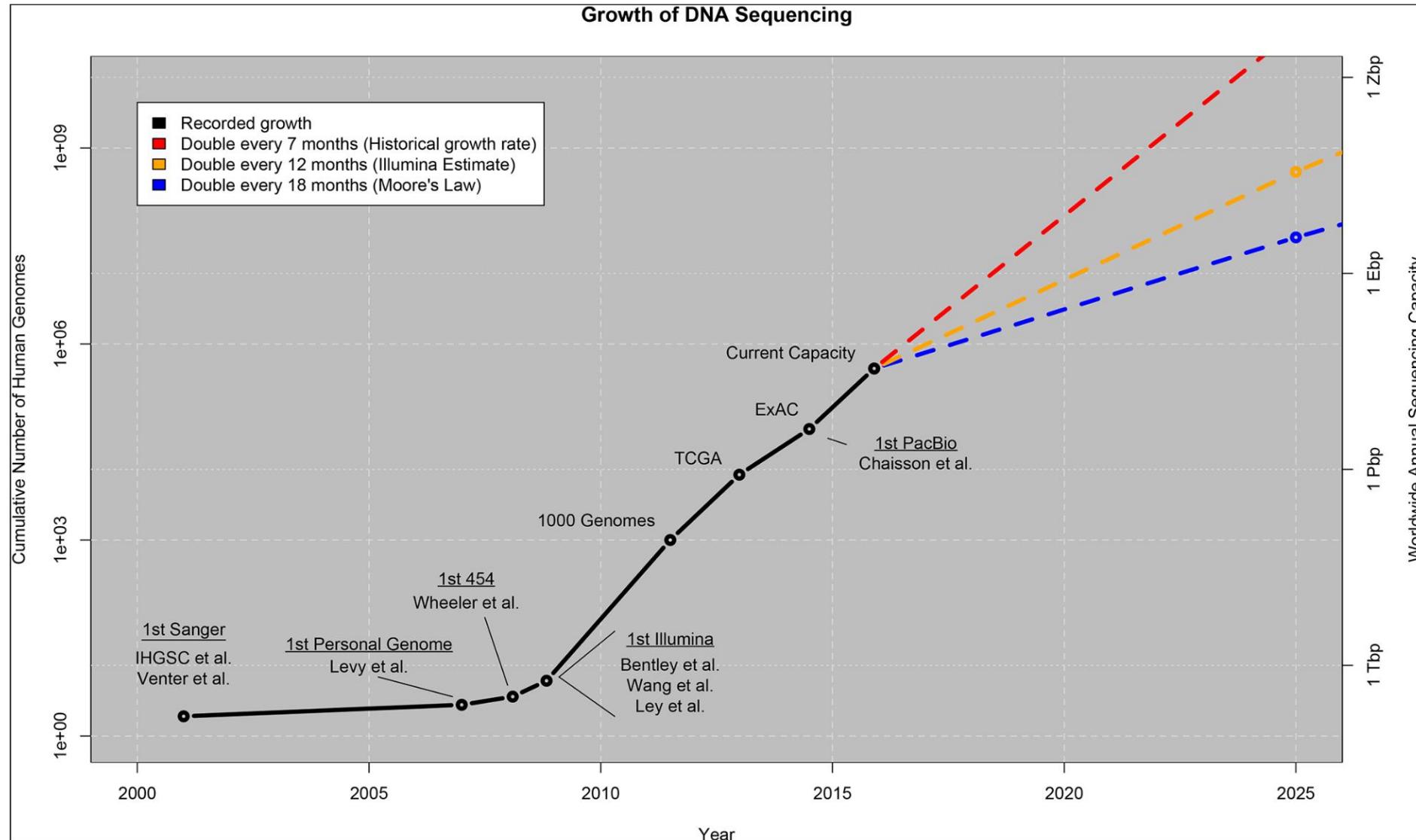
What we have today?



Multidimensional dataset + Cell, Tissue, Disease + Functional annotations = Big Data

Complex Problem + Big Data = Machine Learning

*Goal of machine learning is not to make a perfect guess but an useful guess that can be used in future*

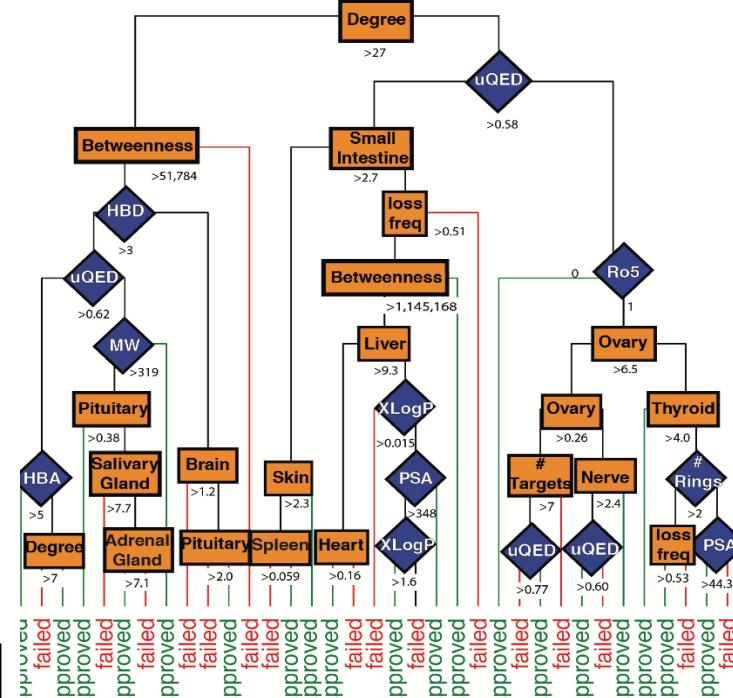


1 zettabyte (ZB) = 1024 EB  
 1 exabyte (EB) = 1024 PB  
 1 petabyte (PB) = 1024 TB  
 1 terabyte (TB) = 1024 GB

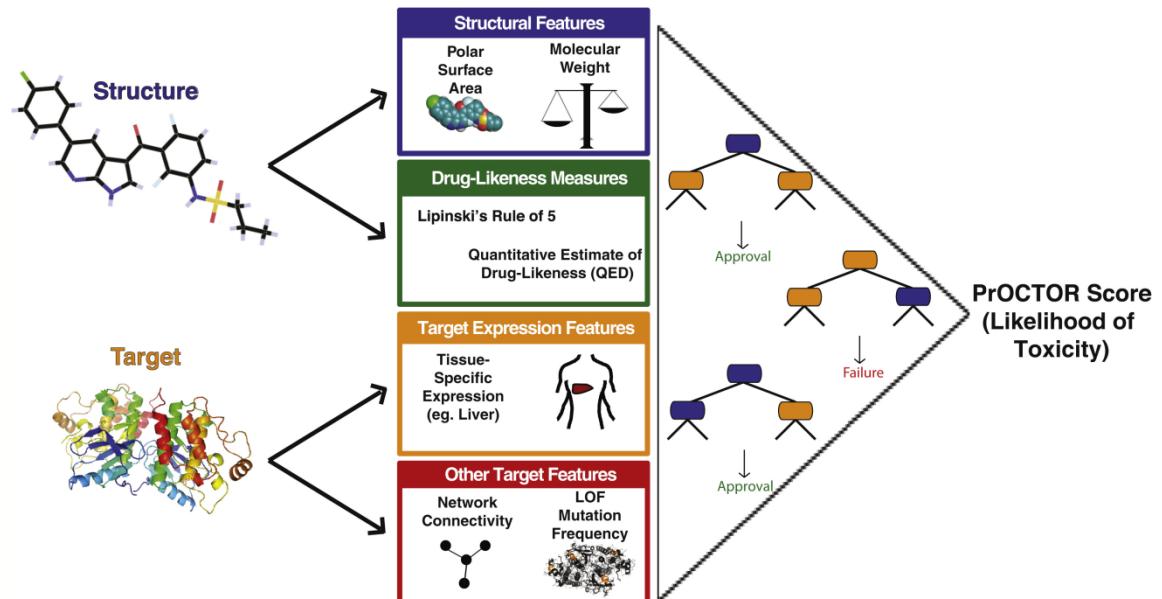
Difficult to identify the unfavourable toxicity properties before conducting clinical trials

## Facts:

1. FDA approved drugs pass Lipinski Ro5 (80.6%) and Ghose Rule (64.9%)
2. FTT drugs pass Lipinski Ro5 (73%) and Ghose Rule (54%)



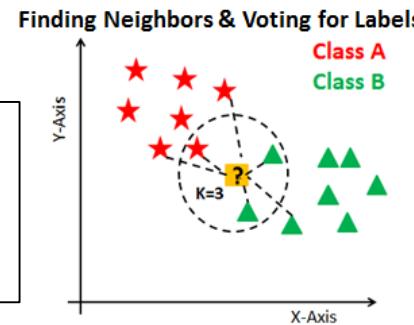
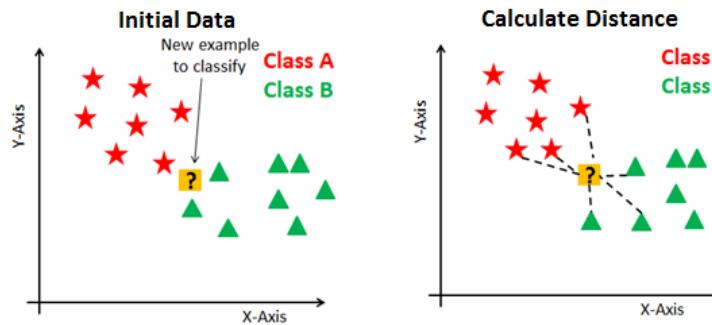
Supervised learning with multiple features can predict the likelihood of success and failures in clinical trials



Predicting the odds of clinical trial outcomes using random forest  
(PrOCTOR) #

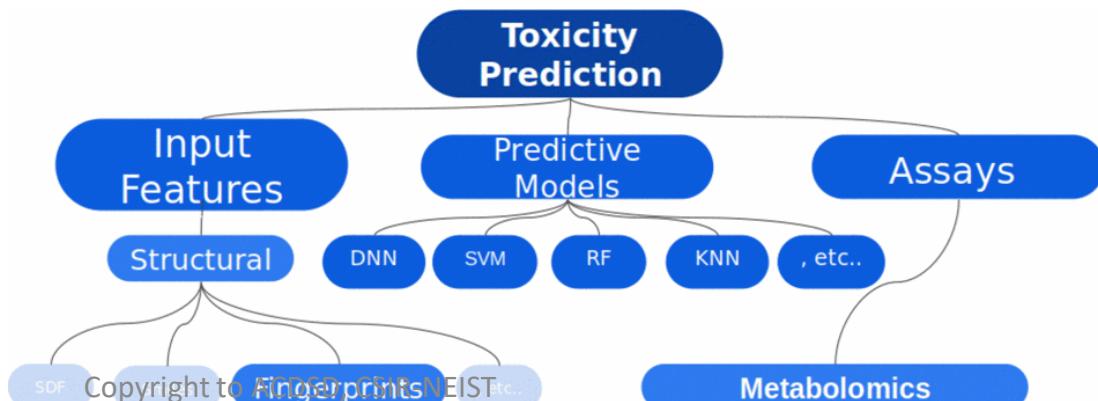
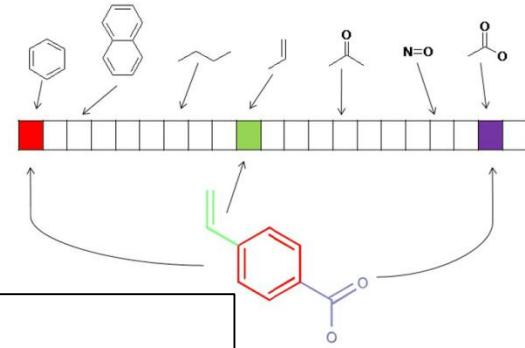
No predefined basis for acceptance or exclusion of a given chemical into a category

**Design a classification model,  $y = f(x)$  where**  
 **$y$ : response variable &  $x$ : fingerprint**



*Class A & class B are acute toxicity based categories*

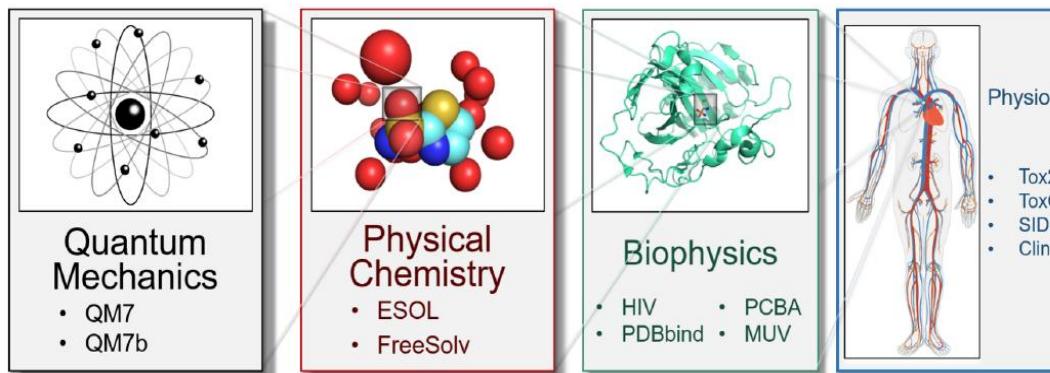
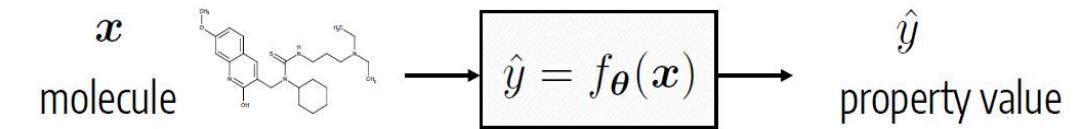
- Step. 1: Label the dataset with category for training
- Step. 2: For new chemical, calculate distance using techniques like “*Jaccard -Tanimoto*”
- Step. 3: Identify closest neighbour using distance matrix
- Step. 4: Out of the  $k$  nearest neighbours, the category having maximum values will be the category of the new chemical
- Step. 5: Optimal value of  $k$  can be found using cross validation



# Machine Learning in Molecular & Material Science

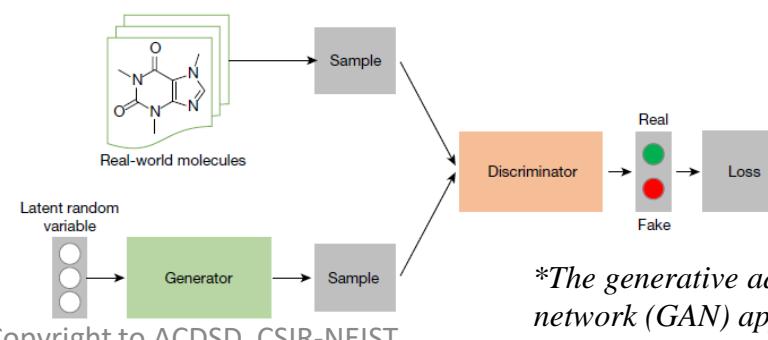
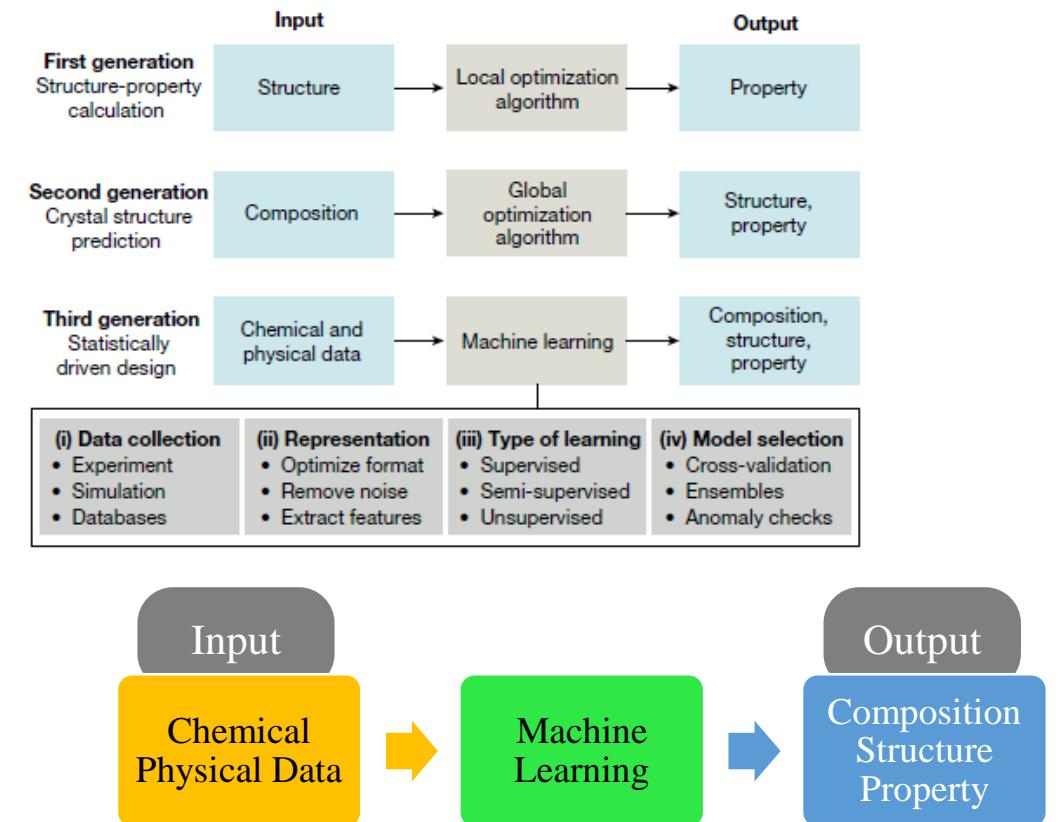
ACDS, CSIR-NEIST

## Quantitative structure–activity/property relationship (QSAR/QSPR)



MoleculeNet<sup>#</sup>: A Benchmark for Molecular Machine Learning

- ❖ Assess the likelihood that a product will crystallize
- ❖ Targeting discovery of new compounds relating to system descriptors

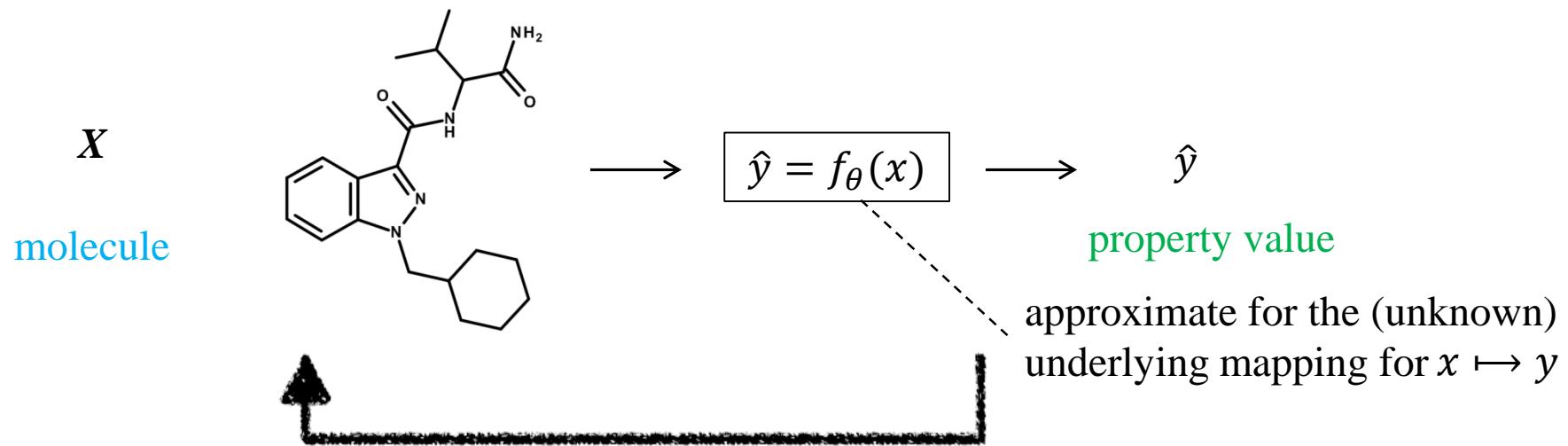


\*The generative adversarial network (GAN) approach to molecular discovery

<sup>#</sup><https://arxiv.org/abs/1703.00564>

<https://github.com/deepchem/deepchem>

\*Machine learning for molecular and materials science, Nature Reviews, 2018



Given  $n$  **input-output** instances as the training data

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

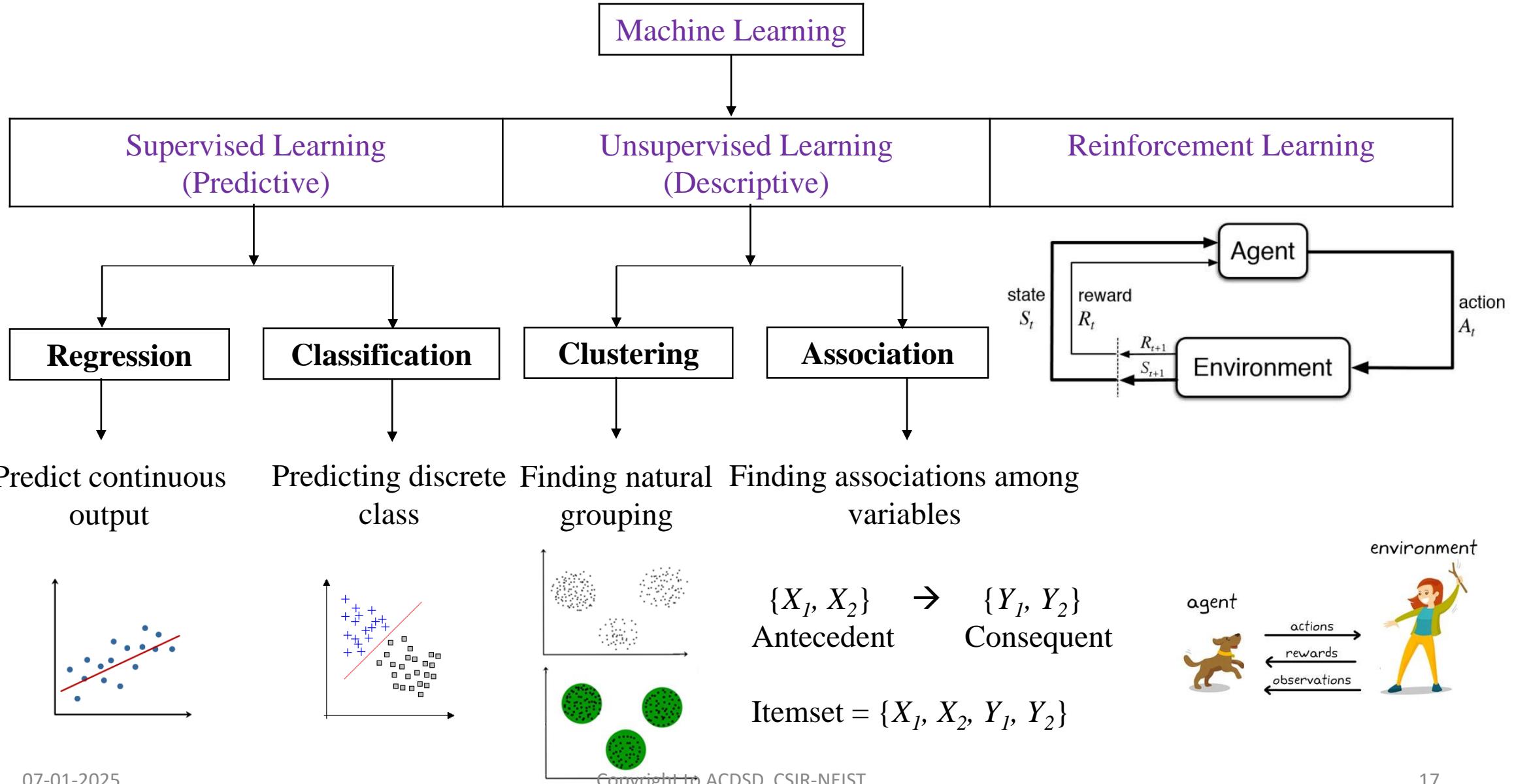
Fit the model  $f_{\theta}$  by tuning  $\Theta$  as

$$\min_{\theta} \sum_{i=1}^n \text{error}(y_i, \hat{y}_i) \quad \text{where} \quad \hat{y}_i = f_{\theta}(x_i)$$

*Note: the error measure (called "loss function" in ML) depends on problems*

# Types of Machine Learning

ACDS, CSIR-NEIST



```
## Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 1          5.1         3.5       1.4       0.2   setosa
## 2          4.9         3.0       1.4       0.2   setosa
## 3          4.7         3.2       1.3       0.2   setosa
## 4          4.6         3.1       1.5       0.2   setosa
## 5          5.0         3.6       1.4       0.2   setosa
## Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 51         7.0         3.2       4.7       1.4 versicolor
## 52         6.4         3.2       4.5       1.5 versicolor
## 53         6.9         3.1       4.9       1.5 versicolor
## 54         5.5         2.3       4.0       1.3 versicolor
## 55         6.5         2.8       4.6       1.5 versicolor
## Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 101        6.3         3.3       6.0       2.5 virginica
## 102        5.8         2.7       5.1       1.9 virginica
## 103        7.1         3.0       5.9       2.1 virginica
## 104        6.3         2.9       5.6       1.8 virginica
## 105        6.5         3.0       5.8       2.2 virginica
```

## Dataset 1: Labelled data

Make classification of all data based on labels

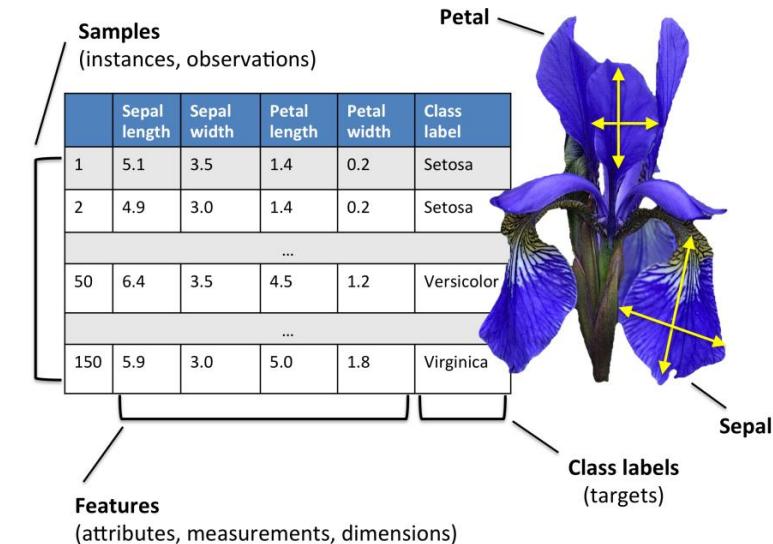


Three classes  
“setosa” or  
“versicolor” or  
“virginia”

New data with sepal length, sepal width, petal length and petal width but no label??



It belongs to  
“setosa” or  
“versicolor” or  
“virginia”



## Supervised Learning

```

## Sepal.Length Sepal.Width Petal.Length Petal.Width
## 1          5.1         3.5       1.4      0.2
## 2          4.9         3.0       1.4      0.2
## 3          4.7         3.2       1.3      0.2
## 4          4.6         3.1       1.5      0.2
## 5          5.0         3.6       1.4      0.2

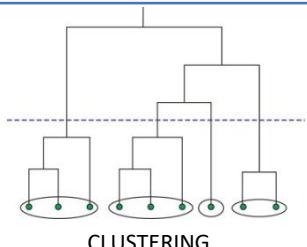
## Sepal.Length Sepal.Width Petal.Length Petal.Width
## 51         7.0         3.2       4.7      1.4
## 52         6.4         3.2       4.5      1.5
## 53         6.9         3.1       4.9      1.5
## 54         5.5         2.3       4.0      1.3
## 55         6.5         2.8       4.6      1.5

## Sepal.Length Sepal.Width Petal.Length Petal.Width
## 101        6.3         3.3       6.0      2.5
## 102        5.8         2.7       5.1      1.9
## 103        7.1         3.0       5.9      2.1
## 104        6.3         2.9       5.6      1.8
## 105        6.5         3.0       5.8      2.2

```

## Dataset 2: Unlabelled data

Cluster the data based on similarity and differences



Clusters of similar ‘petal length’, ‘petal width’ ..

New data with sepal length, sepal width, petal length and petal width but no label??

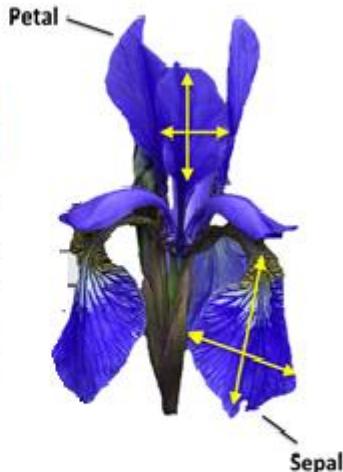


Identify the cluster to which the new sample belongs

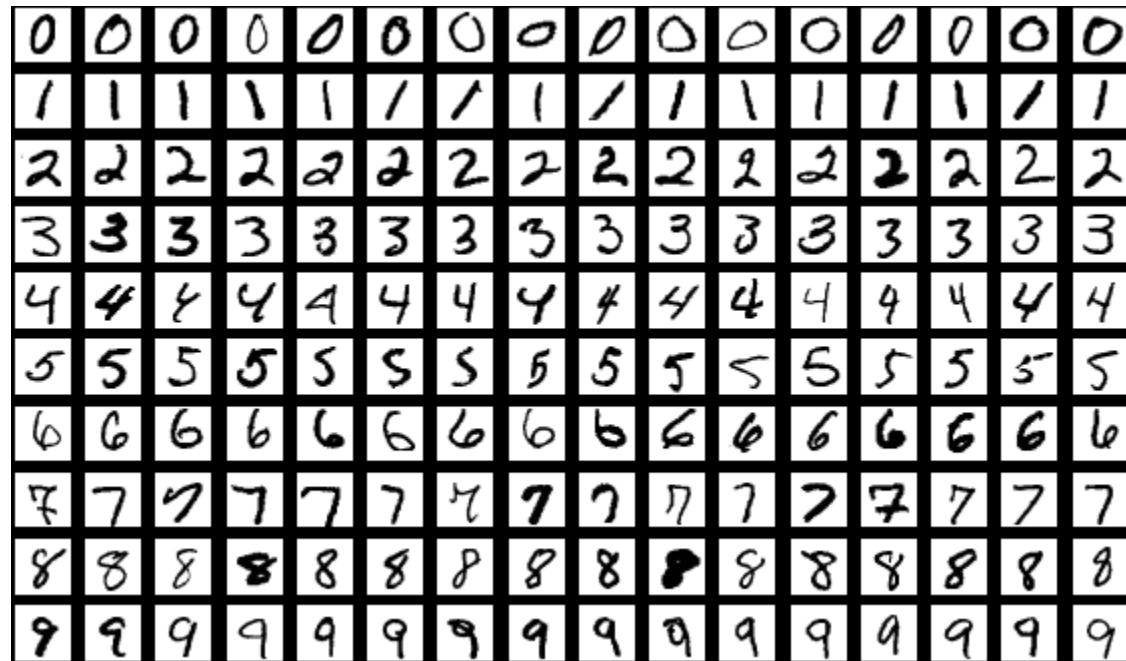
**Samples**  
(instances, observations)

	Sepal length	Sepal width	Petal length	Petal width
1	5.1	3.5	1.4	0.2
2	4.9	3.0	1.4	0.2
..	..	..	..	..
50	6.4	3.5	4.5	1.2
..	..	..	..	..
150	5.9	3.0	5.0	1.8

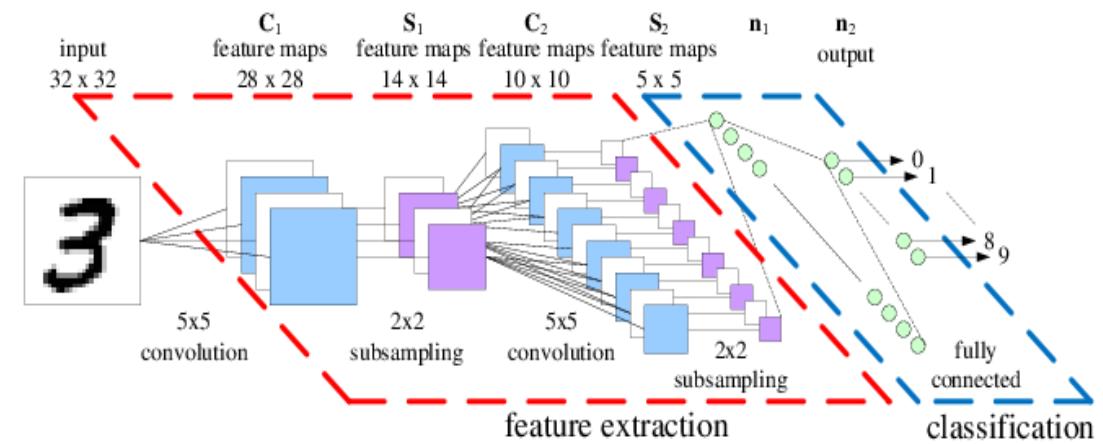
**Features**  
(attributes, measurements, dimensions)



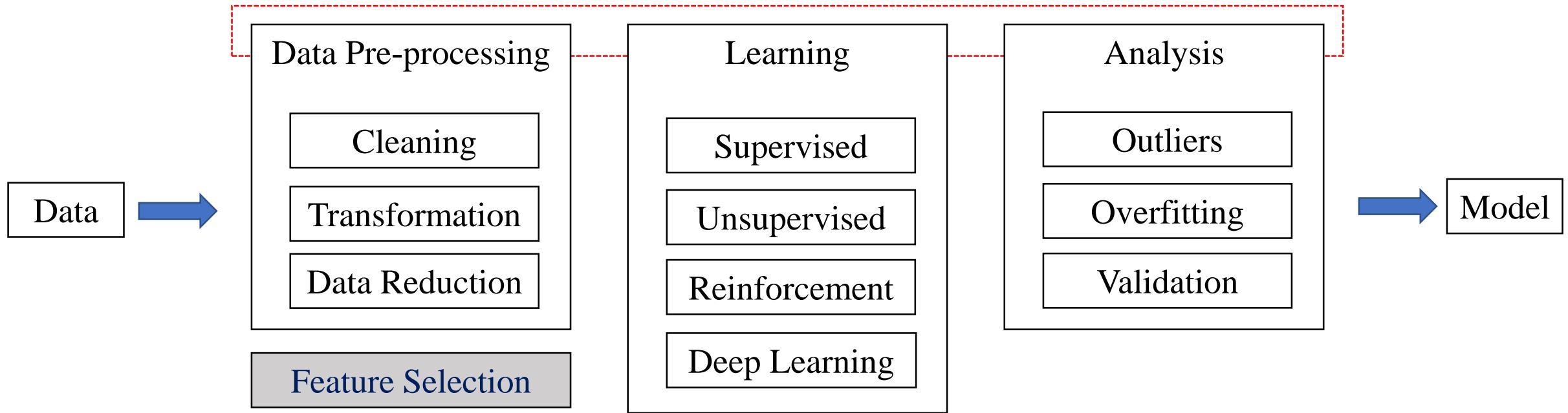
## Unsupervised Learning



**Dataset 3: Unlabelled data**



*Deep Learning*



Data Preprocessing		
Data Cleaning	Data Transformation	Data Reduction
Noisy Data	Discretization	Aggregation
Missing Data	Attribute Selection	Feature Subset Selection
Duplicate Data	Feature Scaling	Dimensionality Reduction

## Noisy values

- Due to noise (Data disrupted)
- Handled by collecting more data, principle component analysis, regularization and cross-validation

## Missing values

- Due to information not collected or attribute not applicable to all cases.
- Handled by either ignoring/dropping or replacing with some values (mean, median, max\_count etc.), predicting missing values, imputation methods

## Duplicate data

- Due to collection of data from multiple sources.
- Handled by deleting the duplicate entry.

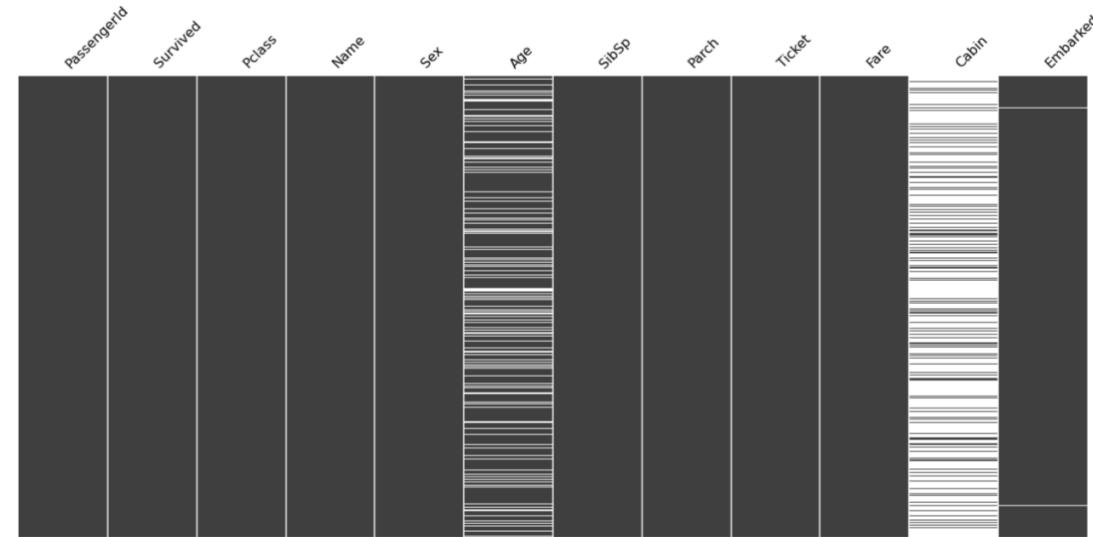
## Noisy Dataset Example

Value	Colour	Class
0.25	Red	Positive
0.25	Red	Negative
0.98	Green	Negative
1.02	Green	Positive
2.05	?	Negative
=	Green	Positive

Attribute noise

Class noise

## Visualization of missing data in a Dataset



Missing data

## Duplicate data in a Dataset

ID	0	1	2	3	4
34	4.9	3.1	1.5	0.1	Iris-Setosa
37	4.9	3.1	1.5	0.1	Iris-Setosa
142	5.8	2.7	5.1	1.9	Iris-Virginica

Duplicate data

- **Discretization:** Replace the raw values of numeric attribute by interval levels or conceptual levels.
- **Attribute Selection:** New attributes are constructed from given set of attributes
- **Feature Scaling:** Scale the data values in a specified range (-1.0 to 1.0 or 0.0 to 1.0)

- Min-max (Normalization):  $v' = \frac{v - min_A}{max_A - min_A} (new\_max_A - new\_min_A) + new\_min_A$

- Z score (standardization):  $v' = \frac{v - mean_A}{stand\_dev_A}$       or       $z_{ij} = \frac{x_{ij} - m_j}{s_j}$

$x_{ij}$  is the value for the  $i^{th}$  sample and  $j^{th}$  feature

$m_j$  is the average of all  $x_{ij}$  for feature  $j$

$s_j$  is the standard deviation of all  $x_{ij}$  over all input samples

Data discretization converts a large number of data values into smaller once, so that data evaluation and data management becomes very easy.

## ***Techniques for discretization:***

**Histogram analysis:** A plot used to present the underlying frequency distribution of a set of continuous data.

**Binning:** A data smoothing technique and its helps to group the a huge number of continuous values into a smaller number of bins.

**Correlation analysis:** Extensively used technique to identifies interesting relationships in data. These relationships help us realize the relevance of attributes with respect to the target class to be predicted.

**Cluster analysis:** Also commonly known as clustering, is the task of grouping with similar objects in one group, commonly called cluster

**Decision tree analysis:** Use a decision tree to identify the optimal number of bins. When the model makes a decision, it assigns an observation for each node.

**Equal width partitioning:** Separating all possible values into ‘n’ number of bins, each having the same width.

$$\text{width} = \frac{(\text{maximum value} - \text{minimum value})}{n} \quad \text{\#where } n \text{ is the number of bins or intervals.}$$

Age	10, 11, 13, 14, 17, 19, 30, 31, 32, 38, 40, 42, 70, 72, 73, 75
-----	--

Table: Before discretization

Attribute	Age	Age	Age
	10, 11, 13, 14, 17, 19,	30, 31, 32, 38, 40, 42	70, 72, 73, 75
After Discretization	Young	Mature	Old

The importance of features (attributes) selection:

- Reduce the cost of learning by reducing the number of attributes.
- Provide better learning performance compared to using full attribute set.

There are two approach for attribute selection.

- Filter approach attempt to assess the merits of attributes from the data, ignoring learning algorithm.
- Wrapper approach the attributes subset selection is done using the learning algorithm as a black box.

Most simple filter approach is ranking of attributes according value chi-square for two way tables. Two way table for this case is confusion matrix.

$$\chi^2 = \sum \frac{(observed\ count - expected\ count)^2}{expected\ count}$$

The expected count in any cell a two-way table is

$$expected\ count = \frac{row\ total \times column\ total}{total\ count}$$

Large value  $\chi^2$  better attribute

		PREDICTED CLASS	
		$D^0$	$BG$
ACTUAL CLASS	$D^0$	TP	FN
	$BG$	FP	TN

Preparing data for processing in machine learning. Normalization brings all data in the data set into a common scale.

## Techniques for normalization

- Z-Score Normalization:  $z = \frac{x - \mu}{\sigma}$  where,  $Z$  = standard score,  $x$  = observed value,  $\mu$  = mean of the sample  $\sigma$  = standard deviation of the sample
- Min-Max Normalization:  $\frac{value - min}{max - min}$ , where max and min are the maximum and minimum values in the dataset
- Decimal scaling Normalization:  $v' = \frac{v}{10^j}$ , where  $j$  is the lowest integer while  $\text{Max}(|v|) < 1$
- Standard Deviation Normalization:  $s_x = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$  where n it the number of data points,  $x_i$  is the data point,  $x'$  mean of  $x_i$

## Worked out example on normalizing the data using Min-Max

```
In [2]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
```

```
In [3]: df = pd.DataFrame({
    'Income': [15000, 1800, 120000, 10000],
    'Age': [25, 18, 42, 51],
    'Department': ['HR', 'Legal', 'Marketing', 'Management']
})
```

```
In [4]: df
```

Out[4]:

	Income	Age	Department
0	15000	25	HR
1	1800	18	Legal
2	120000	42	Marketing
3	10000	51	Management

```
In [7]: df_scaled = df.copy()
col_names = ['Income', 'Age']
features = df_scaled[col_names]
```

```
In [8]: from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
```

```
In [9]: df_scaled[col_names] = scaler.fit_transform(features.values)
```

```
In [10]: df_scaled
```

Out[10]:

	Income	Age	Department
0	0.111675	0.212121	HR
1	0.000000	0.000000	Legal
2	1.000000	0.727273	Marketing
3	0.069374	1.000000	Management

Scaling the data  
without range

```
In [11]: scaler = MinMaxScaler(feature_range=(5, 10))
```

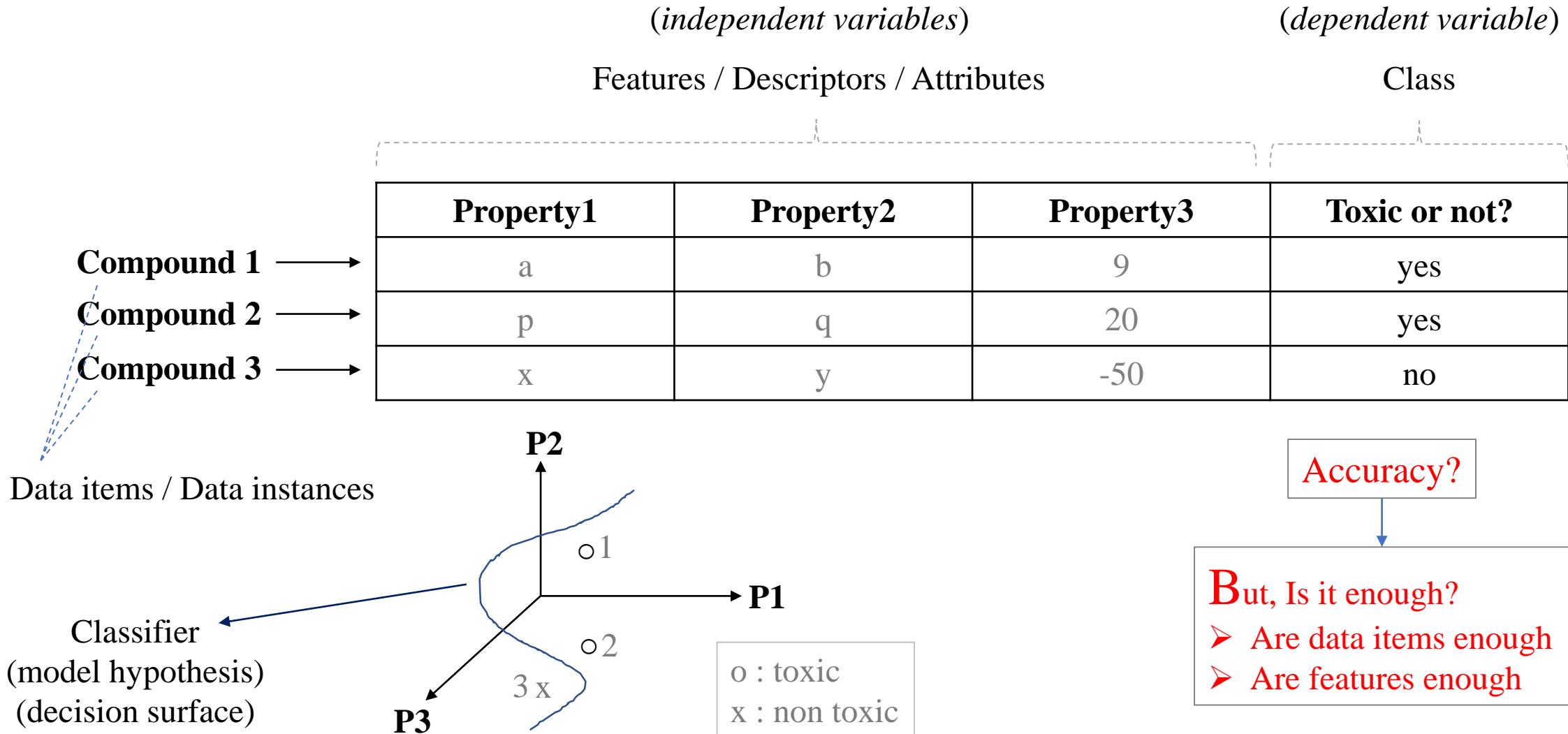
```
df_scaled[col_names] = scaler.fit_transform(features.values)
df_scaled
```

Out[11]:

	Income	Age	Department
0	5.558376	6.060606	HR
1	5.000000	5.000000	Legal
2	10.000000	8.636364	Marketing
3	5.346870	10.000000	Management

Scaling the data  
with range (5-10)

Let us consider the following example – Toxicity prediction (toxic or not?)



Let us include more data items and more features to increase accuracy of the model -

	P1	P2	P3	P4	...	...	...	P97	P98	P99	P100	Toxic or not?
Compound 1 →												yes
Compound 2 →												yes
Compound 3 →												no
...	...											
...	...											
...	...											
Compound 998 →												no
Compound 999 →												no
Compound 1000 →												no

Now, Is it enough?

- Are data items enough
- Are features enough

If yes

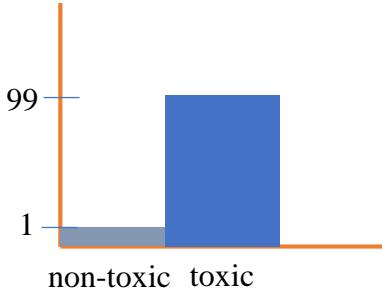
What about performance?

Overfitting - Resampling  
Curse of dimensionality - Dimension Reduction

## Resampling

Dataset Perspective  
(imbalanced dataset)

- Undersampling
- Oversampling
- Random sampling



Model Performance Perspective  
(overfitting)

- Hold out validation
- K-fold cross validation
- Bootstrapping

```
check_data(x) {  
    model();  
}  
  
model() {  
    print("Toxic");  
}
```

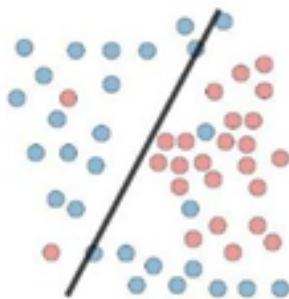
		toxic	non-toxic
toxic	toxic	99	1
	non-toxic	0	0

Undersampling

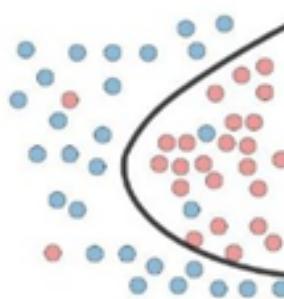
Oversampling

99% accuracy on train data (without doing anything)

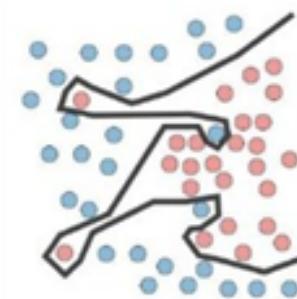
- Generated models have two important aspects to fulfill (learning and generalisation):
    - They must learn decision surfaces correctly (zero error) from training data.
    - They must be able to generalize i.e. deal appropriately with new data.
  - Usually we want our models to learn well and also to generalize well.
- 
- Zero error is possible - training too much!, but leads to overfitting - trade-off between learning and generalization.
  - Overfittings may fit the training data well, even outliers but fail to generalize to new data - loose performance on new data.
  - Allowing not classify all training data totally accurately (allowing misclassification) likely to lead better generalization.
  - Thus non-perfect learning is better in this case!



Underfitting



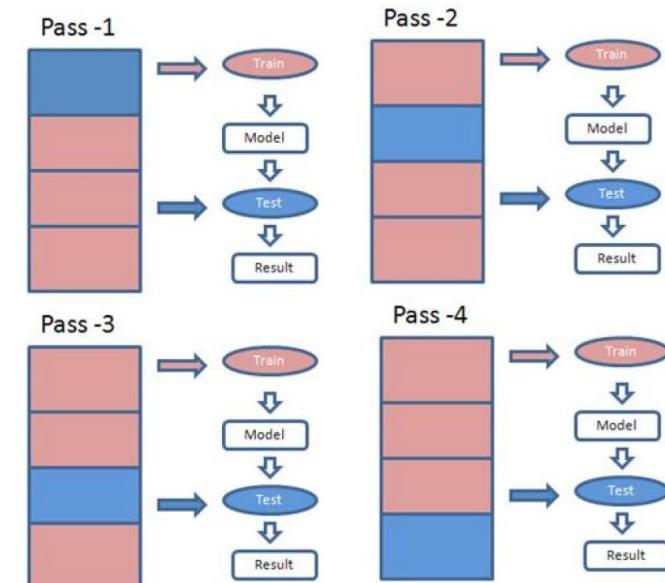
Bestfitting



Overfitting

**Occam's Razor (1300s):** “plurality should not be assumed without necessity” - The simplest model which explains the majority of the data is usually the best.

- Split-sample or hold-out validation (60-20-20)
  - Reserve some data as a "test set", which must not be used during training.
  - Disadvantage is that it reduces the amount of data available for both training and validation.
- K-fold cross-validation (e.g., leave one out)
  - Data is divided into k subsets and train k times, each time leaving one subset out for computing the error.
  - "Crossing" makes an improvement over split-sampling allowing all data used for training.
  - Disadvantage is that the network must be re-trained many times (k times in k-fold crossing).
- Bootstrapping
  - Works on random sub-samples (random shares) chosen from the full data set.
  - Any data item may be selected any number of times for validation.
  - The sub-samples are repeatedly analysed.



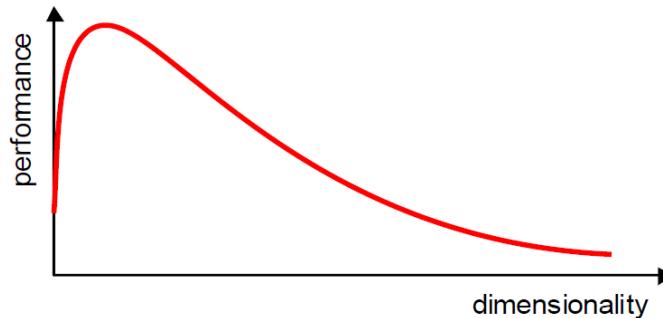
## How long should a model be trained?

- Stop if the error reaches an acceptable level (e.g., 95%)
- Stop if the error fails to improve (has reached a minimum)
- Stop if the rate of improvement drops below a certain level

➤ The objective is to establish a balance between memorization and generalization.

## Curse of dimensionality

- Increasing the number of features will not always improve classification accuracy.
- In practice, the inclusion of more features might actually lead to **worse** performance.



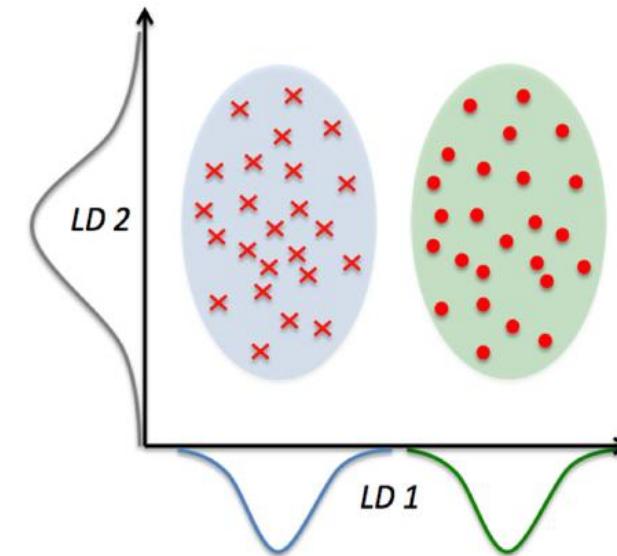
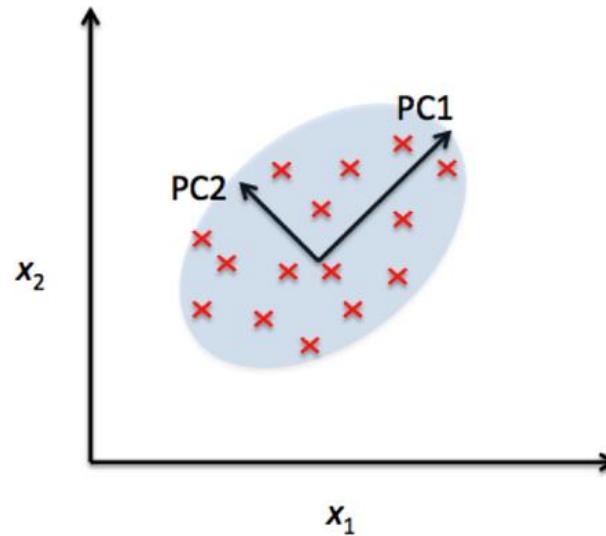
## Dimensionality Reduction

Finds a set of **new** features from **existing** features  
(i.e., through some mapping  $f()$ )

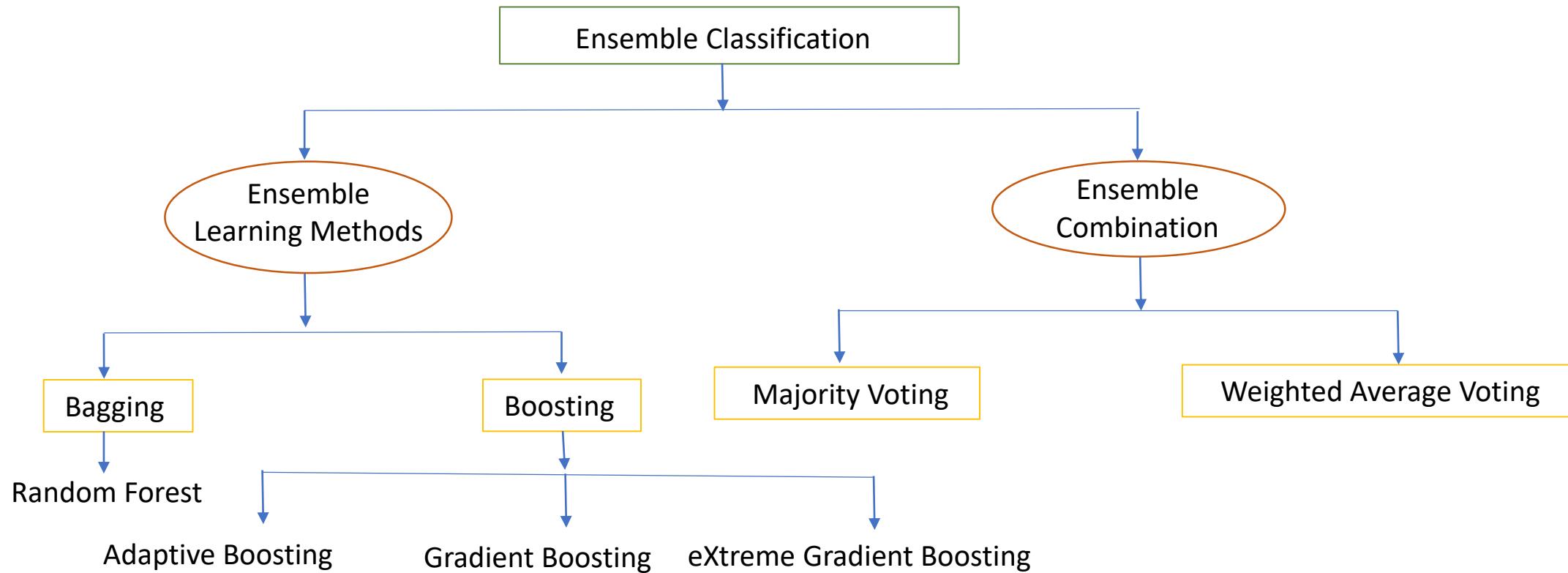
— PCA

— LDA

- Principal Components Analysis (PCA): Seeks a projection **component** that **preserves** as much **information** in the data as possible.
- Linear Discriminant Analysis (LDA): Seeks a projection that **best discriminates** the data.



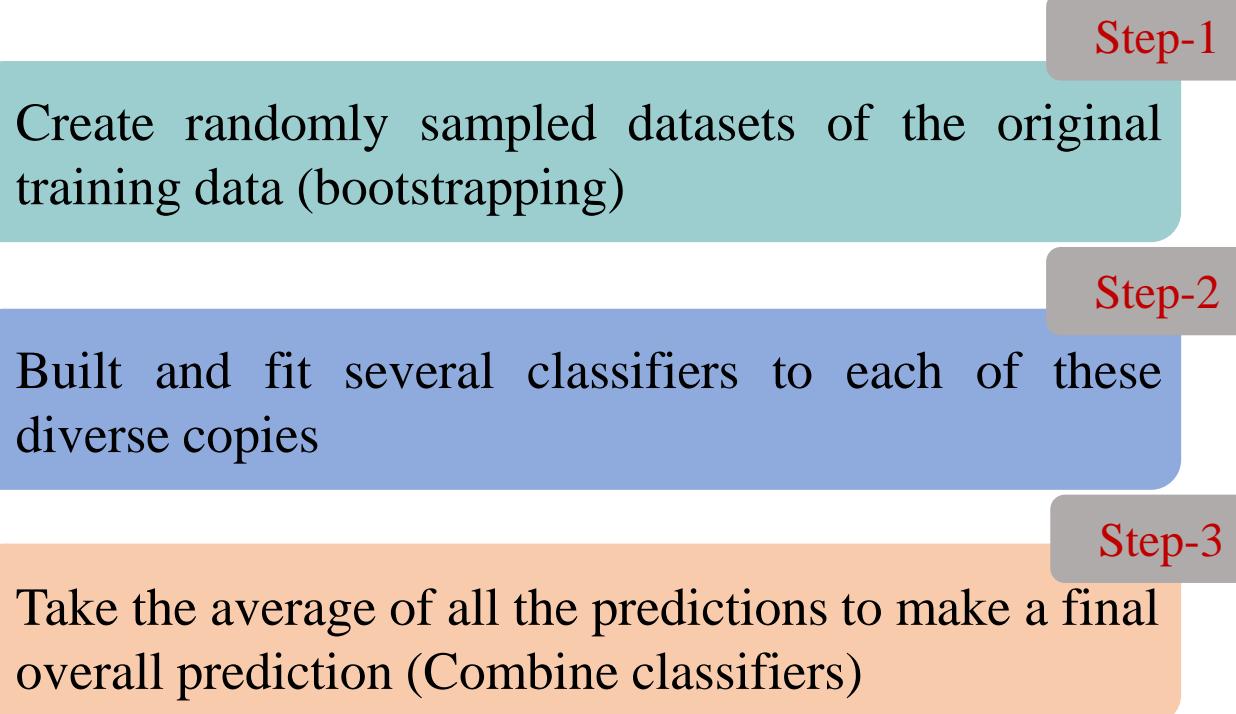
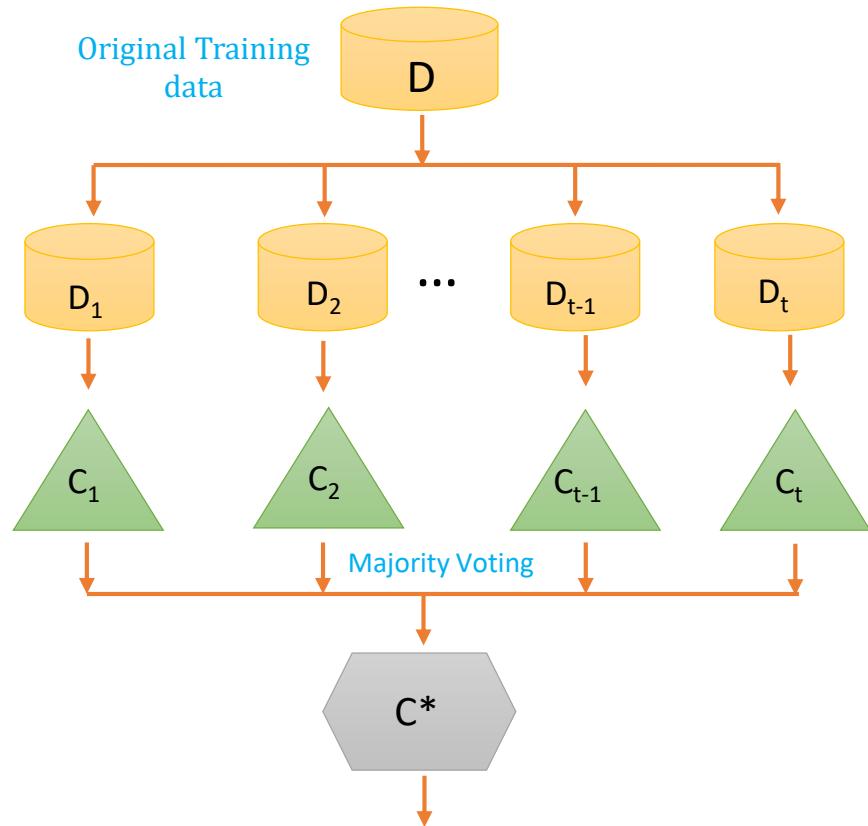
- Ensemble learning is the process by which multiple models, such as classifiers or experts, are strategically generated and combined to solve a particular computational intelligence problem.
- Ensemble learning is primarily used to improve the (classification, prediction, function approximation, etc.) performance of a model, or reduce the likelihood of an unfortunate selection of a poor one.



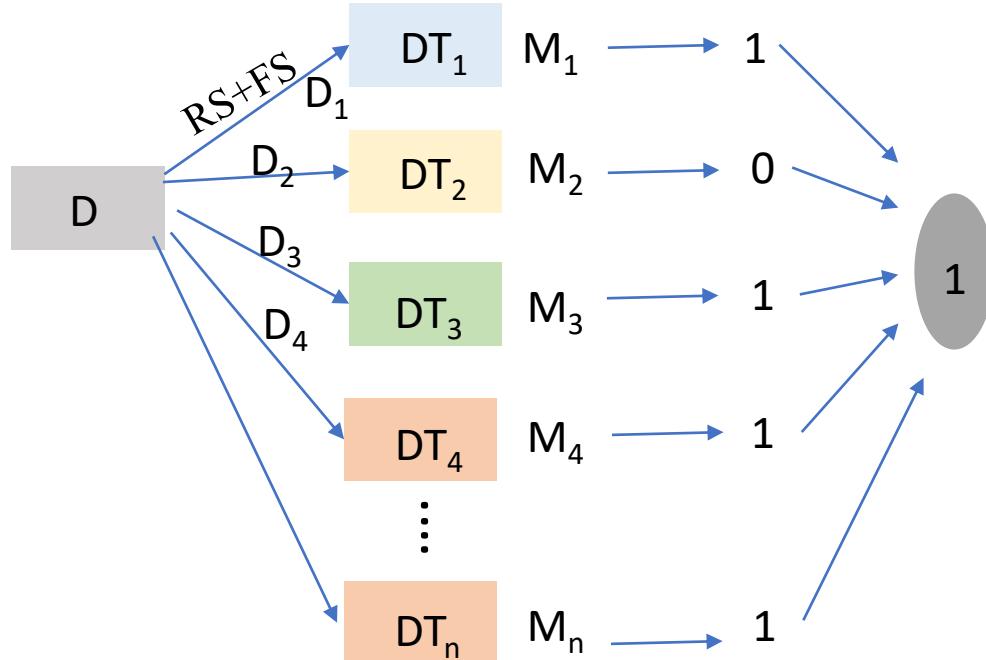
Method	Minimize Bias?	Minimize Variance?	Other Comments
Bagging	Complex model class. (Deep DTs)	Bootstrap aggregation (resampling training data)	Does not work for simple models.
Random Forests	Complex model class. (Deep DTs)	Bootstrap aggregation + bootstrapping features	Only for decision trees.
Gradient Boosting (AdaBoost)	Optimize training performance.	Simple model class. (Shallow DTs)	Determines which model to add at run-time.
Ensemble Selection	Optimize validation performance	Optimize validation performance	Pre-specified dictionary of models learned on training set.

- State-of-the-art prediction performance
  - Won Netflix Challenge
  - Won numerous KDD Cups
  - Industry standard

- Bagging/bootstrap aggregation is a way to decrease the variance in the prediction by generating additional data for training from dataset using combinations with repetitions to produce multi-sets of the original data.
- In general it reduces variance of an estimate by taking mean of multiple estimates.



## Random Forest



- Random forest technique combines various decision trees to produce a more generalized model.
- Random forests are utilized to produce de-correlated decision trees.
- Random forest creates random subsets of the features.
- Smaller trees are built using these subsets, creating tree diversity.
- To overcome overfitting, diverse sets of decision trees are required.

The dataset consisting of :

- Weather information of last 14 days
- Whether match was played or not on that particular day

Now using random forest we need to predict whether the game will happen if the weather condition is

Outlook = Rain

Humidity = High

Wind = Weak

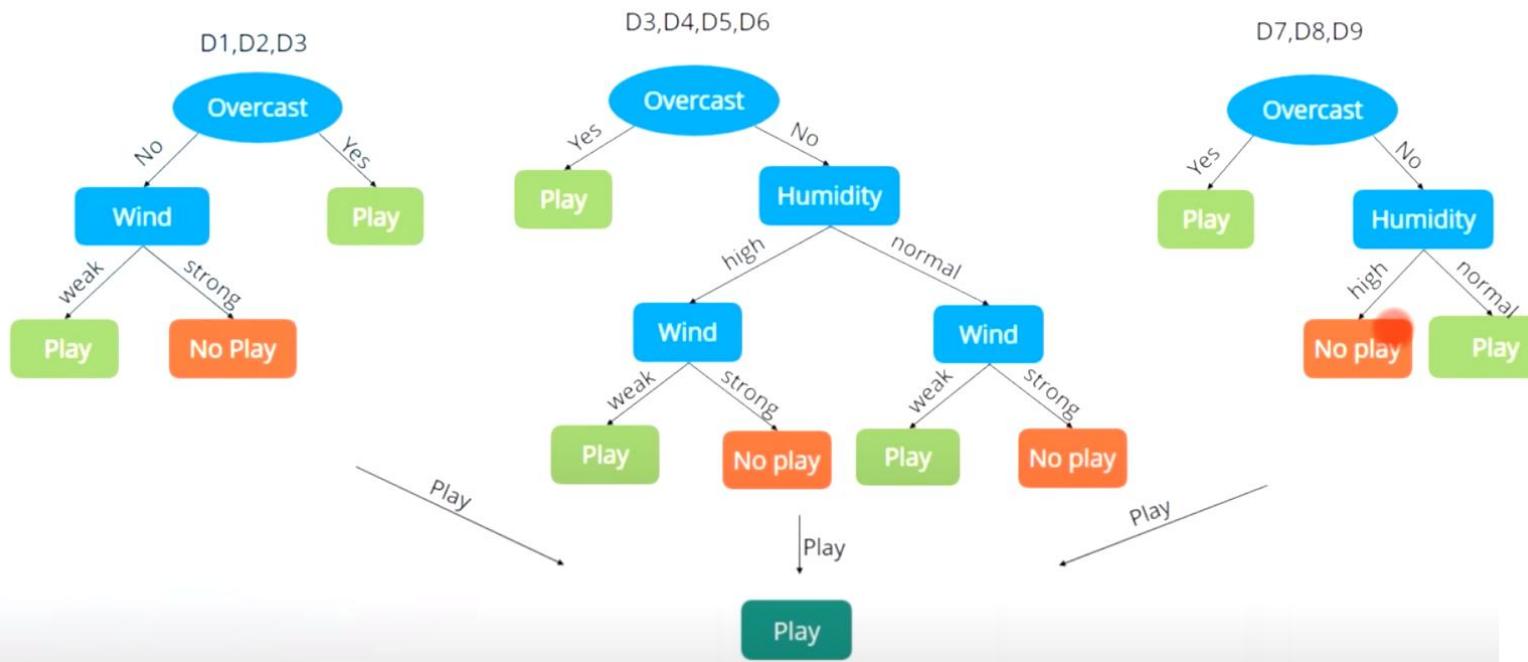
Play = ?

Day	Outlook	Humidity	Wind	Play
D1	Sunny	High	Weak	No
D2	Sunny	High	Strong	No
D3	Overcast	High	Weak	Yes
D4	Rain	High	Weak	Yes
D5	Rain	Normal	Weak	Yes
D6	Rain	Normal	Strong	No
D7	Overcast	Normal	Strong	Yes
D8	Sunny	High	Weak	No
D9	Sunny	Normal	Weak	Yes
D10	Rain	Normal	Weak	Yes
D11	Sunny	Normal	Strong	Yes
D12	Overcast	High	Strong	Yes
D13	Overcast	Normal	Weak	Yes
D14	Rain	High	Strong	No

# How Random Forest Works?

AODS, CSIR-NEST

- Step 1: The first step in random forest is that it will divide the data into smaller subsets.
- Every subset need not be distinct, some subsets may be overlapped.



Day	Outlook	Humidity	Wind	Play
D1	Sunny	High	Weak	No
D2	Sunny	High	Strong	No
D3	Overcast	High	Weak	Yes
D4	Rain	High	Weak	Yes
D5	Rain	Normal	Weak	Yes
D6	Rain	Normal	Strong	No
D7	Overcast	Normal	Strong	Yes
D8	Sunny	High	Weak	No
D9	Sunny	Normal	Weak	Yes
D10	Rain	Normal	Weak	Yes
D11	Sunny	Normal	Strong	Yes
D12	Overcast	High	Strong	Yes
D13	Overcast	Normal	Weak	Yes
D14	Rain	High	Strong	No



Most accurate learning algorithms



Works well for both classification and regression problems



Runs efficiently on large databases



Requires almost no input preparation



Performs implicit feature selection

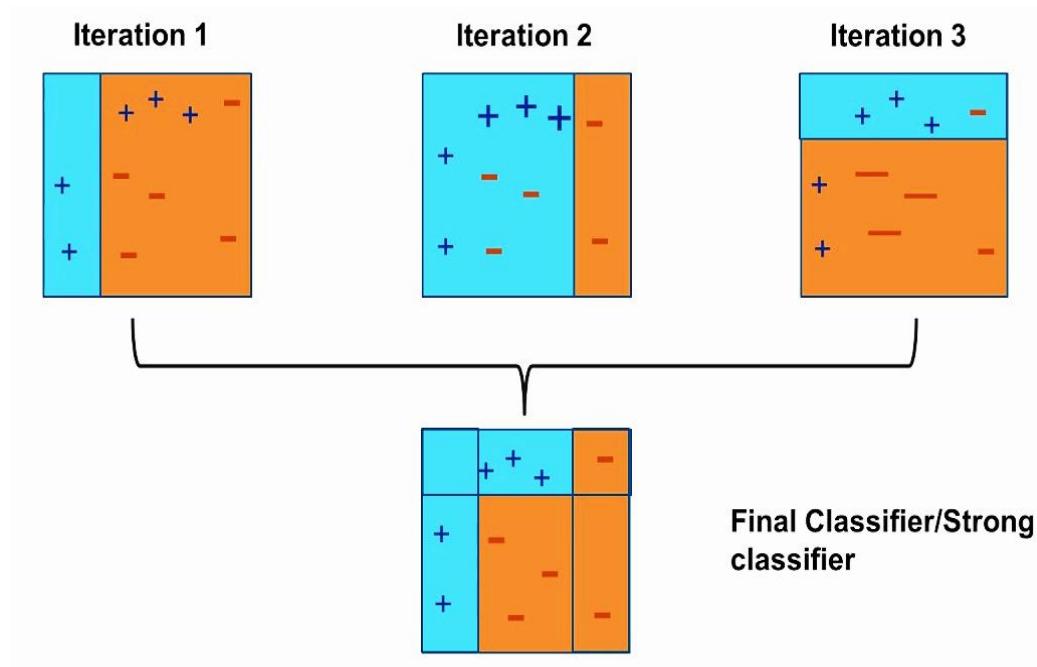


Can be easily grown in parallel



Methods for balancing error in unbalanced data sets

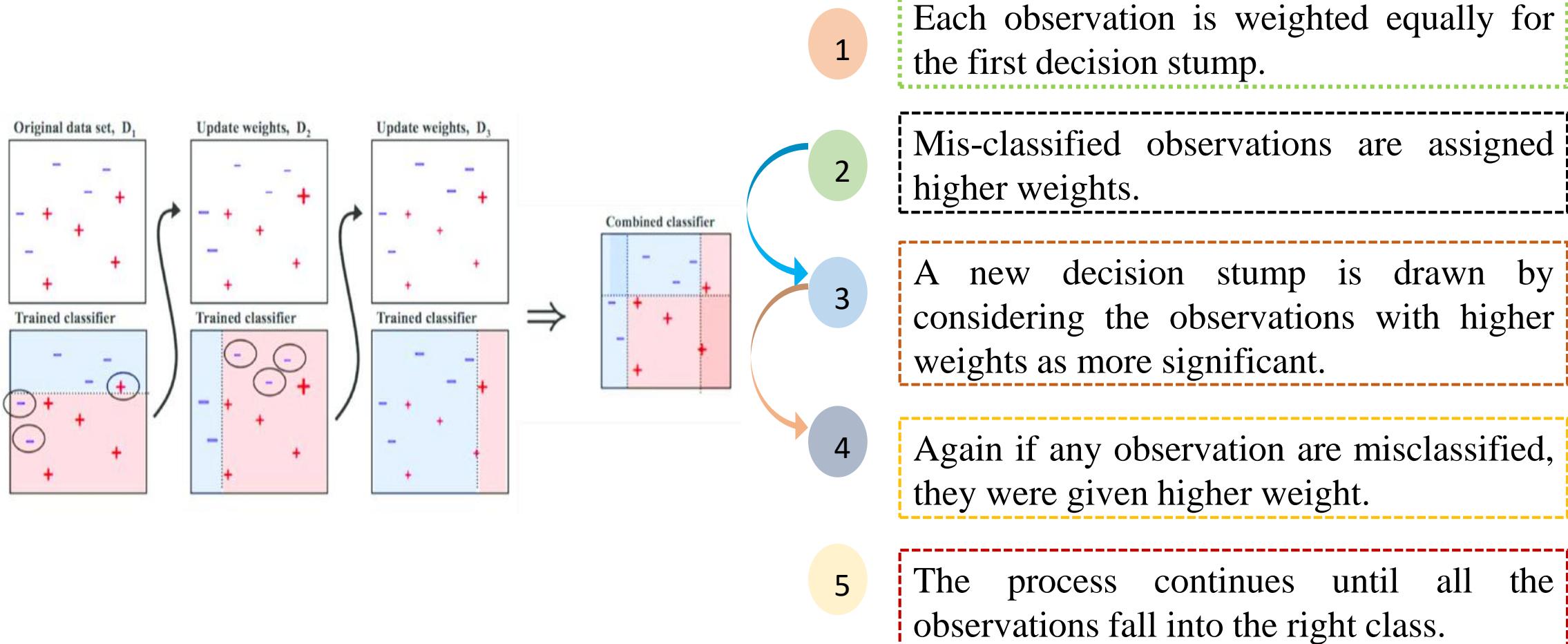
- Boosting is an ensemble learning technique that uses a set of machine learning algorithms in order to convert or combine weak learners to strong learners in order to increase the accuracy of the model.
- So boosting is very effective method in order to increase the efficiency of your model.



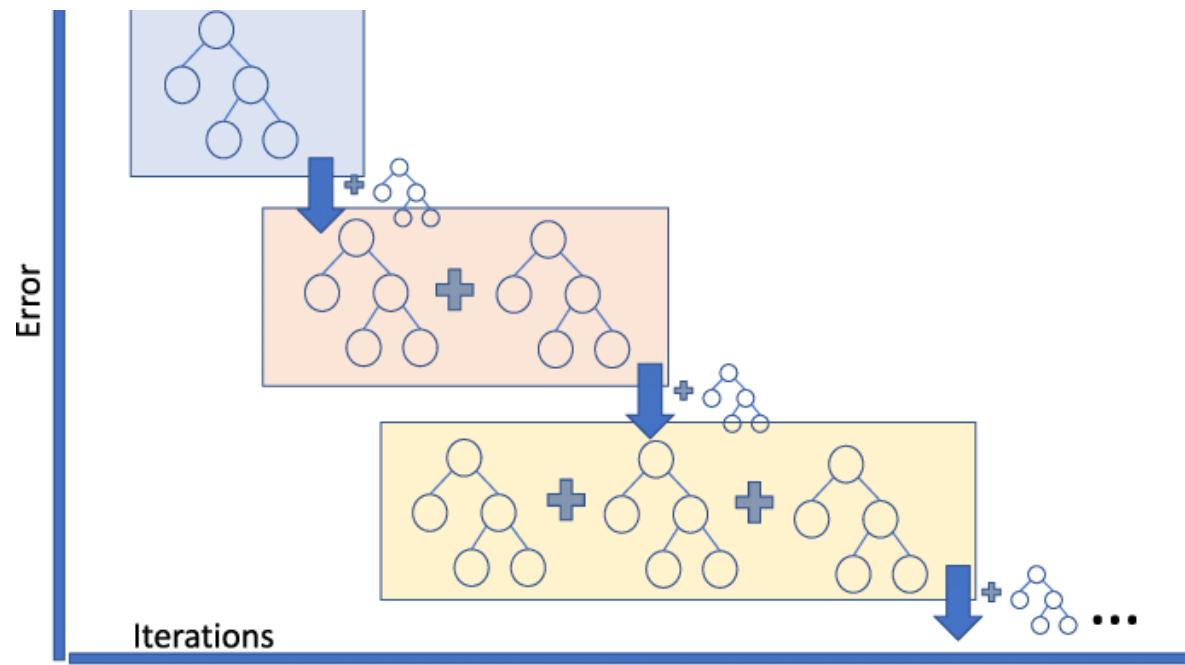
**Step-1:** The base algorithm reads the data and assigns equal weight to each sample observation.

**Step-2:** False predictions are assigned to the next base learner with a higher weightage on these incorrect predictions

**Step-3:** Repeat step 2 until the algorithm can correctly classify the output.



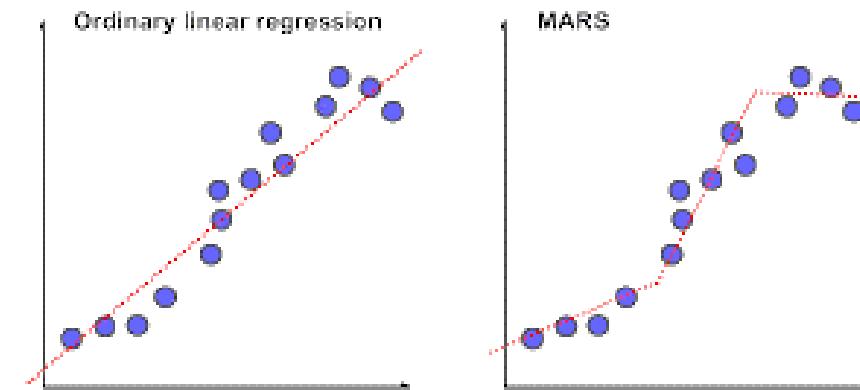
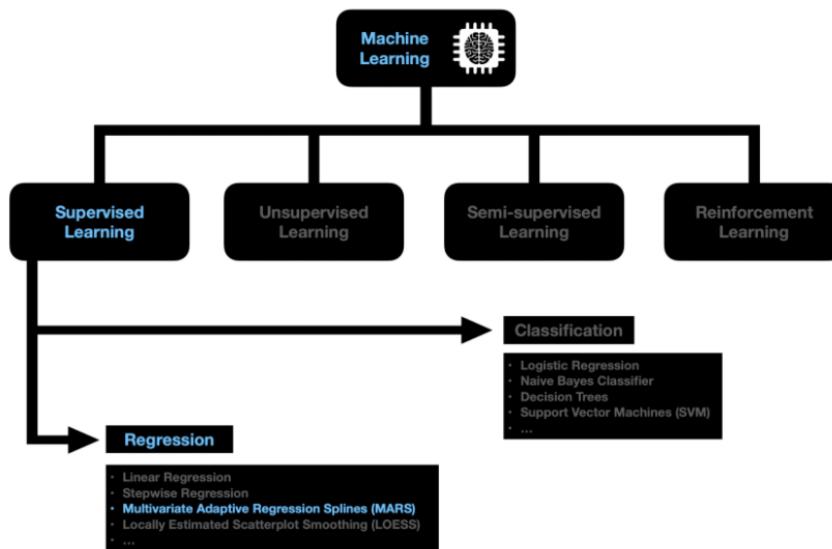
- In Gradient Boosting, base learners are generated sequentially in such a way that the present base learner is always more effective than the previous one.
- The key idea is to set the target outcomes for this next model in order to minimize the error.
- It has three main components:
  - Loss function that needs to be optimized.
  - Weak learner for computing predictions and forming strong learners.
  - An additive Model that will regularize the loss function.



- XGBoost is an advanced version of gradient boosting method that is designed to focus on computational speed and model efficiency.
- It supports parallelization by creating decision trees parallel
- It implements distributed computing methods for evaluating any large or complex method.
- It also uses Out-of-Core computing in order to analyze huge and varied datasets.



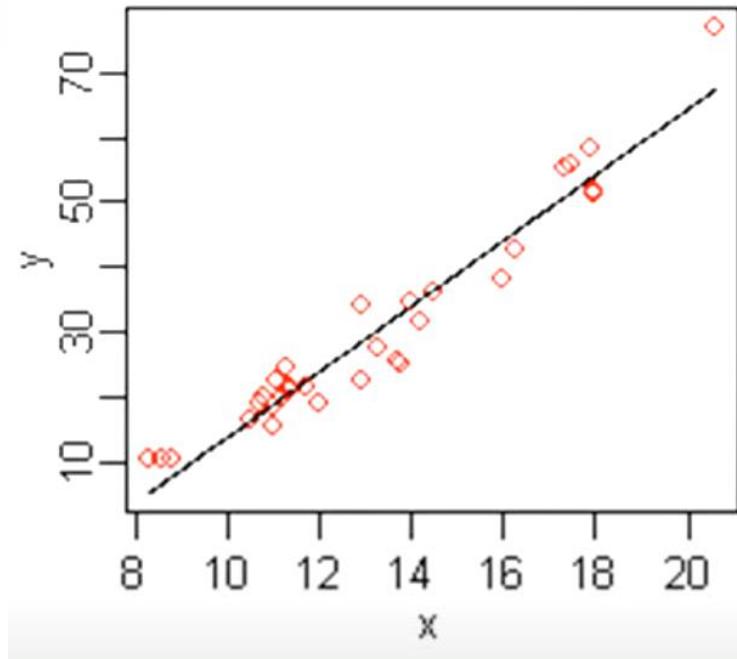
- MARS is a type of ensemble of simple linear functions and can achieve good performance on challenging regression problems with many input variables and complex non-linear relationships.
- It is a non-parametric regression technique and can be seen as an extension of linear models that automatically models non-linearities and interactions between variables.
- The algorithm involves finding a set of simple linear functions that in aggregate result in the best predictive performance.



# Multivariate Adaptive Regression Splines

Example: A set of data with a matrix of input variables  $x$ , and a vector of the observed responses  $y$ , with a response for each row in  $x$ . For example, the data could be: BASICS x y 10.5 16.4 10.7 18.8 10.8 19.7 ... ... 20.6 77.0

X	Y
10.5	16.4
10.7	18.8
10.8	19.7
...	...
20.6	77



$$y = b_0 + b_1 * x_1$$

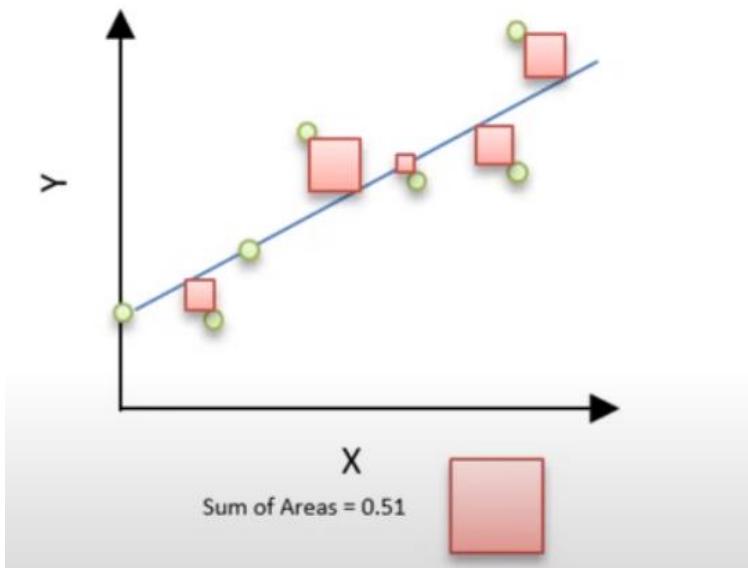
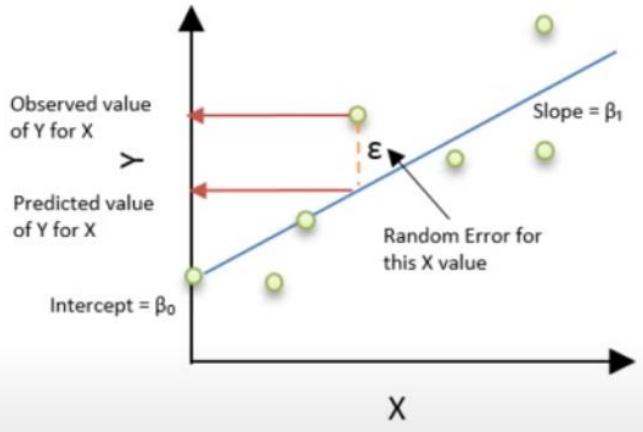
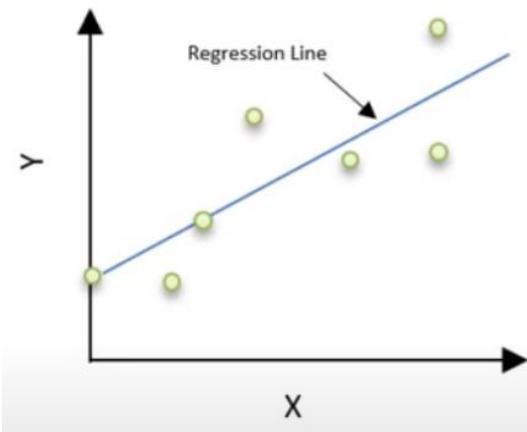
Annotations for the equation:

- Dependent variable: Points to the  $y$  term.
- Independent variable: Points to the  $x_1$  term.
- Constant: Points to the  $b_0$  term.
- Coefficient: Points to the  $b_1$  term.

A linear model for the above data  
is  $\tilde{y} = -37 + 5.1x$

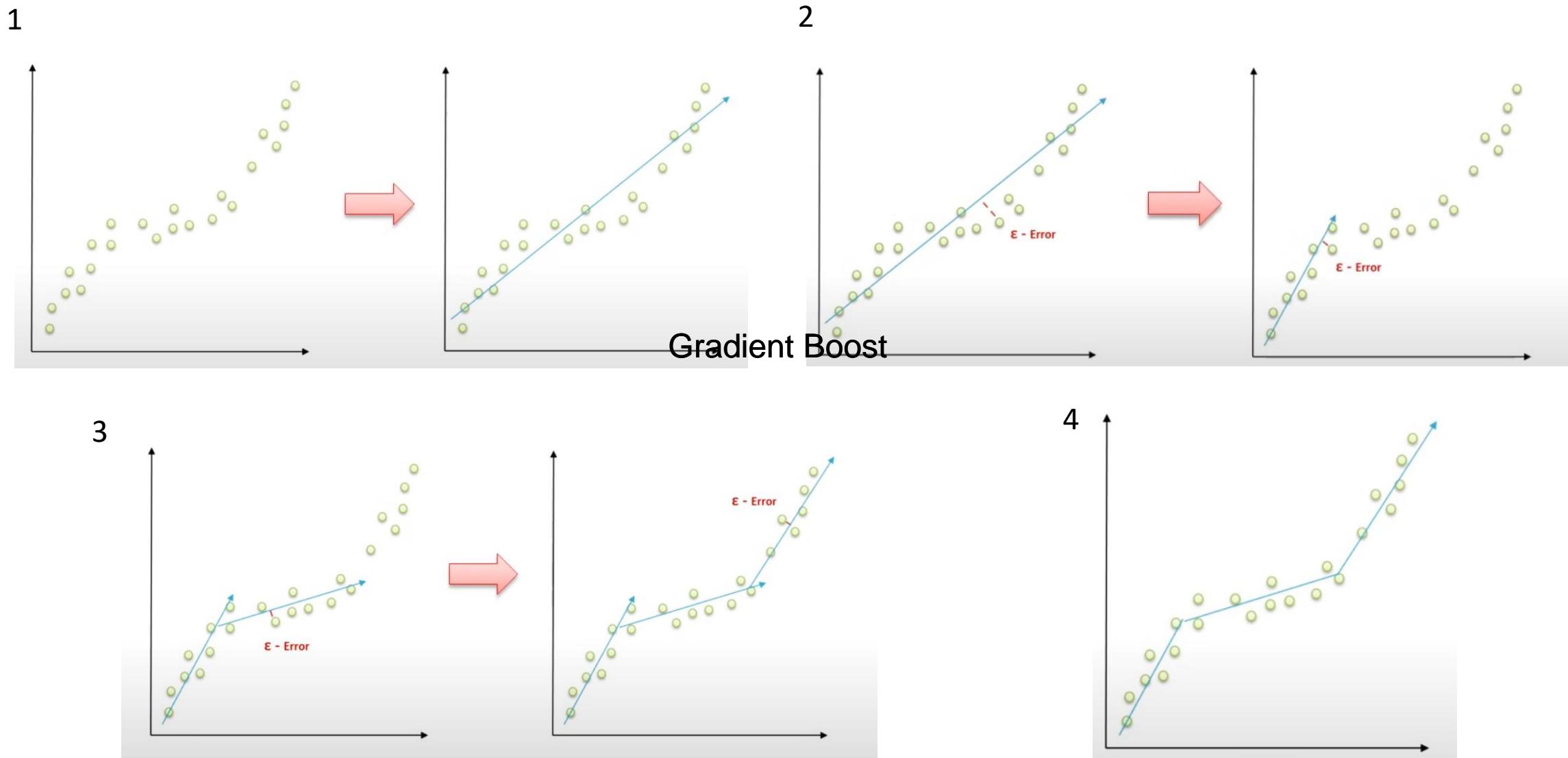
The figure shows a plot of this function: a line giving the predicted  $\tilde{y}$  versus  $x$ , with the original values of  $y$  shown as red dots.

# Multivariate Adaptive Regression Splines

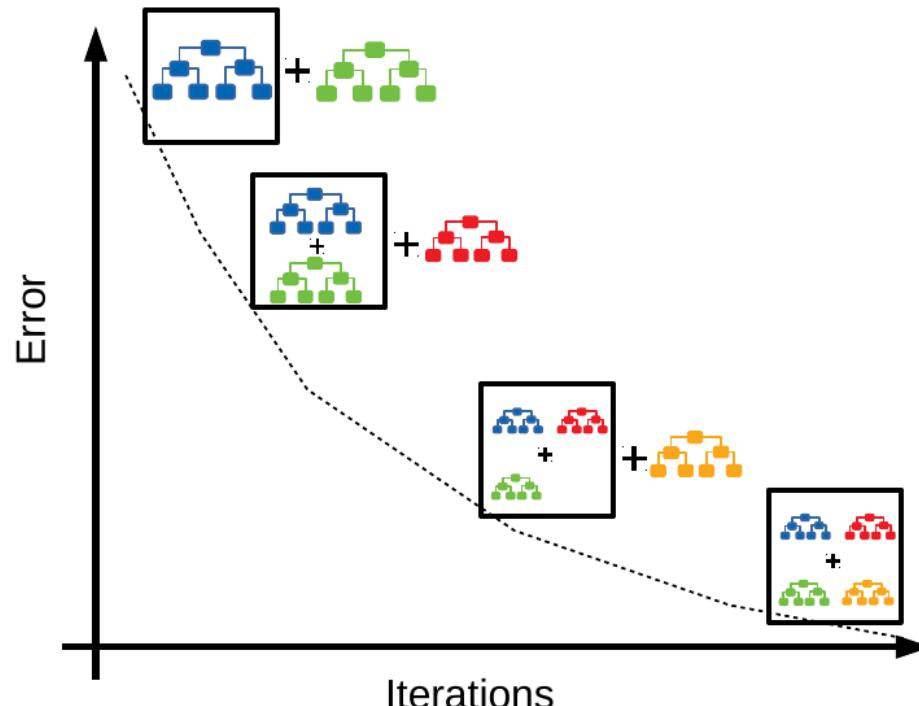


# Multivariate Adaptive Regression Splines

AODS, CSIR-NEST

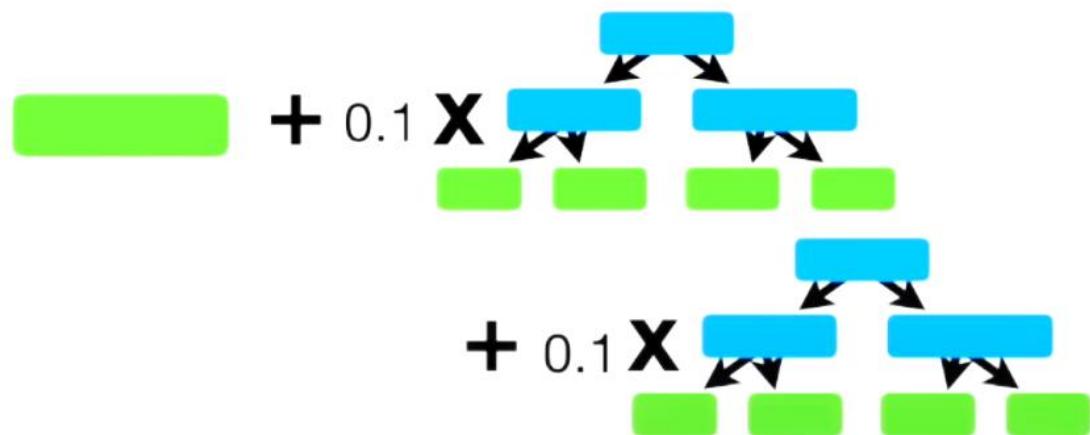


- In gradient boosting machines, or simply, GBMs, the learning procedure consecutively fits new models to provide a more accurate estimate of the response variable.
- A new tree is added, which is trained on labels that result from the errors made by the existing ensemble of trees.
- In gradient boosting, it trains many models sequentially. Each new model gradually minimizes the loss function and needs special attention as it is an error term of the whole system.
- The learning procedure consecutively fit new models to provide a more accurate estimate of the response variable.

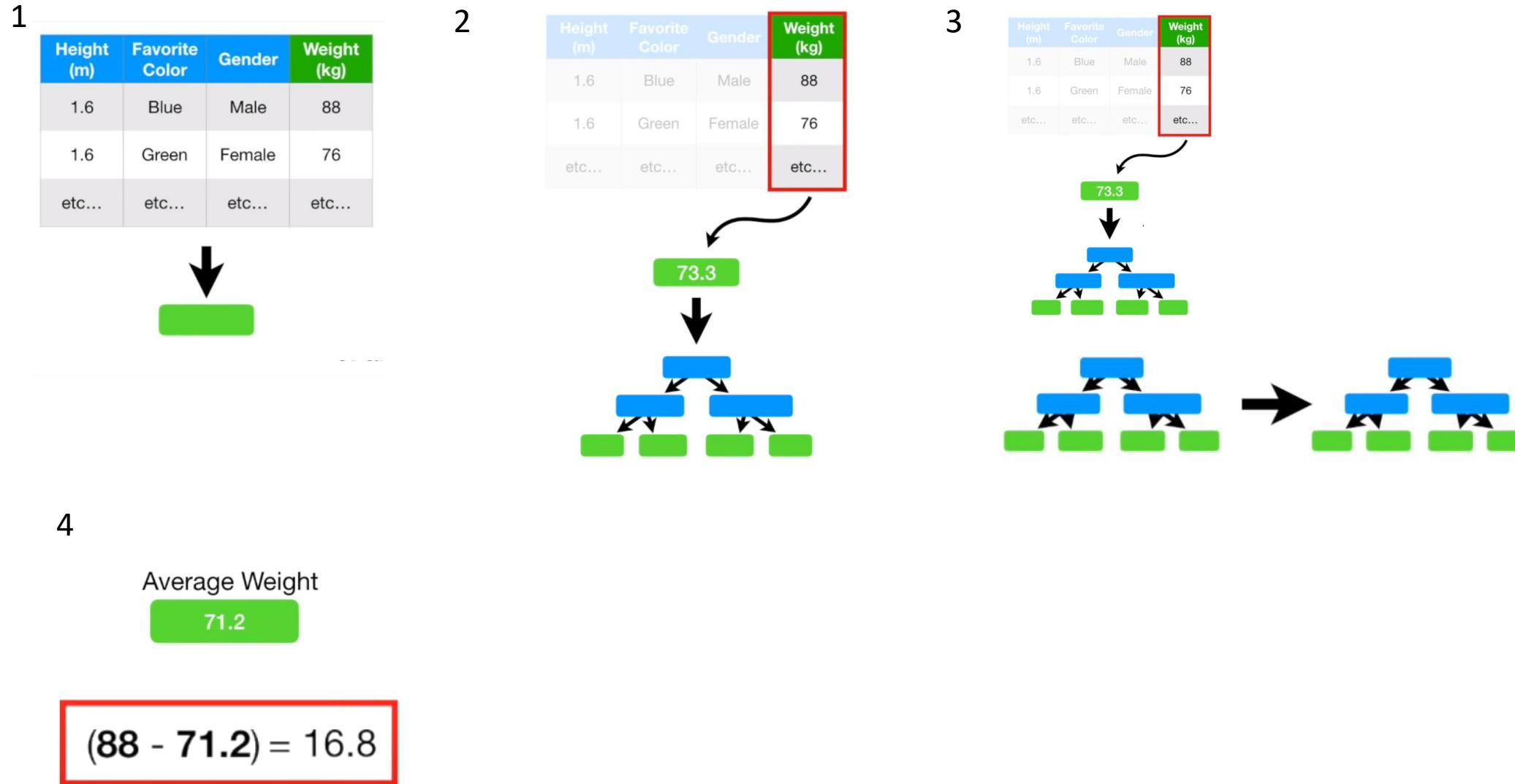


- Type of Problem – You have a set of variables vectors  $x_1$ ,  $x_2$  and  $x_3$ . You need to predict  $y$  which is a continuous variable.

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56
1.8	Red	Male	73
1.5	Green	Male	77
1.4	Blue	Female	57



# Gradient Boosting



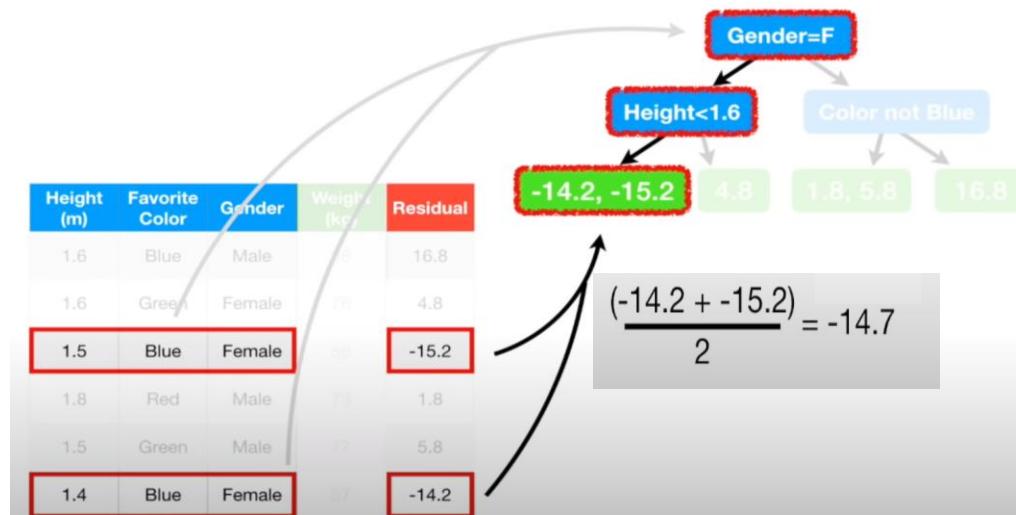
# Gradient Boosting

AODS, CSIR-NEST

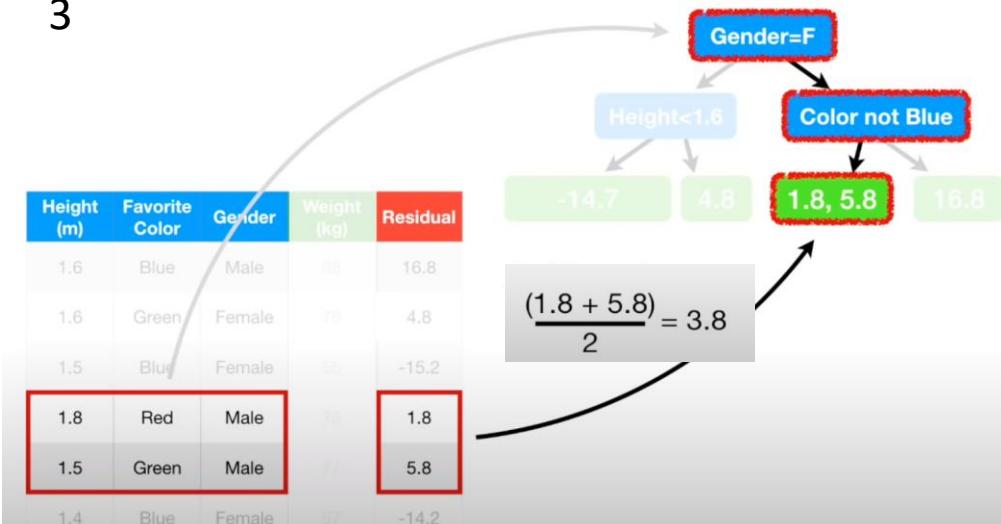
1

Height (m)	Favorite Color	Gender	Weight (kg)	Residual
1.6	Blue	Male	88	16.8
1.6	Green	Female	76	4.8
1.5	Blue	Female	56	-15.2
1.8	Red	Male	73	1.8
1.5	Green	Male	77	5.8
1.4	Blue	Female	57	-14.2

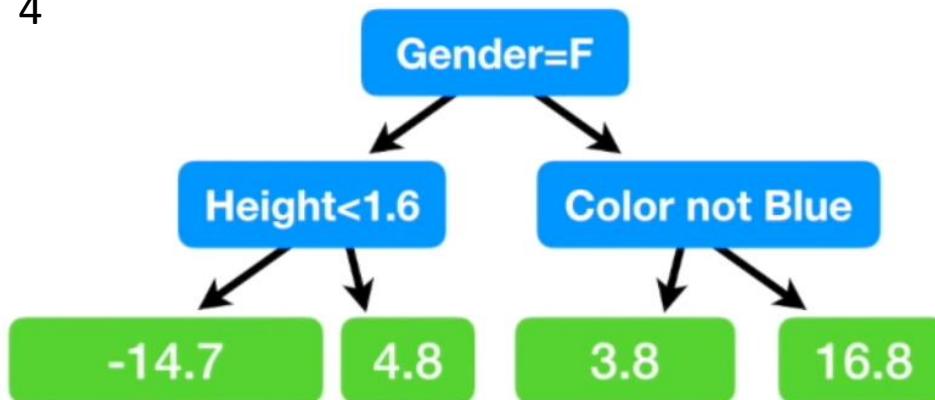
2



3



4



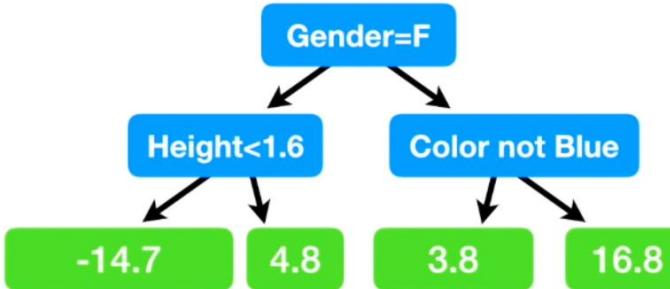
5

Height (m)	Favorite Color	Gender	Weight (kg)	Residual
1.6	Blue	Male	88	16.8
1.6	Green	Female	76	4.8
1.5	Blue	Female	56	-15.2
1.8	Red	Male	73	1.8
1.5	Green	Male	77	5.8
1.4	Blue	Female	57	-14.2

Average Weight

71.2

+



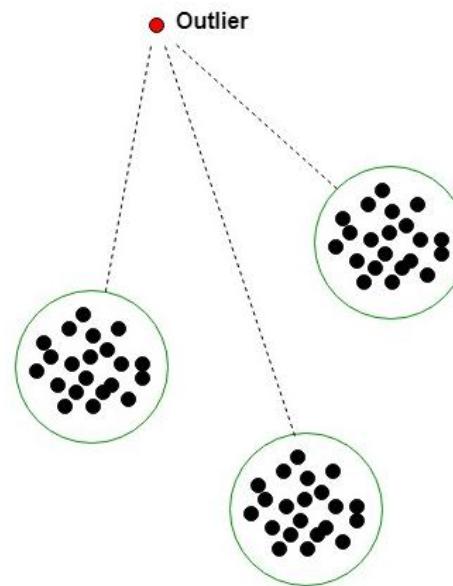
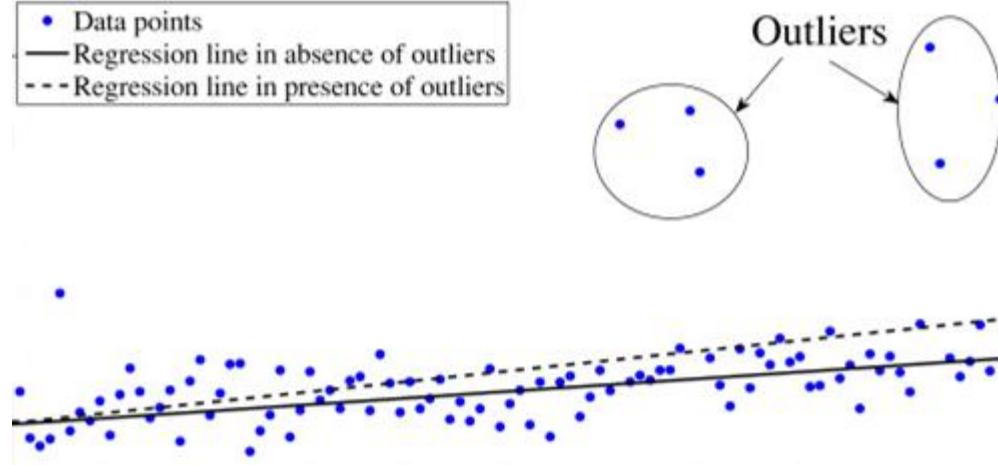
$$\text{Predicted Weight} = 71.2 + 16.8 = 88$$

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88

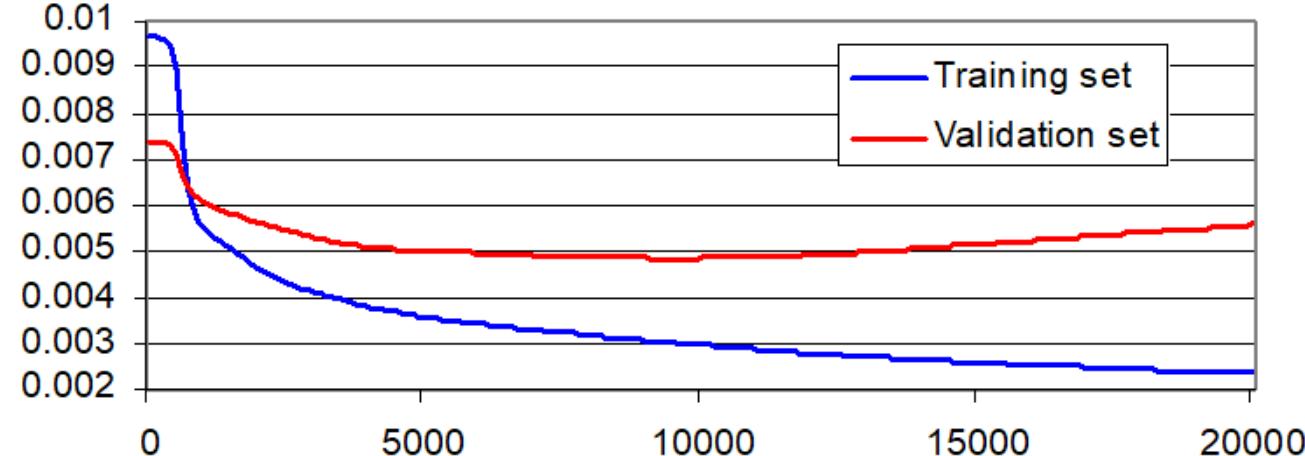
## Methods for estimating generalization errors

- Single-sample statistics
  - Finding error on each data on generalization.
- Split-sample or hold-out validation (60-20-20)
  - Reserve some data as a "test set", which must not be used during training.
  - The test set must represent the cases that the ANN should generalize to.
  - A re-run with the random test set provides an unbiased estimate of the generalization error.
  - Disadvantage is that it reduces the amount of data available for both training and validation.
- Cross-validation (e.g., leave one out)
  - In k-fold cross-validation the data is divided into k subsets, and the network is trained k times, each time leaving one subset out for computing the error.
  - “Crossing” makes an improvement over split-sampling allowing all data used for training.
  - Disadvantage is that the network must be re-trained many times (k times in k-fold crossing).
- Bootstrapping
  - Works on random sub-samples (random shares) chosen from the full data set.
  - Any data item may be selected any number of times for validation.
  - The sub-samples are repeatedly analyzed.

- Data points
- Regression line in absence of outliers
- Regression line in presence of outliers



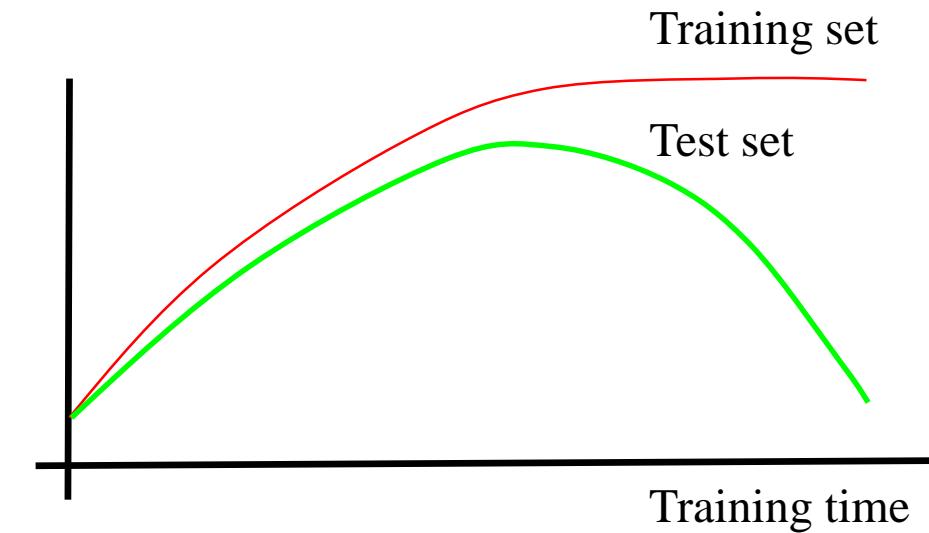
- Due to overtraining (training too much!)
- Overfittings may fit the training data well, even outliers but fail to generalize to new data.
- But loose performance on test sets (new data)



Underfitting  
(high bias)

Number of weight updates

Overfitting  
(high variance)



- A tree that classifies the training data perfectly may not lead to best generalization performance.
  - There may be noise in the training data the tree is fitting.
  - The algorithm might be making decisions based on very little data.

A hypothesis  $h$  is said to overfit the training data

- if there is another hypothesis,  $h'$ , such that  $h$  has smaller error than  $h'$  on the training data but  $h$  has larger error on the test data than  $h'$ .

## Avoiding Overfitting

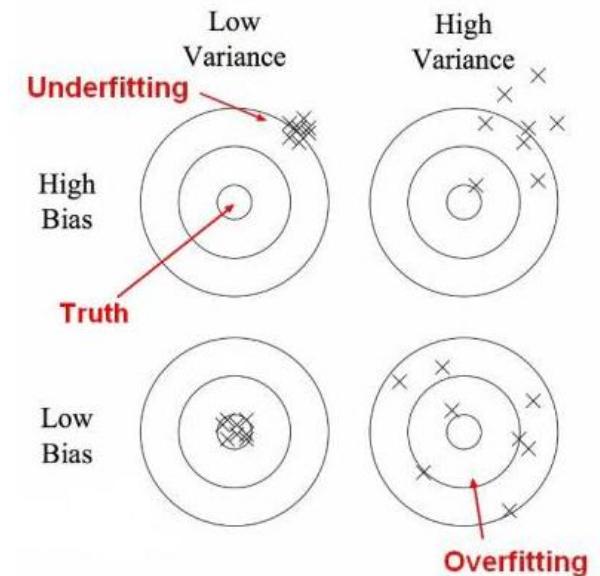
- Two basic approaches
  - Prepruning: Stop growing the tree at some point when there is not enough data.
  - Postpruning: Grow the full tree and then remove nodes not having sufficient evidence.

## Bias

- Bias is how far are the predicted values from the actual values.
- A model having high bias implies that it is too simple and thus underfitting the data.

## Variance

- Occurs when the model performs good on trained dataset but not on test dataset.
- A model having high variance causes overfitting the data.



**Accuracy:** The number of true predictions over total examples.

		ACTUAL CLASS	
		Class=Yes (Positive)	Class>No (Negative)
PREDICTED CLASS	Class=Yes (Positive)	a (TP)	b (FN)
	Class>No (Negative)	c (FP)	d (TN)

$$\text{Accuracy} = \frac{a + d}{a + b + c + d} = \frac{TP + TN}{TP + TN + FP + FN}$$

## Limitation of Accuracy

- Accuracy is **not** a good choice with **unbalanced** classes!
- Consider a 2-class problem
  - Number of Class 0 examples = 9990
  - Number of Class 1 examples = 10
- If model predicts everything to be class 0, accuracy is  $9990/10000 = 99.9\%$ 
  - Accuracy is misleading because model does not detect any class 1 example

**Precision:** Proportion of selected items that are correct.

- Identify only the relevant instances i.e. proportion of relevant instances that actually relevant.

$$\text{Precision} = \frac{TP}{TP+FP}$$

**Recall (Sensitivity):** Proportion of correct items that are selected.

- Identify all the relevant instances in a dataset.

$$\text{Recall} = \frac{TP}{TP+FN}$$

**Note:** Often you have a trade-off between Recall and Precision.

**F measure:** Combined metric of P/R to find optimal blend P/R tradeoff (Weighted Harmonic Mean).

$$F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

**F1 Score:** Balanced F measure with  $\alpha=1/2$  and  $\beta=1$

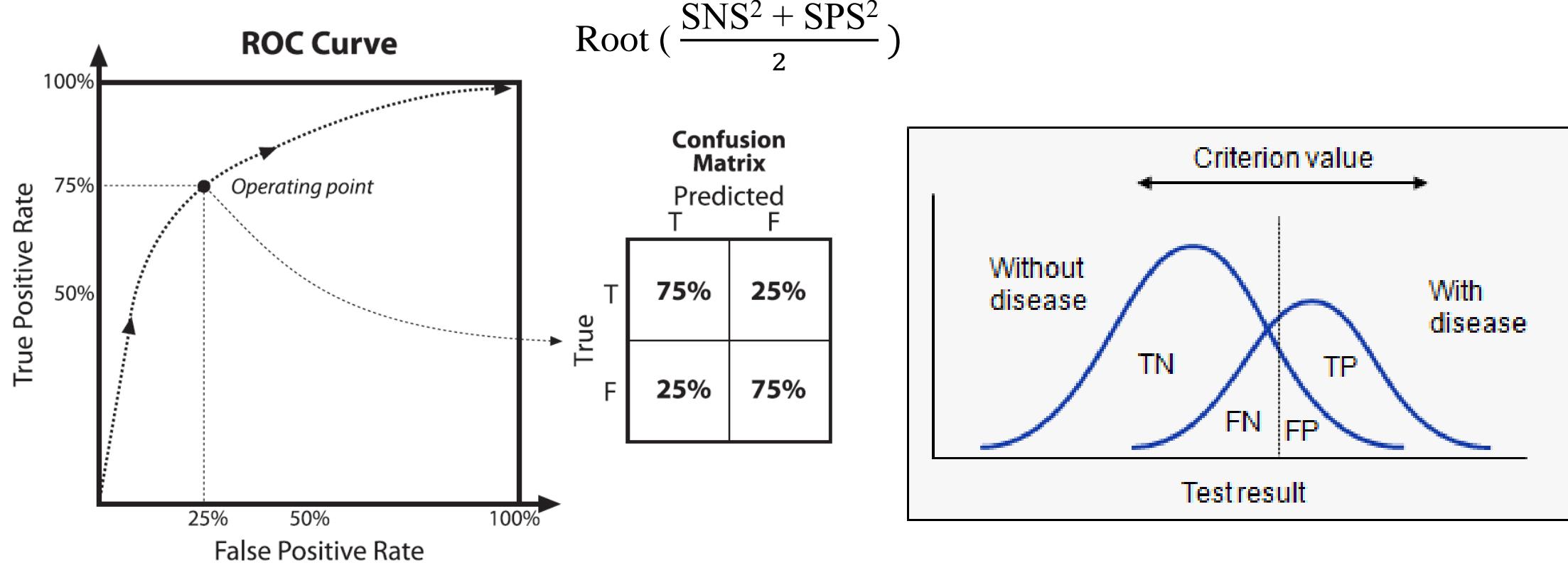
$$\begin{aligned} F_1 &= \left( \frac{\text{recall}^{-1} + \text{precision}^{-1}}{2} \right)^{-1} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}. \\ &= \frac{2TP}{2TP+FN+FP} \end{aligned}$$

**Note:** We use the harmonic mean instead of a simple average because it punishes extreme values.

# Receiver Operating Characteristics

ACDS, CSIR-NEIST

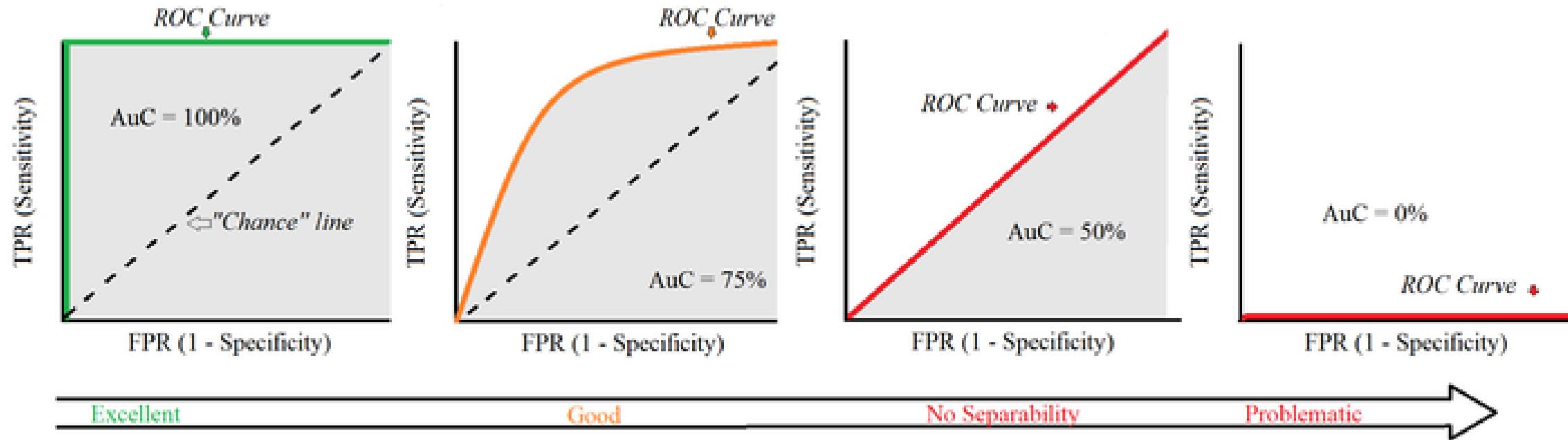
**AUROC (Area Under Receiver Operating Characteristic):** Combined metric on ROC space.



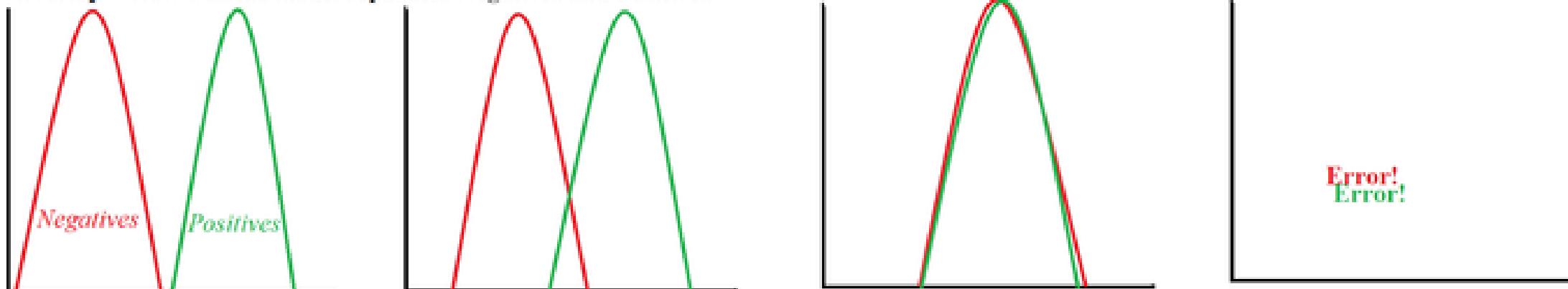
- The ROC graph plots **sensitivity on the y-axis** and **(1-specificity) on the x-axis**.
- We put a point on the graph for each threshold value of the test, plotting the sensitivity and specificity of the test at that value.
- Connecting those points creates a curve - the ROC curve.

# Receiver Operating Characteristics

ACDS, CSIR-NEIST

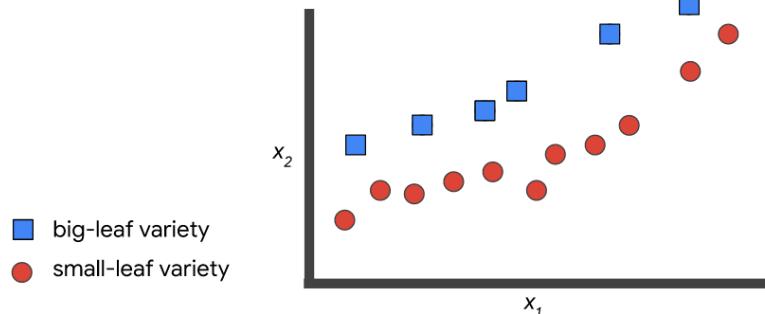


Overlap = How well the model separates Negatives and Positives



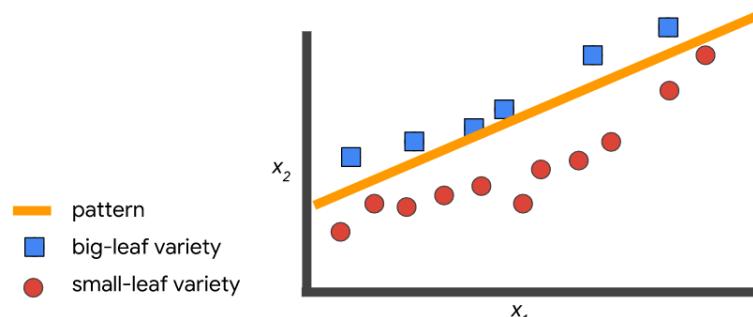
**Question 1.** Suppose you are an amateur botanist determined to differentiate between two species of the Lilliputian plant genus (a completely made-up plant). The two species look pretty similar. Fortunately, a botanist has put together a data set of Lilliputian plants she found in the wild along with their species name. Leaf width and leaf length are the features, while the species is the label. The goal of the data set is to help other botanists answer the question, "Which species is this plant?"

**Solution:**



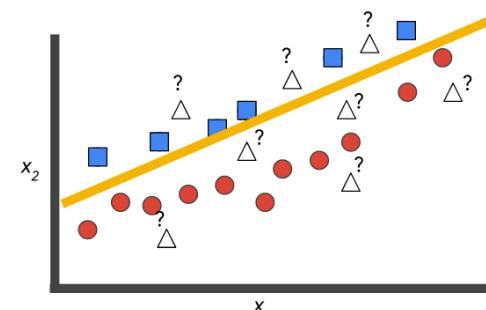
A

Leaf Width	Leaf Length	Species
2.7	4.9	small-leaf
3.2	5.5	big-leaf
2.9	5.1	small-leaf
3.4	6.8	big-leaf



B

- △ unknown
- pattern
- big-leaf variety
- small-leaf variety



C

**Question 2.** A test for a rare disease claims that it will report a positive result for 99.5% of people with the disease, and will report a negative result for 99.9% of those without the disease. We know that the disease is present in the population at 1 in 100,000. Knowing this information, what is the likelihood that an individual who tests positive will actually have the disease?

**Solution:**

This problem provides a simple example of how Bayes' rule can be useful. As before, we begin by establishing our notation,

$P(+\text{test} / +\text{disease})$  --- probability of a positive test result given that the patient has the disease. From the data this probability is 0.995.

$P(-\text{test} / -\text{disease})$  --- probability of a positive test result given that the patient does not have the disease. From the data this probability is 0.999.

$P(+\text{disease})$  --- probability that the patient has the disease given no other information. From the population average, this value is 0.00001.

$P(-\text{disease})$  --- probability that the patient does not have the disease given no other information. Calculated from the population average, this value is 1-0.00001 or 0.99999.

What we want to find is

$P(+\text{disease} / +\text{test}) \rightarrow$  Probability that a patient has the disease and given a positive test result.  
This unknown can be found directly using Bayes' rule

$$P(+\text{disease} / +\text{test}) = [P(+\text{test} / +\text{disease}) \times P(+\text{disease})] / P(+\text{test})$$

To evaluate this expression, we need to know the probability of a positive test, which can be found by marginalizing over all possible disease states (+disease and disease)

$$\begin{aligned} P(+\text{test}) &= P(+\text{test}, \text{disease}) \\ &= P(+\text{test} | \text{disease}) \times P(\text{disease}) \\ &= P(+\text{test} | +\text{disease}) \times P(+\text{disease}) + P(+\text{test} | -\text{disease}) \times P(-\text{disease}) \\ &= P(+\text{test} | +\text{disease}) \times P(+\text{disease}) + (1 - P(-\text{test} | -\text{disease})) \times P(-\text{disease}) \end{aligned}$$

For our given data this works out to be

$$P(+\text{test}) = (0.995) \times (0.00001) + (1 - 0.999) \times (0.99999) = 0.00100994$$

Now, all of the elements are known and can be directly calculated disease states

$$\begin{aligned} P(+\text{disease} | +\text{test}) &= [P(+\text{test} | +\text{disease}) \times P(+\text{disease})] / P(+\text{test}) \\ &= ((0.995) \times (0.00001)) / 0.00100994 \\ &= 0.0099 \end{aligned}$$

Thus, if 100 people get a positive result, only one of these people will really have the disease.

**Question 3.** A researcher hypothesizes that it is possible to detect membrane proteins using the fraction of hydrophobic residues alone. To test this model, the researcher creates a library of 7500 proteins and scores each of these proteins based on their fraction of hydrophobic residues and whether they are membrane proteins. The results of this analysis are shown below. Given this information, we wish to calculate the likelihood that a novel protein that is primarily hydrophobic is also a membrane protein.

**Solution:**

	Majority hydrophobic	Majority hydrophilic
Membrane Bound	2911	961
Cytosolic	713	2915

To solve this problem, we will first summarize our data as a table of all of the possible combinations of possibilities.

	H	NOT H
M	P(H, M)	P(NOT H, M)
NOT M	P(H, NOT M)	P(NOT H, NOT M)

In this table, H identifies hydrophobic proteins and M membrane proteins.

Next we also include the sums of the probabilities in the margins to calculate the marginal probabilities

	H	NOT H	Sum
M	P(H, M)	P(NOT H, M)	P(M)
NOT M	P(H, NOT M)	P(NOT H, NOT M)	P(NOT M)
Sum	P(H)	P(NOT H)	1

The values in this table can be filled directly from the given data. For example,

$$P(H, M) = 2911/7500 = 0.388$$

Note that the sum in the lower right corner must equal one, for it is the sum of all possible outcomes of each variable. Filling in the remaining values, we calculate the following probabilities

	H	NOT H	Sum
M	0.388	0.128	0.516
NOT M	0.095	0.389	0.484
Sum	0.483	0.517	1

By rearranging our definition of conditional probability, we can now answer the question of the likelihood of a novel protein being membrane bound given that it is hydrophobic:

$$P(M | H) = P(H, M) / P(H) = 0.388 / 0.482 = 0.803$$

**Question 4.** Given the following data of transactions at a shop, calculate the support and confidence values of milk → bananas, bananas → milk, milk → chocolate, and chocolate → milk.

**Solution:**

milk → bananas : Support = 2/6, Confidence = 2/4

bananas → milk : Support = 2/6, Confidence = 2/2

milk → chocolate : Support = 3/6, Confidence = 3/4

chocolate → milk : Support = 3/6, Confidence = 3/5

Transaction	Items in basket
1	milk, bananas, chocolate
2	milk, chocolate
3	milk, bananas
4	chocolate
5	chocolate
6	milk, chocolate

$$\text{Support } (S) = \frac{\text{freq}(x, y)}{N}$$

$$\text{Confidence } (C) = \frac{\text{freq}(x, y)}{\text{freq}(x)}$$

Though only half of the people who buy milk buy bananas too, anyone who buys bananas also buys milk.

**Question 5.** Draw a decision tree with the given rules:

R1: IF (age > 38.5) AND (years-in-job > 2.5) THEN  $y = 0.8$

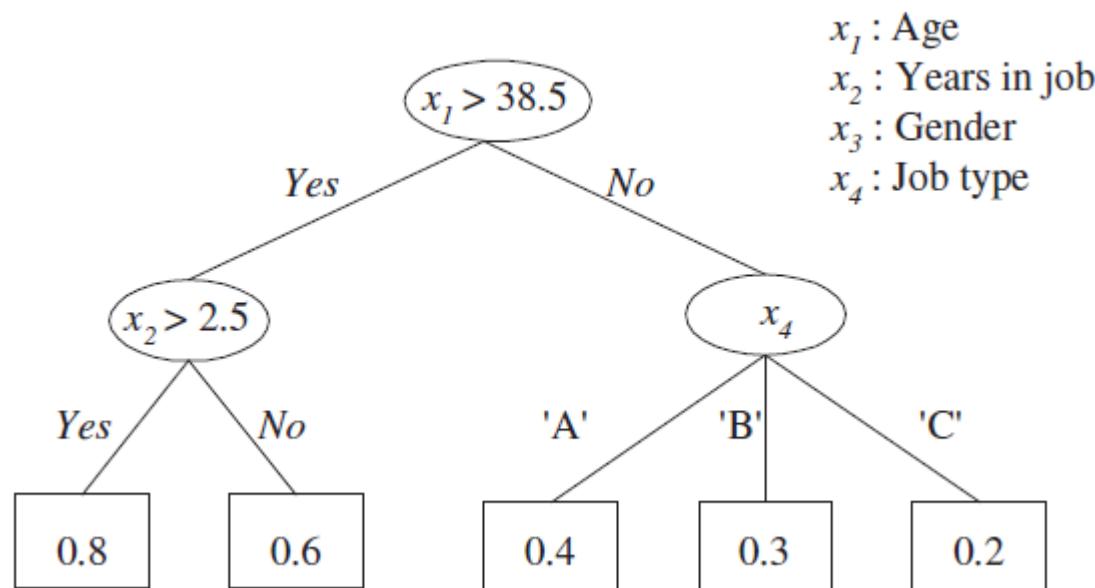
R2: IF (age > 38.5) AND (years-in-job  $\leq 2.5$ ) THEN  $y = 0.6$

R3: IF (age  $\leq 38.5$ ) AND (job-type = 'A') THEN  $y = 0.4$

R4: IF (age  $\leq 38.5$ ) AND (job-type = 'B') THEN  $y = 0.3$

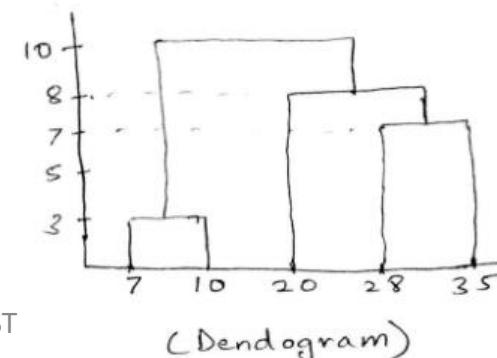
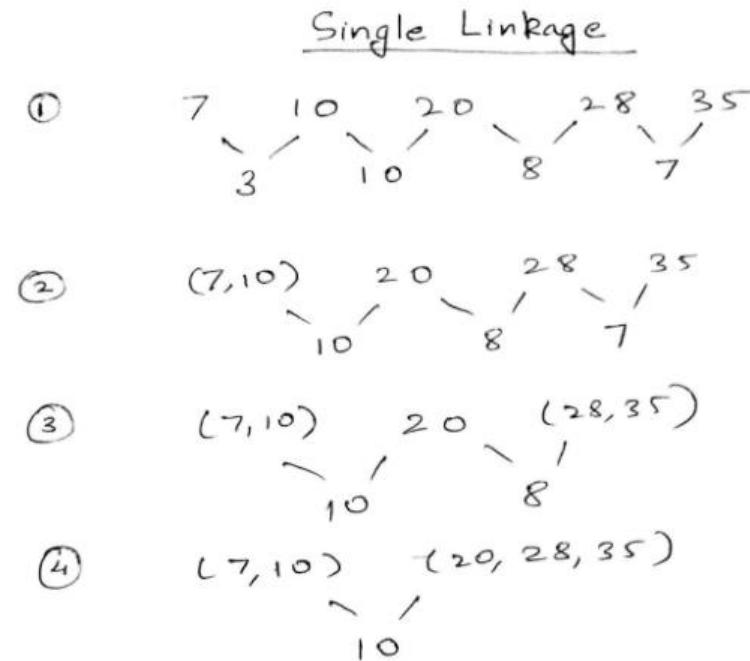
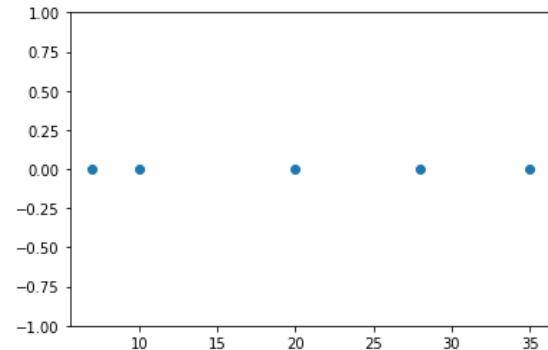
R5: IF (age  $\leq 38.5$ ) AND (job-type = 'C') THEN  $y = 0.2$

**Solution:**



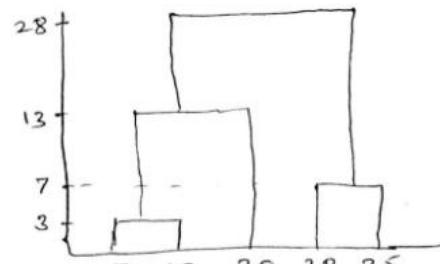
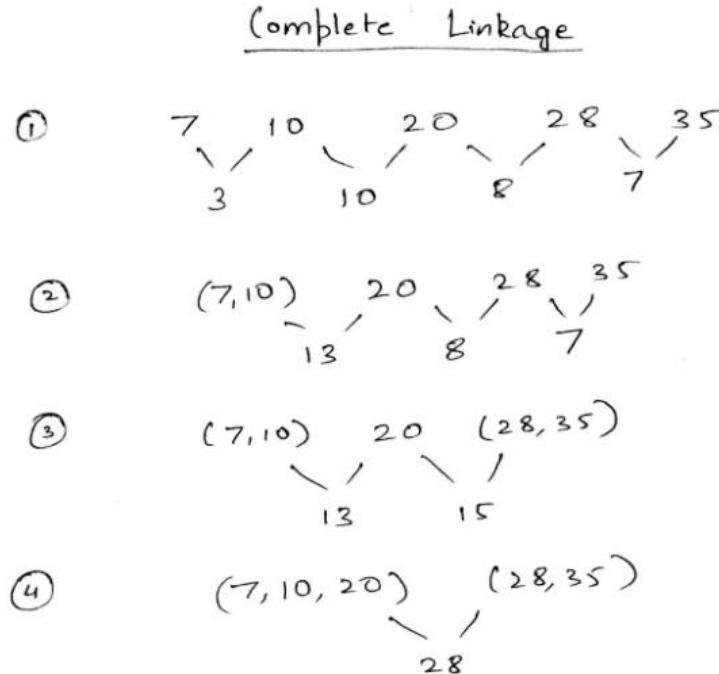
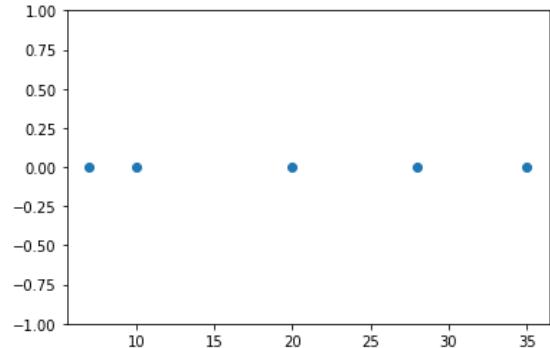
**Question 6.** For the one dimensional data set  $\{7, 10, 20, 28, 35\}$ , perform hierarchical clustering and plot the dendrogram to visualize it using single linkage (*the distance between two clusters is defined as the shortest distance between two points in each cluster*).

**Solution:**



**Question 7.** For the one dimensional data set  $\{7, 10, 20, 28, 35\}$ , perform hierarchical clustering and plot the dendrogram to visualize it using complete linkage (*the distance between two clusters is defined as the largest distance between two points in each cluster*).

**Solution:**



**Question 8.** Consider a n dimensional space of descriptors as shown in the figure.

Find the distance between  $A = (3, 0, 1)$  and  $B = (5, 2, 0)$

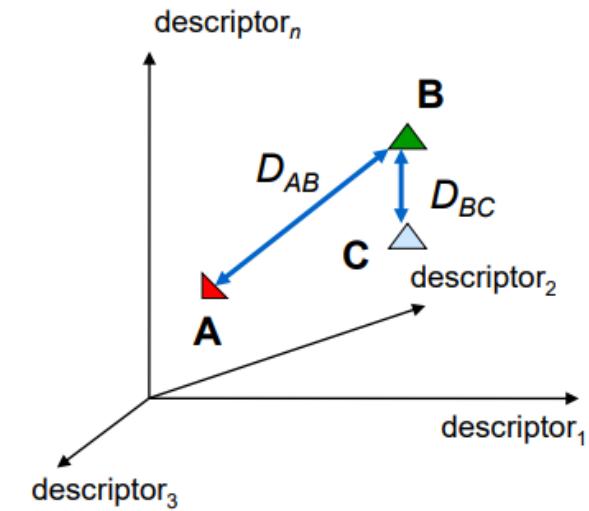
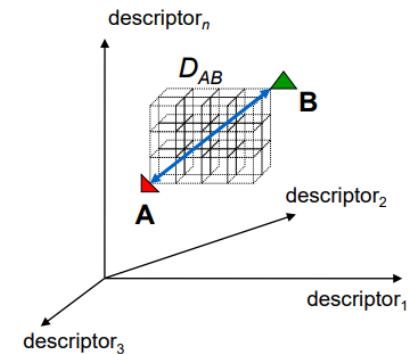
**Solution:** Given two n-dimensional vectors, A and B –

$A = (a_1, a_2, \dots, a_n) - B = (b_1, b_2, \dots, b_n)$ , the Euclidean distance is given by

$$D_{AB} = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$$

$$\mathbf{A} = (3, 0, 1); \mathbf{B} = (5, 2, 0)$$

$$D_{AB} = \sqrt{(3-5)^2 + (0-2)^2 + (1-0)^2} = 3$$



**Question 9.** Four molecules in four-dimensional space - descriptor values listed in the table.

Predict the most similar molecules based on the descriptors.

**Solution:** Using Euclidean distance,

$$D_{AB} = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$$

- $D_{12} = [(2.4250 - 3.4700)^2 + (2 - 2)^2 + (3 - 3)^2 + (2 - 3)^2]^{1/2} = 1.45$
- $D_{13} = [(2.4250 - 0.7090)^2 + (2 - 0)^2 + (3 - 2)^2 + (2 - 1)^2]^{1/2} = 2.99$
- $D_{14} = [(2.4250 - 2.4900)^2 + (2 - 1)^2 + (3 - 2)^2 + (2 - 1)^2]^{1/2} = 1.73$
- $D_{34} = [(0.7090 - 2.4900)^2 + (0 - 1)^2 + (2 - 2)^2 + (1 - 1)^2]^{1/2} = 2.04$

mol	logP(o/w)	b_rotN	a_acc	a_don
1	2.4250	2	3	2
2	3.4700	2	3	3
3	0.7090	0	2	1
4	2.4900	1	2	1

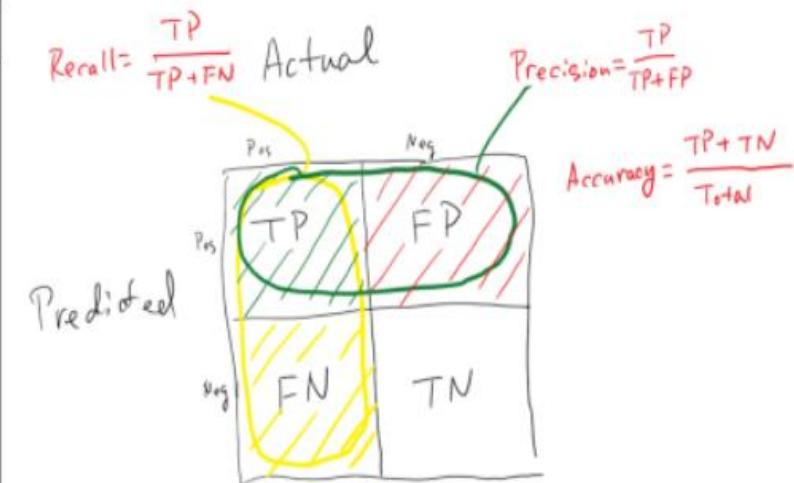
- logP(o/w): index of hydrophobicity
- b\_rotN: number of rotatable bonds
- a\_acc: number of hydrogen acceptors
- a\_don: number of hydrogen donors

**Question 9.** A model has been developed for prediction of  $y$  values. From the given data set for actual value and predicted value of variable  $y$ , build the confusion matrix and find the accuracy, precision and recall of the model.

**Solution:**

$y$	$y_{pred}$	output for threshold 0.6
0	0.5	0
1	0.9	1
0	0.7	1
1	0.7	1
1	0.3	0
0	0.4	0
1	0.5	0

$y$	$y_{pred}$	output for threshold 0.6	Recall	Precision	Accuracy
0	0.5	0	1/2	2/3	4/7
1	0.9	1			
0	0.7	1			
1	0.7	1			
1	0.3	0			
0	0.4	0			
1	0.5	0			



**Question 10.** The confusion matrix for a spam detection model is provided below. Compute the sensitivity and specificity of the model of the model in making the predictions. Also find the precision and recall of the model.

**Solution:**

*Sensitivity* measures how apt the model is to detecting events in the positive class

*Specificity* measures how exact the assignment to the positive class is.

Predicted class POSITIVE (spam ✉ )	Predicted class NEGATIVE (normal 📎 )	
Actual class POSITIVE (spam ✉ )	TRUE POSITIVE (TP) ✉ ✉ 320	FALSE NEGATIVE (FN) ✉ 📎 43
Actual class NEGATIVE (normal 📎 )	FALSE POSITIVE (FP) ✉ ✉ 20	TRUE NEGATIVE (TN) ✉ 📎 538

Predicted class POSITIVE (spam ✉ )	Predicted class NEGATIVE (normal 📎 )	
Actual class POSITIVE (spam ✉ )	TRUE POSITIVE (TP) ✉ ✉ 320	FALSE NEGATIVE (FN) ✉ 📎 43
Actual class NEGATIVE (normal 📎 )	FALSE POSITIVE (FP) ✉ ✉ 20	TRUE NEGATIVE (TN) ✉ 📎 538

Predicted class POSITIVE (spam ✉ )	Predicted class NEGATIVE (normal 📎 )	
Actual class POSITIVE (spam ✉ )	TRUE POSITIVE (TP) ✉ ✉ 320	FALSE NEGATIVE (FN) ✉ 📎 43
Actual class NEGATIVE (normal 📎 )	FALSE POSITIVE (FP) ✉ ✉ 20	TRUE NEGATIVE (TN) ✉ 📎 538

**Question 1.** You want to predict how many people are infected with a contagious virus in times before they show the symptoms, and isolate them from the healthy population. The two values for our target variable would be: Sick and Not Sick. There are 100 observations for which the confusion matrix is given below. Compute the accuracy, precision and recall of the model.

		ACTUAL VALUES	
		POSITIVE	NEGATIVE
PREDICTED VALUES	POSITIVE	560	60
	NEGATIVE	50	330

ID	Actual Sick?	Predicted Sick?	Outcome
1	1	1	1 TP
2	0	0	0 TN
3	0	0	0 TN
4	1	1	1 TP
5	0	0	0 TN
6	0	0	0 TN
7	1	0	0 FP
8	0	1	1 FN
9	0	0	0 TN
10	1	0	0 FP
:	:	:	:
1000	0	0	0 FN

Solution

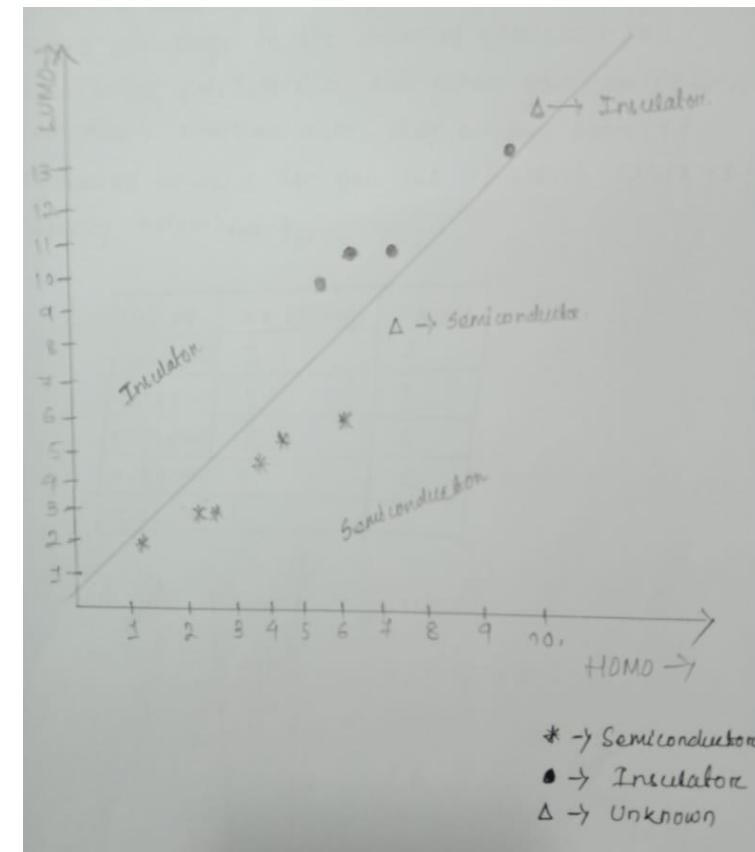
$TP = 560$	$Precision = \frac{TP}{TP+FP}$
$TN = 330$	$= \frac{560}{620} = 0.90$
$FP = 60$	$Recall = \frac{TP}{TP+FN}$
$FN = 50$	$= \frac{560}{610}$

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} = \frac{560+60}{610} = 0.918$$

**Question 2:** Suppose we have data for HOMO & LUMO for several semiconductors & insulators (just for demonstration purposes) and later you are given a material's HOMO-LUMO data, then on the basis of regression model; can you tell to which class it belongs?

HOMO	LUMO	Property
2.3	2.9	Semiconductor
5.4	9.7	Insulator
6.0	5.6	Semiconductor
1.2	1.8	Semiconductor
3.6	4.2	Semiconductor
5.7	10.0	Insulator
4.5	5.1	Semiconductor
8.7	13.0	Insulator
2.1	2.7	Semiconductor
6.7	11.0	Insulator

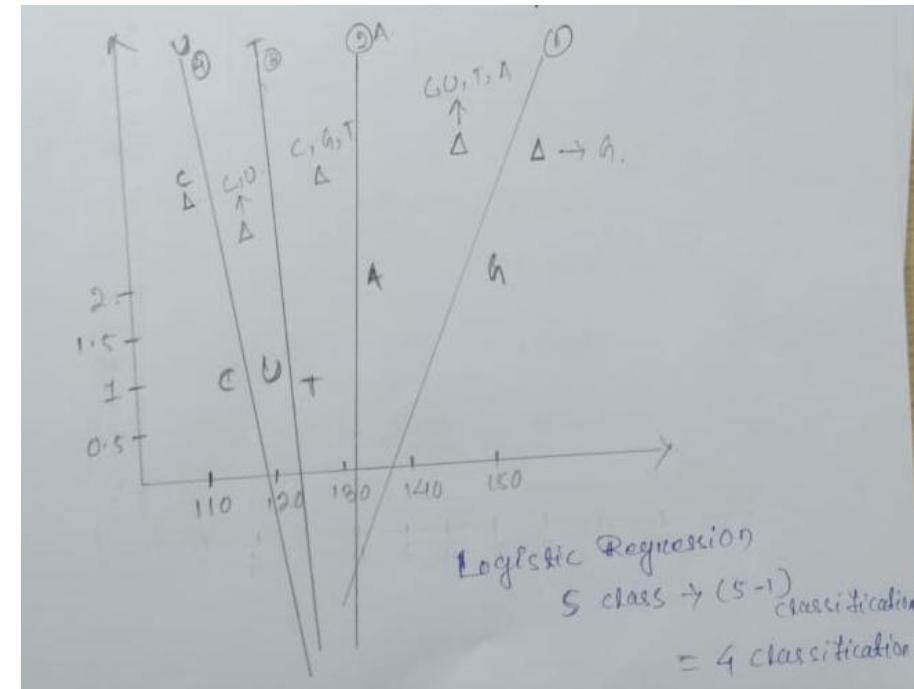
**Solution:**



**Question 3:** Suppose we have data for molecular weight & number of rings in the chemical structure for nucleobases (A, T, G, C, U) and later you are given a structure's similar data, then on the basis of regression model; can you tell to which class it closely resembles?

Mol. Wt.	No. of Rings	Species
135.13 g/mol	2	A
126.11 g/mol	1	T
151.13 g/mol	2	G
111.10 g/mol	1	C
112.08 g/mol	1	U

### Solution:



**Question 4:** The following data for number of students opting for different subject combinations in their final year of graduation; calculate support & confidence values of: P→C; P→M; M→C; B→C; Z→B.

Where, P=Physics, C=Chemistry, M=Mathematics, B=Botany, Z=Zoology

Given that, Support(S) =  $\frac{freq(x,y)}{N}$  & Confidence(C) =  $\frac{freq(x,y)}{freq(x)}$

Solution :-

$$P \rightarrow C : S = \frac{8500}{46000} = 18.47\%$$

$$C = \frac{8500}{18500} = 45.94\%$$

$$M \rightarrow C : S = \frac{9000}{46000} = 19.56\%$$

$$C = \frac{9000}{19000} = 47.36\%$$

$$P \rightarrow M : S = \frac{12000}{46000} = 26.08\%$$

$$C = \frac{12000}{18500} = 64.66\%$$

$$B \rightarrow C : S = \frac{6500}{46000} = 14.13\%$$

$$C = \frac{6500}{14500} = 44.82\%$$

$$Z \rightarrow B : S = \frac{9500}{46000} = 20.65$$

$$C = \frac{9500}{15500} = 61.29\%$$

Number (x10 <sup>3</sup> )	Combination
2	PCM
1.5	ZBC
10	PM
7	MC
6.5	PC
8	ZB
6	ZC
5	BC

**Question 5 :** The following data for number of papers published on interaction studies with Graphene with different dopants; calculate support & confidence values of: Gr→O; Gr→Al; Gr→B; Gr→ Si; Gr→N.

Where, Gr=Graphene; O=oxygen; B=boron; Al=Aluminum; Si=Silicon; Ge=Germanium; N=Nitrogen

Given that, Support(S) =  $\frac{freq(x,y)}{N}$  & Confidence(C) =  $\frac{freq(x,y)}{freq(x)}$

System	No. of Papers
Gr+O	125
Gr+B	251
Gr+Al	165
Gr+Si	115
Gr+Ge	151
Gr+N	100

Solution-5:

Gr → O :  $S = \frac{276}{907} = 30.42\%$   
 $C = \frac{276}{505} = 54.65\%$

Gr → Al :  $S = \frac{165}{907} = 18.19\%$   
 $C = \frac{165}{505} = 32.67\%$

Gr → Si :  $S = \frac{115}{907} = 12.67\%$   
 $C = \frac{115}{505} = 22.77\%$

Gr → N :  $S = \frac{351}{907} = 38.69\%$   
 $C = \frac{351}{505} = 69.50\%$

**Question 6 :** Draw a decision tree with the given data:

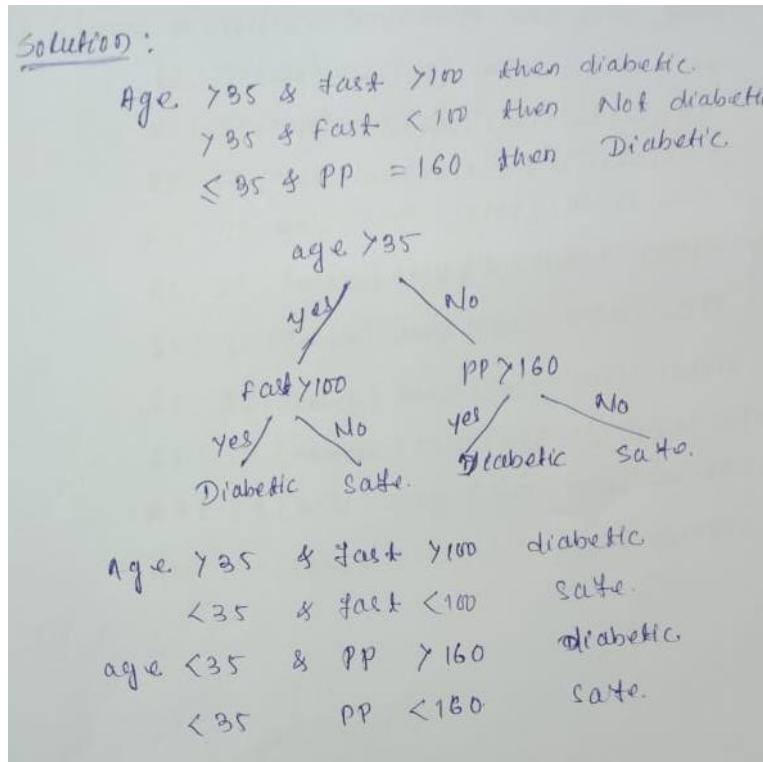
R1: IF( $\text{age} \geq 25$ ) AND ( $\text{fasting} \geq 110\text{mg/dL}$ ) THEN  $y = \text{diabetic}$

R2: IF( $\text{age} \geq 25$ ) AND ( $\text{PP} \geq 200\text{mg/dL}$ ) THEN  $y = \text{diabetic}$

R3: IF( $\text{age} \geq 35$ ) AND ( $\text{fasting} \geq 100\text{mg/dL}$ ) THEN  $y = \text{diabetic}$

R4: IF( $\text{age} \geq 35$ ) AND ( $\text{PP} \geq 160\text{mg/dL}$ ) THEN  $y = \text{diabetic}$

R5: IF( $\text{age} \geq 35$ ) AND ( $\text{RBS} \geq 160\sim 200\text{mg/dL}$ ) THEN  $y = \text{diabetic}$



**Question 7:** A ML model claims that it can positively report accurate predictions for 99.7% and 97.4% negatively. We know that 2 out of 10000 predictions made by the model are not accurate. What is the likelihood that an input file will achieve convergence using the new software?

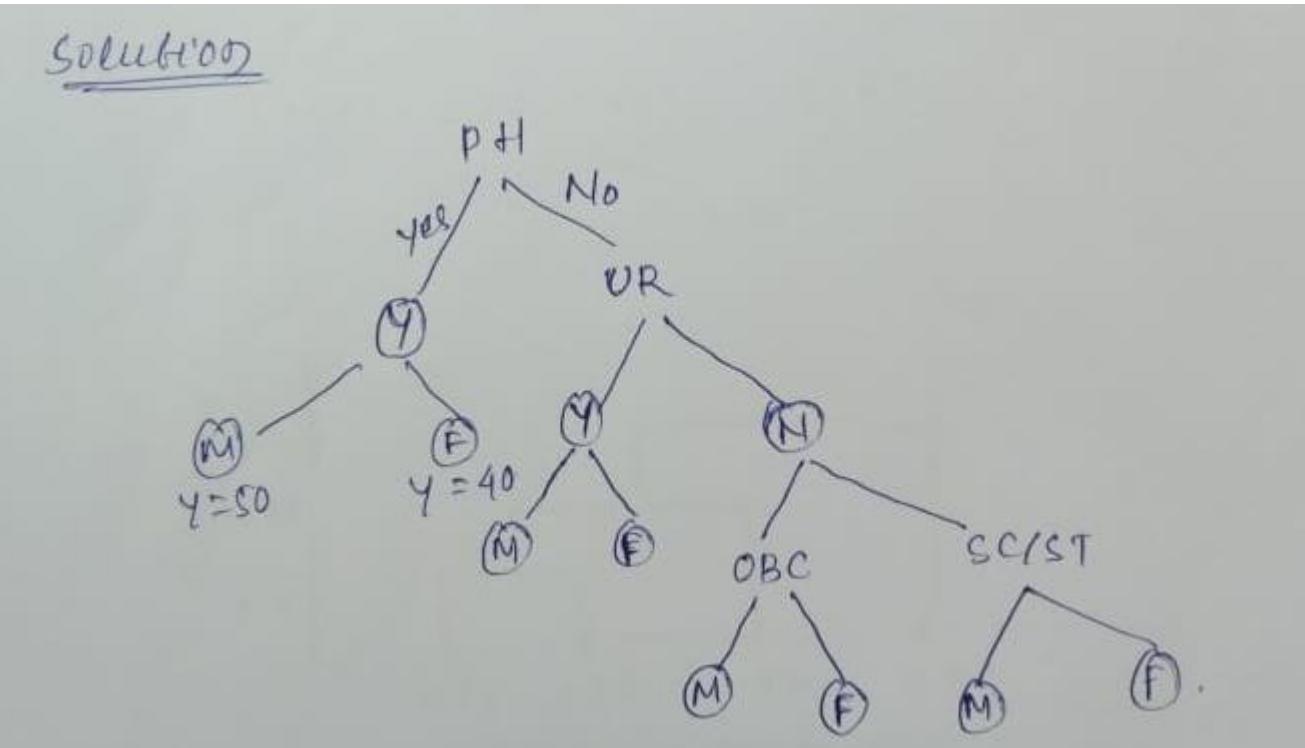
Solution

$$\begin{aligned} P(+|a) &= 0.997 \\ P(-|a) &= 0.003 \\ P(+| \neg a) &= 0.974 \\ P(-| \neg a) &= 0.026 \\ P(+t) &= P(+t, a) = P(+t|a) \cdot P(a) \\ &= P(+t|a) \cdot P(a) + P(+t|\neg a) \cdot P(\neg a) \\ &= P(+t|a) \cdot P(a) + \{1 - P(-t|\neg a)\} \\ &\quad \cdot P(\neg a) \\ &= 0.997 \times 0.00002 + 0.026 \times 0.99998 \\ &= 0.0351 \end{aligned}$$
$$\begin{aligned} P(+a|+t) &= \frac{P(+t|+a) \cdot P(+a)}{P(+t)} \\ &= \frac{0.997 \times 0.00002}{0.0351} \\ &= \frac{0.0001999}{0.0351} \\ &= 0.0056 \\ &= 0.56\% \quad (\text{Ans}) \end{aligned}$$

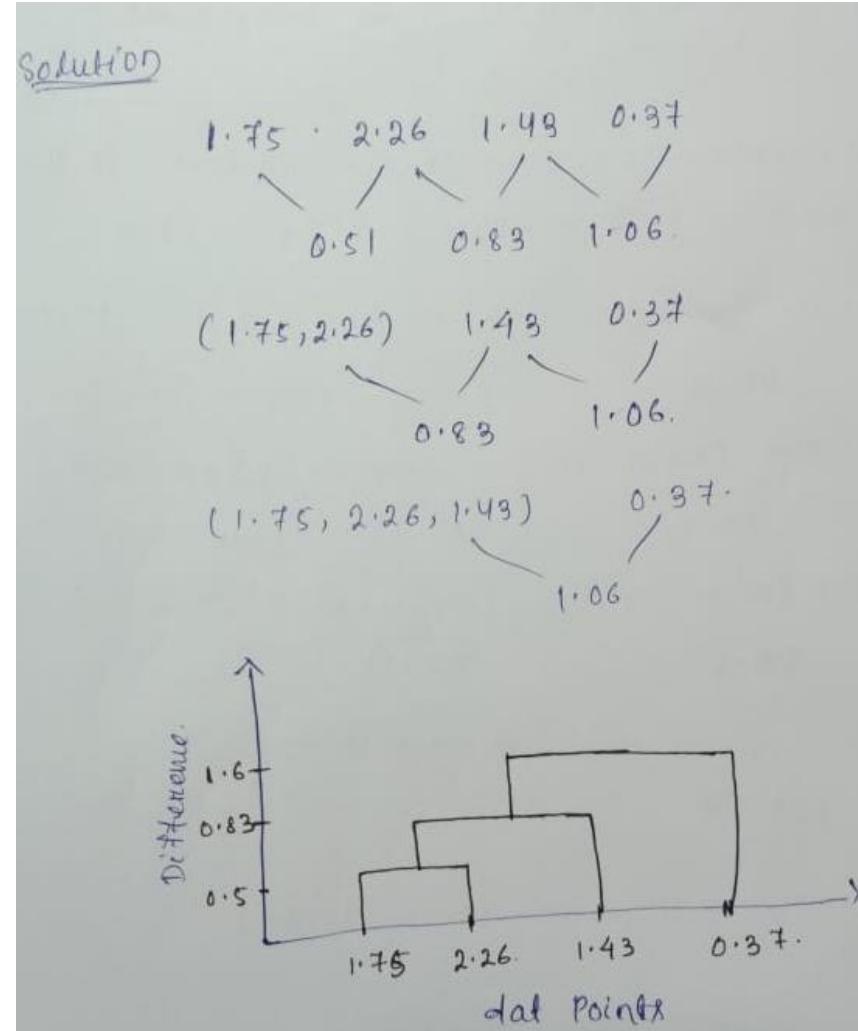
---

**Question 8:** Draw a decision tree with the given data:

- R1: IF(male) AND (UR) THEN cutoff=95
- R2: IF(male) AND (OBC) THEN cutoff=85
- R3: IF(male) AND (EWS) THEN cutoff=85
- R4: IF(male) AND (SC/ST) THEN cutoff=70
- R5: IF(female) AND (UR) THEN cutoff=90
- R6: IF(female) AND (OBC) THEN cutoff=80
- R7: IF(female) AND (EWS) THEN cutoff=80
- R8: IF(female) AND (SC/ST) THEN cutoff=65
- R9: IF(male) AND (PH) THEN cutoff=50
- R10: IF(female) AND (PH) THEN cutoff=40



**Question 9:** For the following one dimensional dataset for the Band Gap of photovoltaic materials { $\text{Cu}_3\text{N}=1.75$ ,  $\text{GaP}=2.26$ ,  $\text{GaAs}=1.43$ ,  $\text{PbS}=0.37$ }, perform hierarchical clustering and plot the dendrogram to visualize it using single linkage.



**Question 10:** For the following one dimensional dataset for the best docked posed energy for a protein ligand system for SARS-COV19 {-11.75, -12.26, -13.43, -15.37, -15.44, -18.42, 13.42, 10.56}, perform hierarchical clustering and plot the dendrogram to visualize it using single linkage.

Solution

-11.75    -12.26    -13.43    -15.37    -15.44    -18.42    10.56.  
 0.51        1.17        1.94        (0.07)        2.02        2.86.

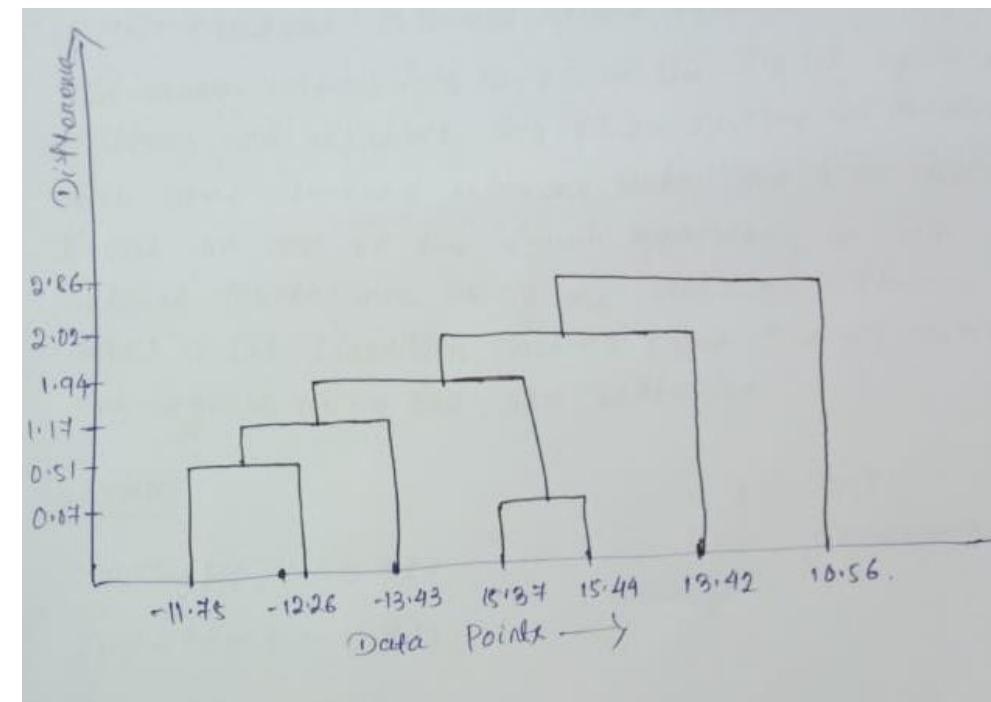
-11.75    -12.26    -13.43    (-15.37, -15.44)    -18.42    -10.56.  
 (0.51)    1.17    1.94                  2.02        2.86.

- (-11.75, -12.26)    -13.43    (-15.37, -15.44)    -18.42    -10.56.  
 (1.17)    1.94                  2.02        2.86.

- (-11.75, -12.26, -13.43)    (-15.37, -15.44)    -18.42    -10.56  
 (1.94)                  2.02        2.86.

- (-11.75, -12.26, -13.43, -15.37, -15.44)    -18.42    -10.56.  
 \* (2.02)                  2.86

- (-11.75, -12.26, -13.43, -15.37, -15.44, -18.42)    -10.56.  
 (2.86)



**Question 11:** A DFT based new software claims that it can report SCF energy values for 98.7% of the input systems positively and negatively for 99.2% systems on comparison with some standard software data. We know that only 5 out of 1000 of the input geometry do not achieve convergence using the standard software. What is the likelihood that an input file will achieve convergence using the new software?

Solution

$P(+t|+c) = 0.987$   
 $P(-t|-c) = 0.992$ .  
 $P(+c) = 0.005$   
 $P(-c) = 0.995$   
 $P(+t) = P(+t, c)$   
 $= P(+t|c) \cdot P(c)$   
 $= P(+t|+c) \cdot P(+c) + P(+t|-c) \cdot P(-c)$   
 $= P(+t|+c) \cdot P(+c) + [1 - P(+t|-c)] \cdot P(-c)$   
 $= 0.987 \times 0.005 + 0.008 \times 0.995$   
 $= 0.012$ .  
 ~~$P(+c|+t) = \frac{0.004}{0.012} = 0.33 = 33\%$ .~~  
 $P(+c|+t) = \frac{P(+t|+c) \cdot P(+c)}{P(+t)}$   
 $= \frac{0.987 \times 0.005}{0.012} = 0.33 = 33\%$ .

- Mitchell, T.M. (1997). *Machine Learning*. McGraw Hill.
  - Covers the field of machine learning, which is the study of algorithms that allow computer programs to automatically improve through experience. ★★★★
- Bishop, C.M. (2006). *Pattern Recognition and Machine Learning*. Springer.
  - This book reflects these developments while providing a grounding in the basic concepts of pattern recognition and machine learning. ★★★★
- Hastie, T., Tibshirani, R. Friedman, J. (2009). *The Elements of Statistical Learning*. Springer.
  - Describes the important ideas in a variety of fields such as medicine, biology, finance, and marketing in a common conceptual framework. While the approach is statistical, the emphasis is on concepts rather than mathematics. ★★★
- James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013). *An Introduction to Statistical Learning*. Springer.
  - Provides tools for Statistical Learning that are essential for practitioners in science, industry and other fields. Analyses and methods are presented in R. Covers topics like linear regression, classification, resampling methods, shrinkage approaches, tree-based methods, support vector machines, and clustering. ★★★

- Zaki, M.J., and Jr., W.B. (2014). *Data Mining and Analysis: Fundamental Concepts and Algorithm*. Cambridge University Press.
- Fundamental algorithms in data mining, automated methods to analyze patterns and models for all kinds of data, with applications ranging from scientific discovery to business intelligence and analytics are covered. Provides a broad yet in-depth overview of data mining, integrating related concepts from machine learning and statistics including exploratory data analysis, pattern mining, clustering, and classification. ★★★★
- Tang, P.N., Steibach, M., and Kumar, V. (2016). *Introduction to Data Mining*. Pearson.
- This book presents fundamental concepts and algorithms for those learning data mining for the first time. Each concept is explored thoroughly and supported with numerous examples. ★★★
- Shwartz, S.S., and Davis, S.B. (2014). *Understanding Machine Learning: From Theory to Algorithm*. Cambridge University Press.
- Introduces machine learning and provides a theoretical account of the fundamentals underlying machine learning and the mathematical derivations that transform these principles into practical algorithms. ★★★★

- Machine Learning | Stanford University

[https://www.youtube.com/watch?v=PPLop4L2eGk&list=PLLssT5z\\_DsKh9vYZkQkYNWcItqhlRJLN](https://www.youtube.com/watch?v=PPLop4L2eGk&list=PLLssT5z_DsKh9vYZkQkYNWcItqhlRJLN)

- Introduction to Machine Learning | University of Toronto

<https://www.youtube.com/watch?v=FvAibtlARQ8&list=PL-Mfq5QSs8iS9XqKuApPE1TSlnZblFHF>

- Machine Learning for Intelligent Systems | Cornell

<https://www.youtube.com/watch?v=MrLPzBxG95I&list=PLl8OHZGYOQ7bkVbuRthEsaLr7bONzbXS>

- Machine Learning Course | Caltech

<https://www.youtube.com/watch?v=mbyG85GZ0PI&list=PLD63A284B7615313A>

# Thank You