



UNSUPERVISED LEARNING AND REINFORCEMENT LEARNING

Introduction to Machine Learning, Artificial Intelligence and Deep Learning

Session 06

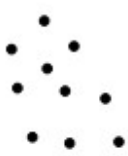
Pramod Sharma
pramod.sharma@prasami.com

2

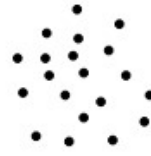
Unsupervised Learning

3

Clusters



Two Clusters



Single Cluster



Not Sure

Unsupervised learning uses procedures that attempt to find natural partitions of patterns.

7/2/2023

pra-sâmi

4

What is Clustering Good for

- ❑ **Market segmentation** - group customers into different market segments
- ❑ **Social network analysis** - Facebook "smartlists"
- ❑ **Organizing computer clusters** and data centers for network layout and location
- ❑ **Astronomical data analysis** - Understanding galaxy formation

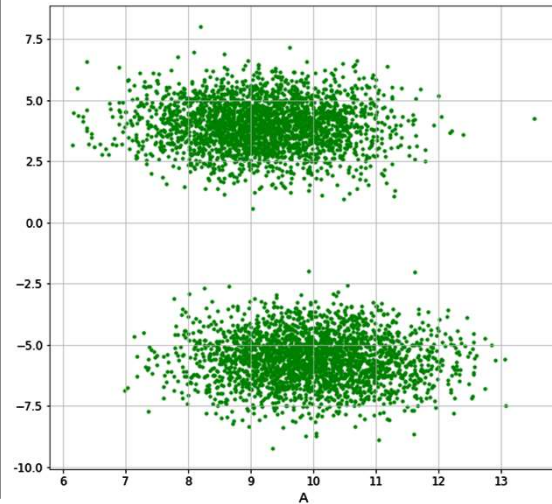
7/2/2023

pra-sâmi

5

K-means Algorithm

- ❑ Randomly allocate two points as the cluster centroids
- ❑ There can be as many cluster centroids as needed (K cluster centroids, in fact)
- ❑ Consider data as shown



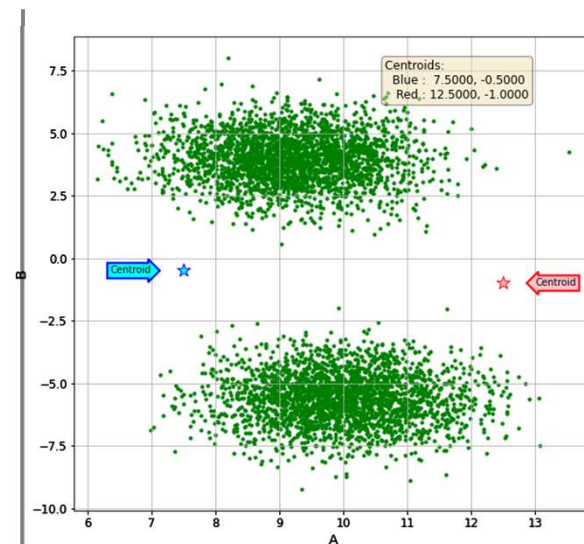
7/2/2023

pra-sâmi

6

Iterations

- ❑ Our data has two groups
- ❑ Consider centroids for two groups arbitrarily
- ❑ Calculate Euclidean distance of each point from the two centroid.
- ❑ Allocate the point to centroid with minimum distance



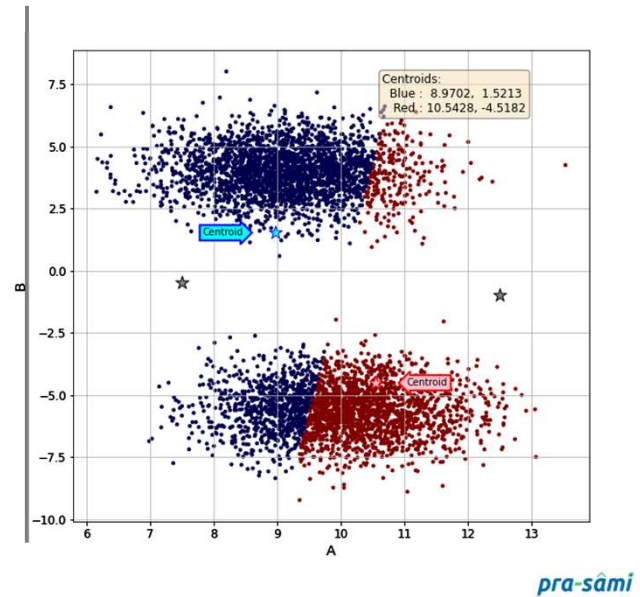
7/2/2023

pra-sâmi

7

Iterations

- ❑ After allotment, calculate coordinates of centroid of each of the group
- ❑ Move the centroid to new location.

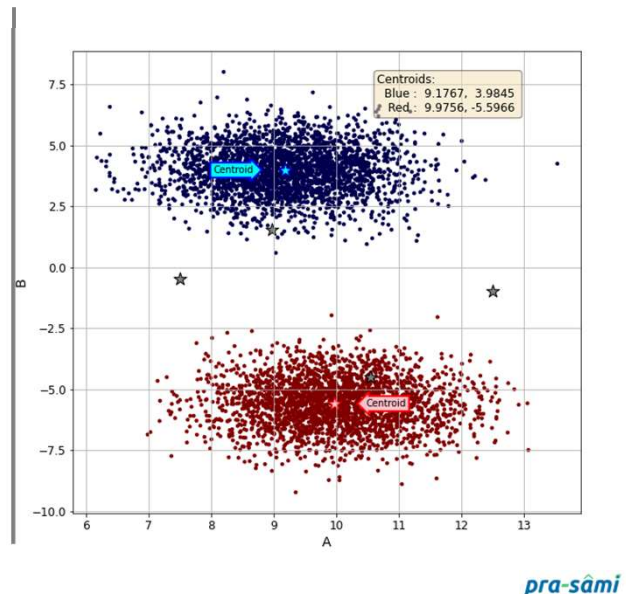


7/2/2023

8

Iterations

- ❑ Once again calculate Euclidean distance of each point from the two centroid
- ❑ Allot the points to closest centroid
- ❑ And we have identified the clusters

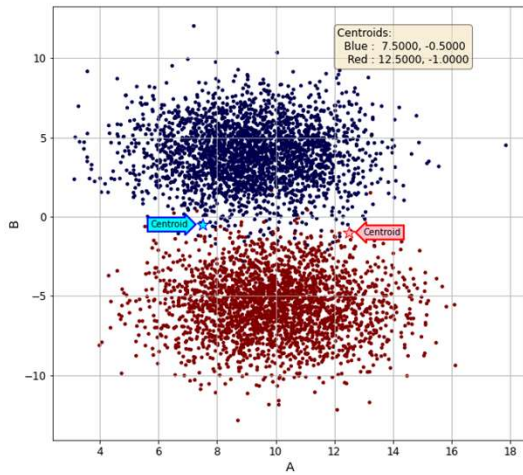


7/2/2023

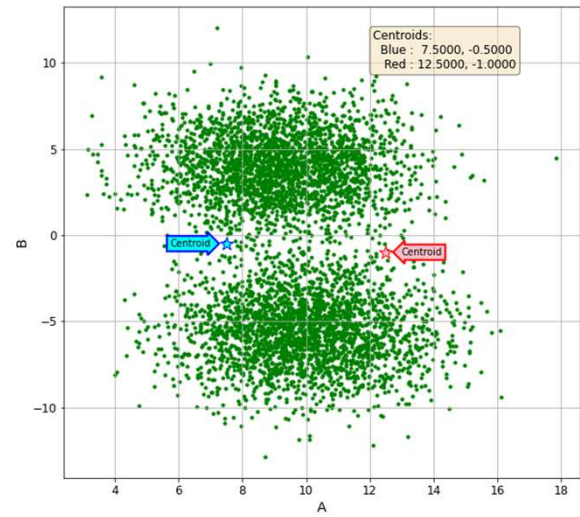
9

Iterations

- Another example with not so obvious separation



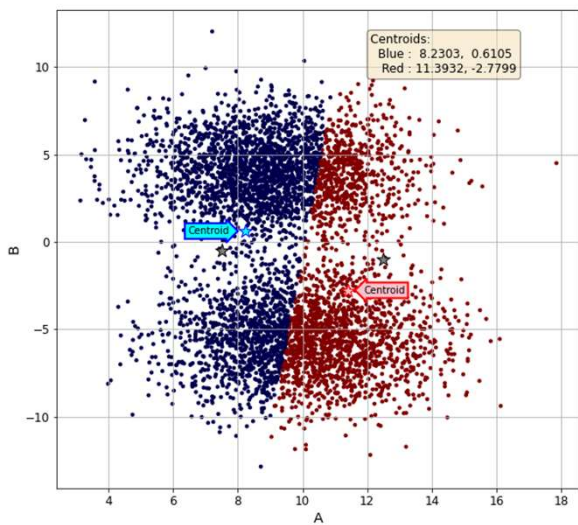
7/2/2023



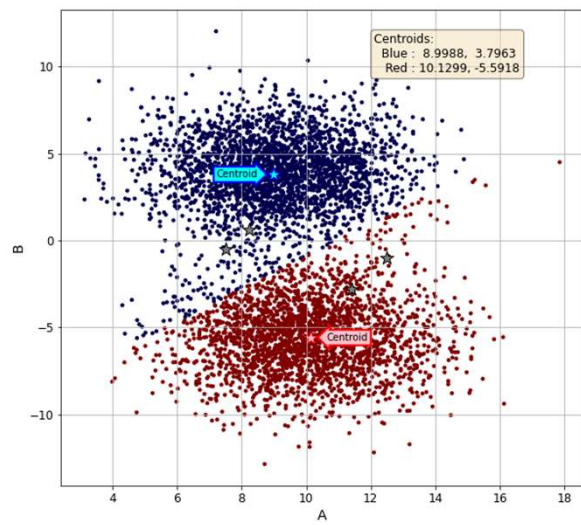
pra-sâmi

10

Iteration 1 and 2



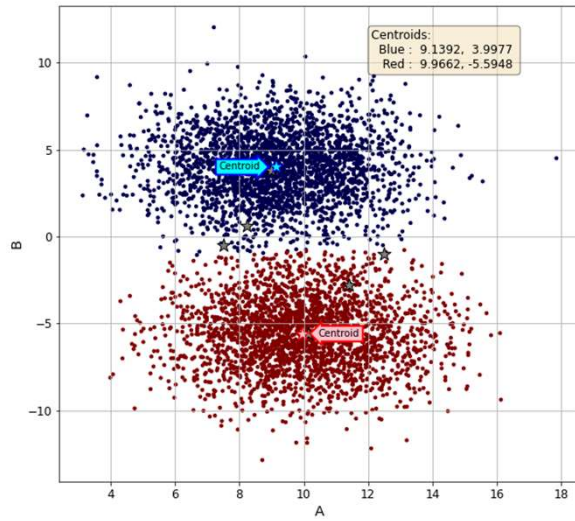
7/2/2023



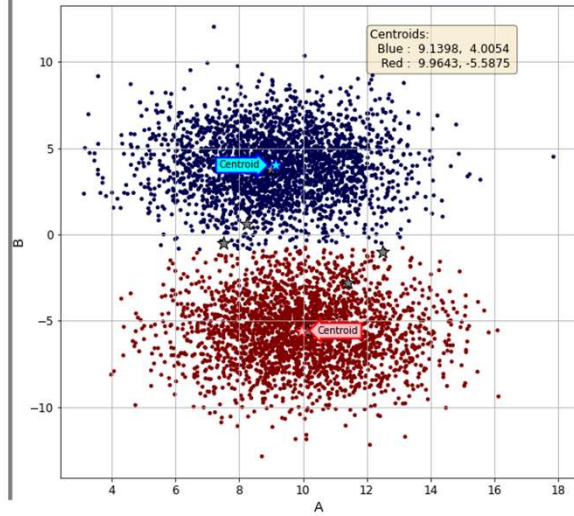
pra-sâmi

11

Iteration 3 and 4



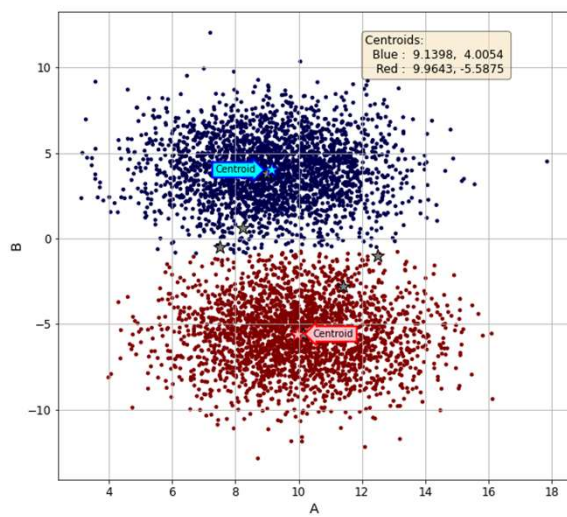
7/2/2023



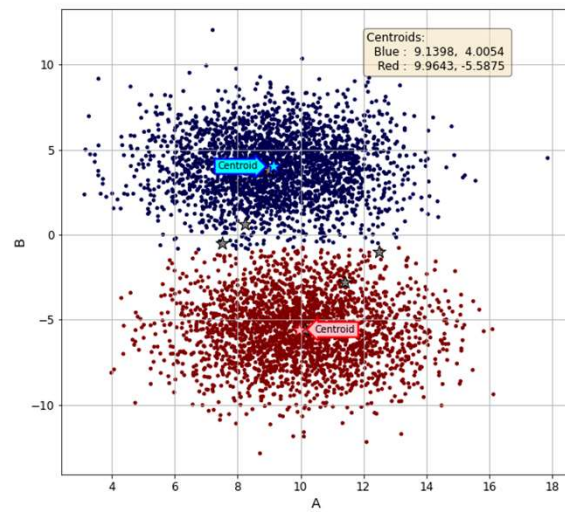
pra-sâmi

12

Iteration s 4 and 5



7/2/2023



pra-sâmi

13

Clustering Distance Measures

7/2/2023

pra-sâmi

14

Dissimilarity or Distance matrix

- ❑ Classification of observations into groups requires
 - ❖ Some methods for computing the **distance** or
 - ❖ The (dis)**similarity** between each pair of observations
- ❑ Distance measures is a critical step in clustering.
- ❑ Wide varieties methods available
- ❑ Method of calculations of the similarity of two elements (x, y) influences the shape of the clusters

7/2/2023

pra-sâmi

15

Methods for Distance Measures

- ❑ The classical methods for distance measures :
 - ❖ Euclidean Distances
 - ❖ Manhattan Distances
 - ❖ Minkowski Distances
 - ❖ Hamming Distance
- ❑ Others:
 - ❖ Pearson correlation distance
 - ❖ Eisen cosine correlation distance (Eisen et al., 1998)
 - ❖ Spearman correlation distance
 - ❖ Kendall correlation distance
- ❑ Pearson correlation analysis is the most commonly used method.
 - ❖ Also known as a parametric correlation which depends on the distribution of the data
- ❑ Kendall and Spearman correlations are non-parametric
 - ❖ Used to perform rank-based correlation analysis

7/2/2023

pra-sâmi

16

Euclidean distance

- ❑ Most common method of distance measurement
- ❑ The distance between two real-valued vectors : $x, y \in \mathbb{R}$
- ❑ Used for calculating the distance between numerical values – Floats or integers.

$$d_{euc}(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

- ❑ If columns have values with differing scales:
 - ❖ normalize or standardize the numerical values to prevent one column prevail over other

7/2/2023

pra-sâmi

17

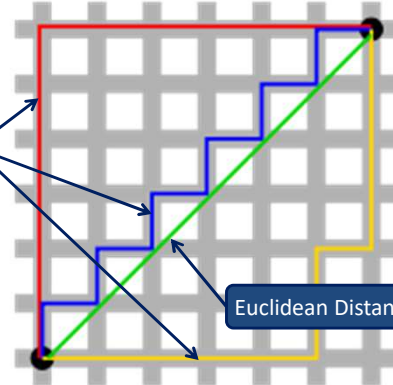
Manhattan Distance

- Names allude to the grid layout of most streets on the island of Manhattan
- Causes the shortest path a car could take between two intersections to have length equal to the intersections' distance
- Manhattan Distance is calculated as follows:

$$d_m(x, y) = \sum_{i=1}^n |x_i - y_i|$$

Manhattan Distance

Euclidean Distance



7/2/2023

pra-sâmi

18

Minkowski Distances

- Both Euclidean and Manhattan Distances are special case of Minkowski Distance
- The formula for Minkowski Distance is given as:

$$d_m(x, y) = (\sum_{i=1}^n (x_i - y_i)^p)^{1/p}$$

- ❖ For $p = 1$ and $p = 2$, it becomes Manhattan and Euclidean Distances

7/2/2023

pra-sâmi

19

Hamming Distance

- ❑ Measures the similarity between two strings of the same length.
- ❑ The number of positions at which the corresponding characters are different

7/2/2023

pra-sâmi

20

Pearson correlation distance

- ❑ A correlation-based distances
- ❑ Widely used for gene expression data analyses
- ❑ Correlation-based distance is defined by subtracting the correlation coefficient from 1
- ❑ Pearson correlation measures the degree of a linear relationship between two profiles
- ❑ It is calculated as follows:

$$d_{\text{cor}}(x, y) = 1 - \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}}$$

7/2/2023

pra-sâmi

21

Clustering Examples

7/2/2023

pra-sâmi

22

News Clustering

The screenshot displays the Google News interface with a focus on COVID-19 news. The left sidebar shows navigation options: Top stories, For you, Following, Saved searches, and a highlighted COVID-19 section. Below this, regional and topic filters are listed: India, World, Your local news, Business, Technology, Entertainment, and Sports. The main content area features a search bar and a 'Local news' section with 'Suggested locations' (Pune, New Delhi, Mumbai, Bengaluru). Four news articles are displayed, each with a title, source, and time: 'Coronavirus in Pune: UK returnee found to have tested positive for COVID-19' (Free Press Journal, 2 hours ago), 'Pune district reports 618 fresh cases, five Covid deaths' (Hindustan Times, 18 hours ago), '3,431 new coronavirus cases in Maharashtra, 71 deaths' (Times of India, 11 hours ago), and 'Pune will have night curfew till Jan 5 amid scare of UK COVID variant | Pune NYOOOZ' (NYOOOZ, 3 days ago). Each article is accompanied by a small thumbnail image.

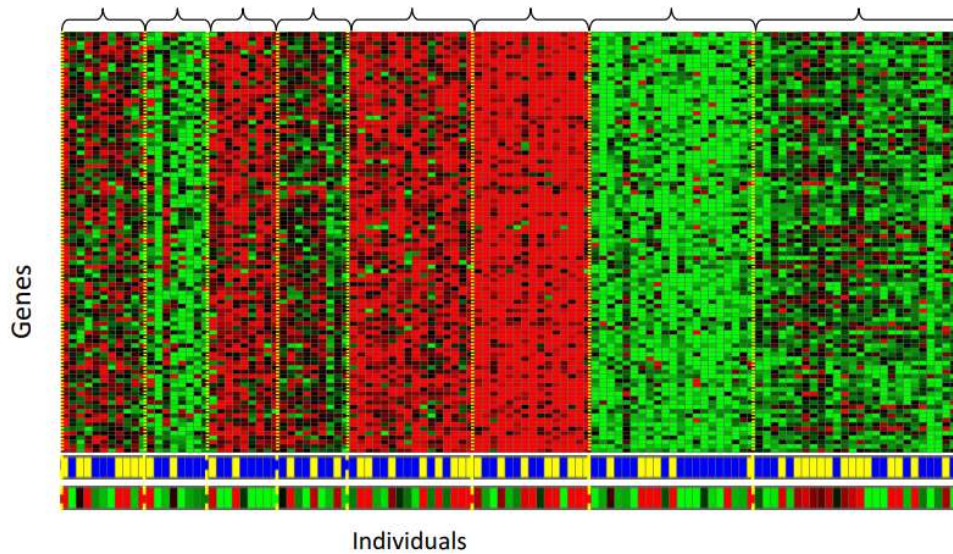
7/2/2023

pra-sâmi

23

Unsupervised Learning

- Genomics application: group individuals by genetic similarity



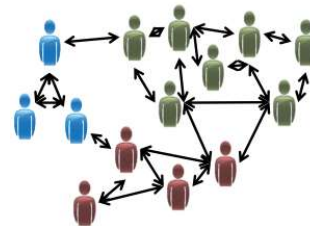
pra-sâmi

24

Unsupervised Learning



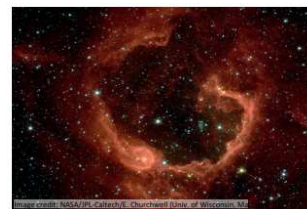
Organize computing clusters



Social network analysis



Market segmentation



Astronomical data analysis

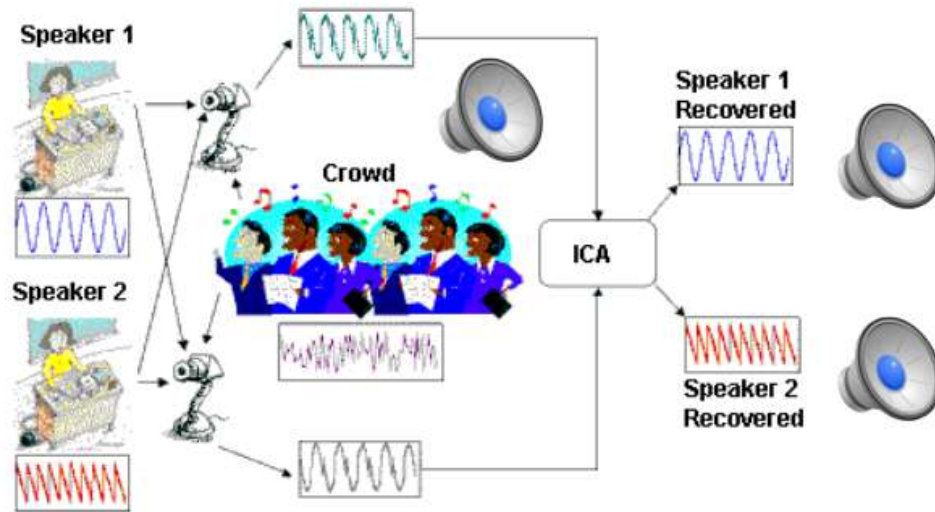
7/2/2023

pra-sâmi

25

Unsupervised Learning

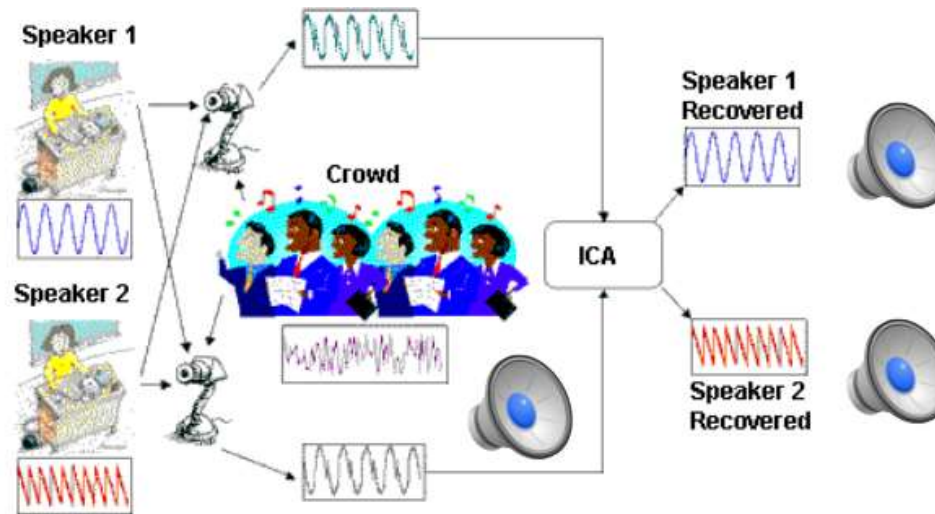
- Independent component analysis – separate a combined signal into its original sources



26

Unsupervised Learning

- Independent component analysis – separate a combined signal into its original sources



27

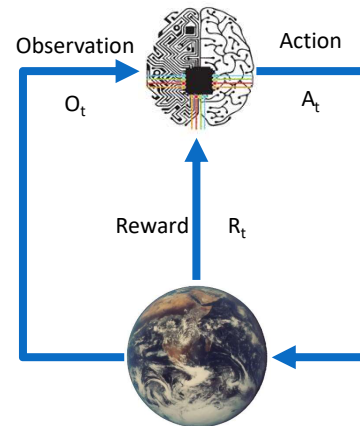
Reinforcement Learning

□ Reinforcement Learning

- ❖ Given a sequence of states and actions with (delayed) rewards, output a policy
- ❖ Policy is a mapping from states \Rightarrow actions that tells you what to do in a given state

□ Examples:

- ❖ Credit assignment problem
- ❖ Game playing
- ❖ Robot in a maze
- ❖ Balance a pole on your hand



7/2/2023

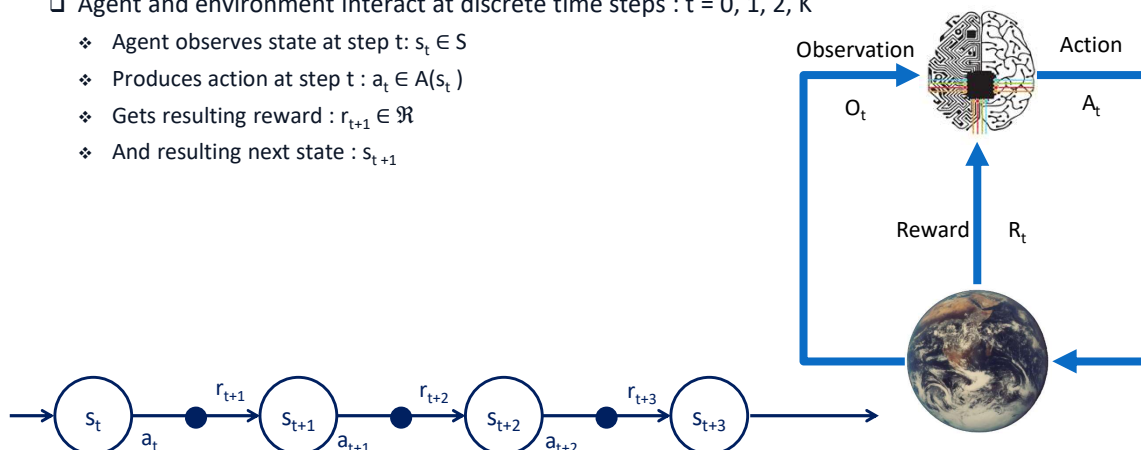
pra-sâmi

28

The Agent-Environment Interface

□ Agent and environment interact at discrete time steps : $t = 0, 1, 2, K$

- ❖ Agent observes state at step t : $s_t \in S$
- ❖ Produces action at step t : $a_t \in A(s_t)$
- ❖ Gets resulting reward : $r_{t+1} \in \mathfrak{R}$
- ❖ And resulting next state : s_{t+1}



7/2/2023

pra-sâmi

29

Examples of Reinforcement Learning

- ❑ Fly stunt maneuvers in a helicopter
- ❑ Defeat the world champion at Backgammon
- ❑ Manage an investment portfolio
- ❑ Control a power station
- ❑ Make a humanoid robot walk
- ❑ Play many different Atari games better than humans

7/2/2023

pra-sâmi

30

Fly stunt maneuvers in a helicopter

Stanford Computer Scientists have developed an artificial intelligence system that enables robotic helicopters to teach themselves to fly difficult stunts by watching other helicopters perform the same maneuvers. The technique is known as "apprenticeship learning." The result is an autonomous helicopter that can fly dazzling stunts on its own.



<https://youtu.be/M-QUkgk3HyE>

7/2/2023

pra-sâmi

31

Game Over: Kasparov and the Machine (trailer)

The documentary "Game Over: Kasparov and the Machine", about the 1997 chess match between Garry Kasparov and Deep Blue, the IBM computer. It also features Yasser Seirawan, Anatoly Karpov and others.



<https://youtu.be/y9UMt-8gfW8>

7/2/2023

pra-sâmi

32

DeepMind Made A Superhuman AI For 57 Atari Games!



<https://youtu.be/dJ4rWhpAGFI>

7/2/2023

pra-sâmi

33

Automated Cargo Wharf Yangshan, Shanghai

- ❑ World's biggest automated cargo wharf, the fourth phase of the Yangshan deep-water port started operation on Sunday.
- ❑ The core technology of the robotic port was developed independently by China.
- ❑ The forth phrase of Yangshan port takes up an area of 2.23 million square meters, whose coastline stretches as long as 2,350 meters. It consists of two 70,000 dead-weight tonnage (DWT) berths and five 50,000 DWT berths.



- ❑ <https://youtu.be/IzOeAGAu60k>

7/2/2023

pra-sâmi

34

Evaluating Machine Learning techniques

7/2/2023

pra-sâmi

35

Common Terms

- ❑ Model
 - ❖ Collection of parameters you are trying to fit
- ❑ Data
- ❑ Instance
- ❑ Feature
 - ❖ Attributes of your data that will be used in prediction
- ❑ Labels
 - ❖ Value you are trying to predict with your model
- ❑ Categorical Data
 - ❖ Nominal : categorically discrete data
 - ❖ Ordinal : quantities that have a natural ordering
 - ❖ Interval : intervals between each value are equally split
 - ❖ Ratio: interval data with a natural zero point.
- ❑ Continuous Data

	A	B	date	gender	grades	tide level	target
0	0.623961	-2.441386	2015-01-01	male	A	3.10 m	0
1	0.262080	-0.951155	2015-01-02	female	D	3.66 m	0
2	-0.653064	0.615920	2015-01-03	male	C	3.47 m	1
3	0.526069	0.421712	2015-01-04	female	D	3.01 m	0
4	-0.472418	1.843539	2015-01-05	female	C	3.32 m	1

Continuous

Time series

Categorical Nominal

Categorical Ordinal

Categorical Ratio

pra-sâmi

7/2/2023

36

How am I Doing?

- ❑ Ever suffered crash and burn in production?
- ❑ No one wants to see a failure in Prod
- ❑ How well a model will generalize to new cases?
 - ❖ Try it out on new cases
- ❑ It's better to test and validate in Dev

7/2/2023

pra-sâmi

37

Train – Test – Validate

- ❑ Split your data into Three sets
 - ❖ The training set (64 %)
 - ❖ The test set (16 %)
 - ❖ The Validation set (20 %)
 - ❑ Train your model using the training set,
 - ❑ Test it using the test set.
 - ❑ The error rate on new cases(Test Set) is called the generalization error (or out-of-sample error),
 - ❑ First indication of how well model will perform
- ❑ Case 1: Training Error High, Test Error High
 - ❖ Model is under fitting
 - ❑ Case 2: Training Error Low, Test Error High
 - ❖ Model is over-fitting the training data.
 - ❑ In train-test cycle, hyper parameters are tuned keeping test data as target
 - ❖ Model is **indirectly exposed** to test set

7/2/2023

pra-sâmi

38

Training Data and Test Data

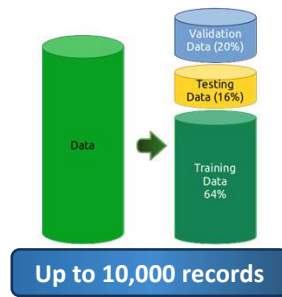
- ❑ Same data cannot be used for Training and Testing
- ❑ Models give better results on seen data,
- ❑ Given huge number of epoch, it may even give perfect score
 - ❖ But would fail to predict anything useful on yet-unseen data
- ❑ In supervised machine learning hold out part of the available data as a test set and another part as validation set
- ❑ Train and tune the model on train + test set,
- ❑ Once tuning is complete, train the model on train + test data and validate on validation data
 - ❖ Remove random seed!
- ❑ Typical cross validation and grid search techniques are also clubbed in this process.
 - ❖ Be very careful not to leak information from test set to train set specially in cross validation
- ❑ Time series: No peeping into future

7/2/2023

pra-sâmi

39

Data Split

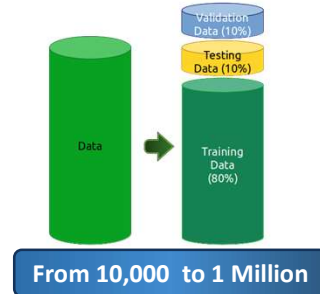
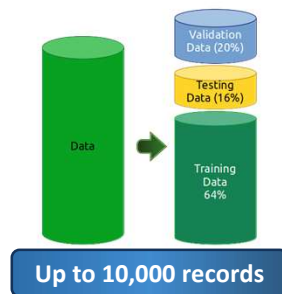


7/2/2023

pra-sâmi

40

Data Split

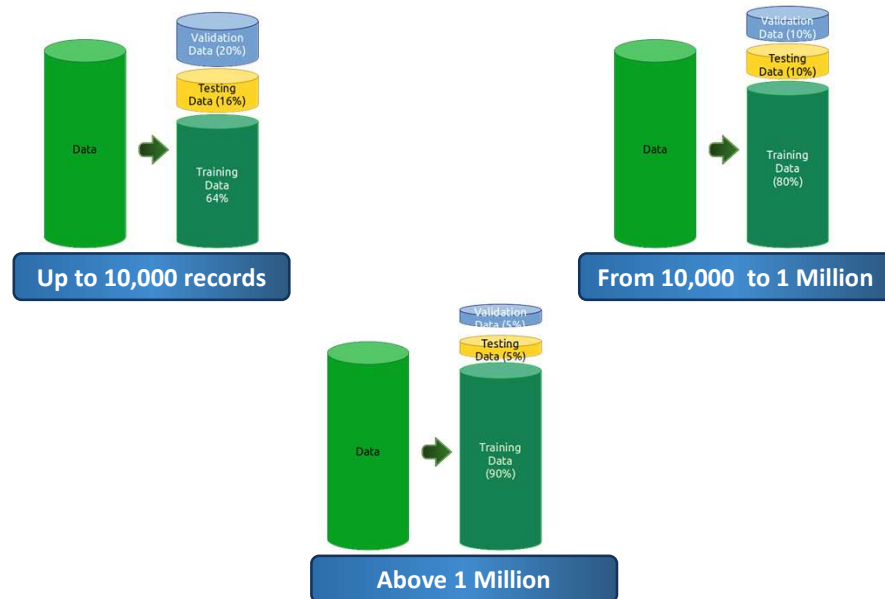


7/2/2023

pra-sâmi

41

Data Split



7/2/2023

pra-sâmi

42

Matching Distribution of Test and Train Data Sets

- ❑ Data Scientists are always looking for data
- ❑ There are plenty of curated datasets available
- ❑ Curated data sets are clean with sufficient features to explain a concept,
 - ❖ May not be suitable for real life application
 - ❖ Real data is dirty
- ❑ Example:
 - ❖ For object identification or image analytics, it is very common to perform web search and download picture
 - ❖ These pictures are of good resolution, taken under favorable lighting condition
 - ❖ After rollout, pictures taken by actual user may not be of similar quality
 - ❖ Resulting in significant drop of accuracy on real data

7/2/2023

pra-sâmi

43

Cross Validation

- ❑ Worried about more than $\frac{1}{3}$ data being used in test-validate
- ❑ The training set is split into complementary subsets,
 - ❖ Each model is trained against a different combination of these subsets
 - ❖ Test against the remaining parts.
 - ❖ Happy with the model type and hyper-parameters
 - ❖ Final model is trained using these hyper-parameters on the full training set,
 - ❖ The generalized error is measured on the validation set.

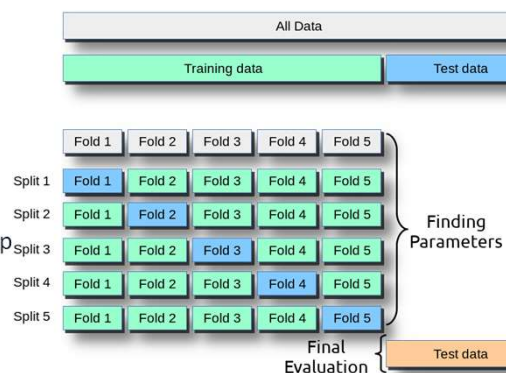
7/2/2023

pra-sâmi

44

Cross-validation : Source Sklearn

- ❑ Partitioning the available data into three sets, drastically reduce the number of samples in each bucket
- ❑ Use cross validation to solve it.
 - ❖ The training set is split into k smaller sets
- ❑ Model is trained using multiple folds as training data;
- ❑ By rotation one fold is kept out
- ❑ Performance measure is average of the values in the loop
- ❑ Approach is computationally expensive,
- ❑ Recommended where data is limited

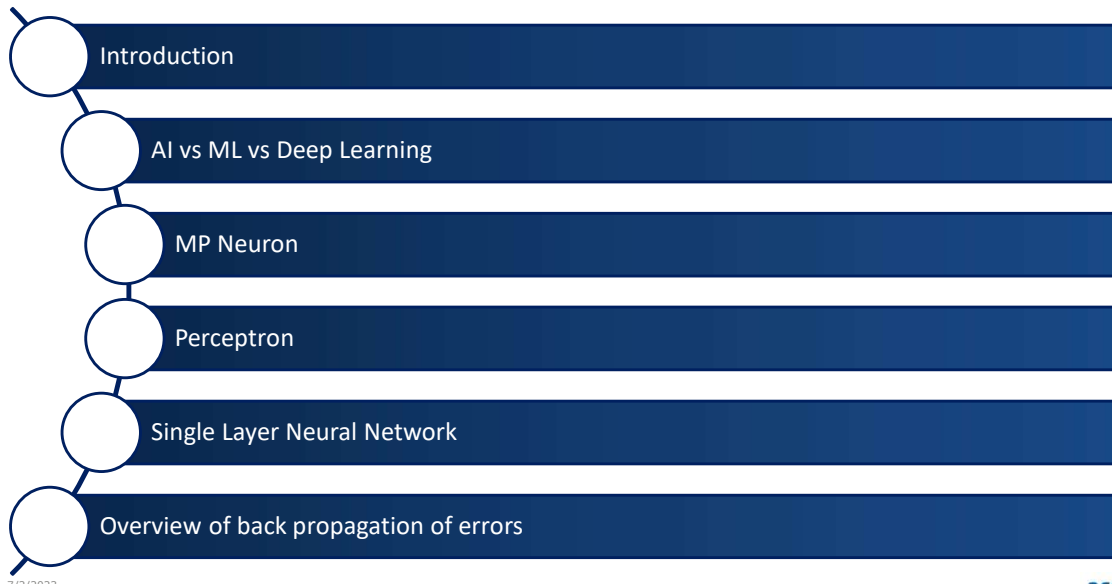


7/2/2023

pra-sâmi

46

Next Session



7/2/2023

pra-sâmi

47



7/2/2023

pra-sâmi