Universiteti i Prishtinës "Hasan Prishtina" Kosovë

# Hyrje në shkencën e të dhënave

Pjesa 10 – Mësimi jo i mbikqyrur (Unsupervised Learning)

Prof. Asoc. Dr. Ermir Rogova

# What is unsupervised learning?

Universiteti i Prishtinës "Hasan Prishtina" Kosovë

- Organize or describe data when no labels are available.
- A major area: clustering

# Agglomerative clustering

- Use any computable cluster similarity measure $sim(C_i, C_j)$ e.g., Euclidean distance, cosine similarity etc.

- For n objects $v_1,..., v_n$, assign each to a singleton cluster $C_i = \{v_i\}$

- Repeat {
  - - identify two most similar clusters $C_j$ and $C_k$ (could be ties-chose one pair)
  - - delete $C_j$ and $C_k$ and add $(C_j \cup C_k)$ to the set of clusters.
  - } until just one cluster.

- Dendrograms diagram the sequence of cluster merges.

Universiteti i Prishtinës "Hasan Prishtina" Kosovë

# Divisive clustering

Universiteti i Prishtinës "Hasan Prishtina" Kosovë

- Put all objects in one cluster
- Repeat until all clusters are singletons {
- - choose a cluster to split based on some criterion.
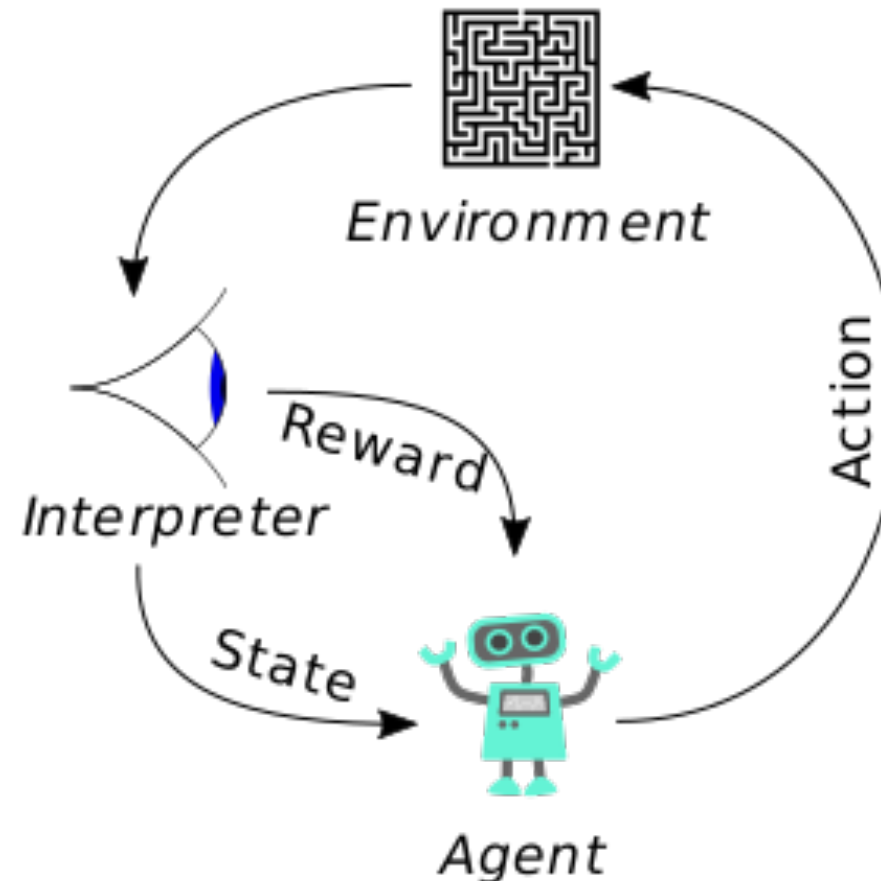- - replace the chosen cluster with sub-clusters.
- }

# K-means algorithm

- Begin with a decision on the value of K = number of clusters.

- Put any initial partition that classifies the data into K clusters. You may assign the training samples randomly or systematically.

- Take each sample in sequence and compute its distance from the centroid of each of the clusters. If a sample is not currently in the cluster with the closest centroid, switch this sample to that cluster and update the centroid of the cluster gaining the new sample and the cluster losing the sample.

- Repeat the above three steps until convergence is achieved.

# Expectation Maximization (EM)

Universiteti i Prishtinës "Hasan Prishtina" Kosovë

- Enables parameter estimation (learning) in probabilistic models with incomplete data.

- Uses Maximum Likelihood Estimation (MLE).

- MLE is a way to assess the quality of a statistical model based on the probability that model assigns to the observed data. The model that has the highest probability of generating the data is the best one.

- EM algorithm is used to find (locally) MLE parameters of a statistical model in cases where the equations cannot be solved directly.

# Introduction to reinforcement learning

Branch of machine learning that attempts to model how **agents** should take **actions** in an **environment** that will maximize some form of cumulative **reward**.

# Summary

- Unsupervised learning methods help with data exploration and identifying useful patterns in the data.

- Clustering is the assignment of a set of observations into subsets (called clusters) so that observations in the same cluster are similar in some sense.

- Clustering approaches:
  - Agglomerative
  - Divisive

Universiteti i Prishtinës "Hasan Prishtina" Kosovë

# Pyetje???