

## SMDA 2018/19 – Exercise 6, Lecture L10 - 12/12/2018

### Exercise 6: Analysis of Human Tumor Microarray dataset – unsupervised learning, clustering with k-means

Please, execute the following tasks and provide answers to the proposed questions.

#### 1. Download the “14-cancer microarray data” from the book website

<https://web.stanford.edu/~hastie/ElemStatLearn/>

- Get informations about the dataset in file 14cancer.info and in Chapter 1 of the book (Hastie et al., 2009)

#### 2. Generate a new Kernel in Kaggle

#### 3. Load the data in Kaggle

- Use for instance the train data
- Load also the labels

#### 4. Use the `sklearn.cluster` module to perform clustering analysis on the dataset. In particular, repeat the analysis proposed in section 14.3.8 of the book (Hastie et al., 2009)

- Start using K-means and then test some other clustering algorithms at your choice
- Cluster the samples (i.e., columns). Each sample has a label (tumor type)
- Do not use the labels in the clustering phase but examine them posthoc to interpret the clusters
- Run k-means with K from 2 to 10 and compare the clusterings in terms of within-sum of squares
- Show the chart of performance depending on K
- Select some K and analyze the clusters as done in the book