# Series 5, April 23rd, 2018
# (Clustering and PCA)

**Solutions will be published on Friday, May 4th 2018.**

### Problem 1 (K-means convergence):

In the K-means clustering algorithm, you are given a set of $n$ points $x_i \in \mathbb{R}^d$, $i \in \{1, \ldots, n\}$ and you want to find the centers of $k$ clusters $\mu = (\mu_1, \ldots, \mu_k)$ by minimizing the average distance from the points to the closest cluster center. Formally, you want to minimize the following loss function

$$L(\mu) = \sum_{i=1}^{n} \min_{j \in \{1,\ldots,k\}} \|x_i - \mu_j\|_2^2.$$

To approximate the solution, we introduce new assignment variables $z_i \in \arg\min_{j \in \{1,\ldots,k\}} \|x_i - \mu_j\|_2^2$ for each data point $x_i$. The K-means algorithm iterates between updating the variables $z_i$ (*assignment step*) and updating the centers $\mu_j = \frac{1}{|\{i : z_i = j\}|} \sum_{i : z_i = j} x_i$ (*refitting step*). The algorithm stops when no change occurs during the *assignment step*.

Show that K-means is guaranteed to converge (to a local optimum). *Hint:* You need to prove that the loss function is guaranteed to decrease monotonically in each iteration until convergence. Prove this separately for the *assignment step* and the *refitting step* .

### Problem 2 (K-medians clustering):

In this exercise, you are asked to derive a new clustering algorithm that would use a different loss function given by

$$L(\mu) = \sum_{i=1}^{n} \min_{j \in \{1,\ldots,k\}} \|x_i - \mu_j\|_1.$$

(a) Find the update steps for both $z_i$ and for $\mu_j$ in this case.

(b) What can you say about the convergence of your algorithm?

(c) In which situation would you prefer to use K-medians clustering instead of K-means clustering?

### Problem 3 (PCA):

Suppose we have a dataset with 4 points:

$$\mathcal{D} = \{(1, 5), (0, 6), (-7, 0), (-6, -1)\}$$

(a) Plot the dataset and try to guess two principal components ($k = 2$).

(b) Compute the empirical covariance matrix, its eigenvalues and eigenvectors. Do the eigenvectors correspond to your guess of principal components? Please do not forget the assumptions of PCA. (The dataset should be centered and we want unit eigenvectors.)