

COMS 4772 Fall 2016 Homework 2
Due Friday, October 28

Instructions:

- The usual homework policies (<http://www.cs.columbia.edu/~djhsu/coms4772-f16/about.html>) are, of course, in effect.
- Using this L^AT_EX template will be helpful for grading purposes.

Problem 1 (25 points). Let \mathbf{X} be a random vector in \mathbb{R}^d whose distribution is a mixture of k spherical Gaussians:

$$\mathbf{X} \sim \pi_1 \mathcal{N}(\boldsymbol{\mu}_1, \sigma_1^2 \mathbf{I}) + \pi_2 \mathcal{N}(\boldsymbol{\mu}_2, \sigma_2^2 \mathbf{I}) + \cdots + \pi_k \mathcal{N}(\boldsymbol{\mu}_k, \sigma_k^2 \mathbf{I}).$$

For any set $C \subset \mathbb{R}^d$, define

$$\text{cost}(C) := \mathbb{E} \left[\min_{\mathbf{y} \in C} \|\mathbf{X} - \mathbf{y}\|_2^2 \right].$$

Let $M := \{\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_k\}$. Prove that if $k < e^{d/2}$, then

$$\text{cost}(M) \leq \frac{1}{1 - \frac{2 \ln(k)}{d}} \cdot \min_{\substack{C \subset \mathbb{R}^d \\ |C| \leq k}} \text{cost}(C).$$

Solution.

□

Problem 2 (25 points). Suppose $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times d}$ each have rank d . Give unambiguous pseudocode for an algorithm that, when given \mathbf{A} and \mathbf{B} as inputs, finds all solutions $\mathbf{v} \in S^{d-1}$ satisfying

$$\exists \lambda \in \mathbb{R} \setminus \{0\} \text{ s.t. } \mathbf{A}^\top \mathbf{A} \mathbf{v} = \lambda \mathbf{B}^\top \mathbf{B} \mathbf{v}.$$

If there is an entire subspace of solutions, the algorithm just needs to return an orthonormal basis for this subspace. Your pseudocode can use things like SVD, Gram-Schmidt, etc. as black-box subroutines. Prove that the algorithm is correct.

Solution.

□

Problem 3 (25 points). Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ be a data matrix whose rows are $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n \in \mathbb{R}^d$. Let $\mathbf{D} \in \mathbb{R}^{n \times n}$ be the matrix whose (i, j) -th entry is the squared Euclidean distance $D_{i,j} = \|\mathbf{a}_i - \mathbf{a}_j\|_2^2$. Suppose you are given the squared Euclidean distance matrix \mathbf{D} as input, and you are asked to recover the set of original points $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$ up to some translation. *You do not have access to the original data matrix \mathbf{A} .*

- (a) Let $\mathbf{s} \in \mathbb{R}^n$ be the vector whose i -th entry is $\|\mathbf{a}_i\|_2^2$. Prove that $\mathbf{D} = \mathbf{s}\mathbf{1}_n^\top - 2\mathbf{A}\mathbf{A}^\top + \mathbf{1}_n\mathbf{s}^\top$, where $\mathbf{1}_n \in \mathbb{R}^n$ is the all-ones vector.
- (b) Let $\mathbf{\Pi} \in \mathbb{R}^{n \times n}$ be the orthogonal projector for the $(n-1)$ -dimensional subspace

$$\{\mathbf{x} \in \mathbb{R}^n : \langle \mathbf{1}_n, \mathbf{x} \rangle = 0\}.$$

Prove that $-(1/2)\mathbf{\Pi D \Pi} = \mathbf{\Pi A A}^\top \mathbf{\Pi}$.

- (c) Explain how to determine points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$ from \mathbf{D} such that:

- $D_{i,j} = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$ for all $i, j \in [n]$; and
- $\sum_{i=1}^n \mathbf{x}_i = \mathbf{0}$.

(You may assume that you are told the original dimension d .)

- (d) *Optional.* Suppose the matrix \mathbf{D} is corrupted (say, because your distance measuring device is imperfect), so the entries no longer correspond to the squared Euclidean distances between the \mathbf{a}_i . Explain how to determine points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^n$ (yes, n and not d) from \mathbf{D} such that:

- $\sum_{i=1}^n \mathbf{x}_i = \mathbf{0}$;
- $\|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \geq D_{i,j}$ for all $i \neq j$; and
- $\max \left\{ \frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{D_{i,j}} : 1 \leq i < j \leq n \right\}$ is as small as possible.

Hint: use semidefinite programming.

Solution.

□

Problem 4 (25 points). Exercise 3.25 from BHK.

Solution.

□