

## HOMWORK SET 5

CSCI5525 Machine Learning: Analysis and Methods (Fall 2024)

**Due** 11:59 pm, Dec 23 2024

**Instruction** Your writeup, either typeset or scanned, should be a single PDF file. For problems requiring coding, organize all codes for all problems into **ONE** Jupyter notebook file (i.e., .ipynb file) with cell execution outputs. Your submission to Gradescope should include the single PDF and the one notebook file—**please DO NOT zip them!** Please assign page(s) to each question to help reduce the TA's navigation time. If your notebook submission does not display properly on Gradescope due to a large file error, please try to remove images and/or figures from your cell outputs and re-upload it. No late submission will be accepted. For each problem, you should acknowledge your collaborators—**including AI tools**, if any.

**About the use of AI tools** You are strongly encouraged to use AI tools—they are becoming our workspace friends, such as ChatGPT (<https://chat.openai.com/>), Claude (<https://claude.ai/chats>), and Github Copilot (<https://github.com/features/copilot>), to help you when trying to solve problems. It takes a bit of practice to ask the right and effective questions/prompts to these tools; we highly recommend that you go through this popular free short course **ChatGPT Prompt Engineering for Developers** offered by <https://learn.deeplearning.ai/> to get started.

**If you use any AI tools for any of the problems, you should include screenshots of your prompting questions and their answers in your writeup.** The answers provided by such AI tools often contain factual errors and reasoning gaps. **So, if you only submit an AI answer with such bugs for any problem, you will obtain a zero score for that problem.** You obtain the scores only when you explain the bugs and also correct them in your own writing. You can also choose not to use any of these AI tools, in which case we will grade based on the efforts you have made.

**Reminder about notations** We will use small letters (e.g.,  $u$ ) for scalars, small boldface letters (e.g.,  $\mathbf{a}$ ) for vectors, and capital boldface letters (e.g.,  $\mathbf{A}$ ) for matrices. For a matrix  $\mathbf{A}$ ,  $\mathbf{a}^i$  (supscripting) means its  $i$ -th row as a *row vector*, and  $\mathbf{a}_j$  (subscripting) means the  $j$ -th column as a column vector, and  $a_{ij}$  means its  $(i, j)$ -th element.  $\mathbb{R}$  is the set of real numbers.  $\mathbb{R}^n$  is the space of  $n$ -dimensional real vectors, and similarly  $\mathbb{R}^{m \times n}$  is the space of  $m \times n$  real matrices. The dotted equal sign  $\doteq$  means defining.

### Problem 1 (Dimension reduction; 7.5/15)

- (a) In our geometric view of PCA, we try to find a best-fitting low-dimensional subspace to our dataset. Mathematically, let  $\mathbf{X} \in \mathbb{R}^{N \times d}$  be our data matrix, and  $\bar{\mathbf{X}}$  the centered version. PCA tries to solve the following problem

$$\min_{\mathbf{U} \in \mathbb{R}^{d \times r}: \mathbf{U}^\top \mathbf{U} = \mathbf{I}} \frac{1}{N} \sum_{i \in [N]} \|\bar{\mathbf{x}}_i - \mathbf{U} \mathbf{U}^\top \bar{\mathbf{x}}_i\|_2^2, \quad (1)$$

where  $r$  is our target dimension for the subspace. We know that the optimal solution is the  $\mathbf{U}$  that collects the top  $r$  right singular vectors of  $\bar{\mathbf{X}}$  (or equivalently, the top  $r$  eigenvectors of  $\frac{1}{N} \bar{\mathbf{X}}^\top \bar{\mathbf{X}}$ ).

Instead of performing the centering step separately, we can also try to fit an  $r$ -dimensional affine subspace directly. Since any affine subspace can be represented as  $\{\mathbf{U} \mathbf{z} + \mathbf{v} : \mathbf{z} \in \mathbb{R}^r\}$  for some orthonormal  $\mathbf{U} \in \mathbb{R}^{d \times r}$  and some  $\mathbf{v}$ , we can formulate the learning problem as

$$\min_{\mathbf{U} \in \mathbb{R}^{d \times r}: \mathbf{U}^\top \mathbf{U} = \mathbf{I}, \mathbf{z}_i \in \mathbb{R}^r \forall i, \mathbf{v} \in \mathbb{R}^d} \frac{1}{N} \sum_{i \in [N]} \|\mathbf{x}_i - \mathbf{v} - \mathbf{U} \mathbf{z}_i\|_2^2. \quad (2)$$

- (i) Prove that for any fixed  $U$ , an optimal solution for Eq. (2) is  $\mathbf{v} = \frac{1}{N} \sum_{i \in [N]} \mathbf{x}_i$  and  $\mathbf{z}_i = U^\top(\mathbf{x}_i - \mathbf{v})$ , i.e., solving Eq. (2) is equivalent to centering as data preprocessing and then solving Eq. (1). (1/15)
- (ii) Consider  $d = 50$ ,  $N = 200$ , and  $r = 10$ . Fix a random seed, and generate data as follows: (1) generate a subspace basis  $\mathbf{A}_0 \in \mathbb{R}^{d \times r}$  with iid  $N(0, 1/r)$  entries; (2) generate a coefficient matrix  $\mathbf{Z}_0 \in \mathbb{R}^{r \times N}$  with iid  $N(0, 1)$  entries; (3) generate the data as  $\mathbf{X} = (\mathbf{A}_0 \mathbf{Z}_0 + \mathbf{N})^\top$ , where  $\mathbf{N} \in \mathbb{R}^{d \times N}$  contains iid  $N(0, 0.05)$  entries. **Implement PCA (including centering, but no variance normalization) on  $\mathbf{X}$  to obtain your estimation of the subspace basis, call it  $\mathbf{A}_1 \in \mathbb{R}^{d \times r}$ .** To compare the two subspaces represented by  $\mathbf{A}_0$  and  $\mathbf{A}_1$ , a reasonable metric is

$$d(\text{span}(\mathbf{A}_1), \text{span}(\mathbf{A}_0)) = \|\mathbf{A}_1 \mathbf{A}_1^\dagger - \mathbf{A}_0 \mathbf{A}_0^\dagger\|_F / \|\mathbf{A}_0 \mathbf{A}_0^\dagger\|_F, \quad (3)$$

where for any matrix  $M$ ,  $M^\dagger$  denotes its pseudoinverse ([https://en.wikipedia.org/wiki/Moore%E2%80%93Penrose\\_inverse](https://en.wikipedia.org/wiki/Moore%E2%80%93Penrose_inverse)). Here,  $\mathbf{A}_1 \mathbf{A}_1^\dagger$  and  $\mathbf{A}_0 \mathbf{A}_0^\dagger$  are the projectors on the subspaces  $\mathbf{A}_1$  and  $\mathbf{A}_0$  span, respectively, and  $\|\mathbf{A}_0 \mathbf{A}_0^\dagger\|_F$  performs proper normalization of the subspace distance. **Report your  $d(\text{span}(\mathbf{A}_1), \text{span}(\mathbf{A}_0))$ ; ideally, it should be close to 0.** (1/15)

- (iii) Linear autoencoders are also known to perform PCA. Consider the following formulation of a linear autoencoder:

$$\min_{\mathbf{A} \in \mathbb{R}^{d \times r}} \frac{1}{N} \sum_{i \in [N]} \|\bar{\mathbf{x}}_i - \mathbf{A} \mathbf{A}^\top \bar{\mathbf{x}}_i\|_2^2, \quad (4)$$

where  $\mathbf{A}^\top$  denotes the linear encoder and  $\mathbf{A}$  denotes the linear decoder, and our encoder and decoder are symmetric with respect to each other (i.e., the transpose of each other). **Implement the gradient descent with line search to solve Eq. (4) with the same data  $\mathbf{X}$  and  $\bar{\mathbf{X}}$  you generate in (ii), and obtain an estimate of the subspace basis  $\mathbf{A}_2$ . Report your  $d(\text{span}(\mathbf{A}_2), \text{span}(\mathbf{A}_0))$ ; is it close to 0?** (1/15)

- (iv) Now keep your  $\mathbf{A}_0$  and  $\mathbf{Z}_0$  in (ii), but change to noise model to Laplace noise ([https://en.wikipedia.org/wiki/Laplace\\_distribution](https://en.wikipedia.org/wiki/Laplace_distribution)), i.e.,  $\mathbf{X}' = (\mathbf{A}_0 \mathbf{Z}_0 + \mathbf{L})^\top$ , where  $\mathbf{L}$  contains iid Laplace(0, 0.15) entries. Due to the long-tailed nature of Laplace noise, the noise magnitudes on certain coordinates can be substantially larger than the rest. In this case, it is appropriate to switch the loss in Eq. (2) to a robust loss, say the  $\ell_1$  loss, yielding:

$$\min_{\mathbf{U} \in \mathbb{R}^{d \times r}: \mathbf{U}^\top \mathbf{U} = \mathbf{I}, \mathbf{z}_i \in \mathbb{R}^r \forall i, \mathbf{v} \in \mathbb{R}^d} \frac{1}{N} \sum_{i \in [N]} \|\mathbf{x}_i - \mathbf{v} - \mathbf{U} \mathbf{z}_i\|_1. \quad (5)$$

However, unlike for Eq. (2) that we can eliminate  $\mathbf{z}_i$ 's and  $\mathbf{v}$  (as proved in (i)), here, we have to solve the three groups of variables altogether. To solve Eq. (5), we can ignore the  $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$ , and consider alternately minimizing over the two groups of variables:  $\mathbf{U}$  vs.  $\mathbf{z}_i$ 's,  $\mathbf{v}$ —when we fix one group, the objective function is convex with respect to the other group. So, we can implement the alternating minimization algorithm (see, e.g., <https://www.mit.edu/~rakhlin/6.883/lectures/lecture07.pdf> for more information), and solve the two convex subproblems with say, the CVXPY package (<https://www.cvxpy.org/>; we have used it in our earlier HW sets). **Implement the suggested algorithm, or whatever numerical algorithm you like to solve Eq. (5) to yield your subspace estimate  $\mathbf{A}_3$ , and report your  $d(\text{span}(\mathbf{A}_3), \text{span}(\mathbf{A}_0))$ . For comparison, also apply the plain PCA on  $\mathbf{X}'$  to produce another subspace estimate  $\mathbf{A}_4$ ; report your  $d(\text{span}(\mathbf{A}_4), \text{span}(\mathbf{A}_0))$ . Do what you observe?** (1.5/15)

- (b) The famous Johnson–Lindenstrauss (J-L) lemma states: for any  $\varepsilon \in (0, 1)$  and any  $N$  set  $S = \{\mathbf{x}_i\}_{i \in [N]} \subset \mathbb{R}^d$ , there exists a matrix  $\mathbf{A} \in \mathbb{R}^{d' \times d}$  for every  $d' \geq C \log N / \varepsilon^2$  ( $C > 0$  is a universal constant) so that

$$(1 - \varepsilon) \|\mathbf{u} - \mathbf{v}\|_2^2 \leq \|\mathbf{A}\mathbf{u} - \mathbf{A}\mathbf{v}\|_2^2 \leq (1 + \varepsilon) \|\mathbf{u} - \mathbf{v}\|_2^2 \quad \forall \mathbf{u}, \mathbf{v} \in S, \quad (6)$$

$$(1 - \varepsilon) \|\mathbf{u} + \mathbf{v}\|_2^2 \leq \|\mathbf{A}\mathbf{u} + \mathbf{A}\mathbf{v}\|_2^2 \leq (1 + \varepsilon) \|\mathbf{u} + \mathbf{v}\|_2^2 \quad \forall \mathbf{u}, \mathbf{v} \in S, \quad (7)$$

where the second line is not typically stated in the J-L lemma, but can be proved easily using a very similar argument to that for the first line. In particular, a random draw iid Gaussian matrix  $\frac{1}{\sqrt{d'}} \mathbf{G}$  satisfies the property with very high probability.

Now, assume that all points inside  $S$  have unit norm, i.e.,  $\|\mathbf{x}_i\|_2 = 1 \quad \forall i \in [N]$ . Show that for any  $d' \geq C \log N / \varepsilon^2$ , we can also construct a matrix  $\mathbf{B}$  so that

$$|\langle \mathbf{B}\mathbf{u}, \mathbf{B}\mathbf{v} \rangle - \langle \mathbf{u}, \mathbf{v} \rangle| \leq \varepsilon \quad \forall \mathbf{u}, \mathbf{v} \in S. \quad (8)$$

(Hint: consider the identity  $4 \langle \mathbf{x}, \mathbf{x}' \rangle = (\|\mathbf{x} + \mathbf{x}'\|_2^2 - \|\mathbf{x} - \mathbf{x}'\|_2^2)$  that holds for arbitrary  $\mathbf{x}$  and  $\mathbf{x}'$  in  $\mathbb{R}^d$ .) (1/15)

- (c) In the ISOMAP method for nonlinear dimension reduction, we need to turn a Euclidean distance matrix into a similarity matrix. Let  $\mathbf{X} \in \mathbb{R}^{N \times d}$  be the data matrix (i.e., we have  $N$  points in  $\mathbb{R}^d$ ), and let  $\mathbf{X}^* = \mathbf{X} - \frac{1}{N} \mathbf{1}\mathbf{1}^\top \mathbf{X}$  be the centered version of  $\mathbf{X}$  ( $\mathbf{1}$ 's are all-one vectors of appropriate dimensions). Moreover, let  $\mathbf{D} \in \mathbb{R}^{N \times N}$  denote the squared Euclidean distance matrix, i.e.,  $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$ , and write  $\mathbf{G} = \mathbf{X}^\top \mathbf{X}$ . Our goal is to write  $\mathbf{G}^* \doteq (\mathbf{X}^*)^\top \mathbf{X}^*$  as a function of  $\mathbf{D}$ .

- (i) Can you represent  $\mathbf{G}^*$  as a function of  $\mathbf{G}$ ? (0.5/15)  
(ii) Combine (i) and the fact that  $g_{ij} = \frac{1}{2} (\|\mathbf{x}_i\|_2^2 + \|\mathbf{x}_j\|_2^2 - d_{ij})$ ,  $\forall i, j$  to show that

$$g_{ij}^* = -\frac{1}{2} \left( d_{ij} - \frac{1}{N} \mathbf{d}^i \mathbf{1} - \frac{1}{N} \mathbf{1}^\top \mathbf{d}_j + \frac{1}{N^2} \mathbf{1}^\top \mathbf{D} \mathbf{1} \right). \quad (9)$$

(1/15)

- (iii) Show that  $\mathbf{G}^* = -\frac{1}{2} \left( \mathbf{I} - \frac{1}{N} \mathbf{1}\mathbf{1}^\top \right) \mathbf{D} \left( \mathbf{I} - \frac{1}{N} \mathbf{1}\mathbf{1}^\top \right)$ . (0.5/15)

## Problem 2 (Clustering; 6/15)

- (a) The classical k-means algorithm can be viewed as an alternating minimization algorithm to solve

$$\min_{\mu_1, \dots, \mu_K, \mathbf{M} \in \{0,1\}^{N \times K}, \mathbf{M}\mathbf{1} = \mathbf{1}} \frac{1}{N} \sum_{i \in [N]} \sum_{k \in [K]} m_{ik} d(\mathbf{x}_i, \mu_k). \quad (10)$$

Here,  $\mathbf{M} \in \{0, 1\}^{N \times K}$  is the cluster-assignment matrix, and  $\forall i, k, m_{ik} = 1$  means  $\mathbf{x}_i$  is assigned to cluster  $k$ . The constraint  $\mathbf{M}\mathbf{1} = \mathbf{1}$  ensures that each point is only assigned to one cluster. Consider the k-medoids variant of k-means, which requires the cluster centroids  $\mu_1, \dots, \mu_K$  come from the dataset  $S = \{\mathbf{x}_i\}_{i \in [N]}$ , not from outside. **Change the formulation in Eq. (10) to reflect this, and also derive an alternating minimization algorithm to solve the k-medoids problem.** (1/15)

- (b) In spectral clustering, the graph is represented by a weight/similarity matrix  $\mathbf{W} \in \mathbb{R}^{N \times N}$  ( $N$  points are assumed), whose entries are large when the connection between the corresponding nodes are strong. A critical quantity is the graph Laplacian:

$$\mathbf{L} = \mathbf{D} - \mathbf{W}, \quad \mathbf{D} \doteq \text{diag}(\mathbf{W}\mathbf{1}) : \text{degree matrix.} \quad (11)$$

- (i) A crucial identity we use to derive the final learning objective in spectral clustering is

$$\forall \mathbf{v}, \quad \mathbf{v}^\top \mathbf{L} \mathbf{v} = \frac{1}{2} \sum_{i \in [N], j \in [N]} w_{ij} (v_i - v_j)^2. \quad (12)$$

Show it! (1/15)

- (ii) Assume that all entries of  $\mathbf{W}$  are nonnegative, i.e.,  $\mathbf{W} \geq 0$ . Prove that  $\mathbf{L}$  is positive semidefinite, and also that 0 is an eigenvalue of  $\mathbf{L}$ . (1/15)
- (iii) Suppose there are exactly  $K \leq N$  connected subgraphs. It can be shown that  $\mathbf{L}$  has exactly  $K$  zero eigenvalues and there are  $K$  eigenvectors exactly indicate the cluster memberships of all data points. If we obtain such eigenvectors, we are done. But, when we have  $K$  equal eigenvalues, all vectors from the corresponding  $K$ -dimensional subspace are valid eigenvectors, so our eigensolver may not return the ones we like. Let  $\mathbf{H}^* \in \mathbb{R}^{N \times K}$  be the ideal eigenvectors, and  $\widehat{\mathbf{H}}$  the one we actually obtain. It must hold that  $\mathbf{H}^* = \widehat{\mathbf{H}} \mathbf{R}$  for an unknown rotation matrix  $\mathbf{R} \in \mathbb{R}^{K \times K}$ , with  $\mathbf{R}^\top \mathbf{R} = \mathbf{Q} \mathbf{R}^\top = \mathbf{I}$ . **Show that for any two vectors  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^K$ ,  $\|(\mathbf{u} - \mathbf{v}) \mathbf{R}\|_2 = \|\mathbf{u} - \mathbf{v}\|_2$ , and thereby explain why the subsequent  $k$ -means step with  $\ell_2$  distance on the rows of  $\widehat{\mathbf{H}}$  makes sense.** (1/15)
- (c) In the mean-shift algorithm for clustering, the underlying density is estimated via kernel density estimation that typically takes the form

$$\widehat{p}(\mathbf{x}) = \frac{1}{N} \sum_{i \in [N]} K(\|\mathbf{x} - \mathbf{x}_i\|_2), \quad (13)$$

where  $K$  is a smoothing kernel that satisfies  $\int_{\mathbf{x}} K(\|\mathbf{x} - \mathbf{x}_i\|_2) d\mathbf{x} = 1$  for all  $i \in [N]$  so that  $\widehat{p}(\mathbf{x})$  is a valid density function. Consider the kernel

$$K(\|\mathbf{x} - \mathbf{x}'\|_2) = c_d \max\left(0, 1 - \|\mathbf{x} - \mathbf{x}'\|_2^2\right), \quad (14)$$

where  $c_d$  is a universal constant to ensure  $\int_{\mathbf{x}} c_d \max\left(0, 1 - \|\mathbf{x} - \mathbf{x}'\|_2^2\right) d\mathbf{x} = 1$ . **Derive the key update formula for mean-shift with this kernel.** (1/15)

- (d) In  $\mathbb{R}^d$ , consider mixture modeling with  $K$  multivariate Bernoulli components, i.e.,

$$p(\mathbf{x}; \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K) = \sum_{k \in [K]} w_k \prod_{j \in [d]} \mu_{kj}^{x_j} (1 - \mu_{kj})^{1-x_j}, \quad (15)$$

i.e., for each mixture component, the  $d$  coordinates are assumed to be independent for simplicity. **Derive an EM style algorithm for fitting this mixture model to a dataset  $S = \{\mathbf{x}_i\}_{i \in [N]}$  using the principle of maximum likelihood estimation.** (1/15)

**Problem 3 (Generative modeling; 1.5/15)** Consider a transformation  $\mathbf{x} = f(\mathbf{z})$  and its inverse  $\mathbf{z} = g(\mathbf{x})$ .

- (a) By differentiating both sides of  $\mathbf{x} = f \circ g(\mathbf{x})$  with respect to  $\mathbf{x}$ , show that  $\mathbf{J}_f(g(\mathbf{x}))\mathbf{J}_g(\mathbf{x}) = \mathbf{I}$ , where  $\mathbf{J}$  denotes Jacobian matrix (refer to our notes on calculus review); (0.5/15)
- (b) Given the density function of  $\mathbf{z}$  as  $p_z(\mathbf{z})$ , the change-of-variable formula reads

$$p_{\mathbf{x}}(\mathbf{x}) = p_z(g(\mathbf{x})) |\det \mathbf{J}_g(\mathbf{x})|. \quad (16)$$

Show that we can also write it as

$$p_{\mathbf{x}}(\mathbf{x}) = p_z(g(\mathbf{x})) |\det \mathbf{J}_f(\mathbf{z})|^{-1}. \quad (17)$$

**(Hint: the determinant of the product of two square matrices is the product of their determinants.)** (0.5/15)

- (c) Assume that the entries inside  $\mathbf{z} \in \mathbb{R}^d$  are iid, and each follows a half-normal distribution ([https://en.wikipedia.org/wiki/Half-normal\\_distribution](https://en.wikipedia.org/wiki/Half-normal_distribution)), i.e.,  $\forall j \in [d], z_j \geq 0$  and  $p(z_i) = \frac{2}{\sqrt{2\pi}} e^{-z_j^2/2}$ . Let  $\mathbf{x} = \mathbf{z}^2$ , where the square is taken elementwise. Derive the probability density function for  $\mathbf{x}$ . (0.5/15)