

Due: Wednesday, May 6 at 11:59 pm

Deliverables:

1. Submit a PDF of your homework, **with an appendix listing all your code**, to the Gradescope assignment entitled “Homework 7 Write-Up”. In addition, please include, as your solutions to each coding problem, the specific subset of code relevant to that part of the problem. You may typeset your homework in LaTeX or Word (submit PDF format, **not** .doc/.docx format) or submit neatly handwritten and scanned solutions. **Please start each question on a new page**. If there are graphs, include those graphs in the correct sections. **Do not** put them in an appendix. We need each solution to be self-contained on pages of its own.
 - In your write-up, please state with whom you worked on the homework.
 - In your write-up, please copy the following statement and sign your signature next to it. (Mac Preview and FoxIt PDF Reader, among others, have tools to let you sign a PDF file.) We want to make it *extra* clear so that no one inadvertently cheats.

“I certify that all solutions are entirely in my own words and that I have not looked at another student’s solutions. I have given credit to all external sources I consulted.”
2. Submit all the code needed to reproduce your results to the Gradescope assignment entitled “Homework 7 Code”. Yes, you must submit your code twice: in your PDF write-up following the directions as described above so the readers can easily read it, and once in compilable/interpretable form so the readers can easily run it. Do **NOT** include any data files we provided. Please include a short file named README listing your name, student ID, and instructions on how to reproduce your results. Please take care that your code doesn’t take up inordinate amounts of time or memory. If your code cannot be executed, your solution cannot be verified.

1 Low-Rank Approximation

Low-rank approximation tries to find an approximation to a given matrix, where the approximation matrix has a lower rank compared to the original matrix. This is useful for mathematical modeling and data compression. Mathematically, given a matrix M , we try to find \hat{M} in the optimization problem

$$\operatorname{argmin}_{\hat{M}} \|M - \hat{M}\|_F \quad \text{subject to} \quad \operatorname{rank}(\hat{M}) \leq k \quad (1)$$

where $\|A\|_F = \sqrt{\sum_i \sum_j a_{ij}^2}$ is the Frobenius norm, i.e., the sum of the squares of all the entries in the matrix, followed by a square root.

This problem can be solved with a singular value decomposition (SVD). Specifically, let $M = U\Sigma V^\top$, where $\Sigma = \operatorname{diag}(\sigma_1, \dots, \sigma_n)$. Then a rank- k approximation of M can be written as $\hat{M} = U\hat{\Sigma}V^\top$, where $\hat{\Sigma} = \operatorname{diag}(\sigma_1, \dots, \sigma_k, 0, \dots, 0)$. In this problem, we aim to perform this approximation method on gray-scale images, which can be thought of as a 2D matrix.

- Using the image `low-rank_data/face.jpg`, perform a rank-5, rank-20, and rank-100 approximation on the image. Show both the original image as well as the low-rank images you obtain in your report.
- Now perform the same rank-5, rank-20, and rank-100 approximation on `low-rank_data/sky.jpg`. Show both the original image as well as the low-rank images you obtain in your report.
- In one plot, plot the Mean Squared Error (MSE) between the rank- k approximation and the original image for both `low-rank_data/face.jpg` and `low-rank_data/sky.jpg`, for k ranging from 1 to 100. Be sure to label each curve in the plot. The MSE between two images $I, J \in \mathbb{R}^{w \times h}$ is

$$\operatorname{MSE}(I, J) = \sum_{i,j} (I_{i,j} - J_{i,j})^2. \quad (2)$$

- Find the lowest-rank approximation for which you begin to have a hard time differentiating the original and the approximated images. (This is a subjective judgement!) Compare your results for the face and the sky image. What are the possible reasons for the difference?

2 Regularized and Kernel k-Means

Recall that in k -means clustering we attempt to minimize the objective

$$\min_{C_1, C_2, \dots, C_k} \sum_{i=1}^k \sum_{X_j \in C_i} \|X_j - \mu_i\|_2^2, \quad \text{where}$$
$$\mu_i = \operatorname{argmin}_{\mu_i \in \mathbb{R}^d} \sum_{X_j \in C_i} \|X_j - \mu_i\|_2^2 = \frac{1}{|C_i|} \sum_{X_j \in C_i} X_j, \quad i = 1, 2, \dots, k.$$

The sample points are $\{X_1, \dots, X_n\}$, where $X_j \in \mathbb{R}^d$. C_i is the set of sample points assigned to cluster i and $|C_i|$ is its cardinality. Each sample point is assigned to exactly one cluster.

- What is the minimum value of the objective when $k = n$ (the number of clusters equals the number of sample points)?

- (b) (Regularized k -means) Suppose we add a regularization term to the above objective. The objective is now

$$\sum_{i=1}^k \left(\lambda \|\mu_i\|_2^2 + \sum_{X_j \in C_i} \|X_j - \mu_i\|_2^2 \right).$$

Show that the optimum of

$$\min_{\mu_i \in \mathbb{R}^d} \lambda \|\mu_i\|_2^2 + \sum_{X_j \in C_i} \|X_j - \mu_i\|_2^2$$

is obtained at

$$\mu_i = \frac{1}{|C_i| + \lambda} \sum_{X_j \in C_i} X_j.$$

- (c) Here is an example where we would want to regularize clusters. Suppose there are n students who live in a \mathbb{R}^2 Euclidean world and who wish to share rides efficiently to Berkeley for their final exam in CS 189. The university permits k vehicles which may be used for shuttling students to the exam location. The students need to figure out k good locations to meet at. The students will then walk to the closest meeting point and then the shuttles will deliver them to the exam location. Let X_j be the location of student j , and let the exam location be at $(0, 0)$. Assume that we can walk/drive as the crow flies, i.e., by taking the shortest path between two points. Write down an appropriate objective function to minimize the total distance that the students and vehicles need to travel. Hint: your result should be similar to the regularized k -means objective, but without the squares.
- (d) (Kernel k -means) Suppose we have a dataset $\{X_i\}_{i=1}^n, X_i \in \mathbb{R}^\ell$ that we want to split into k clusters, i.e., finding the best k -means clustering (without regularization). Furthermore, suppose we know *a priori* that this data is best clustered in an impractically high-dimensional feature space \mathbb{R}^m with an appropriate metric. Fortunately, instead of having to deal with the (implicit) feature map $\Phi : \mathbb{R}^\ell \rightarrow \mathbb{R}^m$ and (implicit) distance metric¹, we have a kernel function $\kappa(X_1, X_2) = \Phi(X_1) \cdot \Phi(X_2)$ that we can compute easily on the raw samples. How should we perform the kernelized counterpart of k -means clustering?

Derive the underlined portion of this algorithm, and show your work in deriving it. The main issue is that although we define the means μ_i in the usual way, we can't ever compute Φ explicitly because it's way too big. Therefore, in the step where we determine which cluster each sample point is assigned to, we must use the kernel function κ to obtain the right result. (Review the lecture on kernels if you don't

¹Just as how the interpretation of kernels in kernelized ridge regression involves an implicit prior/regularizer as well as an implicit feature space, we can think of kernels as generally inducing an implicit distance metric as well. Think of how you would represent the squared distance between two points in terms of pairwise inner products and operations on them.

remember how that's done.)

Algorithm 1: Kernel k -means

Require: Data matrix $X \in \mathbb{R}^{n \times d}$; number of clusters K ; kernel function $\kappa(X_1, X_2)$

Ensure: Cluster $\text{class}(j)$ assigned for each sample point X_j .

function KERNEL-K-MEANS(X, K)

 Randomly initialize $\text{class}(j)$ to be an integer in $1, 2, \dots, K$ for each X_j .

while *not converged* **do**

for $i \leftarrow 1$ **to** K **do**

 Set $S_i \leftarrow \{j \in \{1, 2, \dots, n\} : \text{class}(j) = i\}$.

for $j \leftarrow 1$ **to** n **do**

 Set $\text{class}(j) \leftarrow \arg \min_k$ _____

 Return S_i for $i \leftarrow 1, 2, \dots, K$.

end function

- (e) The expression you derived may have unnecessary terms or redundant kernel computations (especially when you compute $\text{class}(j)$ for every $j \in [1, n]$). Explain how to eliminate them; that is, how to perform the computation quickly without doing irrelevant computations or redoing computations already done.

3 The Training Error of AdaBoost

Recall that in AdaBoost, our input is an $n \times d$ design matrix X with n labels $y_i = \pm 1$, and at the end of iteration T the importance of each sample is reweighted as

$$w_i^{(T+1)} = w_i^{(T)} \exp(-\beta_T y_i G_T(X_i)), \quad \text{where} \quad \beta_T = \frac{1}{2} \ln \left(\frac{1 - \text{err}_T}{\text{err}_T} \right) \quad \text{and} \quad \text{err}_T = \frac{\sum_{y_i \neq G_T(X_i)} w_i^{(T)}}{\sum_{i=1}^n w_i^{(T)}}.$$

Note that err_T is the weighted error rate of the classifier G_T . Recall that $G_T(z)$ is ± 1 for all points z , but the meta-learner has a non-binary decision function $M(z) = \sum_{t=1}^T \beta_t G_t(z)$. To classify a test point z , we calculate $M(z)$ and return its sign.

In this problem we will prove that if every learner G_t achieves 51% accuracy (that is, only slightly above random), AdaBoost will converge to zero training error. (If you get stuck on one part, move on; all five parts below can be done without solving the other parts, and part (e) is the easiest.)

- (a) We want to change the update rule to “normalize” the weights so that each iteration’s weights sum to 1; that is, $\sum_{i=1}^n w_i^{(T+1)} = 1$. That way, we can treat the weights as a discrete probability distribution over the sample points. Hence we rewrite the update rule in the form

$$w_i^{(T+1)} = \frac{w_i^{(T)} \exp(-\beta_T y_i G_T(X_i))}{Z_T} \tag{3}$$

for some scalar Z_T . Show that if $\sum_{i=1}^n w_i^{(T)} = 1$ and $\sum_{i=1}^n w_i^{(T+1)} = 1$, then

$$Z_T = 2 \sqrt{\text{err}_T(1 - \text{err}_T)}. \tag{4}$$

Hint: sum over both sides of (3), then split the right summation into misclassified points and correctly classified.

(b) The initial weights are $w_1^{(1)} = w_2^{(1)} = \dots = w_n^{(1)} = \frac{1}{n}$. Show that

$$w_i^{(T+1)} = \frac{1}{n \prod_{t=1}^T Z_t} e^{-y_i M(X_i)}. \quad (5)$$

(c) Let B (for “bad”) be the number of sample points out of n that the meta-learner classifies incorrectly. Show that

$$\sum_{i=1}^n e^{-y_i M(X_i)} \geq B. \quad (6)$$

Hint: split the summation into misclassified points and correctly classified.

(d) Use the formulas (4), (5), and (6) to show that if $\text{err}_t \leq 0.49$ for every learner G_t , then $B \rightarrow 0$ as $T \rightarrow \infty$.
Hint: (4) implies that every $Z_t < 0.9998$. How can you combine this fact with (5) and (6)?

(e) Explain why AdaBoost with short decision trees is a form of subset selection when the number of features is large.