# Assignment #4
Due: 7:59pm EST, March 19, 2021

# Homework 4: SVM, Clustering, and Ethics

## Introduction

This homework assignment will have you work with SVMs, clustering, and engage with the ethics lecture. We encourage you to read Chapters 5 and 6 of the course textbook.

Please submit the **writeup PDF to the Gradescope assignment 'HW4'**. Remember to assign pages for each question.

Please submit your **LaTeX file and code files to the Gradescope assignment 'HW4 - Supplemental'**.

**Problem 1** (Fitting an SVM by hand, 10pts)

For this problem you will solve an SVM by hand, relying on principled rules and SVM properties. For making plots, however, you are allowed to use a computer or other graphical tools.

Consider a dataset with the following 7 data points each with $x \in \mathbb{R}$ and $y \in \{-1, +1\}$ :

$$\{(x_i, y_i)\}_{i=1}^7 = \{(-3, +1), (-2, +1), (-1, -1), (0, +1), (1, -1), (2, +1), (3, +1)\}$$

Consider mapping these points to 2 dimensions using the feature vector $\phi(x) = (x, -\frac{8}{3}x^2 + \frac{2}{3}x^4)$. The hard margin classifier training problem is:

$$\min_{\mathbf{w}, w_0} \frac{1}{2} \|\mathbf{w}\|_2^2$$
$$\text{s.t.} \quad y_i(\mathbf{w}^\top \phi(x_i) + w_0) \geq 1, \ \forall i \in \{1, \ldots, n\}$$

Make sure to follow the logical structure of the questions below when composing your answers, and to justify each step.

1. Plot the transformed training data in $\mathbb{R}^2$ and draw the optimal decision boundary of the max margin classifier. You can determine this by inspection (i.e. by hand, without actually doing any calculations).

2. What is the value of the margin achieved by the optimal decision boundary found in Part 1?

3. Identify a unit vector that is orthogonal to the decision boundary.

4. Considering the discriminant $h(\phi(x); \mathbf{w}, w_0) = \mathbf{w}^\top \phi(x) + w_0$, give an expression for *all possible* $(\mathbf{w}, w_0)$ that define the optimal decision boundary from 1.1. Justify your answer.

   Hint: The boundary is where the discriminant is equal to 0. Use what you know from 1.1 and 1.3 to solve for $\mathbf{w}$ in terms of $w_0$. (If you solve this problem in this way, then $w_0$ corresponds to your free parameter to describe the set of all possible $(\mathbf{w}, w_0)$.)

5. Consider now the training problem for this dataset. Using your answers so far, what particular solution to $\mathbf{w}$ will be optimal for the optimization problem?

6. What is the corresponding optimal value of $w_0$ for the $\mathbf{w}$ found in Part 5 (use your result from Part 4 as guidance)? Substitute in these optimal values and write out the discriminant function $h(\phi(x); \mathbf{w}, w_0)$ in terms of the variable $x$ .

7. What are the support vectors of the classifier? Confirm that your solution in Part 6 makes the constraints above tight—that is, met with equality—for these support vectors.

# Solution

**Problem 2** (K-Means and HAC, 20pts)

For this problem you will implement K-Means and HAC from scratch to cluster image data. You may use `numpy` but no third-party ML implementations (eg. `scikit-learn`).

We've provided you with a subset of the MNIST dataset, a collection of handwritten digits used as a benchmark for image recognition (learn more at http://yann.lecun.com/exdb/mnist/). MNIST is widely used in supervised learning, and modern algorithms do very well.

You have been given representations of MNIST images, each of which is a $784 \times 1$ greyscale handwritten digit from 0-9. Your job is to implement K-means and HAC on MNIST, and to test whether these relatively simple algorithms can cluster similar-looking images together.

The code in `T4_P2.py` loads the images into your environment into two arrays – `large_dataset`, a 5000x784 array, will be used for K-means, while `small_dataset`, a 300x784 array, will be used for HAC. In your code, you should use the $\ell_2$ norm (i.e. Euclidean distance) as your distance metric.

**Important:** Remember to include all of your plots in your PDF submission!

**Checking your algorithms:** Instead of an Autograder file, we have provided a similar dataset, `P2_Autograder_Data`, and some visualizations, `HAC_visual` and `KMeans_visual`, for how K-means and HAC perform on this data. Run your K-means (with $K = 10$ and `np.random.seed(2)`) and HAC on this second dataset to confirm your answers against the provided visualizations. Do **not** submit the outputs generated from `P2_Autograder_Data`. Load this data with `data = np.load('P2_Autograder_Data.npy')`.

1. Starting at a random initialization and $K = 10$, plot the K-means objective function (the residual sum of squares) as a function of iterations and verify that it never increases.

2. Run K-means for 5 different restarts for different values of $K = 5, 10, 20$. Make a plot of the final K-means objective value after your algorithm converges (y-axis) v. the values of K (x-axis), with each data point having an error bar. To compute these error bars, you will use the 5 final objective values from the restarts for each $K$ to calculate a standard deviation for each $K$.

   How does the final value of the objective function and its standard deviation change with $K$? (Note: Our code takes 10 minutes to run for this part.)

3. For $K = 10$ and for 3 random restarts, show the mean image (aka the centroid) for each cluster. To render an image, use the pyplot `imshow` function. There should be 30 total images. Include all of these images as part of a single plot; your plot must fit on one page.

4. Repeat Part 3, but before running K-means, standardize or center the data such that each pixel has mean 0 and variance 1 (for any pixels with zero variance, simply divide by 1). For $K = 10$ and 3 random restarts, show the mean image (centroid) for each cluster. Again, present the 30 total images in a single plot. Compare to Part 3: How do the centroids visually differ? Why?

5. Implement HAC for min, max, and centroid-based linkages. Fit these models to the `small_dataset`. For each of these 3 linkage criteria, find the mean image for each cluster when using 10 clusters. Display these images (30 total) on a single plot. How do these centroids compare to those found with K-means? **Important Note:** For this part ONLY, you may use `scipy`'s `cdist` function to calculate Euclidean distances between every pair of points in two arrays.

6. For each of the 3 HAC linkages (max/min/centroid), plot "Distance between most recently merged clusters" (y-axis) v. "Total number of merges completed" (x-axis). Does this plot suggest that there are any natural cut points?

7. For each of the max and min HAC linkages, make a plot of "Number of images in cluster" (y-axis) v. "Cluster index" (x-axis) reflecting the assignments during the phase of the algorithm when there were $K = 10$ clusters. Intuitively, what do these plots tell you about the difference between the clusters produced by the max and min linkage criteria?

8. For your K-means with $K = 10$ model and HAC min/max/centroid models using 10 clusters on the `small_dataset` images, use the `seaborn` module's `heatmap` function to plot a confusion matrix of clusters v. actual digits. This is 4 matrices, one per method, each method vs. true labeling. The cell at the $i$th row, $j$th column of your confusion matrix is the number of times that an image with the true label of $j$ appears in cluster $i$. How well do the different approaches match the digits? Is this matching a reasonable evaluation metric for the clustering? Explain why or why not.

**Solution**

**Problem 3** (Ethics Assignment, 15pts)

Read the article "Amazon Doesn't Consider the Race of Its Customers. Should It?". Please write no more than 1 concise paragraph each in response to the below reflection questions. We do not expect you to do any outside research, though we encourage you to connect to lecture materials and the reading for the module where relevant.

1. Some people think that Amazon's process for determining which neighborhoods would receive same-day delivery was wrongfully discriminatory, but others disagree. Based on our definitions and discussions from lecture, do you believe that Amazon's same-day delivery process was wrongfully discriminatory? Explain your reasoning.

2. Basing decisions about how to treat others on social group membership often strikes us as being wrongfully discriminatory. For example, most people would say that refusing to hire someone because they are a woman is wrongful discrimination, at least under normal circumstances.

   However, there are some cases in which some people argue that social group membership *should* be taken into consideration when deciding how to treat others. The title of the article poses the question: Do you think that should Amazon consider the race of its customers in its same-day delivery processes? If so, then how?

3. There are many different technical definitions of fairness in machine learning. In this problem, we'll introduce you to the intuition behind two common definitions and invite you to reflect on their limitations.

   Say that Amazon decides to develop a new algorithm to decide its same-day delivery coverage areas. Given your machine learning expertise, Amazon hires you to help them.

   Assume for simplification that Amazon's same-day delivery coverage algorithm $f$ takes as input $x$, features about an individual Amazon user's account, and outputs binary class label $y$ for whether or not same-day delivery will be offered to user $x$. User $x$ also has (often unobserved) sensitive attributes $a$. In this example, we will assume $a$ is a binary label so that $a = 1$ if the user is not white, 0 otherwise.

   One technical notion of algorithmic fairness is called "group fairness"[a]. An algorithm satisfies group fairness with respect to protected attribute $a$ if it assigns the same proportion of positive labels to the group of white and the group of non-white Amazon users. In other words, if 50% of white users have access to same-day shipping, then 50% of non-white users should have access to same-day shipping too.

   What are some limitations or potential issues that may arise with enforcing this definition of fairness in practice? Are there ways that a classifier that satisfies group fairness may still result in discriminatory outcomes?

4. Another technical notion of algorithmic fairness is called "individual fairness"[b]. An algorithm satisfies individual fairness if for all pairs of users $x_1$ and $x_2$ that are similar *without taking into consideration their race* $a$, the algorithm will assign similar probabilities for the attaining the positive label (roughly, $x_1 \sim x_2 \Rightarrow p(x_1) \sim p(x_2)$)[c]. In other words, if two individuals have almost-identical user profiles, then they should both be eligible for same-day shipping, even if one is white and the other is non-white.

   What are some limitations or potential issues that may arise with enforcing this definition of fairness in practice? Are there ways that a classifier that satisfies individual fairness may still result in discriminatory outcomes?

---

[a] Group fairness is also sometimes referred to as "independence" or "demographic parity". https://fairmlbook.org/classification.html

[b] https://arxiv.org/pdf/1104.3913.pdf

[c] This is an intuitive description of individual fairness (likewise with group fairness above) rather than a precise formalization.

**Problem 4** (Bonus Ethics Assignment, 0pts)

*Estimated total time for completion*: 45 minutes.

In our lecture from class, we discussed philosophical and legal frameworks to examine algorithmic discrimination. But these aren't the only frameworks! A growing body of work in the humanities and social sciences, particularly in feminist studies, critical race studies, sociology, anthropology, and the history of science has emerged to study the social consequences of machine learning.

In this bonus problem, you will be introduced to one framework inspired by scholarship in science and technology studies. To complete the below questions, first watch the 28-minute 2019 NeurIPS talk "The Values of Machine Learning" by Ria Kalluri. Please write no more than 1 paragraph in response to each of the below reflection questions:

1. In their talk, Ria discussed opportunities for shifting power to each of four possible stakeholders: an input source, data source, expert, and decision-maker. Choose one of these stakeholders, and discuss one possible way we as machine learning practitioners and model-builders could shift power to them.

2. What do you think will it take to achieve a world where AI can shift power towards historically marginalized communities? What obstacles or barriers stand in between our current world and your own AI dreams?

## Solution