

CS 446 / ECE 449 — Homework 5

your NetID here

Version 1.0

Instructions.

- Homework is due **Tuesday, April 18, at noon CST**; no late homework accepted.
- Everyone must submit individually on Gradescope under **hw5** and **hw5code**.
- The “written” submission at **hw5** **must be typed**, and submitted in any format Gradescope accepts (to be safe, submit a PDF). You may use L^AT_EX, Markdown, Google Docs, MS Word, whatever you like; but it must be typed!
- When submitting at **hw5**, Gradescope will ask you to select pages for each problem; please do this precisely!
- Please make sure your NetID is clear and large on the first page of the homework.
- Your solution **must** be written in your own words. Please see the course webpage for full academic integrity information. Briefly, you may have high-level discussions with at most 3 classmates, whose NetIDs you should place on the first page of your solutions, and you should cite any external reference you use; despite all this, your solution must be written in your own words.
- We reserve the right to reduce the auto-graded score for **hw5code** if we detect funny business (e.g., your solution lacks any algorithm and hard-codes answers you obtained from someone else, or simply via trial-and-error with the autograder).
- Coding problems come with suggested “library routines”; we include these to reduce your time fishing around APIs, but you are free to use other APIs.
- When submitting to **hw5code**, only upload the two python files **hw5.py** and **hw5_utils.py**. Don’t upload a zip file or additional files.

Version history.

1.0. Initial version.

1. SVM with Biases.

This problem is about SVMs over \mathbb{R}^d with linearly separable data (i.e., the hard margin SVM).

Our formulation of SVM required separators to pass through the origin, which does not provide a geometrically pleasing notion of maximum margin direction.

A first fix is provided by lecture 18: by appending a 1 to the inputs, we obtain the convex program

$$\begin{aligned} \min_{\mathbf{u}} \quad & \frac{1}{2} \|\mathbf{u}\|^2 \\ \text{subject to} \quad & \mathbf{u} \in \mathbb{R}^{d+1} \\ & y_i [\mathbf{x}_i]^\top \mathbf{u} \geq 1 \quad \forall i, \end{aligned}$$

and let $\bar{\mathbf{u}}$ denote the optimal solution to this program.

A second standard fix is to incorporate the bias directly into the optimization problem:

$$\begin{aligned} \min_{\mathbf{v}, b} \quad & \frac{1}{2} \|\mathbf{v}\|^2 \\ \text{subject to} \quad & \mathbf{v} \in \mathbb{R}^d, b \in \mathbb{R} \\ & y_i (\mathbf{v}^\top \mathbf{x}_i + b) \geq 1 \quad \forall i, \end{aligned}$$

and let $(\bar{\mathbf{v}}, \bar{b}) \in \mathbb{R}^d \times \mathbb{R}$ denote an optimal solution to this program. This second version is standard, but we do not use it in lecture for various reasons.

- In lecture, we stated that the first formulation is a *convex program* (formally defined in lecture 19). Show that the second formulation is also a convex program.
- Suppose there is only one datapoint: $\mathbf{x}_1 = \mathbf{e}_1$, the first standard basis vector, with label $y_1 = +1$. The first formulation will have a unique solution $\bar{\mathbf{u}}$, as discussed in lecture. Show that the second formulation does not have a unique solution.
- Let's add another datapoint: $\mathbf{x}_2 = -a\mathbf{e}_1$ for some $a \geq 3$, with label $y_2 = -1$. Now that we have two data points, both of the convex programs now have two constraints. Write out the explicit constraints to the first convex program.
- The optimal solution to the first convex program is $\bar{\mathbf{u}} = \frac{1}{2}\mathbf{e}_1 + \frac{1}{2}\mathbf{e}_{d+1}$, and the optimal solution to the second convex program is $(\bar{\mathbf{v}}, \bar{b}) = (\frac{2}{a+1}\mathbf{e}_1, 1 - \frac{2}{a+1})$ (you do **not** need to prove this). Using this, determine and formally prove the limiting behavior of $\lim_{a \rightarrow \infty} \frac{1}{2} \|\bar{\mathbf{u}}_a\|^2$ and $\lim_{a \rightarrow \infty} \frac{1}{2} \|\bar{\mathbf{v}}_a\|^2$.
- Between the two versions of SVM with bias, which do you prefer? Any answer which contains at least one complete sentence will receive full credit.

Remark: Initially it may have seemed that both optimization problems have the same solutions; the purpose of this problem was to highlight that small differences in machine learning methods can lead to observably different performance.

Solution.

-
-
-
-
-

2. SVM Implementation.

Recall that the dual problem of an SVM is

$$\max_{\alpha \in \mathcal{C}} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j),$$

where the domain $\mathcal{C} = [0, \infty)^n = \{\alpha : \alpha_i \geq 0\}$ for a hard-margin SVM, and $\mathcal{C} = [0, C]^n = \{\alpha : 0 \leq \alpha_i \leq C\}$ for a soft-margin SVM. Equivalently, we can frame this as the minimization problem

$$\min_{\alpha \in \mathcal{C}} f(\alpha) := \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^n \alpha_i.$$

This can be solved by projected gradient descent, which starts from some $\alpha_0 \in \mathcal{C}$ (e.g., $\mathbf{0}$) and updates via

$$\alpha_{t+1} = \Pi_{\mathcal{C}} [\alpha_t - \eta \nabla f(\alpha_t)],$$

where $\Pi_{\mathcal{C}}[\alpha]$ is the *projection* of α onto \mathcal{C} , defined as the closest point to α in \mathcal{C} :

$$\Pi_{\mathcal{C}}[\alpha] := \arg \min_{\alpha' \in \mathcal{C}} \|\alpha' - \alpha\|_2.$$

If \mathcal{C} is convex, the projection is uniquely defined.

- (a) Prove that

$$\left(\Pi_{[0, \infty)^n}[\alpha] \right)_i = \max\{\alpha_i, 0\},$$

and

$$\left(\Pi_{[0, C]^n}[\alpha] \right)_i = \min\{\max\{0, \alpha_i\}, C\}.$$

Hint: Show that the i th component of any other $\alpha' \in \mathcal{C}$ is further from the i th component of α than the i th component of the projection is. Specifically, show that $|\alpha'_i - \alpha_i| \geq |\max\{0, \alpha_i\} - \alpha_i|$ for $\alpha' \in [0, \infty)^n$ and that $|\alpha'_i - \alpha_i| \geq |\min\{\max\{0, \alpha_i\}, C\} - \alpha_i|$ for $\alpha' \in [0, C]^n$.

- (b) Implement an `svm_solver()`, using projected gradient descent formulated as above. Initialize your α to zeros. See the docstrings in `hw5.py` for details.

Remark: Consider using the `.backward()` function in pytorch. However, then you may have to use in-place operations like `clamp_()`, otherwise the gradient information is destroyed.

Library routines: `torch.outer`, `torch.clamp`, `torch.autograd.backward`, `torch.tensor(..., requires_grad=True)`, with `torch.no_grad():`, `torch.tensor.grad.zero_`, `torch.tensor.detach`.

- (c) Implement an `svm_predictor()`, using an optimal dual solution, the training set, and an input. See the docstrings in `hw5.py` for details.

Library routines: `torch.empty`.

- (d) On the area $[-5, 5] \times [-5, 5]$, plot the contour lines of the following kernel SVMs, trained on the XOR data. Different kernels and the XOR data are provided in `hw5_utils.py`. Learning rate 0.1 and 10000 steps should be enough. To draw the contour lines, you can use `hw5_utils.svm_contour()`.

- The polynomial kernel with degree 2.
- The RBF kernel with $\sigma = 1$.
- The RBF kernel with $\sigma = 2$.
- The RBF kernel with $\sigma = 4$.

Include these four plots in your written submission.

Solution.

(a)

(d)

3. Logistic Regression.

Recall the empirical risk $\hat{\mathcal{R}}$ for logistic regression (as presented in lecture 17):

$$\hat{\mathcal{R}}_{\log}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i)).$$

Here you will minimize this risk using gradient descent.

- (a) In your written submission, derive the gradient descent update rule for this empirical risk by taking the gradient. Write your answer in terms of the learning rate η , previous parameters \mathbf{w} , new parameters \mathbf{w}' , number of examples n , and training examples (\mathbf{x}_i, y_i) . Show all of your steps.
- (b) Implement the `logistic()` function in `hw5.py`. You are given as input a training set `X`, training labels `Y`, a learning rate `lr`, and number of gradient updates `num_iter`. Implement gradient descent to find parameters \mathbf{w} that minimize the empirical risk $\hat{\mathcal{R}}_{\log}(\mathbf{w})$. Perform gradient descent for `num_iter` updates with a learning rate of `lr`, initializing $\mathbf{w} = 0$ and returning \mathbf{w} as output. Don't forget to prepend `X` with a column of ones.

Library routines: `torch.matmul` (`@`), `torch.tensor.t`, `torch.exp`.

- (c) Implement the `logistic_vs_ols()` function in `hw5.py`. Use `utils.load_logistic_data()` to generate a training set `X` and training labels `Y`. Run `logistic(X,Y)` from part (b) taking `X` and `Y` as input to obtain parameters \mathbf{w} . Also run `linear_normal(X,Y)` from Homework 4 Problem 2 to obtain parameters \mathbf{w} . Plot the decision boundaries for your logistic regression and least squares models along with the data `X`. Which model appears to classify the data better? Explain why you believe your choice is the better classifier for this problem.

Library routines: `torch.linspace`, `plt.scatter`, `plt.plot`, `plt.show`, `plt.gcf`.

Hints:

- The positive and negative points are guaranteed to be linearly separable (though an algorithm may or may not find the optimal line to separate them).
- The “decision boundary” in the problem description refers to the set of points \mathbf{x} such that $\mathbf{w}^\top \mathbf{x} = 0$ for the chosen predictor. In this case, it suffices to plot the corresponding line.
- In order to make the two models significantly different, we recommend that you train the logistic regression with a large `num_iter` (e.g., 1,000,000 or even larger).

Solution.

- (a)
- (c)

4. Robustness of the Majority Vote Classifier.

The purpose of this problem is to further investigate the behavior of the majority vote classifier using Hoeffding's inequality (*this will be covered in the lecture 22, the Ensemble Methods lecture*). Simplified versions of Hoeffding's inequality are as follows.

Theorem 1. *Given independent random variables (Z_1, \dots, Z_k) with $Z_i \in [0, 1]$,*

$$\Pr \left[\sum_{i=1}^k Z_i \geq \sum_{i=1}^k \mathbb{E}[Z_i] + k\epsilon \right] \leq \exp(-2k\epsilon^2), \quad (1)$$

and

$$\Pr \left[\sum_{i=1}^k Z_i \leq \sum_{i=1}^k \mathbb{E}[Z_i] - k\epsilon \right] \leq \exp(-2k\epsilon^2). \quad (2)$$

In this problem we have an odd number n of classifiers $\{f_1, \dots, f_n\}$ and only consider their behavior on a fixed data example (\mathbf{x}, y) ; by classifier we mean $f_i(\mathbf{x}) \in \{\pm 1\}$. Define the majority vote classifier MAJ as

$$\text{MAJ}(\mathbf{x}; f_1, \dots, f_n) := 2 \cdot \mathbb{1} \left[\sum_{i=1}^n f_i(\mathbf{x}) \geq 0 \right] - 1 = \begin{cases} +1 & \sum_{i=1}^n f_i(\mathbf{x}) > 0, \\ -1 & \sum_{i=1}^n f_i(\mathbf{x}) < 0, \end{cases}$$

where we will not need to worry about ties since n is odd.

To demonstrate the utility of Theorem 1 in analyzing MAJ, suppose that $\Pr[f_i(\mathbf{x}) = y] = p > 1/2$ independently for each i . Then, by defining a random variable $Z_i := \mathbb{1}[f_i(\mathbf{x}) \neq y]$ and noting $\mathbb{E}[Z_i] = 1 - p$,

$$\begin{aligned} \Pr[\text{MAJ}(\mathbf{x}; f_1, \dots, f_n) \neq y] &= \Pr \left[\sum_{i=1}^n \mathbb{1}[f_i(\mathbf{x}) \neq y] \geq \frac{n}{2} \right] \\ &= \Pr \left[\sum_{i=1}^n Z_i \geq n(1-p) - \frac{n}{2} + np \right] \\ &= \Pr \left[\sum_{i=1}^n Z_i \geq n\mathbb{E}[Z_1] + n(p - 1/2) \right] \\ &\leq \exp(-2n(p - 1/2)^2). \end{aligned}$$

The purpose of this problem is to study the behavior of $\text{MAJ}(\mathbf{x})$ when not all of the classifiers $\{f_1, \dots, f_n\}$ are independent.

- (a) Assume n is divisible by 7 and $5n/7$ is odd, and that of the n classifiers $\{f_1, \dots, f_n\}$, now only the first $5n/7$ of them (i.e., $\{f_1, \dots, f_{5n/7}\}$) have independent errors on \mathbf{x} . Specifically, $\Pr[f_i(\mathbf{x}) = y] = p := 4/5$ for classifiers $\{f_1, \dots, f_{5n/7}\}$. By contrast, we make no assumption on the other $2n/7$ classifiers (i.e., $\{f_{5n/7+1}, \dots, f_n\}$) and their errors. Now use Hoeffding's inequality to show that

$$\Pr \left[\sum_{i=1}^{5n/7} \mathbb{1}[f_i(\mathbf{x}) \neq y] \geq \frac{3n}{14} \right] \leq \exp\left(-\frac{n}{70}\right).$$

- (b) Continuing from (a), further show that the majority vote classifier over all n classifiers is still good, specifically showing

$$\Pr[\text{MAJ}(\mathbf{x}; f_1, \dots, f_n) \neq y] \leq \exp\left(-\frac{n}{70}\right).$$

For full points: You need to derive the inequality $\Pr[\text{MAJ}(\mathbf{x}; f_1, \dots, f_n) \neq y] \leq \exp(-n/70)$ rigorously for ANY possible behavior of the $\frac{2n}{7}$ arbitrary classifiers.

- (c) Is the probability of correctly classifying \mathbf{x} reasonably good in part (b) for large n ? Do you have any interesting observations? Any answer which contains at least one complete sentence will receive full credit.

Solution.

- (a)
- (b)
- (c)