

Instructions

Submission: Assignment submission will be via courses.uscdcn.net. By the submission date, there will be a folder set up in which you can submit your files. Please be sure to follow all directions outlined here.

You can submit multiple times, but only *the last submission* counts. That means if you finish some problems and want to submit something first and update later when you finish, that's fine. In fact you are encouraged to do this: that way, if you forget to finish the homework on time or something happens, you still get credit for whatever you have turned in.

Problem sets must be typewritten or neatly handwritten when submitted. In both cases, your submission must be a single PDF. Please also follow the rules below:

- The file should be named as `firstname_lastname_USCID.pdf` (e.g., `Joe.Doe_1234567890.pdf`).
- Do not have any spaces in your file name when uploading it.
- Please include your name and USCID in the header of the report as well.

Total points: 40 points

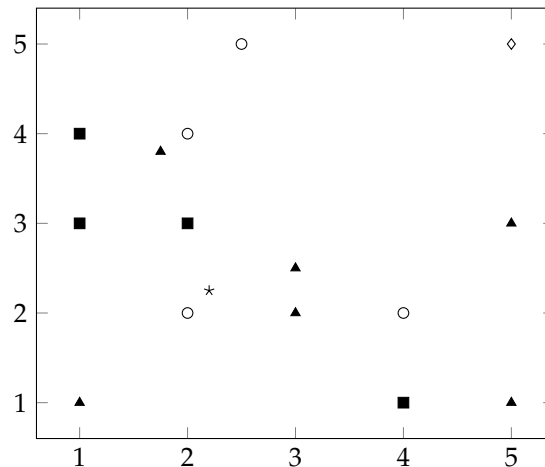
Notes on notation:

- Unless stated otherwise, scalars are denoted by small letter in normal font, vectors are denoted by small letters in bold font and matrices are denoted by capital letters in bold font.
- $\|\cdot\|$ means L2-norm unless specified otherwise i.e. $\|\cdot\| = \|\cdot\|_2$

Problem 1 Nearest Neighbor Classification

(8 points)

For the data given below, squares, triangles, and open circles are three different classes of data in the training set and the diamond (\diamond) and star (*) are test points. We denote the total number of training points as N and consider K-nearest-neighbor (KNN) classifier with L2 distance.



(a) What is the prediction for the test point *star* when $K = 4$? Explain why.

Triangle.

(1 point)

Because the 4 nearest neighbors are: one square, one circle, and two triangles.

(1 point)

(b) What is the *diamond* classified as for $K = N$? Explain why.

Triangle.

(1 point)

When $K = N$ the prediction is always the majority label of the training set. Triangle is the majority in this example.

(1 point)

(c) Suppose one performs leave-one-out validation (that is, N -fold cross validation) to choose the best hyper-parameter K . List the *triangles* that are correctly classified (as a validation point) in this process for the run with $K = 1$.

The two triangles with x -value 3.

(2 points for each)

Problem 2 Linear Regression

(16 points)

2.1 (4 points) In the class we discussed L2 regularized least square solution defined as

$$w_* = \arg \min_{w \in \mathbb{R}^D} \|Xw - y\|_2^2 + \lambda \|w\|_2^2$$

where $X \in \mathbb{R}^{N \times D}$ is the data matrix with each row corresponding to the feature of an example, $y \in \mathbb{R}^N$ is a vector of all the outcomes, $\|\cdot\|_2$ stands for the L2 norm, and λ is the regularization coefficient. In this problem we generalize the L2 regularization by

$$w'_* = \arg \min_{w \in \mathbb{R}^D} \|Xw - y\|_2^2 + w^T M w$$

where $M \in \mathbb{R}^{D \times D}$ is a positive definite matrix (L2 regularization is clearly a special case with M being the identity matrix scaled by λ).

Find the closed form of w'_* .

Setting the gradient $2X^T(Xw - y) + 2Mw$ to be $\mathbf{0}$ and using the fact that M is invertible gives

$$w'_* = (X^T X + M)^{-1} X^T y.$$

Rubrics:

- Write down the correct gradient. (2 points)
- Get the correct final answer. (2 points)

2.2 (12 points) Assume we have a training set $(x_1, y_1), \dots, (x_N, y_N) \in \mathbb{R}^D \times \mathbb{R}$, where each outcome y_n is generated by a probabilistic model $w_*^\top x_n + \epsilon_n$ with ϵ_n being an independent Gaussian noise with zero-mean and variance σ^2 for some $\sigma > 0$. In other words, the probability of seeing any outcome $y \in \mathbb{R}$ given x_n is

$$\Pr(y | x_n; w_*, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y - w_*^\top x_n)^2}{2\sigma^2}\right).$$

1. Assume σ is fixed and given, find the maximum likelihood estimation for w_* . In other words, first write down the probability of seeing the outcomes y_1, \dots, y_N given x_1, \dots, x_N as a function of the value of w_* ; then find the value of w_* that maximizes this probability. You can assume $X^\top X$ is invertible where X is the data matrix as used in Problem 2.1. (6 points)

The probability of seeing the outcomes y_1, \dots, y_N given x_1, \dots, x_N for a linear model w is

$$\mathcal{P}(w) = \prod_{n=1}^N \Pr(y_n | x_n; w, \sigma) = \prod_{n=1}^N \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y_n - w^\top x_n)^2}{2\sigma^2}\right).$$

Taking the negative log, this becomes

$$F(w) = N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - w^\top x_n)^2 = N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} \|Xw - \mathbf{y}\|_2^2.$$

Maximizing \mathcal{P} is the same as minimizing F , which is clearly the same as just minimizing $\sum_{n=1}^N (y_n - w^\top x_n)^2$, the same objective as for least square regression. Therefore the MLE for w_* is exactly the same as the least square solution:

$$w_* = (X^\top X)^{-1} X^\top \mathbf{y}.$$

Rubrics:

- Write down the correct likelihood function. (2 points)
- Any derivation revealing that this is the same as least square regression. (2 points)
- Arrive at the correct final answer. (2 points)

2. Now consider σ as a parameter of the probabilistic model too, that is, the model is specified by both \mathbf{w}_* and σ . Find the maximum likelihood estimation for \mathbf{w}_* and σ . (6 points)

From previous calculation, the MLE for \mathbf{w}_* and σ is the minimizer of the function

$$F(\mathbf{w}, \sigma) = N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2.$$

We first fix σ and minimize over \mathbf{w} , which leads to the same MLE for \mathbf{w}_* :

$$\mathbf{w}_* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

Next we minimize $F(\mathbf{w}_*, \sigma)$ as function of σ by setting the derivative w.r.t. σ to be 0:

$$\frac{\partial F(\mathbf{w}_*, \sigma)}{\partial \sigma} = \frac{N}{\sigma} - \frac{1}{\sigma^3} \|\mathbf{X}\mathbf{w}_* - \mathbf{y}\|_2^2 = 0.$$

Solving for σ gives the MLE estimate:

$$\sigma = \frac{1}{\sqrt{N}} \|\mathbf{X}\mathbf{w}_* - \mathbf{y}\|_2 = \frac{1}{\sqrt{N}} \|\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} - \mathbf{y}\|_2.$$

Rubrics:

- Write down the correct likelihood as a function of both \mathbf{w} and σ . (1 point)
- Arrive at the correct solution for \mathbf{w}_* . (1 point)
- Correct derivative with respect to σ . (2 points)
- Arrive at the correct solution for σ . (2 points)

Problem 3 Convergence of Perceptron Algorithm

(16 points)

In this problem you need to show that when the two classes are linearly separable, the perceptron algorithm will converge. Specifically, for a binary classification dataset of N data points, where every \mathbf{x}_i has a corresponding label $y_i \in \{-1, 1\}$ and is normalized: $\|\mathbf{x}_i\| = \sqrt{\mathbf{x}_i^T \mathbf{x}_i} = 1, \forall i \in \{1, 2, \dots, N\}$, the perceptron algorithm proceeds as below:

Algorithm 1 Perceptron

```
Initialize  $\mathbf{w}_1 = \mathbf{0}$ 
for  $k = 1, 2, \dots$ , do
    Pick a data point  $\mathbf{x}_i$  randomly
    Make a prediction  $y = \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$  using current  $\mathbf{w}_k$ 
    if  $y \neq y_i$  then
         $\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k + y_i \mathbf{x}_i$ 
    else
         $\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k$ 
```

Suppose there exists an unknown hyperplane \mathbf{w}_{opt} such that $\|\mathbf{w}_{opt}\| = 1$ and $y_i \mathbf{w}_{opt}^T \mathbf{x}_i \geq \gamma, \forall i \in [N]$ for some positive value γ . Note that this implies that the linear classifier defined by \mathbf{w}_{opt} perfectly classifies all the N data points.

Following the steps below, you will show that the Perceptron algorithm makes a finite number of mistakes that is at most γ^{-2} , and therefore the algorithm must converge.

3.1 Show that if the algorithm makes a mistake, the update rule moves the weight \mathbf{w}_k towards the direction of the optimal weights \mathbf{w}_{opt} . Specifically, suppose in iteration k we have $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$. Prove

$$\mathbf{w}_{k+1}^T \mathbf{w}_{opt} \geq \mathbf{w}_k^T \mathbf{w}_{opt} + \gamma. \quad (4 \text{ points})$$

$$\mathbf{w}_{k+1} = \mathbf{w}_k + y_i \mathbf{x}_i \quad (2 \text{ points})$$

$$\mathbf{w}_{k+1}^T \mathbf{w}_{opt} = \mathbf{w}_k^T \mathbf{w}_{opt} + y_i \mathbf{x}_i^T \mathbf{w}_{opt} \geq \mathbf{w}_k^T \mathbf{w}_{opt} + \gamma \quad (2 \text{ points})$$

3.2 Show that the length of the weight vector does not increase by a large amount when the algorithm makes a mistake. More specifically, if in iteration k we have $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, then

$$\|\mathbf{w}_{k+1}\|^2 \leq \|\mathbf{w}_k\|^2 + 1. \quad (4 \text{ points})$$

$$\begin{aligned} \|\mathbf{w}_{k+1}\|^2 &= \mathbf{w}_{k+1}^T \mathbf{w}_{k+1} = (\mathbf{w}_k + y_i \mathbf{x}_i)^T (\mathbf{w}_k + y_i \mathbf{x}_i) \\ &= \|\mathbf{w}_k\|^2 + 2y_i \mathbf{w}_k^T \mathbf{x}_i + y_i^2 \mathbf{x}_i^T \mathbf{x}_i \end{aligned} \quad (2 \text{ points})$$

Input \mathbf{x}_i has norm 1 and the algorithm has made a mistake so $y_i \mathbf{w}_k^T \mathbf{x}_i \leq 0$, which implies

$$\|\mathbf{w}_{k+1}\|^2 = \|\mathbf{w}_k\|^2 + 2y_i \mathbf{w}_k^T \mathbf{x}_i + y_i^2 \mathbf{x}_i^T \mathbf{x}_i \leq \|\mathbf{w}_k\|^2 + 1 \quad (2 \text{ points})$$

3.3 Using results from Problem 3.1 and 3.2, show that for any iteration $k + 1$, with M being the total number of mistakes the algorithm has made so far for the first k iterations, we have

$$M\gamma \leq \|\mathbf{w}_{k+1}\| \leq \sqrt{M}$$

Hint: use the fact $\mathbf{a}^T \mathbf{b} \leq \|\mathbf{a}\| \|\mathbf{b}\|$ for any two vectors \mathbf{a} and \mathbf{b} .

(6 points)

By repeatedly applying results from Problem 3.1 we have

$$\mathbf{w}_{k+1}^T \mathbf{w}_{opt} \geq \mathbf{w}_1^T \mathbf{w}_{opt} + M\gamma = M\gamma \quad (2 \text{ points})$$

Using the inequality from the hint, we have

$$M\gamma \leq \mathbf{w}_{k+1}^T \mathbf{w}_{opt} \leq \|\mathbf{w}_{k+1}\| \|\mathbf{w}_{opt}\| = \|\mathbf{w}_{k+1}\| \quad (2 \text{ points})$$

which proves the first inequality. Similarly, using results of Problem 3.2 repeatedly, we have

$$\|\mathbf{w}_{k+1}\|^2 \leq \|\mathbf{w}_1\|^2 + M = M \quad (2 \text{ points})$$

which proves the second inequality.

3.4 Using result of Problem 3.3, conclude $M \leq \gamma^{-2}$.

(2 points)

Simply solving $M\gamma \leq \sqrt{M}$ for M proves $M \leq \gamma^{-2}$.