

# COMS 4771-1 F19 Homework 3 (due November 6, 2019)

## Instructions

Submit your write-up on Gradescope as a neatly typeset (not scanned nor hand-written) PDF document by 2:30 PM of the due date. (Feel free to use Word, L<sup>A</sup>T<sub>E</sub>X, etc.—whatever you like.)

On Gradescope, be sure to select the pages containing your answer for each problem. More details can be found on the [Gradescope Student Workflow help page](#).

(If you don't select pages containing your answer to a problem, you'll receive a zero for that problem.)

Make sure **your name and your UNI** appear prominently on the first page of your write-up.

You are welcome and encouraged to discuss homework assignments with each other in small groups (two to three people). You must **list all discussants in your homework write-up**: also do this prominently on the first page of your write-up.

Remember, discussion about homework assignments may include brainstorming and verbally discussing possible solution approaches, but **must not go as far as one person telling others how to solve a problem**. In addition, **you must write-up your solutions by yourself**, and **you may not look at another student's homework write-up/solutions (whether partial or complete)**. The academic rules of conduct can be found in the [course syllabus](#).

## Source code

Please combine all requested source code files into a **single ZIP file**, along with a plain text file called **README** that contains your name and briefly describes all of the other files in the ZIP file. **Do not include the data files**. Submit this ZIP file on Courseworks.

## Clarity and precision

One of the goals in this class is for you to learn to reason about machine learning problems and algorithms. To demonstrate this reasoning, you must be able to make **clear** and **precise** arguments. A clear and precise argument is not the same as a long, excessively detailed argument. Unnecessary details and irrelevant side-remarks often make an argument less clear. Non-factual statements also detract from the clarity of an argument.

Points may be deducted for answers and arguments that lack sufficient clarity or precision. We will grade your answer/argument based only on a time-economical attempt to understand it.

## Problem 1 (30 points)

In this problem, you will qualitatively compare the inductive biases of different learning methods for linear regression.

Download the synthetic data sets from Courseworks (`hw3p1_train.csv`, `hw3p1_test.csv`) for this problem. In the training data ( $n = 32$  data points), the first column represents the (scalar) input feature, and the second column represents the corresponding label. In the test data ( $m = 2079$  data points), only the input feature is provided. (You'll qualitatively compare the learned functions by plotting their predictions on these test data.)

You will use the feature map  $\varphi: \mathbb{R} \rightarrow \mathbb{R}^{2L}$ , with  $L = 128$ , given by

$$\varphi(x) = [\cos(x) \quad \sin(x) \quad \cos(2x) \quad \sin(2x) \quad \cdots \quad \cos(Lx) \quad \sin(Lx)]^\top,$$

along with the following learning methods.

1. *Minimum Euclidean norm solution.* Given training data  $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R} \times \mathbb{R}$ , return the solution  $\hat{\mathbf{w}}_{\text{Euclid}}$  to the following optimization problem:

$$\begin{aligned} & \min_{\mathbf{w} \in \mathbb{R}^{2L}} \|\mathbf{w}\|_2^2 \\ & \text{subject to } \frac{1}{n} \sum_{i=1}^n (y_i - \varphi(x_i)^\top \mathbf{w}) \varphi(x_i) = \mathbf{0}. \end{aligned}$$

We have seen that this method is a form of *spectral regularization* akin to ridge regression, principal component regression, and gradient descent for the least squares objective.

2. *Minimum weighted Euclidean norm solution.* Given training data  $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R} \times \mathbb{R}$ , return the solution  $\hat{\mathbf{w}}_{\text{weighted}}$  to the following optimization problem:

$$\begin{aligned} & \min_{\mathbf{w} \in \mathbb{R}^{2L}} \sum_{j=1}^{2L} \lceil j/2 \rceil^2 w_j^2 \\ & \text{subject to } \frac{1}{n} \sum_{i=1}^n (y_i - \varphi(x_i)^\top \mathbf{w}) \varphi(x_i) = \mathbf{0}. \end{aligned}$$

Above,  $\lceil r \rceil$  denotes the smallest integer  $z$  such that  $z \geq r$ .

3. *(Weighted) Dantzig selector.* Given training data  $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R} \times \mathbb{R}$ , return the solution  $\hat{\mathbf{w}}_{\text{Dantzig}}$  to the following optimization problem:

$$\begin{aligned} & \min_{\mathbf{w} \in \mathbb{R}^{2L}} \sum_{j=1}^{2L} \lceil j/2 \rceil |w_j| \\ & \text{subject to } \frac{1}{n} \sum_{i=1}^n (y_i - \varphi(x_i)^\top \mathbf{w}) \varphi(x_i) = \mathbf{0}. \end{aligned}$$

The (original) **Dantzig selector** is a close cousin of the Lasso.

### Minimum Euclidean norm solution

First, write a Python function that takes as input  $\mathbf{A} \in \mathbb{R}^{n \times d}$  and  $\mathbf{b} \in \mathbb{R}^n$  and returns the minimum Euclidean norm solution to the system of linear equations  $\mathbf{A}^\top \mathbf{A} \mathbf{w} = \mathbf{A}^\top \mathbf{b}$ . Then use this Python function to compute  $\hat{\mathbf{w}}_{\text{Euclid}}$  for the given training data, and apply the corresponding function  $\hat{f}_{\text{Euclid}}(x) = \varphi(x)^\top \hat{\mathbf{w}}_{\text{Euclid}}$  to the given test data  $\tilde{x}_1, \dots, \tilde{x}_m$ . Plot the points  $(\tilde{x}_i, \hat{f}_{\text{Euclid}}(\tilde{x}_i))$  for  $i = 1, \dots, m$ , with lines connecting consecutive points for illustrative purposes.

## Minimum weighted Euclidean norm solution

Explain (in words and/or pseudocode) how to use the Python function from above to compute  $\hat{\mathbf{w}}_{\text{weighted}}$ . You should be able to do this simply by modifying the input (specifically, the data matrix  $\mathbf{A}$  and label vector  $\mathbf{b}$  that you supply to the Python function to compute  $\hat{\mathbf{w}}_{\text{Euclid}}$ ), calling the Python function, and then suitably modifying the output. Use this approach to compute  $\hat{\mathbf{w}}_{\text{weighted}}$  for the given training data, and apply the corresponding function  $\hat{f}_{\text{weighted}}(x) = \boldsymbol{\varphi}(x)^\top \hat{\mathbf{w}}_{\text{weighted}}$  to the given test data  $\tilde{x}_1, \dots, \tilde{x}_m$ . Plot the points  $(\tilde{x}_i, \hat{f}_{\text{weighted}}(\tilde{x}_i))$  for  $i = 1, \dots, m$ , with lines connecting consecutive points for illustrative purposes.

## Dantzig selector

The (weighted) Dantzig selector seeks the minimum *weighted*  $\ell_1$  norm solution to the normal equations. It can be obtained as the solution to a *linear program*, i.e., an optimization problem where we seek to minimize a linear function of the optimization variables subject to affine equality and inequality constraints on the optimization variables:

$$\begin{aligned} \min_{(\mathbf{v}, \mathbf{w}) \in \mathbb{R}^{2L} \times \mathbb{R}^{2L}} \quad & \sum_{j=1}^{2L} \lceil j/2 \rceil v_j \\ \text{subject to} \quad & \frac{1}{n} \sum_{i=1}^n (y_i - \boldsymbol{\varphi}(x_i)^\top \mathbf{w}) \boldsymbol{\varphi}(x_i) = \mathbf{0} \\ & -v_j + w_j \leq 0 \quad \text{for all } j = 1, \dots, 2L \\ & -v_j - w_j \leq 0 \quad \text{for all } j = 1, \dots, 2L \\ & -v_j \leq 0 \quad \text{for all } j = 1, \dots, 2L. \end{aligned}$$

Here, we have introduced “dummy variables”  $\mathbf{v} \in \mathbb{R}^{2L}$  that are intended to represent the absolute values of (“real”) optimization variables  $\mathbf{w} \in \mathbb{R}^{2L}$ . The linear constraint  $-v_j \leq 0$  ensures that  $v_j$  is non-negative. The pair of linear constraints  $-v_j + w_j \leq 0$  and  $-v_j - w_j \leq 0$  ensures that  $|w_j| \leq v_j$ . Under these constraints, minimizing  $\sum_{j=1}^{2L} \lceil j/2 \rceil v_j$  (a linear function of the optimization variables  $(\mathbf{v}, \mathbf{w})$ ) is the same as minimizing  $\sum_{j=1}^{2L} \lceil j/2 \rceil |w_j|$ , which is what we wanted.

Typical linear programming solvers will require you to supply the following as input:

- a linear objective function (represented as a vector);
- a system of equality constraints (represented by a matrix and a vector);
- a system of inequality constraints (also represented by a matrix and a vector);
- upper- and lower-bounds on the optimization variables (some of which can be  $+\infty$  or  $-\infty$ ).

The linear program described above has  $4L$  optimization variables, so the matrices you supply to the solver will each have  $4L$  columns. There are  $2L$  equality constraints, which correspond to the normal equations. The remaining constraints are inequality constraints, although the ones of the form  $-v_j \leq 0$  can be enforced using the lower-bounds on the variables.

Use a linear programming solver to compute  $\hat{\mathbf{w}}_{\text{Dantzig}}$  for the given training data, and apply the corresponding function  $\hat{f}_{\text{Dantzig}}(x) = \boldsymbol{\varphi}(x)^\top \hat{\mathbf{w}}_{\text{Dantzig}}$  to the given test data  $\tilde{x}_1, \dots, \tilde{x}_m$ . Plot the points  $(\tilde{x}_i, \hat{f}_{\text{weighted}}(\tilde{x}_i))$  for  $i = 1, \dots, m$ , with lines connecting consecutive points for illustrative purposes.

Also, compute the number of non-zero entries in each of  $\hat{\mathbf{w}}_{\text{Euclid}}$ ,  $\hat{\mathbf{w}}_{\text{weighted}}$ , and  $\hat{\mathbf{w}}_{\text{Dantzig}}$ . The weight vector  $\hat{\mathbf{w}}_{\text{Dantzig}}$  should be much sparser than the others.

## What to submit in your write-up

- (a) Brief description of your Python function for computing the minimum Euclidean norm solution to the normal equations.
- (b) Plot of  $(\tilde{x}_i, \hat{f}_{\text{Euclid}}(\tilde{x}_i))$  for  $i = 1, \dots, m$ .
- (c) Brief description of how to use the Python function from Part (a) to compute  $\hat{\mathbf{w}}_{\text{weighted}}$ .
- (d) Plot of  $(\tilde{x}_i, \hat{f}_{\text{weighted}}(\tilde{x}_i))$  for  $i = 1, \dots, m$ .
- (e) Plot of  $(\tilde{x}_i, \hat{f}_{\text{Dantzig}}(\tilde{x}_i))$  for  $i = 1, \dots, m$ .
- (f) Number of non-zero entries in each of  $\hat{\mathbf{w}}_{\text{Euclid}}$ ,  $\hat{\mathbf{w}}_{\text{weighted}}$ , and  $\hat{\mathbf{w}}_{\text{Dantzig}}$ .
- (g) (Optional.) Try to explain the difference between these plots (at least for  $\hat{f}_{\text{Euclid}}$  and  $\hat{f}_{\text{weighted}}$ ). How would you explain the inductive biases of the different methods? Does the Fourier spectrum reveal anything of interest?

*For Parts (a) and (c):*

- First, give a 2-3 sentence high-level overview of the method.
- Then give more details about the subroutines used, their purpose, etc. Feel free to enumerate them in a list.
- Conclude with any citations for external software libraries used.
- For Part (c), also very briefly explain the correctness of your procedure.

*For Parts (b), (d), and (e):*

- Make sure your plots are clearly labeled, and include them directly in your write-up.

Please submit your source code on Courseworks. For this problem, we recommend using the `numpy.linalg` library, as well as `scipy.optimize.linprog` for linear programming. You are welcome to use other software libraries if you like, provided that you include appropriate citations in your write-up.

## Problem 2 (30 points)

In this problem, you will implement a simple learning algorithm for linear classifiers, apply the learning algorithm to a simple data set, and quantitatively and qualitatively study the learned classifiers.

### 20 Newsgroups

Download the “20 Newsgroups data set” `news.mat` from Courseworks. The training feature vectors/labels and test feature vectors/labels are stored as `data/labels` and `testdata/testlabels`. You can load this data in Python as follows:

```
from scipy.io import loadmat
news = loadmat('news.mat')
# news['data'] is the matrix of training feature vectors
# news['labels'] is the vector of training labels
# news['testdata'] is the matrix of test feature vectors
# news['testlabels'] is the vector of test labels
```

Each data point corresponds to a message posted to one of 20 different newsgroups (i.e., message boards). The representation of a message is a (sparse) binary vector in  $\mathcal{X} := \{0, 1\}^d$  (for  $d := 61188$ ) that indicates the words that are present in the message. If the  $j$ -th entry in the vector is 1, it means the message contains the word that is given on the  $j$ -th line of the text file `news.vocab`. The class labels are  $\mathcal{Y} := \{1, \dots, 20\}$ , where the mapping from classes to newsgroups is in the file `news.groups` (which we won't actually need).

### Naive Bayes model

In this problem, you'll develop a classifier based on the *Naive Bayes model*, a particular statistical model for classification data.<sup>1</sup> In this model, the training data are treated as iid random variables  $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$  from a probability distribution  $P_{\theta}$ ; here  $\theta = (\pi, (\mu_y)_{y \in \mathcal{Y}})$  is the parameter vector, with  $\pi = (\pi_y)_{y \in \mathcal{Y}} \in [0, 1]^{|\mathcal{Y}|}$  and  $\mu_y = (\mu_{y,1}, \dots, \mu_{y,d}) \in [0, 1]^d$  for each  $y \in \mathcal{Y}$ . The label  $Y_i$  of the  $i$ -th training example has distribution given by

$$P_{\theta}(Y_i = y) = \pi_y, \quad \text{for all } y \in \mathcal{Y}.$$

Furthermore, for each  $y \in \mathcal{Y}$ , the conditional distribution of the feature vector  $\mathbf{X}_i$  given  $Y_i = y$  is given by

$$P_{\theta}(\mathbf{X}_i = \mathbf{x} \mid Y_i = y) = \prod_{j=1}^d \mu_{y,j}^{x_j} (1 - \mu_{y,j})^{1-x_j}, \quad \text{for all } \mathbf{x} = (x_1, \dots, x_d) \in \mathcal{X}.$$

The parameter vector  $\theta$  must, of course, satisfy  $\sum_{y \in \mathcal{Y}} \pi_y = 1$ .

Aside: The Naive Bayes model is unlike other statistical models we have encountered in this course, in the sense that it explicitly describes the probability distribution of the  $\mathbf{X}_i$ 's (specifically, the conditional distribution of  $\mathbf{X}_i$  given  $Y_i$ ). In contrast, the logistic regression model only describes the conditional distribution of  $Y_i$  given  $\mathbf{X}_i$ .

---

<sup>1</sup>Technically, *Naive Bayes* refers to a broader class of statistical models for classification data where the features are regarded as conditional independent random variables given the class label.

## Parameter estimation for Naive Bayes

Given training data  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \in \mathcal{X} \times \mathcal{Y}$ , the MLE for the parameter vector  $\boldsymbol{\theta}$  is given by

$$\begin{aligned}\hat{\pi}_y &:= \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{y_i=y\}}, \quad \text{for all } y \in \mathcal{Y}; \\ \hat{\mu}_{y,j} &:= \frac{\sum_{i=1}^n \mathbb{1}_{\{y_i=y\}} x_{i,j}}{\sum_{i=1}^n \mathbb{1}_{\{y_i=y\}}}, \quad \text{for all } y \in \mathcal{Y} \text{ and } j = 1, \dots, d;\end{aligned}$$

However, in many applications, it is problematic if some of the parameters turn out to be zero or one (and this will be the case when this model is fit to the 20 Newsgroup data using MLE). So it is common to use a different estimator for the  $\mu_{y,j}$  parameters based on a technique called *Laplace smoothing*:

$$\hat{\mu}_{y,j} := \frac{1 + \sum_{i=1}^n \mathbb{1}_{\{y_i=y\}} x_{i,j}}{2 + \sum_{i=1}^n \mathbb{1}_{\{y_i=y\}}}, \quad \text{for all } y \in \mathcal{Y} \text{ and } j = 1, \dots, d.$$

(We could also use Laplace smoothing to estimate the  $\pi_y$  parameters, but we forego it in this assignment.) Laplace smoothing has a **Bayesian interpretation**.

## Naive Bayes classifiers

Given an estimate  $\hat{\boldsymbol{\theta}}$  of the parameter vector  $\boldsymbol{\theta}$ , we can form a classifier based on the *plug-in principle*. This means we should form the classifier  $\hat{f}: \mathcal{X} \rightarrow \mathcal{Y}$  such that

$$\hat{f}(\mathbf{x}) = \arg \max_{y \in \mathcal{Y}} P_{\hat{\boldsymbol{\theta}}}(Y = y \mid \mathbf{X} = \mathbf{x}).$$

Using Bayes' rule, we have

$$\begin{aligned}P_{\hat{\boldsymbol{\theta}}}(Y = y \mid \mathbf{X} = \mathbf{x}) &= \frac{P_{\hat{\boldsymbol{\theta}}}(Y = y \wedge \mathbf{X} = \mathbf{x})}{P_{\hat{\boldsymbol{\theta}}}(\mathbf{X} = \mathbf{x})} \\ &= \frac{P_{\hat{\boldsymbol{\theta}}}(Y = y) \cdot P_{\hat{\boldsymbol{\theta}}}(\mathbf{X} = \mathbf{x} \mid Y = y)}{P_{\hat{\boldsymbol{\theta}}}(\mathbf{X} = \mathbf{x})} \\ &= \frac{\hat{\pi}_y \cdot \prod_{j=1}^d \hat{\mu}_{y,j}^{x_j} (1 - \hat{\mu}_{y,j})^{1-x_j}}{P_{\hat{\boldsymbol{\theta}}}(\mathbf{X} = \mathbf{x})}.\end{aligned}$$

The denominator  $P_{\hat{\boldsymbol{\theta}}}(\mathbf{X} = \mathbf{x})$  is irrelevant for determining the  $y \in \mathcal{Y}$  for which this conditional probability is highest. Therefore,

$$\begin{aligned}\hat{f}(\mathbf{x}) &= \arg \max_{y \in \mathcal{Y}} P_{\hat{\boldsymbol{\theta}}}(Y = y \mid \mathbf{X} = \mathbf{x}) \\ &= \arg \max_{y \in \mathcal{Y}} \hat{\pi}_y \cdot \prod_{j=1}^d \hat{\mu}_{y,j}^{x_j} (1 - \hat{\mu}_{y,j})^{1-x_j} \\ &= \arg \max_{y \in \mathcal{Y}} \ln \left( \hat{\pi}_y \cdot \prod_{j=1}^d \hat{\mu}_{y,j}^{x_j} (1 - \hat{\mu}_{y,j})^{1-x_j} \right) \\ &= \arg \max_{y \in \mathcal{Y}} \ln \hat{\pi}_y + \sum_{j=1}^d x_j \cdot \ln \hat{\mu}_{y,j} + (1 - x_j) \cdot \ln(1 - \hat{\mu}_{y,j}).\end{aligned}$$

This classifier  $\hat{f}$  is called a Naive Bayes classifier.

## Fitting the Naive Bayes model to the “20 Newsgroups” data

Write code for training and testing a classifier based on the Naive Bayes model described above. Use Laplace smoothing to estimate the  $\mu_{y,j}$  parameters, and MLE for the  $\pi_y$  parameters. You should *not* use or look at any existing implementation (e.g., those that may be provided as library functions). Using your code, train and test a classifier with the data from `news.mat`. Compute the training and test error rates.

## Binary classification and inspection

Consider the *binary* classification problem, where newsgroups  $\{1, 16, 20\}$  comprise the “negative class” (class  $-1$ ), and newsgroups  $\{17, 18, 19\}$  comprise the “positive class” (class  $+1$ ). Newsgroups  $\{1, 16, 20\}$  are “religious” topics, and newsgroups  $\{17, 18, 19\}$  are “political” topics. The training and test data sets for this problem can be downloaded from Courseworks (`news_binary.mat`). Using these data and your codes from the previous part, train and test a Naive Bayes classifier. Compute the training and test error rates.

The binary classifier you learn is ultimately an affine classifier (i.e., a linear classifier with affine feature expansion), which means it has the following form:

$$\mathbf{x} \mapsto \begin{cases} -1 & \text{if } w_0 + \sum_{j=1}^d w_j x_j \leq 0 \\ +1 & \text{if } w_0 + \sum_{j=1}^d w_j x_j > 0 \end{cases}$$

for some real numbers  $w_0, w_1, \dots, w_d$ . Determine the values of these  $w_j$ ’s for your learned classifier. Then, report the vocabulary words whose indices  $j \in \{1, \dots, d\}$  correspond to the 20 largest (i.e., most positive)  $w_j$  value, and also the vocabulary words whose indices  $j \in \{1, \dots, d\}$  correspond to the 20 smallest (i.e., most negative)  $w_j$  value. Don’t report the indices  $j$ ’s, but rather the actual vocabulary words (from `news.vocab`).

## What to submit

- (a) Training and test error rates of the 20-class Naive Bayes classifier  $\hat{f}$  on the full 20 Newsgroup data set.
- (b) Training and test error rates of the binary classifier on the modified 20 Newsgroup data set.
- (c) Two ordered lists (clearly labeled) of 20 words each as described above.

Please submit your source code on Courseworks.

### Problem 3 (30 points)

In this problem, you will implement gradient descent for logistic regression, and apply it to text classification (same data set as in Problem 2).

#### Gradient in logistic regression

Recall the logistic regression objective function, i.e., the empirical logistic loss risk of a linear function (weight vector) on training examples  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^d \times \{-1, +1\}$ :

$$\widehat{\mathcal{R}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell_{\log}(y_i \mathbf{x}_i^\top \mathbf{w})$$

where

$$\begin{aligned} \ell_{\log}(z) &= -\ln \sigma(z) && \text{(logistic loss),} \\ \sigma(z) &= \frac{1}{1 + \exp(-z)} && \text{(sigmoid function).} \end{aligned}$$

The derivative of  $\ell_{\log}$  is given by

$$\begin{aligned} \frac{d\ell_{\log}(z)}{dz} &= -\frac{1}{\sigma(z)} \cdot \frac{d\sigma(z)}{dz} \\ &= -\frac{1}{\sigma(z)} \cdot \sigma(z) \cdot \sigma(-z) \\ &= -\sigma(-z). \end{aligned}$$

Therefore, by linearity and the chain rule, the negative gradient of  $\widehat{\mathcal{R}}$  with respect to  $\mathbf{w}$  is

$$\begin{aligned} -\nabla \widehat{\mathcal{R}}(\mathbf{w}) &= -\frac{1}{n} \sum_{i=1}^n \nabla \ell_{\log}(y_i \mathbf{x}_i^\top \mathbf{w}) \\ &= -\frac{1}{n} \sum_{i=1}^n \left. \frac{d\ell_{\log}(z)}{dz} \right|_{z=y_i \mathbf{x}_i^\top \mathbf{w}} \cdot \nabla (y_i \mathbf{x}_i^\top \mathbf{w}) \\ &= \frac{1}{n} \sum_{i=1}^n \sigma(-y_i \mathbf{x}_i^\top \mathbf{w}) \cdot y_i \mathbf{x}_i. \end{aligned}$$

Now suppose  $\mathbf{A} = [\mathbf{x}_1 | \dots | \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$  and  $\mathbf{b} = [y_1 | \dots | y_n]^\top \in \mathbb{R}^n$ . (Notice that we have omitted the  $1/\sqrt{n}$  scaling that we had for least squares linear regression.) Then the negative gradient of  $\widehat{\mathcal{R}}$  can be written as

$$-\nabla \widehat{\mathcal{R}}(\mathbf{w}) = \frac{1}{n} \mathbf{A}^\top (\mathbf{b} \odot \sigma(-\mathbf{b} \odot (\mathbf{A}\mathbf{w}))),$$

where  $\mathbf{u} \odot \mathbf{v} \in \mathbb{R}^n$  is the coordinate-wise product of vectors  $\mathbf{u} \in \mathbb{R}^n$  and  $\mathbf{v} \in \mathbb{R}^n$ , and  $\sigma(\mathbf{v}) \in \mathbb{R}^n$  is the coordinate-wise application of the sigmoid function to  $\mathbf{v} \in \mathbb{R}^n$ .

#### Gradient descent for logistic regression

Write a Python function executes gradient descent to approximately minimize the empirical logistic loss risk for a given set of training examples for binary classification. The function should take as input the labeled training examples (e.g., as a data matrix  $\mathbf{A} \in \mathbb{R}^{n \times d}$  and label matrix  $\mathbf{b} \in \{-1, +1\}^n$ ), step size parameter  $\eta > 0$ , and the number of iterations  $T$ . It should execute gradient descent on the empirical logistic loss risk objective starting from an initial weight vector of  $\mathbf{w}^{(0)} = \mathbf{0}$ , for  $T$  iterations with step size  $\eta > 0$ , and return the final weight vector  $\mathbf{w}^{(T)}$ .



## Text classification with logistic regression

It turns out the training data for the binary classification task in `news_binary.mat` is linearly separable. Use your gradient descent implementation, with step size  $\eta := 1$ , to find such a linear separator based on the training data. Experimentally determine the number of iterations  $T$  needed so that the *training error rate* (i.e., empirical risk with zero-one loss, not logistic loss) of  $\mathbf{w}^{(T)}$  is zero. What is the test error rate of  $\mathbf{w}^{(T)}$ ?

Report the vocabulary words whose indices  $j \in \{1, \dots, d\}$  correspond to the 20 largest (i.e., most positive)  $w_j^{(T)}$  value, and also the vocabulary words whose indices  $j \in \{1, \dots, d\}$  correspond to the 20 smallest (i.e., most negative)  $w_j^{(T)}$  value. Don't report the indices  $j$ 's, but rather the actual vocabulary words (from `news.vocab`).

## What to submit

- (a) Number of iterations  $T$  needed to reach a training error rate zero.
- (b) Test error rate of  $\mathbf{w}^{(T)}$ .
- (c) Two ordered lists (clearly labeled) of 20 words each as described above.

Please submit your source code on Courseworks.

## Postscript

You should find that the training error rate of the linear classifier obtained via Naive Bayes (in Problem 2) is non-zero, whereas the training error rate for the linear classifier you find using gradient descent in this problem is zero. Moreover, the corresponding test error rates will have a similar comparison. We say that the linear classifier from Naive Bayes is *under-fit* to the data.

Note that both methods are based on some underlying statistical model for the data (the Naive Bayes model and the logistic regression model). Under-fitting is a sign that an underlying statistical model may not be well-suited for the data. [Here is a very good paper that discusses some possible improvements to the Naive Bayes model.](#)

## Problem 4 (10 points)

In this problem, you will discover why the non-zero eigenvalues of  $\mathbf{A}^\top \mathbf{A}$  are the same as those of  $\mathbf{A} \mathbf{A}^\top$ , and then derive the singular value decomposition.

Suppose a matrix  $\mathbf{A} \in \mathbb{R}^{n \times d}$  is given.

- (a) Suppose  $\lambda$  is a non-zero eigenvalue of  $\mathbf{A}^\top \mathbf{A}$  with corresponding eigenvector  $\mathbf{v} \neq \mathbf{0}$ . Prove that  $\lambda$  is an eigenvalue of  $\mathbf{A} \mathbf{A}^\top$ .
- (b) Suppose 0 is an eigenvalue of  $\mathbf{A}^\top \mathbf{A}$  with corresponding eigenvector  $\mathbf{v} \neq \mathbf{0}$ . Prove that  $\mathbf{v}$  is in the null space of  $\mathbf{A}$ .
- (c) Explain why the facts in Parts (a) and (b) imply that  $\mathbf{A}^\top \mathbf{A}$  and  $\mathbf{A} \mathbf{A}^\top$  have the same non-zero eigenvalues.

Now let  $\lambda_1 \geq \dots \geq \lambda_d \geq 0$  be the eigenvalues of  $\mathbf{A}^\top \mathbf{A}$ , and let  $\mathbf{v}_1, \dots, \mathbf{v}_d$  be corresponding orthonormal eigenvectors. Let  $r$  denote the number of  $\lambda_j$ 's that are (strictly) positive. Define, for each  $j = 1, \dots, r$ ,

$$\mathbf{u}_j := \frac{1}{\|\mathbf{A} \mathbf{v}_j\|_2} \mathbf{A} \mathbf{v}_j.$$

- (d) Prove that  $\mathbf{u}_1, \dots, \mathbf{u}_r$  are orthonormal.
- (e) Prove the following statement. For any two matrices  $\mathbf{M} \in \mathbb{R}^{n \times d}$  and  $\mathbf{M}' \in \mathbb{R}^{n \times d}$  and any collection of  $d$  linearly independent vectors  $\mathbf{q}_1, \dots, \mathbf{q}_d \in \mathbb{R}^d$ , if  $\mathbf{M} \mathbf{q}_i = \mathbf{M}' \mathbf{q}_i$  for all  $i = 1, \dots, d$ , then  $\mathbf{M} = \mathbf{M}'$ .
- (f) Use the strategy suggested in the claim from Part (e) to prove that

$$\mathbf{A} = \sum_{j=1}^r \sqrt{\lambda_j} \mathbf{u}_j \mathbf{v}_j^\top.$$

*Note: Answers to most of the parts of this problem can be found in the reading material. It is fine to use the reading material to help you with this problem, although we recommend you attempt it on your own first. In any case, you must write your answers in your own words, without referencing the reading material.*

## Postscript

The decomposition of  $\mathbf{A}$  in Part (f) is called the *singular value decomposition (SVD)* of  $\mathbf{A}$ . The vectors  $\mathbf{u}_1, \dots, \mathbf{u}_r$  are the *left singular vectors* of  $\mathbf{A}$ , the vectors  $\mathbf{v}_1, \dots, \mathbf{v}_r$  are the *right singular vectors* of  $\mathbf{A}$ , and the positive scalars  $\sqrt{\lambda_1}, \dots, \sqrt{\lambda_r}$  are the *singular values* of  $\mathbf{A}$  (usually denoted by  $\sigma_1, \dots, \sigma_r$ ). Every matrix has such a decomposition!

We can “extend” the decomposition beyond the rank of  $\mathbf{A}$  to have  $\min\{n, d\}$  components as follows. We extend the orthonormal bases  $\mathbf{u}_1, \dots, \mathbf{u}_r$  to  $\mathbf{u}_1, \dots, \mathbf{u}_n$  and  $\mathbf{v}_1, \dots, \mathbf{v}_r$  to  $\mathbf{v}_1, \dots, \mathbf{v}_d$  (both via Gram-Schmidt orthogonalization), and then define  $\sigma_j := 0$  for  $j > r$ . Then we can write

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$$

where  $\mathbf{U} = [\mathbf{u}_1 | \dots | \mathbf{u}_n] \in \mathbb{R}^{n \times n}$ ,  $\mathbf{V} = [\mathbf{v}_1 | \dots | \mathbf{v}_d] \in \mathbb{R}^{d \times d}$ , and  $\mathbf{\Sigma} \in \mathbb{R}^{n \times d}$  is the matrix whose  $(j, j)$ -th entry is  $\sigma_j$  for  $j = 1, \dots, \min\{n, d\}$ , and all other entries are zero. Since  $\mathbf{V}$  and  $\mathbf{U}$  are orthogonal matrices, we have  $\mathbf{A}^\top \mathbf{A} = \mathbf{V} (\mathbf{\Sigma}^\top \mathbf{\Sigma}) \mathbf{V}^\top$  and  $\mathbf{A} \mathbf{A}^\top = \mathbf{U} (\mathbf{\Sigma} \mathbf{\Sigma}^\top) \mathbf{U}^\top$ ; this reveals the eigendecompositions of  $\mathbf{A}^\top \mathbf{A}$  and  $\mathbf{A} \mathbf{A}^\top$ , respectively since  $\mathbf{\Sigma}^\top \mathbf{\Sigma}$  and  $\mathbf{\Sigma} \mathbf{\Sigma}^\top$  are diagonal matrices.

The SVD is useful for understanding some regularization methods for linear regression. For any  $\mathbf{b} \in \mathbb{R}^n$ :

- Ridge regression (with hyperparameter  $\lambda > 0$ ):

$$\hat{\mathbf{w}}_{\text{ridge},\lambda} = \sum_{j=1}^r \frac{\sigma_j}{\sigma_j^2 + \lambda} \mathbf{v}_j \mathbf{u}_j^\top \mathbf{b}.$$

- Principal component regression (with hyperparameter  $\lambda > 0$ ):

$$\hat{\mathbf{w}}_{\text{pcr},\lambda} = \sum_{j=1}^r \mathbb{1}_{\{\sigma_j^2 \geq \lambda\}} \cdot \frac{1}{\sigma_j} \mathbf{v}_j \mathbf{u}_j^\top \mathbf{b}.$$

- Minimum Euclidean norm solution to normal equations:

$$\hat{\mathbf{w}}_{\text{min. Euclidean norm}} = \sum_{j=1}^r \frac{1}{\sigma_j} \mathbf{v}_j \mathbf{u}_j^\top \mathbf{b}.$$

This final expression can be written as

$$\hat{\mathbf{w}}_{\text{min. Euclidean norm}} = \mathbf{A}^\dagger \mathbf{b}$$

where

$$\mathbf{A}^\dagger := \sum_{j=1}^r \frac{1}{\sigma_j} \mathbf{v}_j \mathbf{u}_j^\top$$

is called the *Moore-Penrose pseudoinverse* of  $\mathbf{A}$ . (You might want to use this fact on Problem 1. In NumPy, it is `numpy.linalg.pinv`.)

## Problem 5 (10 points)

In this problem, you will practice Bayesian inference with multivariate Gaussian distributions.

Consider the following statistical model for  $n$  scalar observations,  $Y_1, \dots, Y_n$ :

$$\mathbf{Y} := (Y_1, \dots, Y_n) \sim \mathcal{N}(\boldsymbol{\theta}, \mathbf{I}),$$

where  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_n) \in \mathbb{R}^n$  is the parameter vector for the model.

- (a) Give a formula for the MLE of  $\boldsymbol{\theta}$  given data  $y_1, \dots, y_n \in \mathbb{R}$ . Your answer should be given in terms of  $y_1, \dots, y_n$ . *Hint: It is very simple.*
- (b) Let us now take  $\mathcal{N}(\mathbf{0}, \tau \mathbf{I})$  to be a prior distribution for the parameter  $\boldsymbol{\theta}$ , for some fixed positive constant  $\tau > 0$ . Show that, for any data  $y_1, \dots, y_n \in \mathbb{R}$ , the posterior distribution of  $\boldsymbol{\theta}$  given  $\mathbf{Y} = (y_1, \dots, y_n)$  is a multivariate Gaussian distribution, and give explicit formulas for the mean vector and covariance matrix of that multivariate Gaussian distribution in terms of the data  $y_1, \dots, y_n$  and the prior parameter  $\tau$ . *Hint: It suffices to determine the posterior distribution of  $\theta_i$  given  $Y_i = y_i$ .*
- (c) Suppose you take the prior belief from Part (b) very seriously, and instead assume the data is generated as follows. First, we draw  $\boldsymbol{\theta}$  from  $\mathcal{N}(\mathbf{0}, \tau \mathbf{I})$ . Then, conditioned on the value of  $\boldsymbol{\theta}$ , we draw  $\mathbf{Y}$  from  $\mathcal{N}(\boldsymbol{\theta}, \mathbf{I})$ ; i.e.,  $\mathbf{Y} \mid \boldsymbol{\theta} \sim \mathcal{N}(\boldsymbol{\theta}, \mathbf{I})$ . Show that the *marginal distribution* of  $\mathbf{Y}$  is a multivariate Gaussian distribution with mean  $\mathbf{0}$  and covariance matrix  $(1 + \tau)\mathbf{I}$ . *Hint: It suffices to show that the marginal distribution of  $Y_i$  is  $\mathcal{N}(0, 1 + \tau)$ .*

## Postscript

Assume  $n \geq 3$ . It turns out that under the statistical model for  $\mathbf{Y}$  from Part (c), we have

$$\mathbb{E} \left[ \frac{n-2}{\|\mathbf{Y}\|_2^2} \right] = \frac{1}{1+\tau}.$$

The *James-Stein* estimator for  $\boldsymbol{\theta}$  (in the original model  $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\theta}, \mathbf{I})$ ) is given by

$$\left( 1 - \frac{n-2}{\|\mathbf{y}\|_2^2} \right) \mathbf{y}$$

for any  $\mathbf{y} = (y_1, \dots, y_n)$ . (One motivation for this estimator comes from the answer to Part (b).) This estimator is superior to the MLE (i.e., the answer in Part (a)) in a very surprising sense.

## About grading

This assignment will be graded out of 100 points, so you may earn up to 10 points extra credit. (Ignore the placeholder point totals on Gradescope.)