

# CS 229, Fall 2024

## Problem Set #2

Saimai Lau (smlau)

---

**Due Wednesday, October 23 at 11:59 pm on Gradescope.**

**Notes:** (1) These questions require thought, but do not require long answers. Please be as concise as possible.

(2) If you have a question about this homework, we encourage you to post your question on our Ed forum, at <https://edstem.org/us/courses/67631/discussion/>.

(3) If you missed the first lecture or are unfamiliar with the collaboration or honor code policy, please read the policy on the course website before starting work.

(4) For the coding problems, you may not use any libraries except those defined in the provided `environment.yml` file. In particular, ML-specific libraries such as scikit-learn are not permitted.

(5) The due date is Wednesday, October 23 at 11:59 pm. If you submit after Wednesday, October 23 at 11:59 pm, you will begin consuming your late days. The late day policy can be found in the course website: Course Logistics and FAQ.

All students must submit an electronic PDF version of the written question including plots generated from the codes. We highly recommend typesetting your solutions via L<sup>A</sup>T<sub>E</sub>X. All students must also submit a zip file of their source code to Gradescope, which should be created using the `make_zip.py` script. You should make sure to (1) restrict yourself to only using libraries included in the `environment.yml` file, and (2) make sure your code runs without errors. Your submission may be evaluated by the auto-grader using a private test set, or used for verifying the outputs reported in the writeup. Please make sure that your PDF file and zip file are submitted to the corresponding Gradescope assignments respectively. We reserve the right to not give any points to the written solutions if the associated code is not submitted.

**Honor code:** We strongly encourage students to form study groups. Students may discuss and work on homework problems in groups. However, each student must write down the solution independently, and without referring to written notes from the joint session. Each student must understand the solution well enough in order to reconstruct it by him/herself. It is an honor code violation to copy, refer to, or look at written or code solutions from a previous year, including but not limited to: official solutions from a previous year, solutions posted online, and solutions you or someone else may have written up in a previous year. Furthermore, it is an honor code violation to post your assignment solutions online, such as on a public git repo. We run plagiarism-detection software on your code against past solutions as well as student submissions from previous years. Please take the time to familiarize yourself with the Stanford Honor Code<sup>1</sup> and the Stanford Honor Code<sup>2</sup> as it pertains to CS courses.

---

<sup>1</sup><https://communitystandards.stanford.edu/policies-and-guidance/honor-code>

<sup>2</sup><https://web.stanford.edu/class/archive/cs/cs106b/cs106b.1164/handouts/honor-code.pdf>

### 1. [12 points] Logistic Regression: Training stability

In this problem, we will be delving deeper into the workings of logistic regression. The goal of this problem is to help you develop your skills debugging machine learning algorithms (which can be very different from debugging software in general).

We have provided an implementation of logistic regression in `src/stability/stability.py`, and two labeled datasets  $A$  and  $B$  in `src/stability/ds1_a.csv` and `src/stability/ds1_b.csv`.

Please do not modify the code for the logistic regression training algorithm for this problem. First, run the given logistic regression code to train two different models on  $A$  and  $B$ . You can run the code by simply executing `python stability.py` in the `src/stability` directory.

- (a) [2 points] What is the most notable difference in training the logistic regression model on datasets  $A$  and  $B$ ?

**Answer:** The training on dataset  $A$  converged in 30370 iterations, whereas the training on dataset  $B$  does not converge.

- (b) [5 points] Investigate why the training procedure behaves unexpectedly on dataset  $B$ , but not on  $A$ . Provide hard evidence (in the form of math, code, plots, etc.) to corroborate your hypothesis for the misbehavior. Remember, you should address why your explanation does *not* apply to  $A$ .

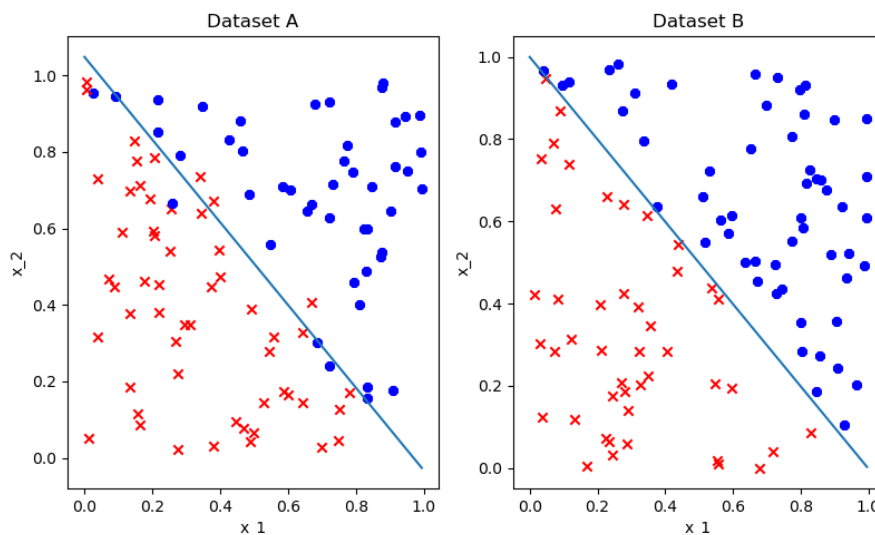
**Hint:** The issue is not a numerical rounding or over/underflow error.

**Answer:** Logistic regression is based on optimizing the hypotheses

$$h_{\theta}(x) = \frac{1}{1 + e^{\theta^T x}}$$

where  $\|\theta\|_2$  determines how quickly the prediction change between 0 and 1 when varying  $x$ . So, if the dataset is perfectly separable, the best transition would be an immediate jump from 0 to 1 across the boundary, which corresponds to  $\|\theta\| \rightarrow \infty$

Now, plotting the datasets and decision boundary  $\theta^T x = 0$ ,



As shown, dataset  $B$  can be perfectly separated by a line  $\theta^T x = 0$ , so the ideal  $\|\theta\|_2 \rightarrow \infty$ , it cannot converge. Whereas dataset  $A$  is not separable, so there is a finite  $\theta$  to converge to.

- (c) [5 points] For each of these possible modifications, state whether or not it would lead to the provided training algorithm converging on datasets such as  $B$ . Justify your answers.
- i. Using a different constant learning rate.
  - ii. Decreasing the learning rate over time (e.g. scaling the initial learning rate by  $1/t^2$ , where  $t$  is the number of gradient descent iterations thus far).
  - iii. Linear scaling of the input features.
  - iv. Adding a regularization term  $\|\theta\|_2^2$  to the loss function.
  - v. Adding zero-mean Gaussian noise to the training data or labels.

**Answer:**

- i. No, since this doesn't change the fact that there isn't a finite  $\theta_{best}$
- ii. No, since this doesn't change the fact that there isn't a finite  $\theta_{best}$
- iii. No, since the dataset would still be separable after linearly scaling.
- iv. Yes, since that would stop  $\|\theta\|_2$  from going to  $\infty$ , and the algorithm would converge when a balance is reached.
- v. Yes, since the dataset might no longer be separable after adding the noise.

## 2. [18 points] Decision Trees and Gini Loss

When growing a decision tree, we split the input space in a greedy, top-down, recursive manner. Given a parent region  $R_p$ , we can choose a split  $s_p(j, t)$  which yields two child regions  $R_1 = \{X \mid x_j < t, X \in R_p\}$  and  $R_2 = \{X \mid x_j \geq t, X \in R_p\}$ . Assuming we have defined a per region loss  $L(R)$ , at each branch we select the split that minimizes the weighted loss of the children:

$$\min_{j,t} \frac{|R_1|L(R_1) + |R_2|L(R_2)}{|R_1| + |R_2|}$$

When performing classification, a commonly used loss is the Gini loss, defined for the K-class classification problem as:

$$G(R_m) = G(\vec{p}_m) = \sum_{k=1}^K p_{mk}(1 - p_{mk})$$

Where  $\vec{p}_m = [p_{m1} \ p_{m2} \ \dots \ p_{mK}]$  and  $p_{mk}$  is the proportion of examples of class  $k$  that are present in region  $R_m$ . However, we are oftentimes more interested in optimizing the final misclassification loss:

$$M(R_m) = M(\vec{p}_m) = 1 - \max_k p_{mk} \quad (1)$$

For the problems below, assume we are dealing with binary classification and that there are no degenerate cases where positive and negative datapoints overlap in the feature space.

- (a) [5 points] Show that for any given split, the weighted Gini loss of the children can not exceed that of the parent. (**Hint:** first show that the Gini loss is strictly concave. And then use the fact that G is strictly concave meaning:

$$\forall p_1 \neq p_2, \forall t \in (0, 1) : G(tp_1 + (1-t)p_2) > tG(p_1) + (1-t)G(p_2)$$

**Answer:**

$$\forall k \in \{1, \dots, K\}, \frac{\partial^2 G(\vec{p}_m)}{\partial p_{mk}^2} = -2 < 0 \text{ and } \forall k \neq r, \frac{\partial^2 G(\vec{p}_m)}{\partial p_{mk} \partial p_{mr}} = 0$$

$\therefore$  The Gini loss is strictly concave.

i.e.  $\forall p_1 \neq p_2, \forall t \in [0, 1] : G(tp_1 + (1-t)p_2) \geq tG(p_1) + (1-t)G(p_2)$

Now, for any given split,

weighted Gini loss of the children =  $\frac{|R_1|G(R_1) + |R_2|G(R_2)}{|R_1| + |R_2|} = \frac{|R_1|}{|R_p|}G(\vec{p}_1) + \frac{|R_2|}{|R_p|}G(\vec{p}_2),$

where  $\frac{|R_1|}{|R_p|} \in [0, 1]$  and  $\frac{|R_1|}{|R_p|} + \frac{|R_2|}{|R_p|} = 1$

$\therefore$  Gini loss of the parent =  $G(\vec{p}_p)$

$$\begin{aligned} &= G\left(\frac{|R_1|}{|R_p|}\vec{p}_1 + \frac{|R_2|}{|R_p|}\vec{p}_2\right) \\ &\geq \frac{|R_1|}{|R_p|}G(\vec{p}_1) + \frac{|R_2|}{|R_p|}G(\vec{p}_2) = \text{Weighted Gini loss of the children} \end{aligned}$$

Hence, (Gini loss of the parent)  $\geq$  (weighted Gini loss of the children) for any given split.

- (b) [5 points] List out the cases where Gini loss will stay the same after a split. Show why these do not violate the strong concavity of the Gini loss. Briefly explain why these cases do not prevent a fully grown tree from achieving zero Gini loss. (**Hint:** Recall the definition of strict concavity).

**Answer:**

1.  $R_1 = R_p$  and  $|R_2| = 0$ , then  $\frac{|R_1|}{|R_p|} = 1 \notin (0, 1)$ , does not contradict with the strong concavity.
2.  $R_2 = R_p$  and  $|R_1| = 0$ , then  $\frac{|R_2|}{|R_p|} = 1 \notin (0, 1)$ , does not contradict with the strong concavity.

Further splits down the tree can still decrease the Gini loss.

- (c) [4 points] If instead we use misclassification loss, what additional case causes the loss to stay the same after a split? Show why this is (hint: you may find it useful to define  $N_m = |R_m|$  and  $N_{mk}$  as the number of examples of class  $k$  present in  $R_m$ ).

**Answer:**

Define  $N_m = |R_m|$  and  $N_{mk}$  as the number of examples of class  $k$  present in  $R_m$ .

Additional case:  $\min_k(N_{pk}) = \min_k(N_{1k}) + \min_k(N_{2k})$ . i.e. Both children regions have the same dominating class as the parent, since binary.

Then, parent misclassification loss =  $\frac{\min_k(N_{pk})}{N_p} = \frac{\min_k(N_{1k}) + \min_k(N_{2k})}{N_1 + N_2} = \text{weighted children misclassification loss}$

- (d) [4 points] Consider a training set  $X$ . In bootstrap sampling, each time we draw a random sample  $Z$  of size  $N$  from the training data and obtain  $Z_1, Z_2, \dots, Z_B$  after  $B$  times, i.e. we generate  $B$  different bootstrapped training data sets. If we apply bagging to regression trees, each time a tree  $T_i (i = 1, 2, \dots, B)$  is grown based on the bootstrapped data  $Z_i$ , and we average all the predictions to get:

$$\hat{T}(x) = \frac{1}{B} \sum_{i=1}^B T_i(x)$$

Now, if  $T_1, T_2, \dots, T_B$  is independent from each other, but each has the same variance  $\sigma^2$ , the variance of the average  $\hat{T}$  is  $\sigma^2/B$ . However, in practice, the bagged trees could be similar to each other, resulting in correlated predictions. Assume  $T_1, T_2, \dots, T_B$  still share the same variance  $\sigma^2$ , but have a positive pair-wise correlation  $\rho$ . We define the correlation between two random variables as:

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$$

Thus, we have  $\rho = \text{Corr}(T_i(x), T_j(x)), i \neq j$ .

Show that in this case, the variance of the average is given by:

$$\text{Var}\left(\frac{1}{B} \sum_{i=1}^B T_i(x)\right) = \rho\sigma^2 + \frac{1-\rho}{B}\sigma^2$$

**Answer:**

$$\begin{aligned}
\text{Var}\left(\frac{1}{B} \sum_{i=1}^B T_i(x)\right) &= \frac{1}{B^2} \left( \sum_{i=1}^B \sum_{j=1}^B \text{Cov}(T_i(x), T_j(x)) \right) \\
&= \frac{1}{B^2} \left( B \cdot \text{Var}(T_1(x)) + \sum_{i=1, j=1, i \neq j}^B \text{Cov}(T_i(x), T_j(x)) \right) \\
&= \frac{1}{B^2} \left( B\sigma^2 + \sum_{i=1, j=1, i \neq j}^B \rho\sigma\sigma \right) \\
&= \frac{1}{B^2} (B\sigma^2 + B(B-1)\rho\sigma^2) \\
&= \frac{\sigma^2 + B\rho\sigma^2 - \rho\sigma^2}{B} \\
&= \rho\sigma^2 + \frac{1-\rho}{B}\sigma^2
\end{aligned}$$

### 3. [20 points] AdaBoost Performance

We learned about boosting in lecture. Statistician Kevin Murphy claims that “It can be shown that, as long as each base learner has an accuracy that is better than chance (even on the weighted dataset), then the final ensemble of classifiers will have higher accuracy than any given component.” We will now verify this in the AdaBoost framework.

- (a) [3 points] Given a set of  $n$  observations  $(x_i, y_i)$  where  $y_i$  is the label  $y_i \in \{-1, 1\}$ , let  $f_t(x)$  be the weak classifier at step  $t$  and let  $\hat{w}_t$  be its weight. First we note that the final classifier after  $T$  steps is defined as

$$F(x) = \text{sign} \left\{ \sum_{t=1}^T \hat{w}_t f_t(x) \right\} = \text{sign}\{f(x)\},$$

where

$$f(x) = \sum_{t=1}^T \hat{w}_t f_t(x).$$

We can assume that  $f(x)$  is never exactly zero.

Show that

$$\varepsilon_{\text{training}} := \frac{1}{n} \sum_{i=1}^n 1_{\{F(x_i) \neq y_i\}} \leq \frac{1}{n} \sum_{i=1}^n \exp(-f(x_i)y_i),$$

where  $1_{\{F(x_i) \neq y_i\}}$  is 1 if  $F(x_i) \neq y_i$  and 0 otherwise.

**Answer:**

For any  $i = 1, \dots, n$ ,

if  $F(x_i) = y_i$ ,

$$1_{\{F(x_i) \neq y_i\}} = 0 < \exp(-f(x_i)y_i) \text{ since } \exp() > 0 \text{ always true.}$$

if  $F(x_i) \neq y_i$ ,

$$-f(x_i)y_i > 0, \implies \exp(-f(x_i)y_i) > \exp(0) = 1$$

$$\therefore 1_{\{F(x_i) \neq y_i\}} = 1 < \exp(-f(x_i)y_i)$$

Hence,  $\frac{1}{n} \sum_{i=1}^n 1_{\{F(x_i) \neq y_i\}} \leq \frac{1}{n} \sum_{i=1}^n \exp(-f(x_i)y_i)$

- (b) [8 points] The weight for each data point  $i$  at step  $t+1$  can be defined recursively by

$$\alpha_{i,(t+1)} = \frac{\alpha_{i,t} \exp(-\hat{w}_t f_t(x_i)y_i)}{Z_t},$$

where  $Z_t$  is a normalizing constant ensuring the weights sum to 1

$$Z_t = \sum_{i=1}^n \alpha_{i,t} \exp(-\hat{w}_t f_t(x_i)y_i).$$

Show that

$$\frac{1}{n} \sum_{i=1}^n \exp(-f(x_i)y_i) = \prod_{t=1}^T Z_t.$$

**Answer:**

$$\begin{aligned}
\alpha_{i,1} &= \frac{1}{n} \\
\alpha_{i,2} &= \frac{\frac{1}{n} \exp(-\hat{w}_1 f_1(x_i) y_i)}{Z_1} \\
\alpha_{i,3} &= \frac{\frac{1}{n} \exp(-\hat{w}_1 f_1(x_i) y_i) \exp(-\hat{w}_2 f_2(x_i) y_i)}{Z_1 Z_2} \\
&\vdots \\
\alpha_{i,T} &= \frac{\frac{1}{n} \prod_{t=1}^{T-1} \exp(-\hat{w}_t f_t(x_i) y_i)}{\prod_{t=1}^{T-1} Z_t} \\
\therefore Z_T &= \sum_{i=1}^n \alpha_{i,T} \exp(-\hat{w}_T f_T(x_i) y_i) \\
&= \sum_{i=1}^n \frac{\frac{1}{n} \prod_{t=1}^T \exp(-\hat{w}_t f_t(x_i) y_i)}{\prod_{t=1}^{T-1} Z_t} \\
\therefore \prod_{t=1}^T Z_t &= \frac{1}{n} \sum_{i=1}^n \prod_{t=1}^T \exp(-\hat{w}_t f_t(x_i) y_i) \\
&= \frac{1}{n} \sum_{i=1}^n \exp\left(-\sum_{t=1}^T \hat{w}_t f_t(x_i) y_i\right) \\
&= \frac{1}{n} \sum_{i=1}^n \exp(-f(x_i) y_i)
\end{aligned}$$

- (c) [9 points] We showed above that training error is bounded above by  $\prod_{t=1}^T Z_t$ . At step  $t$  the values  $Z_1, Z_2, \dots, Z_{t-1}$  are already fixed therefore at step  $t$  we can choose  $\alpha_t$  to minimize  $Z_t$ . Let

$$\varepsilon_t = \sum_{i=1}^n \alpha_{i,t} 1_{\{f_t(x_i) \neq y_i\}}$$

be the weighted training error for the weak classifier  $f_t(x)$ . Then we can re-write the formula for  $Z_t$  as

$$Z_t = (1 - \varepsilon_t) \exp(-\hat{w}_t) + \varepsilon_t \exp(\hat{w}_t).$$

- (i) [3 points] First find the value of  $\hat{w}_t$  that minimizes  $Z_t$ . Then show that the corresponding optimal value is

$$Z_t^{\text{opt}} = 2\sqrt{\varepsilon_t(1 - \varepsilon_t)}.$$

- (ii) [3 points] Assume we choose  $Z_t$  this way. Then re-write  $\varepsilon_t = 1/2 - \gamma_t$ , where  $\gamma_t > 0$  implies better than random and  $\gamma_t < 0$  implies worse than random. Then show that

$$Z_t \leq \exp(-2\gamma_t^2).$$

(You may want to use the fact that  $\log(1 - x) \leq -x$  for  $0 \leq x < 1$ .)



- (iii) [3 points] Finally, show that if each classifier is better than random, i.e.,  $\gamma_t > \gamma$  for all  $t$  and  $\gamma > 0$ , then

$$\epsilon_{\text{training}} \leq \exp(-2T\gamma^2),$$

which shows that the training error can be made arbitrarily small with enough steps.

**Answer:**

(i)

$$\begin{aligned} \frac{\partial Z_t}{\partial \hat{w}_t} &= -(1 - \epsilon_t) \exp(-\hat{w}_t) + \epsilon_t \exp(\hat{w}_t) \\ \therefore (1 - \epsilon_t) \exp(-\hat{w}_t^{\text{opt}}) &= \epsilon_t \exp(\hat{w}_t^{\text{opt}}) \\ \exp(\hat{w}_t^{\text{opt}}) &= \sqrt{\frac{(1 - \epsilon_t)}{\epsilon_t}} \\ \hat{w}_t^{\text{opt}} &= \frac{\log(1 - \epsilon_t) - \log(\epsilon_t)}{2} \end{aligned}$$

$$\begin{aligned} \therefore Z_t^{\text{opt}} &= (1 - \epsilon_t) \sqrt{\frac{\epsilon_t}{(1 - \epsilon_t)}} + \epsilon_t \sqrt{\frac{(1 - \epsilon_t)}{\epsilon_t}} \\ &= 2\sqrt{\epsilon_t(1 - \epsilon_t)} \end{aligned}$$

(ii)

$$\begin{aligned} Z_t &= 2\sqrt{(1/2 - \gamma_t)(1 - (1/2 - \gamma_t))} \\ &= 2\sqrt{\frac{1}{4} - \gamma_t^2} \\ &= \exp(\log(2) + \frac{1}{2} \log(1 - 4\gamma_t^2) - \frac{1}{2} \log(4)) \\ &= \exp(\frac{1}{2} \log(1 - 4\gamma_t^2)) \\ &= \sqrt{\exp(\log(1 - 4\gamma_t^2))} \\ &\leq \sqrt{\exp(-4\gamma_t^2)} \\ &= \exp(-2\gamma_t^2) \end{aligned}$$

(iii)

$$\begin{aligned} \epsilon_{\text{training}} &\leq \prod_{t=1}^T Z_t \\ &\leq \prod_{t=1}^T \exp(-2\gamma_t^2) \\ &\leq \prod_{t=1}^T \exp(-2\gamma^2) \\ &= \exp(-2T\gamma^2) \end{aligned}$$

#### 4. [20 points] Spam classification

In this problem, we will use the naive Bayes algorithm to build a spam classifier.

In recent years, spam on electronic media has been a growing concern. Here, we'll build a classifier to distinguish between real messages, and spam messages. For this class, we will be building a classifier to detect SMS spam messages. We will be using an SMS spam dataset developed by Tiago A. Almeida and José María Gómez Hidalgo which is publicly available on <http://www.dt.fee.unicamp.br/~tiago/smsspamcollection><sup>3</sup>

We have split this dataset into training and testing sets and have included them in this assignment as `src/spam/spam_train.tsv` and `src/spam/spam_test.tsv`. See `src/spam/spam_readme.txt` for more details about this dataset. Please refrain from redistributing these dataset files. The goal of this assignment is to build a classifier from scratch that can tell the difference the spam and non-spam messages using the text of the SMS message.

- (a) [5 points] Implement code for processing the the spam messages into numpy arrays that can be fed into machine learning models. Do this by completing the `get_words`, `create_dictionary`, and `transform_text` functions within our provided `src/spam.py`. Do note the corresponding comments for each function for instructions on what specific processing is required.

The provided code will then run your functions and save the resulting dictionary into `spam_dictionary` and a sample of the resulting training matrix into `spam_sample_train_matrix`. In your writeup, report the vocabular size after the pre-processing step. You do not need to include any other output for this subquestion.

**Answer:** Size of dictionary: 1758

- (b) [10 points] In this question you are going to implement a naive Bayes classifier for spam classification with **multinomial event model** and Laplace smoothing.

Code your implementation by completing the `fit_naive_bayes_model` and `predict_from_naive_bayes_model` functions in `src/spam/spam.py`.

Now `src/spam/spam.py` should be able to train a Naive Bayes model, compute your prediction accuracy and then save your resulting predictions to `spam_naive_bayes_predictions`.

In your writeup, report the accuracy of the trained model on the **test set**.

**Remark.** If you implement naive Bayes the straightforward way, you will find that the computed  $p(x|y) = \prod_i p(x_i|y)$  often equals zero. This is because  $p(x|y)$ , which is the product of many numbers less than one, is a very small number. The standard computer representation of real numbers cannot handle numbers that are too small, and instead rounds them off to zero. (This is called “underflow.”) You’ll have to find a way to compute Naive Bayes’ predicted class labels without explicitly representing very small numbers such as  $p(x|y)$ . [**Hint:** Think about using logarithms.]

**Answer:** Naive Bayes had an accuracy of 0.978494623655914 on the testing set

- (c) [5 points] Intuitively, some tokens may be particularly indicative of an SMS being in a particular class. We can try to get an informal sense of how indicative token  $i$  is for the SPAM class by looking at:

$$\log \frac{p(x_j = i \mid y = 1)}{p(x_j = i \mid y = 0)} = \log \left( \frac{P(\text{token } i \mid \text{email is SPAM})}{P(\text{token } i \mid \text{email is NOTSPAM})} \right).$$

<sup>3</sup>Almeida, T.A., Gómez Hidalgo, J.M., Yamakami, A. Contributions to the Study of SMS Spam Filtering: New Collection and Results. Proceedings of the 2011 ACM Symposium on Document Engineering (DOCENG’11), Mountain View, CA, USA, 2011.

Complete the `get_top_five_naive_bayes_words` function within the provided code using the above formula in order to obtain the 5 most indicative tokens. Report the top five words in your writeup.

**Answer:**

The top 5 indicative words for Naive Bayes are: ['claim', 'won', 'prize', 'tone', 'urgent!']

### 5. [25 points] Linear Classifiers (logistic regression and GDA)

In this problem, we cover two probabilistic linear classifiers we have covered in class so far. First, a discriminative linear classifier: logistic regression. Second, a generative linear classifier: Gaussian discriminant analysis (GDA). Both of the algorithms find a linear decision boundary that separates the data into two classes, but make different assumptions. Our goal in this problem is to get a deeper understanding of the similarities and differences (and, strengths and weaknesses) of these two algorithms.

For this problem, we will consider two datasets, along with starter codes provided in the following files:

- `src/linearclass/ds1_{train,valid}.csv`
- `src/linearclass/ds2_{train,valid}.csv`
- `src/linearclass/logreg.py`
- `src/linearclass/gda.py`

Each file contains  $n$  examples, one example  $(x^{(i)}, y^{(i)})$  per row. In particular, the  $i$ -th row contains columns  $x_1^{(i)} \in \mathbb{R}$ ,  $x_2^{(i)} \in \mathbb{R}$ , and  $y^{(i)} \in \{0, 1\}$ . In the subproblems that follow, we will investigate using logistic regression and Gaussian discriminant analysis (GDA) to perform binary classification on these two datasets.

Typically, a trained model is evaluated by its performance on the validation dataset. The validation dataset is a set of examples drawn from the same (or a similar) distribution as the training data. Intuitively, this is because we need the trained model to correctly predict the label for not only the training data, but also new samples from the same distribution.

- (a) [5 points] Recall that in GDA we model the joint distribution of  $(x, y)$  by the following equations:

$$p(y) = \begin{cases} \phi & \text{if } y = 1 \\ 1 - \phi & \text{if } y = 0 \end{cases} \quad (2)$$

$$p(x|y=0) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1}(x - \mu_0)\right) \quad (3)$$

$$p(x|y=1) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)\right),$$

where  $\phi$ ,  $\mu_0$ ,  $\mu_1$ , and  $\Sigma$  are the parameters of our model.

Suppose we have already fit  $\phi$ ,  $\mu_0$ ,  $\mu_1$ , and  $\Sigma$ , and now want to predict  $y$  given a new point  $x$ . To show that GDA results in a classifier that has a linear decision boundary, show the posterior distribution can be written as

$$p(y=1 | x; \phi, \mu_0, \mu_1, \Sigma) = \frac{1}{1 + \exp(-(\theta^T x + \theta_0))},$$

where  $\theta \in \mathbb{R}^d$  and  $\theta_0 \in \mathbb{R}$  are appropriate functions of  $\phi$ ,  $\Sigma$ ,  $\mu_0$ , and  $\mu_1$ . State the value of  $\theta$  and  $\theta_0$  as a function of  $\phi, \mu_0, \mu_1, \Sigma$  explicitly.

**Answer:**

$$\begin{aligned}
 p(y = 1 \mid x) &= \frac{p(x \mid y = 1)p(y = 1)}{p(x \mid y = 1)p(y = 1) + p(x \mid y = 0)p(y = 0)} \\
 &= \frac{\exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)\right) \phi}{\exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)\right) \phi + \exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1}(x - \mu_0)\right) (1 - \phi)} \\
 &= \frac{1}{1 + \exp\left(-\frac{1}{2}((x - \mu_0)^T \Sigma^{-1}(x - \mu_0) - (x - \mu_1)^T \Sigma^{-1}(x - \mu_1))\right) \frac{1 - \phi}{\phi}} \\
 &= \frac{1}{1 + \exp\left(-\frac{1}{2} \sum_{i,j=1}^m \left\{ \Sigma_{ij}^{-1}[(x - \mu_0)_i(x - \mu_0)_j - (x - \mu_1)_i(x - \mu_1)_j] \right\} + \log\left(\frac{1 - \phi}{\phi}\right)\right)}
 \end{aligned}$$

where,

$$\begin{aligned}
 &\sum_{i,j=1}^m \left\{ \Sigma_{ij}^{-1}[(x - \mu_0)_i(x - \mu_0)_j - (x - \mu_1)_i(x - \mu_1)_j] \right\} \\
 &= \sum_{i,j=1}^m \left\{ \Sigma_{ij}^{-1}[x_i x_j - x_i \mu_{0j} - x_j \mu_{0i} + \mu_{0i} \mu_{0j} - x_i x_j + x_i \mu_{1j} + x_j \mu_{1i}] - \mu_{1i} \mu_{1j} \right\} \\
 &= \sum_{i,j=1}^m \left\{ \Sigma_{ij}^{-1}[x_i(\mu_{1j} - \mu_{0j}) + x_j(\mu_{1i} - \mu_{0i}) + \mu_{0i} \mu_{0j} - \mu_{1i} \mu_{1j}] \right\} \\
 &= \mu_0^T \Sigma^{-1} \mu_0 + \mu_1^T \Sigma^{-1} \mu_1 + (\mu_1 - \mu_0)^T \Sigma^{-1} x + (\mu_1 - \mu_0)^T (\Sigma^{-1})^T x \\
 &= \mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1 + \left( (\Sigma^{-1})^T (\mu_1 - \mu_0) + \Sigma^{-1}(\mu_1 - \mu_0) \right)^T x
 \end{aligned}$$

$$\begin{aligned}
 \therefore p(y = 1 \mid x) &= \frac{1}{1 + \exp\left(-\left\{ \frac{1}{2} \left( (\Sigma^{-1})^T (\mu_1 - \mu_0) + \Sigma^{-1}(\mu_1 - \mu_0) \right)^T x + \left( \frac{1}{2}(\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1) - \log\left(\frac{1 - \phi}{\phi}\right) \right) \right\}\right)} \\
 i.e. \begin{cases} \theta = \Sigma^{-1}(\mu_1 - \mu_0) \\ \theta_0 = \frac{1}{2} (\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1) - \log\left(\frac{1 - \phi}{\phi}\right) \end{cases}
 \end{aligned}$$

- (b) [7 points] Given the dataset, we claim that the maximum likelihood estimates of the parameters are given by

$$\phi = \frac{1}{n} \sum_{i=1}^n 1\{y^{(i)} = 1\} \quad (4)$$

$$\mu_0 = \frac{\sum_{i=1}^n 1\{y^{(i)} = 0\} x^{(i)}}{\sum_{i=1}^n 1\{y^{(i)} = 0\}} \quad (5)$$

$$\mu_1 = \frac{\sum_{i=1}^n 1\{y^{(i)} = 1\} x^{(i)}}{\sum_{i=1}^n 1\{y^{(i)} = 1\}} \quad (6)$$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T$$

The log-likelihood of the data is

$$\begin{aligned}\ell(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^n p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \\ &= \log \prod_{i=1}^n p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi).\end{aligned}\tag{7}$$

By maximizing  $\ell$  with respect to the four parameters, prove that the maximum likelihood estimates of  $\phi$ ,  $\mu_0$ ,  $\mu_1$ , and  $\Sigma$  are indeed as given in the formulas above. (You may assume that there is at least one positive and one negative example, so that the denominators in the definitions of  $\mu_0$  and  $\mu_1$  above are non-zero.)

**Answer:**

$$\begin{aligned}\ell(\phi, \mu_0, \mu_1, \Sigma) &= \sum_{i=1}^n [\log p(x^{(i)} | y^{(i)}) + \log p(y^{(i)})] \\ &= \sum_{i=1}^n \left\{ y^{(i)} [\log p(x^{(i)} | y^{(i)} = 1) + \log \phi] + (1 - y^{(i)}) [\log p(x^{(i)} | y^{(i)} = 0) + \log(1 - \phi)] \right\}\end{aligned}$$

$$\begin{aligned}\frac{\partial \ell}{\partial \phi} &= \sum_{i=1}^n \left( \frac{y^{(i)}}{\phi} - \frac{1 - y^{(i)}}{1 - \phi} \right) \\ &= \frac{1}{\phi} \sum_{i=1}^n 1\{y^{(i)} = 1\} - \frac{1}{1 - \phi} \left( n - \sum_{i=1}^n 1\{y^{(i)} = 1\} \right)\end{aligned}$$

$$\therefore \text{Optimal condition: } \frac{1}{\phi} \sum_{i=1}^n 1\{y^{(i)} = 1\} = \frac{1}{1 - \phi} \left( n - \sum_{i=1}^n 1\{y^{(i)} = 1\} \right)$$

$$\left( \frac{1}{\phi} - 1 + 1 \right) \sum_{i=1}^n 1\{y^{(i)} = 1\} = n$$

$$\phi = \frac{1}{n} \sum_{i=1}^n 1\{y^{(i)} = 1\}$$

$$\begin{aligned}\nabla_{\mu_0} \ell &= \sum_{i=1}^n \frac{1 - y^{(i)}}{p(x^{(i)} | y^{(i)} = 0)} \nabla_{\mu_0} p(x^{(i)} | y^{(i)} = 0) \\ &= \sum_{i=1}^n (1 - y^{(i)}) [\Sigma^{-1} (x^{(i)} - \mu_0)]\end{aligned}$$

$$\therefore \text{Optimal condition: } \sum_{i=1}^n 1\{y^{(i)} = 0\} [(x^{(i)} - \mu_0)] = 0$$

$$\sum_{i=1}^n 1\{y^{(i)} = 0\} x^{(i)} = \mu_0 \sum_{i=1}^n 1\{y^{(i)} = 0\}$$

$$\mu_0 = \frac{\sum_{i=1}^n 1\{y^{(i)} = 0\} x^{(i)}}{\sum_{i=1}^n 1\{y^{(i)} = 0\}}$$

$$\begin{aligned}\nabla_{\mu_1} \ell &= \sum_{i=1}^n \frac{y^{(i)}}{p(x^{(i)} | y^{(i)} = 1)} \nabla_{\mu_1} p(x^{(i)} | y^{(i)} = 1) \\ &= \sum_{i=1}^n (y^{(i)}) [\Sigma^{-1} (x^{(i)} - \mu_1)]\end{aligned}$$

$$\therefore \text{Optimal condition: } \sum_{i=1}^n 1\{y^{(i)} = 1\} [(x^{(i)} - \mu_1)] = 0$$

$$\begin{aligned}\sum_{i=1}^n 1\{y^{(i)} = 1\} x^{(i)} &= \mu_1 \sum_{i=1}^n 1\{y^{(i)} = 1\} \\ \mu_1 &= \frac{\sum_{i=1}^n 1\{y^{(i)} = 1\} x^{(i)}}{\sum_{i=1}^n 1\{y^{(i)} = 1\}}\end{aligned}$$

$$\frac{\partial \ell}{\partial \Sigma_{jk}} = \sum_{i=1}^n \left\{ \frac{y^{(i)}}{p(x^{(i)} | y^{(i)} = 1)} \frac{\partial}{\partial \Sigma_{jk}} p(x^{(i)} | y^{(i)} = 1) + \frac{1 - y^{(i)}}{p(x^{(i)} | y^{(i)} = 0)} \frac{\partial}{\partial \Sigma_{jk}} p(x^{(i)} | y^{(i)} = 0) \right\}$$

$$\begin{aligned}\text{where } \frac{\partial p(x^{(i)} | y^{(i)} = 0)}{\partial \Sigma_{jk}} &= p(x^{(i)} | y^{(i)} = 0) \left\{ -\frac{1}{2|\Sigma|} |\Sigma| \text{tr} \left( \Sigma^{-1} \frac{\partial \Sigma}{\partial \Sigma_{jk}} \right) + \frac{1}{2} (x^{(i)} - \mu_0)^T \Sigma^{-1} \frac{\partial \Sigma}{\partial \Sigma_{jk}} \Sigma^{-1} (x^{(i)} - \mu_0) \right\} \\ &= -\frac{1}{2} p(x^{(i)} | y^{(i)} = 0) \left\{ \Sigma_{kj}^{-1} - [(\Sigma^{-1})^T (x^{(i)} - \mu_0)]_j [\Sigma^{-1} (x^{(i)} - \mu_0)]_k \right\} \\ &= -\frac{1}{2} p(x^{(i)} | y^{(i)} = 0) \left\{ \Sigma_{kj}^{-1} - [\Sigma^{-1} (x^{(i)} - \mu_0)]_k [(\Sigma^{-1})^T (x^{(i)} - \mu_0)]_j \right\}\end{aligned}$$

$$\text{Similarly, } \frac{\partial p(x^{(i)} | y^{(i)} = 1)}{\partial \Sigma_{jk}} = -\frac{1}{2} p(x^{(i)} | y^{(i)} = 1) \left\{ \Sigma_{kj}^{-1} - [\Sigma^{-1} (x^{(i)} - \mu_1)]_k [(x^{(i)} - \mu_1)^T \Sigma^{-1}]_j \right\}$$

$$\therefore \frac{\partial \ell}{\partial \Sigma_{jk}} = \sum_{i=1}^n -\frac{1}{2} \left\{ \Sigma_{kj}^{-1} - [\Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}})]_k [(\Sigma^{-1})^T (x^{(i)} - \mu_{y^{(i)}})]_j \right\}$$

$$\begin{aligned}\therefore \text{Optimal condition: } \sum_{i=1}^n \Sigma_{kj}^{-1} &= \sum_{i=1}^n [\Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}})]_k [(\Sigma^{-1})^T (x^{(i)} - \mu_{y^{(i)}})]_j \\ n \Sigma^{-1} &= \sum_{i=1}^n \Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}}) (x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} \\ \Sigma &= \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu_{y^{(i)}}) (x^{(i)} - \mu_{y^{(i)}})^T\end{aligned}$$

(c) [5 points] **Coding problem.** In `src/linearclass/gda.py`, fill in the code to calculate  $\phi$ ,  $\mu_0$ ,  $\mu_1$ , and  $\Sigma$ , use these parameters to derive  $\theta$ , and use the resulting GDA model to make predictions on the validation set. Make sure to write your model's predictions on the validation set to the file specified in the code.

In addition, use the provided code in `src/linearclass/logreg.py` to fit a logistic regression and make predictions on the validation set.

Include two plots of the **validation data** with  $x_1$  on the horizontal axis and  $x_2$  on the vertical axis. To visualize the two classes, use a different symbol for examples  $x^{(i)}$  with

$y^{(i)} = 0$  than for those with  $y^{(i)} = 1$ . In the first figure, plot the decision boundary found by GDA (i.e. line corresponding to  $p(y|x) = 0.5$ ). In the second figure, plot the decision boundary found by logistic regression.

**Answer:**

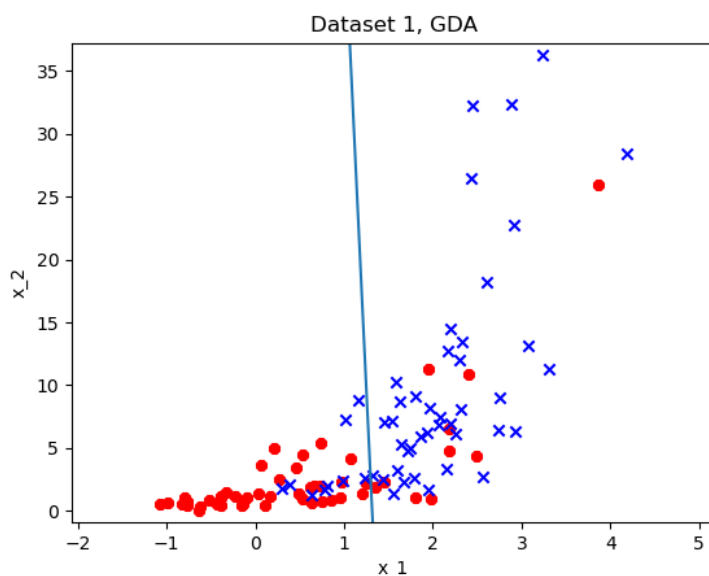


Figure 1: GDA, dataset 1

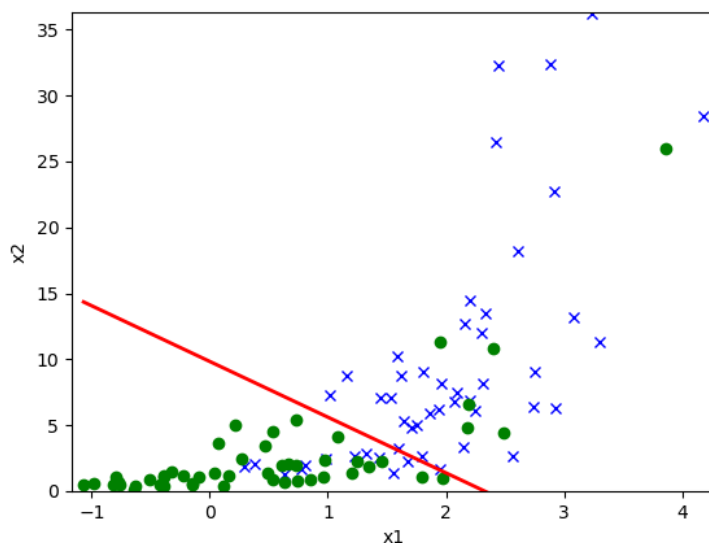


Figure 2: LR, dataset 1

- (d) [2 points] For Dataset 1, compare the validation set plots obtained in part (c) from logistic regression and GDA respectively, and briefly comment on your observation in a couple of



lines.

**Answer:** With LR LR, the decision boundary cut looks clean, the major points group of each class is separated. Whereas with GDA, the cut looks rough even though it is somewhat dividing the classes.

- (e) [5 points] Repeat the steps in part (c) for Dataset 2. Create similar plots on the **validation set** of Dataset 2 and include those plots in your writeup.

On which dataset does GDA seem to perform worse than logistic regression? Why might this be the case?

**Answer:**

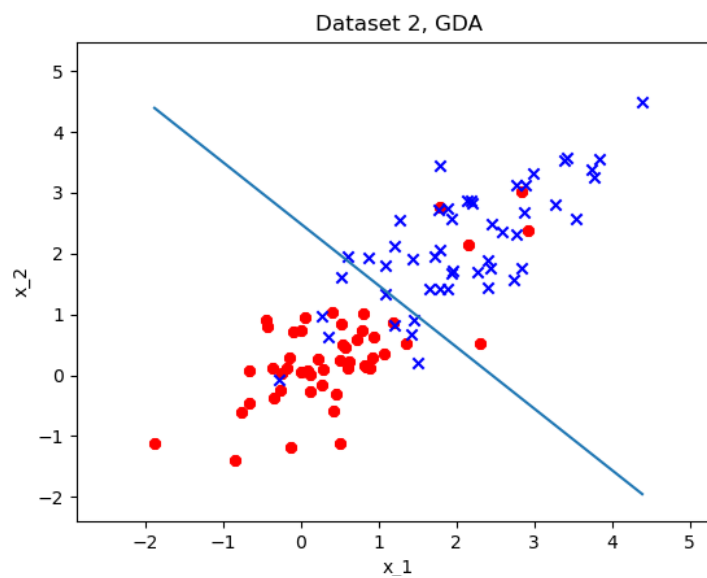


Figure 3: GDA, dataset 2

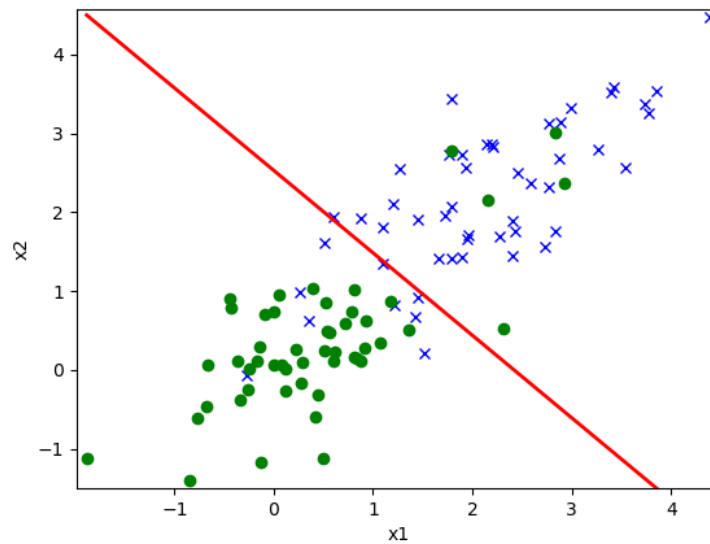


Figure 4: LR, dataset 2

GDA seems to perform worse than LR on dataset 1, since GDA assumes Gaussian distribution, it only performs well when data points of the same class stay together roughly about some center point, like in dataset 2.

- (f) [1 points] For the dataset where GDA performed worse in part (e), can you find a transformation of the  $x^{(i)}$ 's such that GDA performs significantly better? What might this transformation be?

**Answer:** Do a linear scaling transformation, especially scaling down  $x_2$ , to make the crosses closer together. That way, it will be closer to normally distributed.