

# Homework #2

CSE 546: Machine Learning

Prof. Kevin Jamieson

Due: 11/1/2018 11:59 PM

## 1 A Taste of Learning Theory

1. [2 points] Let  $X$  be a random feature vector in  $\mathbb{R}^d$  and  $Y$  be a random label in  $\{1, \dots, K\}$  for some  $K \in \mathbb{N}$  drawn from a joint distribution  $P_{XY}$ . A *randomized classifier*  $\delta(x)$  takes an  $x \in \mathbb{R}^d$  as input and outputs an element  $y \in \{1, \dots, K\}$  with probability  $\alpha(x, y) := P(\delta(x) = y)$ , where  $\sum_{y=1}^K \alpha(x, y) = 1$  for all  $x$ . For any classifier  $\delta$  we define the *risk* as

$$R(\delta) := \mathbb{E}_{X,Y,\delta}[\mathbf{1}\{\delta(X) \neq Y\}]$$

where  $\mathbf{1}(\mathcal{E})$  is the indicator function for the event  $\mathcal{E}$  (the function takes the value 1 if  $\mathcal{E}$  occurs and 0 otherwise). We say a classifier  $\delta$  is *deterministic* if  $\alpha(x, y) \in \{0, 1\}$  for all  $x, y$ . A Bayes classifier is defined as  $\delta_* \in \arg \inf_{\delta} R(\delta)$  where the infimum is taken over all randomized classifiers (note it may not be unique). For an arbitrary  $P_{XY}$ , characterize the set of all Bayes Classifiers (i.e., for a given  $x$  what are permissible Bayes classifiers in terms of  $\alpha(x, y)$ ). Then propose a deterministic decision rule that is a Bayes classifier.

2. [8 points] For  $i = 1, \dots, n$  let  $(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY}$  where  $y_i \in \{-1, 1\}$  and  $x_i$  lives in some set  $\mathcal{X}$  ( $x_i$  is not necessarily a vector). The 0/1 loss, or *risk*, for a deterministic classifier  $f : \mathcal{X} \rightarrow \{-1, 1\}$  is defined as:

$$R(f) = \mathbb{E}_{XY}[\mathbf{1}(f(X) \neq Y)]$$

where  $\mathbf{1}(\mathcal{E})$  is the indicator function for the event  $\mathcal{E}$  (the function takes the value 1 if  $\mathcal{E}$  occurs and 0 otherwise). The expectation is with respect to the underlying distribution  $P_{XY}$  on  $(X, Y)$ . Unfortunately, we don't know  $P_{XY}$  exactly, but we do have our i.i.d. samples  $\{(x_i, y_i)\}_{i=1}^n$  drawn from it. Define the *empirical risk* as

$$\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(f(x_i) \neq y_i)$$

which is just an empirical estimate of our loss. Recall Hoeffding's inequality: if  $Z_1, \dots, Z_m$  are i.i.d. random variables on the interval  $[a, b]$  where  $\mathbb{E}[Z_i] = \mu$  and  $\hat{\mu} = \frac{1}{m} \sum_{i=1}^m Z_i$ , then

$$\mathbb{P}(|\mu - \hat{\mu}| \geq \epsilon) \leq 2 \exp\left(-\frac{2m\epsilon^2}{(b-a)^2}\right).$$

- a. Suppose a domain expert who has not seen your freshly drawn data  $\{(x_i, y_i)\}_{i=1}^n$  gives you a classifying rule  $\tilde{f} : \mathcal{X} \rightarrow \{-1, 1\}$ . We wish to estimate the true risk  $R(\tilde{f})$  of this classifying rule using the empirical risk  $\hat{R}_n(\tilde{f})$ . Justify the use of Hoeffding's inequality, and use it to find a value of  $A$  such that

$$\mathbb{P}(|\hat{R}_n(\tilde{f}) - R(\tilde{f})| \leq A) \geq 1 - \delta.$$

- b. After seeing the confidence interval, the domain expert is not satisfied with her classifying rule  $\tilde{f}$ . Because she understands the underlying process that generated the data, she proposes a finite set of alternative classifying rules  $\mathcal{F} = \{f_1, \dots, f_k\}$  that may fit the data better. The "best in class" function  $f^*$  is:

$$f^* = \arg \min_{f \in \mathcal{F}} R(f).$$

Note that  $R(f^*)$  is the best true loss we could hope to achieve using functions in  $\mathcal{F}$ . If we replace  $\tilde{f}$  with  $f^*$  in part a., does the same confidence interval hold? Why or why not?

- c. Of course, determining  $f^*$  would require exact knowledge of  $R(f)$ , which requires exact knowledge of  $P_{XY}$ , which we don't have. As an alternative, you propose to use the *empirical risk minimizer* (ERM):

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \hat{R}_n(f).$$

If we replace  $\tilde{f}$  with  $\hat{f}$  in part a., does the same confidence interval hold? Why or why not?

- d. Come up with a **simple** example (i.e., give a  $P_{XY}, \mathcal{F}$ ) where for any number of observations  $n$  we have  $\hat{R}_n(\hat{f}) = 0$  and  $R(\hat{f}) = 1/2$  where  $\hat{f} = \arg \min_{f \in \mathcal{F}} \hat{R}_n(f)$ .
- e. Provide a confidence interval that simultaneously holds for the losses of all  $f \in \mathcal{F}$ , with probability of error  $\delta$ . In other words, provide a value  $B$  so that:

$$\mathbb{P}(\text{for all } f \in \mathcal{F}, |\hat{R}_n(f) - R(f)| \leq B) \geq 1 - \delta$$

Show your steps. (Hint: for events  $\mathcal{E}_1, \mathcal{E}_2, \dots, \mathcal{E}_k$ , the “union bound” states that  $\mathbb{P}(\mathcal{E}_1 \text{ or } \mathcal{E}_2 \text{ or } \dots \text{ or } \mathcal{E}_k) = \mathbb{P}(\bigcup_{i=1}^k \mathcal{E}_i) \leq \sum_{i=1}^k \mathbb{P}(\mathcal{E}_i)$ .)

- f. Noting that  $\hat{f} \in \mathcal{F}$ , provide a confidence interval for the loss of  $\hat{f}$  that holds with probability greater than  $1 - \delta$ . In other words, we are seeking a bound on  $|R(\hat{f}) - \hat{R}_n(\hat{f})|$  that holds with probability greater than  $1 - \delta$ .
- g. Provide a bound on how close your loss, using the ERM, is to the best possible loss (in the hypothesis space). Specifically, provide a value  $C$  so that the following holds with probability greater than  $1 - \delta$ ,

$$R(\hat{f}) - R(f^*) \leq C$$

The quantity  $R(\hat{f}) - R(f^*)$  is often referred to as the *excess risk*.

- h. Fix an  $f \in \mathcal{F}$  and suppose  $R(f) > \epsilon$ . Show that  $\mathbb{P}(\hat{R}_n(f) = 0) \leq (1 - \epsilon)^n \leq e^{-n\epsilon}$ . Leverage this insight to show that with probability at least  $1 - \delta$

$$\hat{R}_n(\hat{f}) = 0 \implies R(\hat{f}) - R(f^*) \leq \frac{\log(|\mathcal{F}|/\delta)}{n}$$

where  $\hat{f} = \arg \min_{f \in \mathcal{F}} \hat{R}_n(f)$ .

- i. Let us understand how large a hypothesis class we can utilize, i.e. how large  $|\mathcal{F}|$  can be? Suppose we know we will be provided with a training set of size  $n$ , and, before we look at our training data, we choose the a hypothesis class  $\mathcal{F}$  as a function of the sample size  $n$ . As we get more data, we would expect that we can utilize a larger hypothesis class. Let us examine this more quantitatively. We can think of learning being possible if our regret tends to 0 as  $n$  becomes large (with a probability of error less than  $\delta$ ). Let us determine if learning is possible in each of the following cases. For the following cases, does the above suggest that we are able to learn, and if so, what is our excess risk as a function of  $n$  in big-O notation?

(a)  $|\mathcal{F}|$  is a constant.

(b)  $|\mathcal{F}| = n^p$  for some constant  $p$ .

(c)  $|\mathcal{F}| = \exp(\sqrt{n})$ .

(d)  $|\mathcal{F}| = \exp(10n)$ .

Your answers to part i are one interpretation as to when and why learning complex models is possible (from a statistical perspective). Of particular interest is when  $\mathcal{F}$  is allowed to be infinite in which more advanced tools are necessary. To provide intuition we note that  $n$  points in  $d$ -dimensions can be labeled by a hyperplane  $w$  in at most  $O(n^d)$  ways (i.e.,  $|\{\text{sign}(Xw) : w \in \mathbb{R}^d\}| = O(n^d)$  for fixed  $X \in \mathbb{R}^{n \times d}$ ). While more advanced tools (e.g., VC dimension, covering numbers, etc.) are necessary to formally prove the learning rates for infinite classes, thinking about this number of sign patterns and how this relates to  $|\mathcal{F}|$  is instructive.

---

**Algorithm 1:** Coordinate Descent Algorithm for Lasso

---

```
while not converged do
     $b \leftarrow \frac{1}{n} \sum_{i=1}^n \left( y_i - \sum_{j=1}^d w_j x_{i,j} \right)$ 
    for  $k \in \{1, 2, \dots, d\}$  do
         $a_k \leftarrow 2 \sum_{i=1}^n x_{i,k}^2$ 
         $c_k \leftarrow 2 \sum_{i=1}^n x_{i,k} \left( y_i - (b + \sum_{j \neq k} w_j x_{i,j}) \right)$ 
         $w_k \leftarrow \begin{cases} (c_k + \lambda)/a_k & c_k < -\lambda \\ 0 & c_k \in [-\lambda, \lambda] \\ (c_k - \lambda)/a_k & c_k > \lambda \end{cases}$ 
    end
end
```

---

## 2 Programming: Lasso

Given  $\lambda > 0$  and data  $(x_1, y_1), \dots, (x_n, y_n)$ , the Lasso is the problem of solving

$$\arg \min_{w \in \mathbb{R}^d, b \in \mathbb{R}} \sum_{i=1}^n (x_i^T w + b - y_i)^2 + \lambda \sum_{j=1}^d |w_j| \quad (1)$$

$\lambda$  is a regularization tuning parameter. For the programming part of this homework, you are required to implement the coordinate descent method of Algorithm 1 that can solve the Lasso problem.

You may use common computing packages (such as NumPy or SciPy), but do not use an existing Lasso solver (e.g., of Sci-Kit Learn).

Before you get started, here are some hints that you may find helpful:

- For loops can be slow whereas vector/matrix computation in Numpy is very optimized, exploit this as much as possible.
- As a sanity check, ensure the objective value is nonincreasing with each step.
- It is up to you to decide on a suitable stopping condition. A common criteria is to stop when no element of  $w$  changes by more than some small  $\delta$  during an iteration. If you need your algorithm to run faster, an easy place to start is to loosen this condition.
- You will need to solve the Lasso on the same dataset for many values of  $\lambda$ . This is called a regularization path. One way to do this efficiently is to start at a large  $\lambda$ , and then for each consecutive solution, initialize the algorithm with the previous solution, decreasing  $\lambda$  by a constant ratio until finished.
- The smallest value of  $\lambda$  for which the solution  $\hat{w}$  is entirely zero is given by

$$\lambda_{max} = \max_{k=1, \dots, d} 2 \left| \sum_{i=1}^n x_{i,k} \left( y_i - \left( \frac{1}{n} \sum_{j=1}^n y_j \right) \right) \right|$$

This is helpful for choosing the first  $\lambda$  in a regularization path.

3. [5 points] We will first try out your solver with some synthetic data. A benefit of the Lasso is that if we believe many features are irrelevant for predicting  $y$ , the Lasso can be used to enforce a sparse solution, effectively differentiating between the relevant and irrelevant features. Suppose that  $x \in \mathbb{R}^d, y \in \mathbb{R}, k < d$ , and pairs of data  $(x_i, y_i)$  for  $i = 1, \dots, n$  are generated independently according to the model  $y_i = w^T x_i + \epsilon_i$  where

$$w_j = \begin{cases} j/k & \text{if } j \in \{1, \dots, k\} \\ 0 & \text{otherwise} \end{cases}$$

where  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$  is some Gaussian noise (in the model above  $b = 0$ ). Note that since  $k < d$ , the features  $k + 1$  through  $d$  are unnecessary (and potentially even harmful) for predicting  $y$ .

With this model in mind, let  $n = 500, d = 1000, k = 100$ , and  $\sigma = 1$ . Generate some data by drawing each  $x_i \sim \mathcal{N}(0, I)$  so that  $x_i \in \mathbb{R}^d$  with  $y_i$  generated as specified above.

- With your synthetic data, solve multiple Lasso problems on a regularization path, starting at  $\lambda_{max}$  where 0 features are selected and decreasing  $\lambda$  by a constant ratio (e.g., 1.5) until nearly all the features are chosen. In plot 1, plot the number of non-zeros as a function of  $\lambda$  on the x-axis (Tip: use `plt.xscale('log')`).
- For each value of  $\lambda$  tried, record values for false discovery rate (FDR) (number of incorrect nonzeros in  $\hat{w}$ /total number of nonzeros in  $\hat{w}$ ) and true positive rate (TPR) (number of correct nonzeros in  $\hat{w}/k$ ). In plot 2, plot these values with the x-axis as FDR, and the y-axis as TPR and note that in an ideal situation we would have an (FDR,TPR) pair in the upper left corner, but that can always trivially achieve  $(0,0)$  and  $(\frac{d-k}{d}, 1)$

Comment on the effect of  $\lambda$  in these two plots.

4. [5 points] We'll now put the Lasso to work on some real data from Yelp (an old Kaggle competition, see <http://www.kaggle.com/c/yelp-recruiting> for background, but get the data from the class website).

For this competition, the task is to predict the number of useful upvotes a particular review will receive. One of the most important requirements for learning great models is creating great features. We can use our Lasso solver for this as follows. First, generate a large amount of features from the data, even if many of them are likely unnecessary. Afterward, use the Lasso to reduce the number of features to a more reasonable amount.

Yelp provides a variety of data, such as the review's text, date, and restaurant, as well as data pertaining to each business, user, and check-ins. We have preprocessed this data for you into the following files:

<code>upvote_data.csv</code>	Each row is a review, each column is a real-valued feature ( $n \times d$ )
<code>upvote_labels.txt</code>	Each row is the number of useful vote counts for that review ( $n \times 1$ )
<code>upvote_features.txt</code>	Names of each feature for interpreting results ( $d \times 1$ )

To get you started, the Python following code should load the data:

```
import numpy as np
# Load a csv of floats:
X = np.genfromtxt("upvote_data.csv", delimiter=",")
# Load a text file of integers:
y = np.loadtxt("upvote_labels.txt", dtype=np.int)
# Load a text file of strings:
featureNames = open("upvote_features.txt").read().splitlines()
```

Use the first 4000 samples for training, the next 1000 samples for validation, and the remaining samples for testing. Each sample is a  $(x_i, y_i)$  pair. As a pre-processing step, take the square root of each  $y_i$  so that  $y_i \mapsto \sqrt{y_i}$ . This rescaling ameliorates for outliers.

- Solve lasso to predict the number of useful votes a Yelp review will receive. Starting at  $\lambda_{max}$ , run Lasso on the training set, decreasing  $\lambda$  using previous solutions as initial conditions to each problem. Stop when you have considered enough  $\lambda$ 's that, based on validation error, you can choose a good solution with confidence (for instance, when validation error begins increasing by a lot). Plot the squared error on the training and validation data versus  $\lambda$ . On a different plot, plot the number of nonzeros in each solution versus  $\lambda$ . (Tip: use `plt.xscale('log')` and `plt.gca().invert_xaxis()`)
- Find the  $\lambda$  that achieves best validation performance, and test your model on the remaining set of test data. What is the train, val, and test error for this choice?
- Inspect your solution and take a look at the 10 features with weights largest in magnitude. List the names of these features and their weights, and comment on if the weights generally make sense intuitively. As you use a larger  $\lambda$  so that fewer features are selected, they may make more sense.

## 2.1 Programming: Binary Logistic Regression

5. [5 points] Let us again consider the MNIST dataset, but now just binary classification, specifically, recognizing if a digit is a 2 or 7. Here, let  $Y = 1$  for all the 7's digits in the dataset, and use  $Y = -1$  for 2. We will use regularized logistic regression. Given a binary classification dataset  $\{(x_i, y_i)\}_{i=1}^n$  for  $x_i \in \mathbb{R}^d$  and  $y_i \in \{-1, 1\}$  we showed in class that the regularized negative log likelihood objective function can be written as

$$J(w, b) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i(b + x_i^T w))) + \lambda \|w\|_2^2$$

Note that the offset term  $b$  is not regularized. For all experiments, use  $\lambda = 10^{-1}$ . Let  $\mu_i(w, b) = \frac{1}{1 + \exp(-y_i(b + x_i^T w))}$ .

- a. Derive the gradients  $\nabla_w J(w, b)$ ,  $\nabla_b J(w, b)$  and Hessians  $\nabla_w^2 J(w, b)$ ,  $\nabla_b^2 J(w, b)$  and give your answers in terms of  $\mu_i(w, b)$  (your answers should not contain exponentials).
- b. Implement gradient descent with an initial iterate of all zeros. Try several values of step sizes to find one that appears to make  $J(w, b)$  on the training set converge the fastest. Run until you feel you are near to convergence.
  - (a) For both the training set and the test, plot  $J(w, b)$  as a function of the iteration number (and show both curves on the same plot).
  - (b) For both the training set and the test, classify the points according to the rule  $\text{sign}(b + x_i^T w)$  and plot the misclassification error as a function of the iteration number (and show both curves on the same plot).

Note that you are only optimizing on the training set. The  $J(w, b)$  and misclassification error plots should be on separate plots.

- c. Repeat (b) using stochastic gradient descent with batch size of 1. Note, the expected gradient with respect to the random selection should be equal to the gradient found in part (a). Take careful note of how to scale the regularizer.
- d. Repeat (b) using stochastic gradient descent with batch size of 100. That is, instead of approximating the gradient with a single example, use 100.
- e. Repeat (b) using Newton's method.