

**Machine Learning and Large Scale Data Analysis**

Assignment 3 Due: Thursday, April 26, 2018 at 1:30 pm.

Please hand in this Homework in 5 files:

1. A pdf of the theoretical homework (problems 1-3).
2. A pdf of your jupyter notebook for problem 4.
3. The ipynb file for problem 4.
4. A pdf of your jupyter notebook for problem 5.
5. The ipynb file for problem 5.

Read this article and be prepared to discuss in class. What problem is it related to?

<https://arstechnica.com/information-technology/2016/02/the-nsas-skynet-program-may-be-killing-thousands-of-innocent-people/>

1. *Bounds on the error probability.* (5 points) If  $a$  and  $b$  are non-negative, show that  $\min(a, b) \leq \sqrt{ab}$ . Use this to show that the error rate for a two category Bayes classifier must satisfy

$$P(\text{error}) \leq \sqrt{P(Y=1)P(Y=2)} \int \sqrt{f(x|Y=1)f(x|Y=2)} dx \leq \frac{1}{2} \int \sqrt{f(x|Y=1)f(x|Y=2)} dx.$$

2. *Bayes rule, variances, and priors.* (10 points) Let  $f(x|Y=k) \sim N(\mu_k, \sigma^2)$ ,  $k = 1, 2$  be a two category one-dimensional problem with  $P(Y=1) = \pi_1$ ,  $P(Y=2) = \pi_2$ .

- (a) Compute the decision boundary assuming  $\mu_1 > \mu_2$  and compute the Bayes loss.
- (b) Write the probability of error using the CDF of the normal distribution. Show that the error goes to zero as  $\sigma \rightarrow 0$ .
- (c) For fixed  $\sigma$  what happens to the decision boundary as  $\pi_1 \rightarrow 0$ ? For small  $\pi_1$  what is a very simple classification rule that can guarantee low error rate?
- (d) Assume a classification problem with  $K$  classes. Let  $L_{k,l}$  be the cost of choosing class  $l$  when class  $k$  is true. Let  $p(x, y)$  be the joint distribution of  $X, Y$ . The expected loss for a decision function  $h$  is

$$L(h) = \sum_{k=1}^K \sum_{l=1}^K \int_{h(x)=l} p(x, k) L_{k,l} dx.$$

Show that the decision rule with lowest loss is given by the generalize Bayes rule:

$$h_B(x) = \operatorname{argmin}_{j=1, \dots, K} \sum_{k=1}^K L_{k,j} p(k|x).$$

Hint: Use the loss of an individual example defined as  $L(h, x) = \sum_{k=1}^K L_{k,h(x)} p(k|x)$ .

- (e) One way to avoid the degenerate situation from item ?? is to change the cost function. Set  $L_{2,1}$  to be the cost of choosing class 1 when class 2 is true, and  $L_{1,2}$  the cost of choosing class 2 when class 1 is true. Write the expected loss in this situation. Recompute the decision boundary. What values on  $L_{i,j}$  would you assign to remedy the problem of small  $\pi_1$ .

### 3. Logistic regression (10 points)

- (a) Give a detailed derivation of the Newton algorithm for logistic regression and show that each iteration corresponds to solving a weighted least square problem. Hint: Compute the gradient  $\nabla L$  and the Hessian  $H$  of the loss  $L$  at  $\theta$ . Given current  $\theta$  write the Newton iteration for  $\theta^{new} - \theta$ . Note that  $H\theta$  simplifies.
- (b) Assume the data are perfectly linearly separable, i.e. there exist  $\theta$  such that  $x_i^t \theta < 0$  if  $y_i = 0$  and  $x_i^t \theta > 0$  if  $y_i = 1$ . Show that the conditional maximum likelihood estimator for the logistic regression model does not exist. Comment on the behavior of the iteratively reweighted least squares algorithm. Hint: If  $\theta$  is a perfect separator then  $\alpha\theta$  is also for any  $\alpha > 0$ . It may be easier to work with the likelihood instead of the log-likelihood.

### LSDA

### 4. Sparse coding of natural images and digits (35 points)

In this problem you will implement the sparse coding procedure as described in class on image patches of size 12x12. This was proposed as a possible computational mechanism underlying the evolution of neural representations in the visual cortex of mammals.<sup>1</sup> You will use the actual images used in this landmark paper.

To run the sparse coding algorithm over the images, we have provided a function that selects random patches. This can be run on the Olshausen-Field images using the following code

```
import scipy.io
%matplotlib inline
import matplotlib.pyplot as plt
import random
import numpy
plt.imshow(images[:, :, 0], cmap='gray')
data =
scipy.io.loadmat('/project/cmssc25025/sparsecoding/IMAGES_RAW.mat')
images = data['IMAGESr']
# images is a 3D array of size [512,512,10]
# where the 10 images are of size 512 x 512
```

<sup>1</sup>B. Olshausen and D. Field, "Emergence of simple-cell receptive field properties by learning a sparse code for natural images," *Nature* 381, 607–609, 1996.

```
#Show the first image.
plt.imshow(images[:, :, 0], cmap=?gray?)

# Function to sample image patches from the large images.
def sample_random_square_patches(image, num, width):
    patches = numpy.zeros([width,width,num]);
    for k in range(num):
        i, j = random.sample(range(image.shape[0]-width), 2)
        patches[:, :, k] = image[i:i+width, j:j+width]
    return patches
```

Run the sparse coding scheme over the images.

- Use the class `sklearn.linear_model.Lasso` for the coefficient estimation step.
- We want to fully exploit the vector computation properties of numpy. **Make sure there are no ‘loops’ inside the main loop that iterates over steps in the SGD. All computations should be done with matrix operations.** For the Lasso step, you can have the `fit` function of the `Lasso` class fit all training points at once. Compare the time it takes to do this to a loop calling `Lasso` on each training point separately.
- Monitor the convergence of the SGD algorithm by checking the change in the codebook. How long does it take to converge? Experiment with a step size  $\eta$  for your algorithm, constant or decreasing. Display the codebook after initialization, after convergence, and at several intermediate stages. Comment on your results. Are they consistent with the results presented in paper?
- Show reconstructions of patches images using the sparse representation.

## LSDA

### 5. Stochastic gradient descent on beer reviews (40 points)

In this problem, you will work with 2,924,163 beer reviews from <https://www.ratebeer.com>. One entry contains the text of the consumer’s review together with ratings of *appearance*, *aroma*, *palate*, *style*, *taste*, and an *overall* score from 0 to 20. The data also include the name, id and brewer of the beer. You will give reviews binary label: 1 for a (*positive*) review with overall score at least 14, and otherwise 0 for a (*negative*) review.

#### Part 1: Data inspection.

To warm up, check the mean, median and standard deviation of the overall ratings for each beer and brewer. Do you think people have similar taste?

#### Part 2: Sentiment analysis.

Your first task to classify the sentiment of reviews from the text. Text can have neutral, mixed, sarcastic, or otherwise ambiguous sentiment. *Sentiment analysis*<sup>2</sup> can be challenging

<sup>2</sup>[http://en.wikipedia.org/wiki/Sentiment\\_analysis](http://en.wikipedia.org/wiki/Sentiment_analysis)

appearance	aroma	beer_id	beer_name	brewer	overall	palate	review	review_id	style	taste
4.0	6.0	45842	John Harvards Sim...	3084	13.0	3.0	On tap at the Spr...	0	17	6.0
4.0	6.0	45842	John Harvards Sim...	3084	13.0	4.0	On tap at the Joh...	1	17	7.0
4.0	5.0	95213	John Harvards Cri...	3084	14.0	3.0	UPDATED: FEB 19, ...	2	33	6.0
2.0	4.0	65957	John Harvards Fan...	3084	8.0	2.0	On tap the Spring...	3	33	4.0
5.0	8.0	41336	John Harvards Van...	3084	16.0	4.0	Springfield, PA 1...	5	58	7.0
4.0	5.0	80424	John Harvards Ame...	3084	12.0	3.0	On tap at the Spr...	6	73	6.0
2.0	6.0	51269	John Harvards Gra...	3084	14.0	3.0	Sampled @ the Spr...	7	12	7.0
4.0	8.0	51269	John Harvards Gra...	3084	16.0	3.0	Springfield... Po...	8	12	7.0
3.0	8.0	51269	John Harvards Gra...	3084	17.0	4.0	UPDATED: FEB 19, ...	9	12	8.0
4.0	4.0	73033	John Harvards Bel...	3084	11.0	2.0	UPDATED: FEB 19, ...	10	62	5.0
3.0	5.0	56415	John Harvards Cas...	3084	14.0	3.0	UPDATED: FEB 19, ...	11	24	7.0
3.0	6.0	68538	John Harvards Yin...	3084	12.0	3.0	On tap at Springf...	13	50	6.0
3.0	5.0	68538	John Harvards Yin...	3084	9.0	2.0	On tap at the Spr...	14	50	5.0
4.0	8.0	21566	Barley Island Bar...	1786	15.0	4.0	Handbottled from ...	15	66	8.0
4.0	8.0	21566	Barley Island Bar...	1786	15.0	4.0	On tap at the Gre...	16	66	8.0
3.0	7.0	21566	Barley Island Bar...	1786	14.0	3.0	UPDATED: JUL 7, 2...	17	66	6.0
4.0	5.0	21566	Barley Island Bar...	1786	13.0	4.0	On cask at BI - A...	18	66	4.0
3.0	7.0	21566	Barley Island Bar...	1786	13.0	3.0	Name: BA Count Ho...	20	66	7.0
3.0	7.0	21566	Barley Island Bar...	1786	14.0	3.0	Aroma is sweet bo...	21	66	7.0
3.0	6.0	21566	Barley Island Bar...	1786	13.0	3.0	Cask at GTMW. Po...	22	66	7.0

```
df = spark.read.json('/project/cmsc25025/beer_review/labeled.json')
df.show()
```

even for people.

Here is one example review that you will be working on:

I got this one in Indianapolis. The body was dark brown. The aroma was caramel and chocolate with some pancakes. The taste was pancakes with coffee and chocolate. Well made, but a little boring.

Other reviews seem easier to classify according to positive or negative sentiment:

I was surprised by this one. A really nice local sour beer. Pours a great looking brown with a slight red hue. Not much head, but a good amount of lacing. Aroma is nutty, with cherries and a tartness to it. Very earthy. Flavor is nice and sour, lots of red fruits and nuts, with just a hint of yeast. The wood really comes into play here, and works nicely with the other elements. This one is tasty for sure.

After loading the labeled data, proceed as follows.

(a) *Generating features.* You need to represent text reviews in terms of a vector of features (covariates). One simple but effective representation is to use membership in a fixed vocabulary. Suppose the vocabulary contains  $p$  words. For a given review, you normalize the text, and separate it into space-delimited tokens. For each of the tokens, if it is in the dictionary you have a one for the corresponding word in the feature vector, and you ignore it otherwise.

For example, suppose the review is

"I like this beer! Thanks for sharing! Yay!"

and after normalization you process this into the list

```
["i", "like", "this", "beer", "thanks", "for", "sharing", "yay"].
```

If "i", "this", "for" and "yay" are not in your dictionary, but the other words are, you have four active features

```
["like", "beer", "thanks", "sharing"].
```

So, this review would be a  $p$ -dimensional vector where four features are set to 1, and the rest are set to 0. You should drop a review if none of its words is in the dictionary.

We have processed a few vocabularies for you. The file `vocab_50.json` contains `word:id` pairs for 30,009 words. Those words appear at least 50 times across all the reviews. Similarly for `vocab_{10,20,30}.json`. The most common 50 words are thrown out. You are more than welcome to use your own vocabulary. To load our vocabulary, run

```
with open('/project/cmsc25025/beer_review/vocab_50.json', 'r') as f:
    vocab = json.load(f)
```

Note that your feature vectors are going to be very sparse, as most reviews have very few words. In order to speed up the computation, using a sparse vector representation for high dimensional data is crucial. We recommend `scipy.sparse.csr_matrix`<sup>3</sup> since it's easy to construct from scratch. To construct these features for every document You will have to loop over every document and find which vocabulary words are in the document. Try your function on a few documents to see how long it takes. It could take a few minutes to process the whole data set. Split your data into training, validation and testing sets with proportion 0.7:0.15:0.15.

(b) *Logistic regression using Newton's method.* Train an  $\ell_2$ -regularized logistic regression classifier using the `sklearn.linear_model.LogisticRegression` class. To select the regularization parameter  $C = 1/\lambda$ , you should try different values on the validation set. Pick the best. How long does it take to train?

Do the same thing using the `LinearSVC` class in `sklearn.svm`. Use `loss='hinge'`. Compare the results of the logistic loss to the hinge loss. Is there a difference?

(c) *Stochastic gradient descent.* Your next job is to train an  $\ell_2$ -regularized logistic regression classifier using stochastic gradient descent. Recall the SGD framework that was covered in class using minibatches.

- i. Initialize the model with  $\beta = 0$  (uniform).
- ii. Randomly split the training data into mini-batches. Make one pass of the data, processing one mini-batch in every iteration. This is called one training epoch.
- iii. Repeat the last step a few times.

---

<sup>3</sup>[https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.csr\\_matrix.html](https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.csr_matrix.html)

Remember the role that the learning rate plays in SGD. Tiny learning rates lead to slow convergence, while large rates can impede convergence—it might make the objective value oscillate. Try fixed vs. decreasing learning rates. You can also try adaptive learning rates<sup>4</sup> using Adagrad or Adam.

To select the regularization parameter  $\lambda$ , you should try different values on the validation set (see the next question for the proportion). Select those that yield the smallest validation error. You might want to do this on a smaller training set than the full one. You can do this for the minibatch size as well. Once you've chosen the values of these parameters train the model. Plot the error rate and negative log-likelihood for both training and validation set as a function of iterations. Finally report the error on the test set. How does it compare in terms of time and error rate to the Logistic regression function in the previous item.

### *Part 3: Scores versus text.*

In addition to text reviews, the users also scored *appearance*, *aroma*, *palate*, *style*, *taste* of a beer. In this problem, you will check whether those scores could reflect people's opinion better than text.

Train another logistic regression model using those features. You should use the same SGD algorithm as before. Compare the model using score features with that using review text. Again, use the validation set to tune the regularization parameters, and retrain the model on the union of training and validation set. Finally, compute the prediction error on the testing set.

Which model predicts better? Is the representation you constructed for text more powerful, or are the scores? Why? Comment on your findings and discuss your thinking.

---

<sup>4</sup><http://sebastianruder.com/optimizing-gradient-descent/index.html>