# HOMEWORK 1 MLE, MAP ESTIMATES; LINEAR AND LOGISTIC REGRESSION

CMU 10-701: MACHINE LEARNING (SPRING 2017)

OUT: Jan 31

DUE: Feb 10, 11:59 PM

### START HERE: Instructions

- Collaboration policy: Collaboration on solving the homework is allowed, after you have thought about the problems on your own. It is also OK to get clarification (but not solutions) from books or online resources, again after you have thought about the problems on your own. There are two requirements: first, cite your collaborators fully and completely (e.g., "Jane explained to me what is asked in Question 3.4"). Second, write your solution *independently*: close the book and all of your notes, and send collaborators out of the room, so that the solution comes from you only.
- Submitting your work: Assignments should be submitted as PDFs using Gradescope unless explicitly stated otherwise. Each derivation/proof should be completed on a separate page. Submissions can be handwritten, but should be labeled and clearly legible. Else, submissions can be written in LaTeX. Upon submission, label each question using the template provided by Gradescope.
- Programming: All programming portions of the assignments should be submitted to Gradescope as well. We will not be using this for autograding, but rather for plagiarism detection and to make it simpler to submit code. You may use any language which you like to submit unless the problem states otherwise. There is a separate 'programming assignment' that should allow you to upload your code easily. Code should be uploaded to this separate programming assignment, while visualizations and written answers should still be submitted within the primary Gradescope assignment. In your code, please make it clear in the comments which are the primary functions to compute the answers to each question.

## Part A: Multiple Choice Questions [7 Points] (Yiting)

- There will be only one right answer. Please explain your choice in one or two sentences.
- 1. [4 Points] For each case listed below, what type of machine learning problem does it belong to?
  - (a) Advertisement selection system, which can predict the probability whether a customer will click on an ad or not based on the search history
  - (b) U.S post offices use a system to automatically recognize handwriting on the envelope
  - (c) Reduce dimensionality using principal components analysis (PCA)
  - (d) Trading companies try to predict future stock market based on current market conditions
  - (e) Repair a digital image that has been partially damaged
  - A. Supervised learning: Classification
  - B. Supervised learning: Regression
  - C. Unsupervised learning
- 2. [1 Point] For four statements below, which one is wrong?
  - A. In maximum a posterior (MAP) estimate, data overwhelms the prior if we have enough data
  - B. There are no parameters in non-parametirc models
  - C.  $P(X \cap Y \cap Z) = P(Z|X \cap Y)P(Y|X)P(X)$
  - D. Compared with parametric models, non-parameter models are flexible, since they don't make strong assumptions
- 3. [1 Point] There are about 12% people in U.S. having breast cancer during their lifetime. One patient has a positive result for the medical test. Suppose the sensitivity of this test is 90%, meaning the test will be positive with probability 0.9 if one really has cancer. The false positive is likely to be 2%. Then what is the probability this patient actually having cancer based on Bayes Theorem?
  - A. 90%
- B. 86%
- C. 12%
- D. 43%
- 4. [1 Point] What is the most suitable error function for gradient descent using logistic regression?
  - A. The negative log-likelihood function
  - B. The number of mistakes
  - C. The squared error
  - D. The log-likelihood function

# Part B, Problem 1: Bias-Variance Decomposition [20 Points] (Calvin and Hao)

Consider a p-dimensional vector  $\mathbf{x} \in \mathbb{R}^p$  drawn from a Gaussian distribution with an identity covariance matrix  $\mathbf{\Sigma} = \mathbf{I}_p$  and an unknown mean  $\boldsymbol{\mu}$ , i.e.  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{I}_p)$ . Our goal is to evaluate the effectiveness of an estimator  $\hat{\boldsymbol{\mu}} = \boldsymbol{f}(\mathbf{x})$  of the mean from only a single sample (i.e. n = 1) by measuring its mean squared error  $\mathbb{E}[\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}\|^2]$ , where  $\|\cdot\|^2$  is the squared Euclidean norm and the expectation is taken over the data generating distribution.

Note that for any estimator  $\hat{\theta}$  of a parameter vector  $\theta$ , its mean squared error can be decomposed as:

$$\mathbb{E}\left[\left\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\right\|^2\right] = \left\|\mathrm{Bias}[\hat{\boldsymbol{\theta}}]\right\|^2 + \mathrm{Trace}(\mathrm{Var}[\hat{\boldsymbol{\theta}}]), \text{ where:}$$

$$\operatorname{Bias}[\hat{\boldsymbol{\theta}}] = \mathbb{E}[\hat{\boldsymbol{\theta}}] - \boldsymbol{\theta} \quad \text{and} \quad (\operatorname{Var}[\hat{\boldsymbol{\theta}}])_{j,j} = \operatorname{Var}[\hat{\boldsymbol{\theta}}_j] = \mathbb{E}[(\hat{\boldsymbol{\theta}}_j - \mathbb{E}[\hat{\boldsymbol{\theta}}_j])^2]$$

Here,  $\operatorname{Trace}(\cdot)$  denotes the sum of the diagonal elements of a square matrix,  $(\operatorname{Var}[\hat{\theta}])_{j,j}$  denotes the jth diagonal element of  $\operatorname{Var}[\hat{\theta}]$ , and  $\hat{\theta}_j$  denotes the jth element of  $\hat{\theta}$ .

1. [4 Points] Derive the maximum likelihood estimator:

$$\hat{\boldsymbol{\mu}}_{\text{MLE}} = \arg \max_{\boldsymbol{\mu}} P(\boldsymbol{x}; \boldsymbol{\mu}).$$

What is its mean squared error?

2. [4 Points] Derive the  $\ell_2$ -regularized maximum likelihood estimator:

$$\hat{\boldsymbol{\mu}}_{\text{RMLE}} = \arg \max_{\boldsymbol{\mu}} \log P(\boldsymbol{x}; \boldsymbol{\mu}) - \lambda \|\boldsymbol{\mu}\|^{2}.$$

What is its mean squared error?

3. [4 Points] Consider an estimator of the form  $\hat{\mu}_{\text{SCALE}} = cx$  where  $c \in \mathbb{R}$  is a constant scaling factor. Find the value  $c^*$  that minimizes its mean squared error:

$$c^* = \arg\min_{c} \mathbb{E}[\|c\boldsymbol{x} - \boldsymbol{\mu}\|^2].$$

What is the corresponding minimum mean squared error?

4. Consider the James-Stein estimator:

$$\hat{\boldsymbol{\mu}}_{\mathrm{JS}} = \left(1 - \frac{p-2}{\left\|\boldsymbol{x}\right\|^2}\right) \boldsymbol{x}.$$

Note that  $\hat{\boldsymbol{\mu}}_{JS}$  can be written as  $\boldsymbol{x} - \boldsymbol{g}(\boldsymbol{x})$  where  $\boldsymbol{g}(\boldsymbol{x}) = \frac{p-2}{\|\boldsymbol{x}\|^2} \boldsymbol{x}$ . This allows us to separate the mean squared error into three parts:

$$\mathbb{E}[\|\hat{\boldsymbol{\mu}}_{\mathrm{JS}} - \boldsymbol{\mu}\|^{2}] = \mathbb{E}[\|\boldsymbol{x} - \boldsymbol{g}(\boldsymbol{x}) - \boldsymbol{\mu}\|^{2}]$$

$$= \mathbb{E}[\boldsymbol{x}^{\mathsf{T}}\boldsymbol{x} - 2\boldsymbol{x}^{\mathsf{T}}\boldsymbol{\mu} + \boldsymbol{\mu}^{\mathsf{T}}\boldsymbol{\mu} + \boldsymbol{g}(\boldsymbol{x})^{\mathsf{T}}\boldsymbol{g}(\boldsymbol{x}) - 2\boldsymbol{x}^{\mathsf{T}}\boldsymbol{g}(\boldsymbol{x}) + 2\boldsymbol{\mu}^{\mathsf{T}}\boldsymbol{g}(\boldsymbol{x})]$$

$$= \mathbb{E}[\|\boldsymbol{x} - \boldsymbol{\mu}\|^{2}] + \mathbb{E}[\|\boldsymbol{g}(\boldsymbol{x})\|^{2}] - 2\mathbb{E}[(\boldsymbol{x} - \boldsymbol{\mu})^{\mathsf{T}}\boldsymbol{g}(\boldsymbol{x})]$$

Furthermore, from Stein's lemma, we know that:

$$\mathbb{E}[(\boldsymbol{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{g}(\boldsymbol{x})] = \mathbb{E}\Big[\sum_{j=1}^p \frac{\partial}{\partial x_j} g_j(\boldsymbol{x})\Big]$$

(a) [1 Point] Find  $\mathbb{E}[\|\boldsymbol{x} - \boldsymbol{\mu}\|^2]$ .

- (b) [1 Point] Find  $\mathbb{E}[\|g(x)\|^2]$ . (Hint: your answer will include  $\mathbb{E}[\|x\|^{-2}]$ .)
- (c) [1 Point] Show that:

$$\frac{\partial}{\partial x_j} g_j(\boldsymbol{x}) = \frac{\left\|\boldsymbol{x}\right\|^2 - 2x_j^2}{\left\|\boldsymbol{x}\right\|^4}$$

where  $x_j$  is the jth element of  $\boldsymbol{x}$  and  $g_j(\boldsymbol{x})$  is the jth element of  $\boldsymbol{g}(\boldsymbol{x})$ .

- (d) [1 Point] What is the resulting mean squared error. (Hint: your answer will include  $\mathbb{E}[\|x\|^{-2}]$ .)
- 5. [4 Points] Qualitatively compare these estimators, noting any similarities between them. How does regularization affect an estimator's bias and variance? Which estimator would you choose to approximate  $\mu$  from real data about which you have no prior knowledge? How does the data dimensionality p affect your answer?

# Part B, Problem 2: Linear Regression [18 Points] (Adams and Weixiang)

Suppose we observe N data pairs  $\{(x_i, y_i)\}_{i=1}^N$ , where  $y_i$  is generated by the following rule:

$$y_i = x_i^\mathsf{T} \beta + \epsilon_i,$$

where  $x_i, \beta \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$ , and  $\epsilon_i$  is an i.i.d random noise drawn from the Gaussian Distribution:

$$\epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

with a known constant  $\sigma$ . We further denote  $Y = [y_1, y_2, ..., y_N]^\mathsf{T}$  and  $X = [x_1, x_2, ..., x_N]^\mathsf{T}$ .

Now, we are interested in estimating  $\beta$  from the observed data.

- 1. [3 Points] Derive the likelihood function  $\mathcal{L}(\beta)$ .
- 2. [5 Points] Show that the MLE estimator  $\hat{\beta}_{ml}$  of  $\beta$  is equivalent to the solution of the following linear regression problem:

$$\min_{\beta} \frac{1}{2} \|Y - X\beta\|_2^2 \tag{1}$$

3. [5 Points] Now we suppose  $\beta$  is not a deterministic parameter, but a random variable having a Gaussian prior distribution:

$$p(\beta) \sim \mathcal{N}(0, \frac{\sigma^2}{2\lambda}I),$$

where I is a  $d \times d$  identity matrix and  $\lambda > 0$  is a known parameter. Show that the MAP estimation  $\hat{\beta}_{\text{map}}$  of  $\beta$  is equivalent to the solution of the following ridge regression problem:

$$\min_{\beta} \frac{1}{2} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \tag{2}$$

4. [5 Points] Refer to the closed form solutions of (1) and (2) in the lecture slides, what might be an issue of  $\hat{\beta}_{ml}$  if  $d \gg N$ ? How can  $\hat{\beta}_{map}$  possibly address it?

# Part B, Problem 3: MLE, MAP and Logistic Regression [30 Points] (Prakhar and Yichong)

We learnt about Maximum Likelihood estimation in class. For a fixed set of data and underlying statistical model, the method of maximum likelihood selects the set of values of the model parameters that maximises the likelihood function.

In this problem, we will look at two different ways of estimating parameters in a probability distribution. Suppose we observe n iid random variables  $X_1, \ldots, X_n$ , drawn from a distribution with parameter  $\theta$ . That is, for each  $X_i$  and a natural number k,

$$P(X_i = k) = (1 - \theta)^k \theta$$

Given some observed values of  $X_1$  to  $X_n$ , we want to estimate the value of  $\theta$ .

## 3.1. Maximum Likelihood Estimation [9 Points]

The first kind of estimator for  $\theta$  we will consider is the Maximum Likelihood Estimator (MLE). The probability of observing given data is called the likelihood of the data, and the function that gives the likelihood for a given parameter  $\hat{\theta}$  (which may or may not be equal to the true parameter  $\theta$ ) is called the likelihood function, written as  $L(\hat{\theta})$ . When we use MLE, we estimate  $\theta$  by choosing the  $\hat{\theta}$  that maximizes the likelihood.

$$\hat{\theta}^{\text{MLE}} = \operatorname*{argmax}_{\hat{\theta}} L(\hat{\theta})$$

It is often convenient to deal with the log-likelihood  $(\ell(\hat{\theta}) = \log L(\hat{\theta}))$  instead, and since log is an increasing function, the argmax also applies in the log space:

$$\hat{\theta}^{\text{MLE}} = \operatorname*{argmax}_{\hat{\theta}} \ell(\hat{\theta})$$

- 1. [4 Points] Given a dataset  $\mathcal{D}$ , containing observations  $\{X_1 = k_1, X_2 = k_2, \dots, X_n = k_n\}$ , write an expression for  $\ell(\hat{\theta})$  as a function of  $\mathcal{D}$  and  $\hat{\theta}$ . How does the order of the variables affect the function?
- 2. [5 Points] Derive an expression for the maximum likelihood estimate.

## 3.2. Maximum a Posteriori Estimation [11 Points]

Now we assume that we have some prior knowledge about the true parameter  $\theta$ . We express it by treating  $\theta$  itself as a random variable and definining a prior probability distribution over it. Precisely, we suppose that the data  $X_1, \ldots, X_n$  are drawn as follows:

- $\theta$  is drawn from the prior probability distribution
- Then  $X_1, \ldots, X_n$  are drawn independently from a Geometric distribution with  $\theta$  as the parameter.

Now both  $X_i$  and  $\theta$  are random variables, and they have a joint probability distribution. We now estimate  $\theta$  as follows

$$\hat{\theta}^{\text{MAP}} = \operatorname*{argmax}_{\hat{\theta}} P(\theta = \hat{\theta} | X_1, \dots, X_n)$$

This is called Maximum a Posteriori (MAP) estimation. Using Bayes rule, we can rewrite the posterior probability as follows.

$$P(\theta = \hat{\theta}) = \frac{P(X_i, \dots, X_n | \theta = \hat{\theta}) P(\theta = \hat{\theta})}{P(X_1, \dots, X_n)}$$

Applying this to the MAP estimate, we get the following expression. Notice that we can ignore the denominator since it is not a function of  $\hat{\theta}$ .

$$\begin{split} \hat{\theta}^{\text{MAP}} &= \operatorname*{argmax}_{\hat{\theta}} P(X_1, \dots, X_n | \theta = \hat{\theta}) P(\theta = \hat{\theta}) \\ &= \operatorname*{argmax}_{\hat{\theta}} L(\hat{\theta}) P(\theta = \hat{\theta}) \\ &= \operatorname*{argmax}_{\hat{\theta}} \left( \ell(\hat{\theta}) + \log P(\theta = \hat{\theta}) \right) \end{split}$$

Thus, the MAP estimator maximizes the sum of the log-likelihood and the log-probability of the prior distribution on  $\theta$ . When the prior is a continuous distribution with density function p, we have

$$\hat{\theta}^{\text{MAP}} = \operatorname*{argmax}_{\hat{\theta}} \left( \ell(\hat{\theta}) + \log p(\hat{\theta}) \right)$$

For this problem, we will use the Beta distribution (a popular choice when the data distribution is Geometric or Bernoulli) as the prior, and the density function is given by

$$p(\hat{\theta}) = \frac{\hat{\theta}^{\alpha - 1} (1 - \hat{\theta})^{\beta - 1}}{B(\alpha, \beta)}$$

where  $B(\alpha, \beta)$  is the beta function.

- 4. [5 Points] Derive a close form expression for the maximum a posteriori estimate. (hint: If  $x^*$  maximizes f,  $f'(x^*) = 0$ ).
- 5. [3 Points] Is the bias of Maximum Likelihood Estimate (MLE) typically greater than or equal to the bias of Maximum A Posteriori (MAP) estimate? (Explain your answer in a sentence)
- 6. [3 Points] What can you say about the value of Maximum Likelihood Estimate (MLE) as compared to the value of Maximum A Posteriori (MAP) estimate with a uniform prior? Why?

## 3.3. Logistic Regression [10 Points]

In class, we will learn about MLE of parameters in logistic regression. For a given data  $x \in \mathbb{R}^p$ , the probability of Y being 1 in logistic regression is

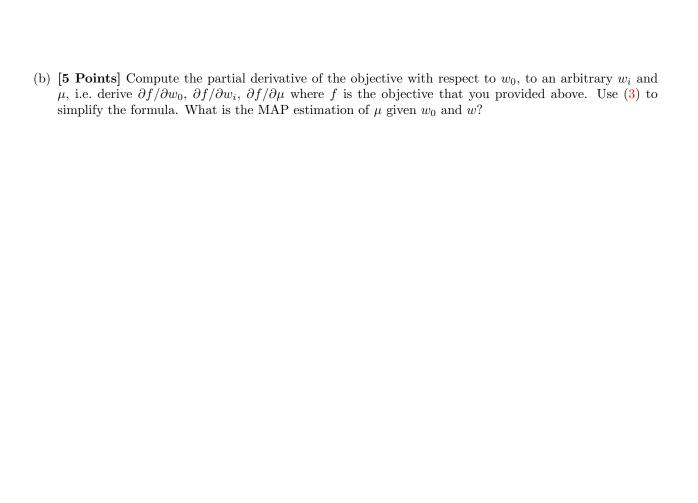
$$P(Y=1|X=x) = \frac{\exp(w_0 + x^T w)}{1 + \exp(w_0 + x^T w)},$$
(3)

where  $w_0$  and  $w = (w_1, w_2, ..., w_p)^T$  are model parameters. In this problem, we consider the maximum a posteriori setting, where we put a Gaussian prior on the parameters:

$$w_i \sim \mathcal{N}(\mu, 1)$$

for i = 0, 1, 2, ..., p.

(a) [5 Points] Choose a conjugate prior for Gaussian on  $\mu$  (choose any higher parameters as you want to ease the computation). Assuming you are given a dataset with n training examples and p features, write down a formula for the conditional log posterior likelihood of the training data in terms of the the class labels  $y^{(i)}$ , the features  $x_1^{(i)}, \ldots, x_p^{(i)}$ , and the parameters  $w_0, w_1, \ldots, w_p$ , where the superscript (i) denotes the sample index. This will be your objective function for gradient ascent.



## Part C: Programming Exercise [25 Points] (Dan and Danish)

Note: Your code for all of the programming exercises should be submitted to Gradescope. There is a separate 'programming assignment' that should allow you to upload your code easily. Code should be uploaded to this separate programming assignment, while visualizations and written answers should still be submitted within the primary Gradescope assignment. In your code, Please make it clear in the comments which are the primary functions to compute the answers to each question.

### Exploring The Effect of Priors in Batting Average Estimation: Dan [10 Points]

In this problem, you will explore how prior knowledge can effect your estimates of batting averages.

#### **Dataset**

In this problem, we have generated data for 5000 fictional baseball players. The data is divided into 3 parts - 'pre\_season.txt', 'mid\_season.txt', and 'end\_season.txt'. Each of these files has 3 columns: the id for the player (an integer), the number of at\_bats for the player (an at-bat is an opportunity to get a hit), and the number of hits the player got during those at-bats. The data files can be loaded using the provided load\_data function in hw1\_baseball.py. The batting average for a player can be computed by dividing the number of hits by the number of at\_bats.

### Maximum Likelihood Estimator [3 Points]

Assume for the moment that you only have access to the data in 'mid\_season.txt'. Midway through the season, you would like to estimate the end of season batting averages for all 5000 players. Write a function to compute the maximum likelihood estimate of the batting average for all 5000 players. Make sure to turn in your code.

### Maximum a Posteriori Estimator [3 Points]

Unsatisfied with the MLE estimates, you decide that you would like to use the pre-season statististics of the players as a prior on what their in-season batting averages will be. Write a function to compute the maximum a posteriori estimate of the batting average for all 5000 players. Briefly describe how you choose to incorporate prior information. Make sure to turn in your code.

#### Visualize Your Estimates [4 Points]

Compute the actual batting averages from 'end\_season.txt' (do not include statistics from the other files in these actual averages) and compare your estimates of the batting average to these estimates. Use the provided visualize function in hw1\_baseball.py to visualize and compare your MLE and MAP estimators. Make sure to turn in your visualizations.

- Does the MLE estimator appear to fail in certain cases? Why?
- Does the MAP estimator appear to fail in certain cases? Why?
- What conclusions do you draw from this experiment?

Note: The data files for this subproblem, and the following subproblem can be found here.

## Logistic Regression on Movie Review Dataset: Danish [15 Points]

In this problem, you will explore logistic regression to classify movie reviews into two classes - positive & negative. The dataset to be used is IMDB Large Movie Review dataset (Maas et. al, ACL 2011). The datafiles are present in the link shared above.

#### Details about dataset

The dataset comprises of two folders: 'train' and 'test', and each of these in turn contain two subfolders pos & neg. Each file in these subfolders is a unique review. In total, we have 25K training reviews (12.5K positive, and remaining 12.5K negative). The test folder too has 12.5K positive and 12.5K negative reviews. For our task, we will use bag of word representation.

#### Exercises

For this exercise, we will directly use Logistic Regression library from sklearn.linear\_models (feel free to use an equivalent library in any language of your choice). We will experiment with different values of  $C \in \{0.001, 0.01, 0.1, 1, 10, 100\}$ . Here, C is the inverse of regularization constant. We will also closely study the learnt parameter/weight/coefficient vector.

- [4 Points] Plot train and test accuracy for varying values of C. First plot should contain both train and test accuracy vs C with L2 regularizer (penalty) and the second plot should employ L1 regularizer (penalty). What do you observe in the two plots? Which value of C is optimum in these two cases?
- [2 Points] While using L2 regularizer, and different values of C, plot the L2 norm of weight vector vs C. What do you observe?
- [2 Points] While using L1 regularizer, and different values of C, plot the L1 norm of weight vector vs C. What do you observe?
- [2 Points] Study how sparsity (i.e percentage of zero elements in a vector) of the parameter vector changes with different values of C. In one plot, depict two curves one for L1 regularizer and the other one for L2 regularizer. Jot down your observations.

Now we will try to visualize the basis of the classification! One way to do so is to look at the weight vector and analyze the top (least) K values.

- [3 Points] While using L2 regularizer, and the optimum value of C (with respect to test accuracy), which 5 words correspond to the largest weight indices in the learnt weight vector? Which 5 words correspond to the least weight indices in the learnt parameter vector?
- [2 Points] While using L2 regularizer, and the optimum value of C (with respect to test accuracy), which review is predicted positive with highest probability? Similarly, which review is predicted negative with highest certainty?

### Implementation Details

• Please ignore all the words in the file 'stopwords.txt'. The vocabulary must be constructed by splitting the raw text only with whitespace characters, and converting them into lowercase.

- Since the vocabulary size will be close to 100K, and there are 25K reviews in training and test set, the bag of word matrix will be somewhat large, hence its best to use sparse formats wherever necessary.
- You may choose to ignore the words that occur in the test set but never show up in the training set.

Note: For the entire programming exercise, please turn in your code in a single zipped folder that might contain multiple source code files.