

Feel free to work with other students, but make sure you write up the homework and code on your own (no copying homework *or* code; no pair programming). Feel free to ask students or instructors for help debugging code or whatever else, though.

The starter files can be found under the Resource tab on course website. The graphs for problem 3 generated by the sample solution could be found in the corresponding zipfile. These graphs only serve as references to your implementation. You should generate your own graphs for submission. Please print out all the graphs generated by your own code and submit them together with the written part, and make sure you upload the code to your Github repository.

**1 (Murphy 8.3)** Gradient and Hessian of the log-likelihood for logistic regression.

(a) Let  $\sigma(x) = \frac{1}{1+e^{-x}}$  be the sigmoid function. Show that

$$\sigma'(x) = \sigma(x) [1 - \sigma(x)] .$$

(b) Using the previous result and the chain rule of calculus, derive an expression for the gradient of the log likelihood for logistic regression.

(c) The Hessian can be written as  $\mathbf{H} = \mathbf{X}^\top \mathbf{S} \mathbf{X}$  where  $\mathbf{S} = \text{diag}(\mu_1(1 - \mu_1), \dots, \mu_n(1 - \mu_n))$ . Derive this and show that  $\mathbf{H} \succeq 0$  ( $A \succeq 0$  means that  $A$  is positive semidefinite).

*Hint:* Use the **negative** log-likelihood of logistic regression for this problem.

a)

$$\begin{aligned} \sigma'(x) &= \frac{d}{dx} \left( \frac{1}{1 + e^{-x}} \right) \\ &= \frac{d}{dx} (1 + e^{-x})^{-1} \\ &= e^{-x} (1 + e^{-x})^{-2} \\ &= \left( \frac{1}{1 + e^{-x}} \right) \left( \frac{e^{-x}}{1 + e^{-x}} \right) \\ &= \sigma(x) \left( \frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \\ &= \sigma(x) \left( 1 - \frac{1}{1 + e^{-x}} \right) \\ &= \sigma(x) [1 - \sigma(x)] \end{aligned}$$

b)

Logistic regression expression for gradient of log-likelihood is :

$$n\ell\ell(\boldsymbol{\theta}) = - \sum_i y_i \log \sigma(\boldsymbol{\theta}^\top \mathbf{x}_i) + (1 - y_i) \log (1 - \sigma(\boldsymbol{\theta}^\top \mathbf{x}_i))$$

Now just take the gradients using theta

$$\begin{aligned} \nabla_{\boldsymbol{\theta}^{n\ell\ell}(\boldsymbol{\theta})} &= - \sum_i y_i \frac{1}{\sigma(\boldsymbol{\theta}^\top \mathbf{x}_i)} \sigma'(\boldsymbol{\theta}^\top \mathbf{x}_i) + (1 - y_i) \frac{1}{1 - \sigma(\boldsymbol{\theta}^\top \mathbf{x}_i)} (-\sigma'(\boldsymbol{\theta}^\top \mathbf{x}_i)) \\ &= - \sum_i y_i (1 - \sigma(\boldsymbol{\theta}^\top \mathbf{x}_i)) \mathbf{x}_i - (1 - y_i) \sigma(\boldsymbol{\theta}^\top \mathbf{x}_i) \mathbf{x}_i \\ &= - \sum_i y_i \mathbf{x}_i - y_i \sigma(\boldsymbol{\theta}^\top \mathbf{x}_i) \mathbf{x}_i - \sigma(\boldsymbol{\theta}^\top \mathbf{x}_i) \mathbf{x}_i + y_i \sigma(\boldsymbol{\theta}^\top \mathbf{x}_i) \mathbf{x}_i \\ &= \sum_i (\sigma(\boldsymbol{\theta}^\top \mathbf{x}_i) - y_i) \mathbf{x}_i \\ &= \mathbf{X}^\top (\sigma(\boldsymbol{\theta}^\top \mathbf{x}_i) - \mathbf{y}) \end{aligned}$$

Each  $\mathbf{x}_i$  is the  $i$  th column of  $\mathbf{X}$  transposed.

C) Using prev results, find Hessian matrix

$$\begin{aligned} \mathbf{H}_{\boldsymbol{\theta}} &= \nabla_{\boldsymbol{\theta}} (\nabla_{\boldsymbol{\theta}} n\ell\ell(\boldsymbol{\theta}))^\top = \nabla_{\boldsymbol{\theta}} [\mathbf{X}^\top (\boldsymbol{\mu} - \mathbf{y})]^\top \\ &= \nabla_{\boldsymbol{\theta}} (\boldsymbol{\mu}^\top \mathbf{X} - \mathbf{y}^\top \mathbf{X}) \\ &= \nabla_{\boldsymbol{\theta}} \boldsymbol{\mu}^\top \mathbf{X} = \nabla_{\boldsymbol{\theta}} \sigma(\mathbf{X}\boldsymbol{\theta})^\top \mathbf{X} \\ &= \mathbf{X}^\top \text{diag}(\boldsymbol{\mu}(1 - \boldsymbol{\mu})) \mathbf{X} \\ &= \mathbf{X}^\top \mathbf{S} \mathbf{X} \end{aligned}$$

To prove this is semi-definite, we to prove that  $\mathbf{S}$ 's eigenvalues are nonnegative.  $\mathbf{S}$  is a diagonal matrix, meaning its eigenvalues are its diagonal values. We can show this by showing  $\mu_i (1 - \mu_i) \geq 0$ . We can assert that:  $\mu_i (1 - \mu_i) = \sigma(\boldsymbol{\theta}^\top \mathbf{x}_i) (1 - \sigma(\boldsymbol{\theta}^\top \mathbf{x}_i))$ .  $\sigma(\cdot)(1 - \sigma(\cdot))$  is always nonnegative since  $\sigma(\cdot)$  is between 0 and 1. This shows that the Hessian matrix is positive semi-definite.

■

**2 (Murphy 2.11)** Derive the normalization constant ( $Z$ ) for a one dimensional zero-mean Gaussian

$$\mathbb{P}(x; \sigma^2) = \frac{1}{Z} \exp\left(-\frac{x^2}{2\sigma^2}\right)$$

such that  $\mathbb{P}(x; \sigma^2)$  becomes a valid density.

$$\int_{\mathbb{R}} \mathbb{P}(x; \sigma^2) dx = \int_{\mathbb{R}} \frac{1}{Z} \exp\left(-\frac{x^2}{2\sigma^2}\right) dx = \frac{1}{Z} \int_{\mathbb{R}} \exp\left(-\frac{x^2}{2\sigma^2}\right) dx = 1$$

which shows

$$Z = \int_{\mathbb{R}} \exp\left(-\frac{x^2}{2\sigma^2}\right) dx$$

Consider  $Z^2$

$$\begin{aligned} Z^2 &= \int_{\mathbb{R}} \exp\left(-\frac{x^2}{2\sigma^2}\right) dx \int_{\mathbb{R}} \exp\left(-\frac{y^2}{2\sigma^2}\right) dy \\ &= \iint_{\mathbb{R}^2} \exp\left(-\frac{x^2 + y^2}{2\sigma^2}\right) dx dy \\ &= \int_0^\infty \int_0^{2\pi} \exp\left(-\frac{r^2}{2\sigma^2}\right) r d\theta dr \\ &= 2\pi \int_0^\infty \exp\left(-\frac{r^2}{2\sigma^2}\right) r dr \\ &= 2\pi \left(-\sigma^2\right) \int_0^\infty \exp\left(-\frac{r^2}{2\sigma^2}\right) \left(-\frac{r}{\sigma^2}\right) dr \\ &= -2\pi\sigma^2 \exp\left(-\frac{r^2}{2\sigma^2}\right) \Big|_0^\infty \\ &= -2\pi\sigma^2(0 - 1) \\ &= 2\pi\sigma^2 \end{aligned}$$

$$Z^2 = 2\pi\sigma^2$$

$$Z = \sqrt{2\pi\sigma^2} = \sqrt{2\pi}\sigma$$

Probability density function integrate to 1 so becomes a valid density.

■

**3 (regression).** In this problem, we will use the online news popularity dataset to set up a model for linear regression. In the starter code, we have already parsed the data for you. However, you might need internet connection to access the data and therefore successfully run the starter code.

We split the csv file into a training and test set with the first two thirds of the data in the training set and the rest for testing. Of the testing data, we split the first half into a ‘validation set’ (used to optimize hyperparameters while leaving your testing data pristine) and the remaining half as your test set. We will use this data for the remainder of the problem. The goal of this data is to predict the **log** number of shares a news article will have given the other features.

- (a) **(math)** Show that the maximum a posteriori problem for linear regression with a zero-mean Gaussian prior  $\mathbb{P}(\mathbf{w}) = \prod_j \mathcal{N}(w_j|0, \tau^2)$  on the weights,

$$\arg \max_{\mathbf{w}} \sum_{i=1}^N \log \mathcal{N}(y_i | w_0 + \mathbf{w}^\top \mathbf{x}_i, \sigma^2) + \sum_{j=1}^D \log \mathcal{N}(w_j | 0, \tau^2)$$

is equivalent to the ridge regression problem

$$\arg \min \frac{1}{N} \sum_{i=1}^N (y_i - (w_0 + \mathbf{w}^\top \mathbf{x}_i))^2 + \lambda \|\mathbf{w}\|_2^2$$

with  $\lambda = \sigma^2 / \tau^2$ .

- (b) **(math)** Find a closed form solution  $\mathbf{x}^*$  to the ridge regression problem:

$$\text{minimize: } \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 + \|\mathbf{\Gamma}\mathbf{x}\|_2^2.$$

- (c) **(implementation)** Attempt to predict the log shares using ridge regression from the previous problem solution. Make sure you include a bias term and *don't regularize the bias term*. Find the optimal regularization parameter  $\lambda$  from the validation set. Plot both  $\lambda$  versus the validation RMSE (you should have tried at least 150 parameter settings randomly chosen between 0.0 and 150.0 because the dataset is small) and  $\lambda$  versus  $\|\boldsymbol{\theta}^*\|_2$  where  $\boldsymbol{\theta}$  is your weight vector. What is the final RMSE on the test set with the optimal  $\lambda^*$ ?

(continued on the following pages)

■

**3 (continued)**

- (d) **(math)** Consider regularized linear regression where we pull the basis term out of the feature vectors. That is, instead of computing  $\hat{\mathbf{y}} = \boldsymbol{\theta}^\top \mathbf{x}$  with  $\mathbf{x}_0 = 1$ , we compute  $\hat{\mathbf{y}} = \boldsymbol{\theta}^\top \mathbf{x} + b$ . This corresponds to solving the optimization problem

$$\text{minimize: } \|\mathbf{A}\mathbf{x} + b\mathbf{1} - \mathbf{y}\|_2^2 + \|\Gamma\mathbf{x}\|_2^2.$$

Solve for the optimal  $\mathbf{x}^*$  explicitly. Use this close form to compute the bias term for the previous problem (with the same regularization strategy). Make sure it is the same.

- (e) **(implementation)** We can also compute the solution to the least squares problem using gradient descent. Consider the same bias-relocated objective

$$\text{minimize: } f = \|\mathbf{A}\mathbf{x} + b\mathbf{1} - \mathbf{y}\|_2^2 + \|\Gamma\mathbf{x}\|_2^2.$$

Compute the gradients and run gradient descent. Plot the  $\ell_2$  norm between the optimal  $(\mathbf{x}^*, b^*)$  vector you computed in closed form and the iterates generated by gradient descent. Hint: your plot should move down and to the left and approach zero as the number of iterations increases. If it doesn't, try decreasing the learning rate.

a)

$$\arg \max_{\mathbf{w}} \sum_{i=1}^N \log \mathcal{N}(y_i | w_0 + \mathbf{w}^\top \mathbf{x}_i, \sigma^2) + \sum_{j=1}^D \log \mathcal{N}(w_j | 0, \tau^2)$$

After applying the probability distribution  $\mathcal{N}(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$ , we get

$$\begin{aligned} & \arg \max_{\mathbf{w}} \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - w_0 - \mathbf{w}^\top \mathbf{x}_i)^2}{2\sigma^2}\right) + \sum_{j=1}^D \log \frac{1}{\sqrt{2\pi}\tau} \exp\left(-\frac{w_j^2}{2\tau^2}\right) \\ &= \arg \max_{\mathbf{w}} \sum_{i=1}^N \left( -\frac{(y_i - w_0 - \mathbf{w}^\top \mathbf{x}_i)^2}{2\sigma^2} - \log \sqrt{2\pi}\sigma \right) + \sum_{j=1}^D \left( -\frac{w_j^2}{2\tau^2} - \log \sqrt{2\pi}\tau \right) \\ &= \arg \max_{\mathbf{w}} - \left( (N + D) \log \sqrt{2\pi}\sigma + \sum_{i=1}^N \frac{(y_i - w_0 - \mathbf{w}^\top \mathbf{x}_i)^2}{2\sigma^2} + \sum_{j=1}^D \frac{w_j^2}{2\tau^2} \right) \end{aligned}$$

To maximize function, we minimize its negative.

$$\arg \min_{\mathbf{w}} \sum_{i=1}^N \left( y_i - w_0 - \mathbf{w}^\top \mathbf{x}_i \right)^2 + \frac{\sigma^2}{\tau^2} \sum_{j=1}^D w_j^2$$

Replace  $\lambda$  with  $\sigma^2/\tau^2$

$$\begin{aligned} & \arg \min_{\mathbf{w}} \sum_{i=1}^N \left( y_i - w_0 - \mathbf{w}^\top \mathbf{x}_i \right)^2 + \lambda \sum_{j=1}^D w_j^2 \\ &= \arg \min_{\mathbf{w}} \sum_{i=1}^N \left( y_i - w_0 - \mathbf{w}^\top \mathbf{x}_i \right)^2 + \lambda \|\mathbf{w}\|_2^2 \end{aligned}$$

b) Find 0 value of gradient of  $f$  with respect to  $\mathbf{x}$  in order to minimize  $f = \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 + \|\Gamma\mathbf{x}\|_2^2$

$$\begin{aligned} \nabla_{\mathbf{x}} f &= \nabla_{\mathbf{x}} \left[ (\mathbf{A}\mathbf{x} - \mathbf{b})^\top (\mathbf{A}\mathbf{x} - \mathbf{b}) + (\Gamma\mathbf{x})^\top (\Gamma\mathbf{x}) \right] \\ &= \nabla_{\mathbf{x}} \left[ \left( \mathbf{x}^\top \mathbf{A}^\top - \mathbf{b}^\top \right) (\mathbf{A}\mathbf{x} - \mathbf{b}) + \mathbf{x}^\top \Gamma^\top \Gamma \mathbf{x} \right] \\ &= \nabla_{\mathbf{x}} \left[ \mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} - 2\mathbf{x}^\top \mathbf{A}^\top \mathbf{b} + \mathbf{b}^\top \mathbf{b} + \mathbf{x}^\top \Gamma^\top \Gamma \mathbf{x} \right] \\ &= 2\mathbf{A}^\top \mathbf{A} \mathbf{x} - 2\mathbf{A}^\top \mathbf{b} + 2\Gamma^\top \Gamma \mathbf{x} \\ \nabla_{\mathbf{x}} f(0) &= \left( \mathbf{A}^\top \mathbf{A} + \Gamma^\top \Gamma \right) \mathbf{x} = \mathbf{A}^\top \mathbf{b} \end{aligned}$$

Closed form solution:  $\mathbf{x}^* = (\mathbf{A}^\top \mathbf{A} + \Gamma^\top \Gamma)^{-1} \mathbf{A}^\top \mathbf{b}$

c)

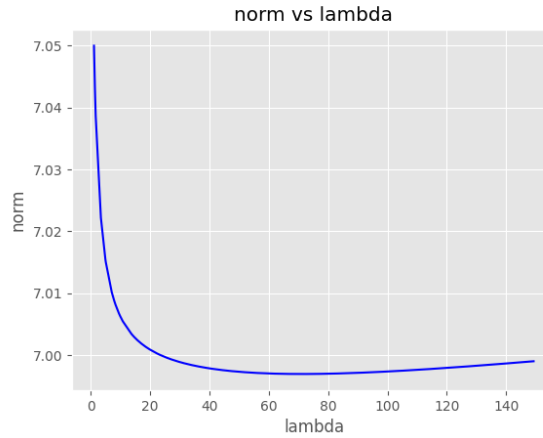


Figure 1: Norm vs Lambda

The optimal regularization parameter is 8.9946. The RMSE on the validation set with the optimal regularization parameter is 0.8341. The RMSE on the test set with the optimal regularization parameter is 0.8628.

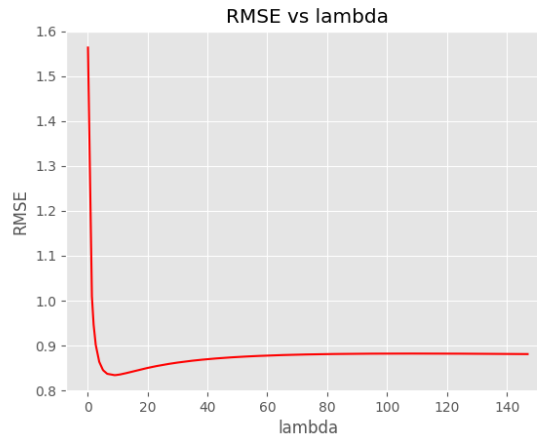


Figure 2: RMSE vs Lambda

d)

$$\begin{aligned}
 f &= \|A\mathbf{x} + b\mathbf{1} - \mathbf{y}\|_2^2 + \|\Gamma\mathbf{x}\|_2^2 \\
 &= (A\mathbf{x} + b\mathbf{1} - \mathbf{y})^\top (A\mathbf{x} + b\mathbf{1} - \mathbf{y}) + (\Gamma\mathbf{x})^\top (\Gamma\mathbf{x}) \\
 &= (\mathbf{x}^\top A^\top + b\mathbf{1}^\top - \mathbf{y}^\top) (A\mathbf{x} + b\mathbf{1} - \mathbf{y}) + \mathbf{x}^\top \Gamma^\top \Gamma \mathbf{x} \\
 &= \mathbf{x}^\top A^\top A \mathbf{x} + 2b\mathbf{1}^\top A \mathbf{x} - 2\mathbf{y}^\top A \mathbf{x} - 2b\mathbf{1}^\top \mathbf{y} + b^2 n + \mathbf{y}^\top \mathbf{y} + \mathbf{x}^\top \Gamma^\top \Gamma \mathbf{x}
 \end{aligned}$$

We know at Optimality is when  $\nabla_{\mathbf{x}} f = 0$ , so find when  $\nabla_{\mathbf{b}} f$  is 0.

$$\begin{aligned}
 \nabla_{\mathbf{x}} f &= 2A^\top A \mathbf{x} + 2bA^\top \mathbf{1} - 2A^\top \mathbf{y} + 2\Gamma^\top \Gamma \mathbf{x} \\
 \nabla_{\mathbf{b}} f &= 2\mathbf{1}^\top A \mathbf{x} - 2\mathbf{1}^\top \mathbf{y} + 2bn \\
 0 &= 2\mathbf{1}^\top A \mathbf{x} - 2\mathbf{1}^\top \mathbf{y} + 2bn
 \end{aligned}$$

$$b^* = \frac{\mathbf{1}^\top (\mathbf{y} - A\mathbf{x})}{n}$$

Shows bias term is the mean value of output when  $\mathbf{x} = 0$ .

Plugging  $b^*$  back to equation to solve for  $\mathbf{x}^*$ .

$$\begin{aligned}
& \left( A^\top A + \Gamma^\top \Gamma \right) \mathbf{x} + \left( \frac{\mathbf{1}^\top (\mathbf{y} - A\mathbf{x})}{n} \right) A^\top \mathbf{1} - A^\top \mathbf{y} = 0 \\
& \left( A^\top A + \Gamma^\top \Gamma \right) \mathbf{x} + \frac{1}{n} A^\top \mathbf{1} \mathbf{1}^\top \mathbf{y} - \frac{1}{n} A^\top \mathbf{1} \mathbf{1}^\top A \mathbf{x} - A^\top \mathbf{y} = 0 \\
& \left[ A^\top A + \Gamma^\top \Gamma - \frac{1}{n} A^\top \mathbf{1} \mathbf{1}^\top A \right] \mathbf{x} = A^\top \mathbf{y} - \frac{1}{n} A^\top \mathbf{1} \mathbf{1}^\top \mathbf{y} \\
& \left[ A^\top \left( \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \right) A + \Gamma^\top \Gamma \right] \mathbf{x} = A^\top \left( \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \right) \mathbf{y} \\
& \mathbf{x}^* = \left[ A^\top \left( \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \right) A + \Gamma^\top \Gamma \right]^{-1} A^\top \left( \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \right) \mathbf{y}
\end{aligned}$$

where  $\mathbf{I}$  is the identity matrix,  $\mathbf{1}$  is a vector of all ones, and  $\mathbf{y} \in \mathbf{R}^n$ .

Bias from sample code and part(c) are extremely similar

Difference in bias is 1.9661E-11

Difference in weights is 2.3683E-10

e)

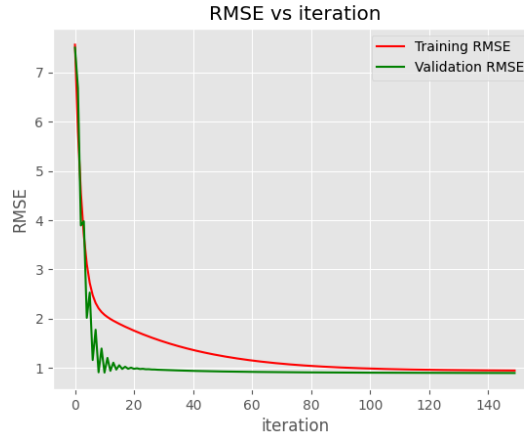


Figure 3: Convergence Plot

Difference in bias is 1.5386E-01

Difference in weights is 7.9580E-01

■