

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 $\mathbf{H} = \mathbf{X}^T \mathbf{S} \mathbf{X}$ 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. Using chain rule,

$$\begin{aligned}\sigma'(x) &= -(1 + e^{-x})^{-2}(-e^{-x}) \\ &= \frac{1}{(1 + e^{-x})} \frac{e^{-x}}{1 + e^{-x}} \\ &= \sigma(x)(1 - \sigma(x))\end{aligned}$$

- b. The negative log likelihood equation is

$$NLL(\mathbf{w}) = -\sum_{i=1}^N [y_i \log \sigma(\mathbf{w}^T x_i) + (1 - y_i) \log(1 - \sigma(\mathbf{w}^T x_i))]$$

Taking the derivative and using the result from a), we get

$$\begin{aligned}\nabla NLL(\mathbf{w}) &= -\sum_{i=1}^N \left[y_i \frac{d}{d\mathbf{w}} (\log \sigma(\mathbf{w}^T x_i)) + \frac{d}{d\mathbf{w}} (1 - y_i) \log(1 - \sigma(\mathbf{w}^T x_i)) \right] \\ &= -\sum_{i=1}^N \left[y_i \frac{1}{\sigma(\mathbf{w}^T x_i)} \sigma'(\mathbf{w}^T x_i) + (1 - y_i) \frac{1}{1 - \sigma(\mathbf{w}^T x_i)} (-\sigma'(\mathbf{w}^T x_i)) \right] \\ &= -\sum_{i=1}^N [y_i(1 - \sigma(\mathbf{w}^T x_i))x_i - (1 - y_i)\sigma(\mathbf{w}^T x_i)x_i] \\ &= -\sum_{i=1}^N [y_i x_i - y_i \sigma(\mathbf{w}^T x_i)x_i - \sigma(\mathbf{w}^T x_i)x_i + y_i \sigma(\mathbf{w}^T x_i)x_i] \\ &= -\sum_{i=1}^N (y_i x_i - \sigma(\mathbf{w}^T x_i)x_i)\end{aligned}$$

Let $\mu_i = \sigma(\mathbf{w}^T x_i)$. Then,

$$\begin{aligned}\nabla NLL(\mathbf{w}) &= -\sum_{i=1}^N (y_i x_i - \mu_i x_i) \\ &= X^T(\boldsymbol{\mu} - \mathbf{y})\end{aligned}$$

c. In order to find the Hessian matrix, we need to take the derivative from the gradient found in (b):

$$\begin{aligned}H &= \nabla(\nabla NLL\mathbf{w})^T \\ &= \nabla(X^T \boldsymbol{\mu} - X^T \mathbf{y})^T\end{aligned}$$

Using the matrix transpose distribution rule, $H = \nabla(\boldsymbol{\mu}^T X - \mathbf{y}^T X)$. If we substitute σ back in, since the second term is not dependent on \mathbf{w} , we can disregard it:

$$\begin{aligned}H &= \nabla \boldsymbol{\mu}^T X \\ &= \nabla(\sigma(\mathbf{w}^T X))^T X \\ &= \nabla \sigma(X \mathbf{w})^T X \\ &= X^T \text{diag}(\boldsymbol{\mu}(1 - \boldsymbol{\mu})) X \\ &= X^T S X\end{aligned}$$

By definition, H is positive semi-definite if all of its eigenvalues are non-negative. Since H can be diagonalized as $X^T S X$, the diagonal entries of S are its eigenvalues and we need to show that they are non-negative. Since $\mu_i = \frac{1}{(1+e^{-\mathbf{w}^T x_i})}$, $\mu_i > 0$ and $\mu_i < 1$, thus $1 - \mu_i > 0$. Therefore, the diagonal entries, or in other words the eigenvalues of H , $\mu_i(1 - \mu_i) > 0$, so H must be positive semidefinite.

■

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.

In order for the distribution to be a valid density, $\int_{-\infty}^{\infty} \frac{1}{Z} e^{-\frac{x^2}{2\sigma^2}} dx = 1$. Rewriting this, we get $\int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} dx = Z$. Since this form cannot be integrated easily, the textbook hints on computing Z^2 instead, using polar coordinate.

$$\begin{aligned} Z^2 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-\frac{x^2+y^2}{2\sigma^2}} dx dy \\ &= \int_0^{2\pi} \int_0^{\infty} e^{-\frac{x^2+y^2}{2\sigma^2}} dr d\theta \\ &= \int_0^{2\pi} \int_0^{\infty} e^{-\frac{r^2}{2\sigma^2}} r dr d\theta \\ &= \int_0^{2\pi} \int_0^{\infty} \sigma^2 e^{-u} du d\theta \\ &= \int_0^{2\pi} \sigma^2 d\theta \\ &= 2\pi\sigma^2 \end{aligned}$$

Thus, $Z = \sqrt{2\pi}\sigma$

■

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:

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

- (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 λ from the validation set. Plot both λ 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 λ versus $\|\theta^*\|_2$ where θ is your weight vector. What is the final RMSE on the test set with the optimal λ^* ?

(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 $\hat{y} = \theta^\top x$ with $x_0 = 1$, we compute $\hat{y} = \theta^\top x + b$. This corresponds to solving the optimization problem

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

Solve for the optimal 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

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

Compute the gradients and run gradient descent. Plot the ℓ_2 norm between the optimal (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) For a Gaussian distribution, $N(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$. Thus, if we substitute in the MAP problem, we get

$$\begin{aligned} & \underset{w}{\operatorname{argmax}} \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - w_0 - \vec{w}^\top \vec{x}_i)^2}{2\sigma^2}} + \sum_{j=1}^D \log \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{w_j^2}{2\sigma^2}} \\ &= \underset{w}{\operatorname{argmax}} - \sum_{i=1}^N \left(\frac{(y_i - w_0 - \vec{w}^\top \vec{x}_i)^2}{2\sigma^2} + \log \frac{1}{\sqrt{2\pi}\sigma} \right) - \sum_{j=1}^D \left(\frac{w_j^2}{2\sigma^2} + \log \frac{1}{\sqrt{2\pi}\sigma} \right) \\ &= \underset{w}{\operatorname{argmax}} - \left[\underbrace{\log \frac{1}{\sqrt{2\pi}\sigma} \cdot (N+D)}_{\text{constants}} + \sum_{i=1}^N \frac{(y_i - w_0 - \vec{w}^\top \vec{x}_i)^2}{2\sigma^2} + \sum_{j=1}^D \frac{w_j^2}{2\sigma^2} \right] \end{aligned}$$

Since these are constants, and argmax answers where the max occurs and not what the max is, they can be ignored. Scaling can also be ignored thus, it becomes

$$\underset{w}{\operatorname{argmax}} - \left(\sum_{i=1}^N (y_i - w_0 - \vec{w}^\top \vec{x}_i)^2 + \frac{\sigma^2}{2} \sum_{j=1}^D w_j^2 \right)$$

Negative of $\operatorname{argmax} = \operatorname{argmin}$.

$$\begin{aligned} &= \underset{w}{\operatorname{argmin}} \left(\sum_{i=1}^N (y_i - w_0 - \vec{w}^\top \vec{x}_i)^2 + \frac{\sigma^2}{2} \sum_{j=1}^D w_j^2 \right) \\ &= \underset{w}{\operatorname{argmin}} \sum_{i=1}^N (y_i - w_0 - \vec{w}^\top \vec{x}_i)^2 + \frac{\sigma^2}{2} ||w||_2^2 \end{aligned}$$

Note that $\|x\|_p = (\sum_i |x_i|^p)^{1/p}$

Continue Problem 3

b) Suppose $f(\vec{x}) = \|A\vec{x} - \vec{b}\|_2^2 + \|\Gamma\vec{x}\|_2^2$

In order to minimize $f(\vec{x})$, we need to find its derivative:

$$\begin{aligned}\nabla f(\vec{x}) &= \nabla [(A\vec{x} - \vec{b})^T (A\vec{x} - \vec{b}) + (\Gamma\vec{x})^T (\Gamma\vec{x})] \quad \text{(distribute } T\text{)} \\ &= \nabla [(\vec{x}^T A^T - \vec{b}^T)(A\vec{x} - \vec{b}) + \vec{x}^T \Gamma^T \Gamma \vec{x}] \\ &= \nabla (\vec{x}^T A^T A \vec{x} - \vec{x}^T A^T \vec{b} - \underbrace{\vec{b}^T A \vec{x}}_{\text{Note that } \vec{x}^T A^T \vec{b} = \vec{b}^T A \vec{x}} + \vec{b}^T \vec{b} + \vec{x}^T \Gamma^T \Gamma \vec{x}) \\ &= \nabla (\vec{x}^T A^T A \vec{x} - 2\vec{x}^T A^T \vec{b} + \vec{b}^T \vec{b} + \vec{x}^T \Gamma^T \Gamma \vec{x}) \\ &= 2A^T A \vec{x} - 2A^T \vec{b} + 2\Gamma^T \Gamma \vec{x}\end{aligned}$$

Solve for \vec{x} by setting $\nabla_x f = 0$:

$$0 = 2(A^T A \vec{x} - A^T \vec{b} + \Gamma^T \Gamma \vec{x})$$

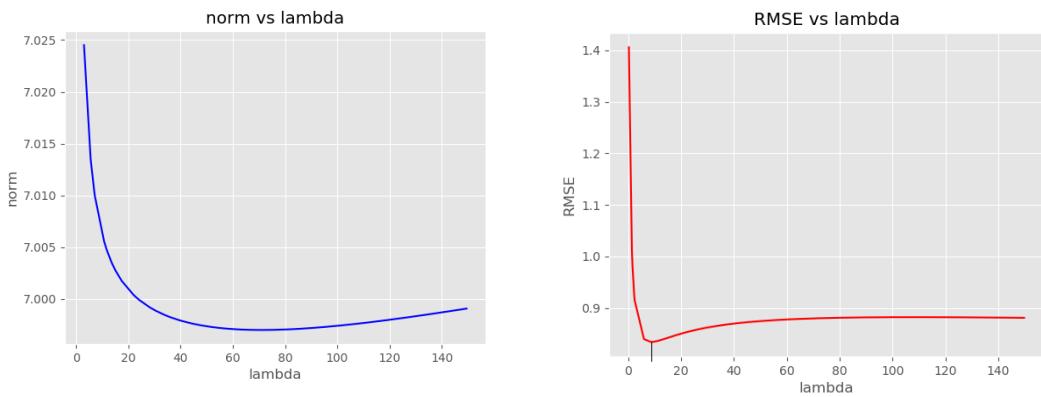
$$A^T \vec{b} = A^T A \vec{x} + \Gamma^T \Gamma \vec{x}$$

$$A^T \vec{b} = (A^T A + \Gamma^T \Gamma) \vec{x}$$

$$\boxed{\vec{x} = (A^T A + \Gamma^T \Gamma)^{-1} A^T \vec{b}}$$

(continued part 3 on next page)

c)



The optimal RMSE on validation set is 0.8340.

The optimal RMSE on test set is 0.8628.

$$\vec{x}^* = 9.0262$$

d) $f(\vec{x}) = \|A\vec{x} + b - \vec{y}\|_2^2 + \|\Gamma\vec{x}\|^2$

Norm 2 can be re-written into the following:

$$\begin{aligned} f(\vec{x}) &= (A\vec{x} + b - \vec{y})^T (A\vec{x} + b - \vec{y}) + (\Gamma\vec{x})^T (\Gamma\vec{x}) \\ &= (\vec{x}^T A^T + \vec{y}^T \cdot b^T - \vec{y}^T) (A\vec{x} + b - \vec{y}) + (\vec{x}^T \Gamma^T) (\Gamma\vec{x}) \\ &= \cancel{\vec{x}^T A^T \vec{x}} + \cancel{\vec{x}^T A^T b^T} - \cancel{\vec{x}^T A^T \vec{y}} + \cancel{\vec{y}^T b^T A \vec{x}} + \cancel{\vec{y}^T b^T b^T} - \cancel{\vec{y}^T b^T \vec{y}} \\ &\quad - \cancel{\vec{y}^T A \vec{x}} - \cancel{\vec{y}^T b^T} + \cancel{\vec{y}^T \vec{y}} + \cancel{\vec{x}^T \Gamma^T \Gamma \vec{x}} \\ &= \vec{x}^T A^T \vec{x} + 2\vec{x}^T A^T b^T - 2\vec{x}^T A^T \vec{y} - 2\vec{y}^T b^T \vec{y} + \vec{y}^T b^T b^T \\ &\quad + \vec{y}^T \vec{y} + \vec{x}^T \Gamma^T \Gamma \vec{x} \end{aligned}$$

To find the optimal b , we need to set $0 = \nabla_b f(\vec{x})$

$$\begin{aligned} 0 &= \nabla_b (\vec{x}^T A^T \vec{x} + 2\vec{x}^T A^T b^T - 2\vec{x}^T A^T \vec{y} - 2\vec{y}^T b^T \vec{y} + \vec{y}^T b^T b^T \\ &\quad + \vec{y}^T \vec{y} + \vec{x}^T \Gamma^T \Gamma \vec{x}) \end{aligned} \quad \begin{matrix} \curvearrowright & \vec{y} \text{ has dimension } 1 \times n \\ & \text{so } \Gamma^T b^T \vec{y} = b^T \vec{y} \end{matrix}$$

Solving for b , we get

$$b = \frac{\Gamma^T \vec{y} - \vec{x}^T A^T \vec{x}}{n} = \frac{\Gamma^T (\vec{y} - A\vec{x})}{n} \quad (1)$$

To find the optimal \vec{x} , we need to set $0 = \nabla_{\vec{x}} f(\vec{x})$

$$\begin{aligned} 0 &= \nabla_{\vec{x}} f(\vec{x}) \\ &= \nabla_{\vec{x}} (\vec{x}^T A^T \vec{x} + 2\vec{x}^T A^T b^T - 2\vec{x}^T A^T \vec{y} - 2\vec{y}^T b^T \vec{y} + \vec{y}^T b^T b^T \\ &\quad + \vec{y}^T \vec{y} + \vec{x}^T \Gamma^T \Gamma \vec{x}) \end{aligned} \quad \begin{matrix} & \text{not dependent on } \vec{x} \end{matrix}$$

(Continue)

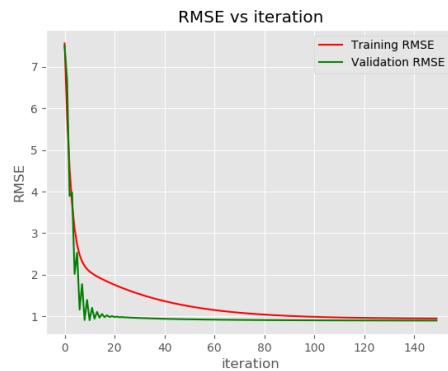
If we substitute in the optimal b ,

$$\begin{aligned} 0 &= A^T A \vec{x} + A^T \left[\frac{T^T(I - A\vec{x})}{n} \right] \vec{i} - A^T \vec{y} + T^T T \vec{x} \\ 0 &= (A^T A - \frac{A^T T^T A \vec{i}}{n} + T^T T) \vec{x} + \frac{A^T T^T \vec{y} \vec{i}}{n} - A^T \vec{y} \\ \vec{x} &= \left[A^T (I - \frac{\vec{i} \vec{i}^T}{n}) A + T^T T \right]^{-1} A^T (I - \frac{\vec{i} \vec{i}^T}{n}) \vec{y} \quad (2) \end{aligned}$$

Substituting (2) in code, we get:

difference in bias = $4.6495 \times 10^{-10} \approx 0$ ⇒ same as previous
 difference in weights = $6.2312 \times 10^{-10} \approx 0$

- c) difference in bias = 1.5386×10^{-1}
 difference in weights = 7.9639×10^{-1}



∴

$$\begin{aligned} \text{Since } \vec{x}^T A^T b \vec{1} &= \begin{bmatrix} (x_1 a_{11} + x_2 a_{21} + \dots) & (x_1 a_{12} + x_2 a_{22} + \dots) & \dots \end{bmatrix} \begin{bmatrix} b \\ b \\ \vdots \end{bmatrix} \\ &= (x_1 a_{11} + x_2 a_{21} + \dots) b + (x_1 a_{12} + x_2 a_{22} + \dots) b + \dots (x_1 a_{1n} + x_2 a_{2n} + \dots) b \\ &= [b \quad b \quad \dots] \begin{bmatrix} (x_1 a_{11} + x_2 a_{21} + \dots) \\ (x_1 a_{12} + x_2 a_{22} + \dots) \\ \vdots \end{bmatrix} \\ &= \vec{1} b \begin{bmatrix} a_{11} & a_{21} & \dots \\ a_{12} & a_{22} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix} \\ &= \vec{1} b A^T \vec{x} \end{aligned}$$