# Math of Big Data, Summer 2018
## Prof: Gu

| | |
|---|---|
| **Name:** | **Forest Kobayashi** |
| **HW #:** | **2** |
| **Day:** | Mon.    Tue.    (Wed.)    Thu.    Fri. |
| **Date:** | **05/16/2018** |

| No. | Points | Acknowledgments |
|---|---|---|
| 1 | | Murphy, Stackexchange (thread about Hessian) |
| 2 | | |
| 3 | | |
| **Total** | | |

This Assignment is (check one):

☒ On Time      ☐ Late, without deduction      ☐ Late, with deduction

**Comments:** 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.

# Problem 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 $\mathbf{S} = \mathrm{diag}(\mu_1(1 - \mu_1), \ldots, \mu_n(1 - \mu_n))$. 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.

---

## Solution:

(a) Note that

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

(b) We have

$$\mathrm{NLL}(\mathbf{w}) = -\sum_{i=1}^{N} y_i \log\left(h(\mathbf{x}_i)\right) + (1 - y_i) \log(1 - h(\mathbf{x}_i))$$

where $h(\mathbf{x}_i) = \sigma(\boldsymbol{\theta}^\top \mathbf{x}_i)$. Hence

$$
\begin{aligned}
\nabla_{\boldsymbol{\theta}} \mathrm{NLL}(\mathbf{w}) &= -\sum_{i=1}^{N} \frac{\partial}{\partial \theta_i} \left[y_i \log\left(h(\mathbf{x}_i)\right) + (1 - y_i) \log(1 - h(\mathbf{x}_i))\right] \\
&= -\sum_{i=1}^{N} \left[\frac{y_i}{h(\mathbf{x}_i)} \frac{\partial h(\mathbf{x}_i)}{\partial \boldsymbol{\theta}} - \frac{1 - y_i}{1 - h(\mathbf{x}_i)} h'(\mathbf{x}_i) \frac{\partial h(\mathbf{x}_i)}{\partial \boldsymbol{\theta}}\right]
\end{aligned}
$$

by the previous result, we have

$$\frac{\partial h(\mathbf{x}_i)}{\partial \boldsymbol{\theta}} = h(\mathbf{x}_i)(1 - h(\mathbf{x}_i))$$

hence

$$\nabla_{\boldsymbol{\theta}} \mathrm{NLL}(\mathbf{w}) = -\sum_{i=1}^{N} [y_i(1 - h(\mathbf{x}_i))\mathbf{x}_i - (1 - y_i)h(\mathbf{x}_i)\mathbf{x}_i]$$

$$= -\sum_{i=1}^{N} \left[ \mathbf{x}_i (y_i(1 - h(\mathbf{x}_i)) - (1 - y_i)h(\mathbf{x}_i)) \right]$$

$$= -\sum_{i=1}^{N} \left[ \mathbf{x}_i (y_i - y_i h(\mathbf{x}_i) - h(\mathbf{x}_i) + y_i h(\mathbf{x}_i)) \right]$$

$$= -\sum_{i=1}^{N} \left[ \mathbf{x}_i (y_i - h(\mathbf{x}_i)) \right]$$

$$= \sum_{i=1}^{N} \left[ \mathbf{x}_i (h(\mathbf{x}_i) - y_i) \right]$$

checking with the solutions, we see that it'd be a good idea to define $\mu_i = \sigma(\boldsymbol{\theta}^\top \mathbf{x}_i)$. Then

$$\nabla_{\boldsymbol{\theta}} \mathrm{NLL}(\mathbf{w}) = \sum_{i=1}^{N} \mathbf{x}_i(\mu_i - y_i)$$

$$= \mathbf{X}^\top (\boldsymbol{\mu} - \mathbf{y})$$

(c) Let $\mathbf{J}$ denote the Jacobian matrix. Then

$$\mathbf{H}(\mathrm{NLL}(\mathbf{w})) = \mathbf{J}(\nabla_{\boldsymbol{\theta}} \mathrm{NLL}(\mathbf{w}))^\top$$

$$= \mathbf{J}(\mathbf{X}^\top (\boldsymbol{\mu} - \mathbf{y}))^\top$$

$$= \mathbf{J}(\mathbf{X}^\top \boldsymbol{\mu})^\top$$

$$= \sum_{i=1}^{N} \left( \frac{\partial \mu_i}{\partial \boldsymbol{\theta}} \right) \mathbf{x}_i^\top$$

$$= \sum_{i=1}^{N} \left( \mathbf{x}_i (\mu_i(1 - \mu_i)) \mathbf{x}_i^\top \right)$$

$$= \mathbf{X}^\top \mathrm{diag}(\mu_1(1 - \mu_1), \ldots, \mu_n(1 - \mu_n)) \mathbf{X}$$

$$= \mathbf{X}^\top \mathbf{S} \mathbf{X}$$

Note that by the definition of $\sigma(x)$, $\forall i \in [N]$, $0 < \mu_i < 1$. Hence $1 - \mu_i > 0$. Thus $\forall i \in [N]$, $\mu_i(1 - \mu_i) > 0$. Now, observe that $\mathbf{S}$ is the diagonal matrix of eigenvalues of $\mathbf{H}$, which are all of the form $\mu_i(1 - \mu_i)$. Thus, all the eigenvalues of $\mathbf{H}$ are nonnegative, so $\mathbf{H}$ is positive semidefinite. ∎

## Problem 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.

---

## Solution:

For $\mathbb{P}(x; \sigma^2)$ to be normalized, it must integrate to 1 over $\mathbb{R}$. Hence, we must have

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

because the antiderivative is not expressable in terms of elementary functions, we pull out some trickery. Squaring $Z$,

$$\begin{aligned}
Z^2 &= \left(\int_{-\infty}^{\infty} \exp\left(-\frac{x^2}{2\sigma^2}\right) dx\right)^2 \\
&= \int_{-\infty}^{\infty} \exp\left(-\frac{x_0^2}{2\sigma^2}\right) dx_0 \int_{-\infty}^{\infty} \exp\left(-\frac{x_1^2}{2\sigma^2}\right) dx_1 \\
&= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp\left(-\frac{x_0^2 + x_1^2}{2\sigma^2}\right) dx_0\, dx_1
\end{aligned}$$

We want to convert to polar coordinates. Consider the transformation $x_0 \mapsto r\cos(\theta)$; $x_1 \mapsto r\sin(\theta)$. Then $x_0^2 + y_0^2 = r^2(\cos^2(\theta) + \sin^2(\theta)) = r^2$, and by the power of Math 60 (and google), we recall that $dx\,dy \mapsto r\,dr\,d\theta$. Hence, our integral becomes

$$\begin{aligned}
Z^2 &= \int_0^{2\pi} \int_0^{\infty} r \exp\left(-\frac{r}{2\sigma^2} dr\, d\theta\right) \\
&= \int_0^{2\pi} \left(-\sigma^2 \exp\left(-\frac{r^2}{2\sigma^2}\right)\Big|_0^{\infty}\right) d\theta \\
&= \int_0^{2\pi} \sigma^2\, d\theta \\
&= 2\pi\sigma^2
\end{aligned}$$

hence

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

as desired.                                                                                     ■

# Problem 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

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

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

(b) (**math**) Find a closed form solution $\mathbf{x}^\star$ 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 $\lambda$ 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 $\|\boldsymbol{\theta}^\star\|_2$ where $\boldsymbol{\theta}$ is your weight vector. What is the final RMSE on the test set with the optimal $\lambda^\star$?

(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

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

Solve for the optimal $\mathbf{x}^\star$ 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

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

Compute the gradients and run gradient descent. Plot the $\ell_2$ norm between the optimal $(\mathbf{x}^\star, b^\star)$ 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.

## Solution:

(a) Substituting in our expression for the standard normal probability distribution, we have

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

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

The constants don't matter (independent of $\mathbf{w}$), so we ignore them.

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

multiplying out by $2\sigma^2$, we obtain

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

$$= -\left( \sum_{i=1}^{N} \left[ \left(y_i - (w_0 + \mathbf{w}^\top \mathbf{x}_i)\right)^2 \right] + \lambda \sum_{j=1}^{D} \left[ w_j^2 \right] \right)$$

by definition of $\ell_2$ norm, this becomes

$$= -\left( \sum_{i=1}^{N} \left[ \left(y_i - (w_0 + \mathbf{w}^\top \mathbf{x}_i)\right)^2 \right] + \lambda \|\mathbf{w}\|_2^2 \right)$$

$\text{MAP}(\mathbf{x})$ is maximized when the part inside the parens is minimized, hence we seek to solve the following estimation problem:

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

Supposedly, I'm missing a factor of $1/N$. But, without redefining $\lambda$, I don't see how I could obtain one.

(b) First, we prove a small lemma.

**Lemma 3.1.** *Let* $\mathbf{x} \in \mathbb{R}^n$*, and* $A \in \mathcal{L}(\mathbb{R}^n)$*. Then* $\nabla_{\mathbf{x}}(\mathbf{x}^\top A\mathbf{x}) = \mathbf{x}^\top(A + A^\top)$*.*

*Proof.* By the product rule, we have

$$\nabla_{\mathbf{x}}(\mathbf{x}^\top A\mathbf{x}) = \left(\nabla_{\mathbf{x}}(\mathbf{x}^\top)\right)A\mathbf{x} + \mathbf{x}^\top A(\nabla_{\mathbf{x}}(\mathbf{x}))$$

because these are all $1 \times 1$ matrices, they are symmetric, hence

$$= \left(\nabla_{\mathbf{x}}(\mathbf{x}^\top)\right)A\mathbf{x} + \left(\mathbf{x}^\top A(\nabla_{\mathbf{x}}(\mathbf{x}))\right)^\top$$
$$= \left(\nabla_{\mathbf{x}}(\mathbf{x}^\top)\right)A\mathbf{x} + \left(\nabla_{\mathbf{x}}(\mathbf{x})\right)^\top A^\top \mathbf{x}$$
$$= \left(\nabla_{\mathbf{x}}(\mathbf{x}^\top)\right)A\mathbf{x} + \left(\nabla_{\mathbf{x}}(\mathbf{x}^\top)\right)A^\top \mathbf{x}$$
$$= \nabla_{\mathbf{x}}(\mathbf{x}^\top)\left(A + A^\top\right)\mathbf{x}$$
$$= 1 \cdot \left(A + A^\top\right)\mathbf{x}$$

as desired.                                                                                            ∎

Now, we proceed to the main problem. By the definition of $\ell_2$ norm, we have

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

$$= \mathbf{x}^\top A^\top A \mathbf{x} - \mathbf{x}^\top A^\top \mathbf{b} - \mathbf{b}^\top A \mathbf{x} + \mathbf{x}^\top \Gamma^\top \Gamma \mathbf{x} + \mathbf{b}^\top \mathbf{b}$$

we want to minimize $f(\mathbf{x})$. Hence, we set the gradient to 0. Note that the $\mathbf{b}^\top \mathbf{b}$ goes away.

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \nabla_{\mathbf{x}} \big( \mathbf{x}^\top A^\top A \mathbf{x} - \mathbf{x}^\top A^\top \mathbf{b} - \mathbf{b}^\top A \mathbf{x} + \mathbf{x}^\top \Gamma^\top \Gamma \mathbf{x} \big)$$
$$= \nabla_{\mathbf{x}} \Big( \mathbf{x}^\top \big( A^\top A + \Gamma^\top \Gamma \big) \mathbf{x} - \big( \mathbf{x}^\top A^\top \mathbf{b} + \big( \mathbf{x}^\top A^\top \mathbf{b} \big)^\top \big) \Big)$$

since these will all be $1 \times 1$ matrices, they'll all be symmetric, hence $\mathbf{x}^\top A^\top \mathbf{b} = (\mathbf{x}^\top A^\top \mathbf{b})^\top$. Thus,

$$= \nabla_{\mathbf{x}} \big( \mathbf{x}^\top \big( A^\top A + \Gamma^\top \Gamma \big) \mathbf{x} - 2 \mathbf{x}^\top A^\top \mathbf{b} \big)$$

by the lemma,

$$= \big( A^\top A + \Gamma^\top \Gamma + A A^\top + \Gamma \Gamma^\top \big) \mathbf{x} - 2 A^\top \mathbf{b}$$
$$= 2 ( A^\top A + \Gamma^\top \Gamma) \mathbf{x} - 2 A^\top \mathbf{b}$$
$$= 0$$

hence,

$$2 ( A^\top A + \Gamma^\top \Gamma) \mathbf{x} = 2 A^\top \mathbf{b}$$
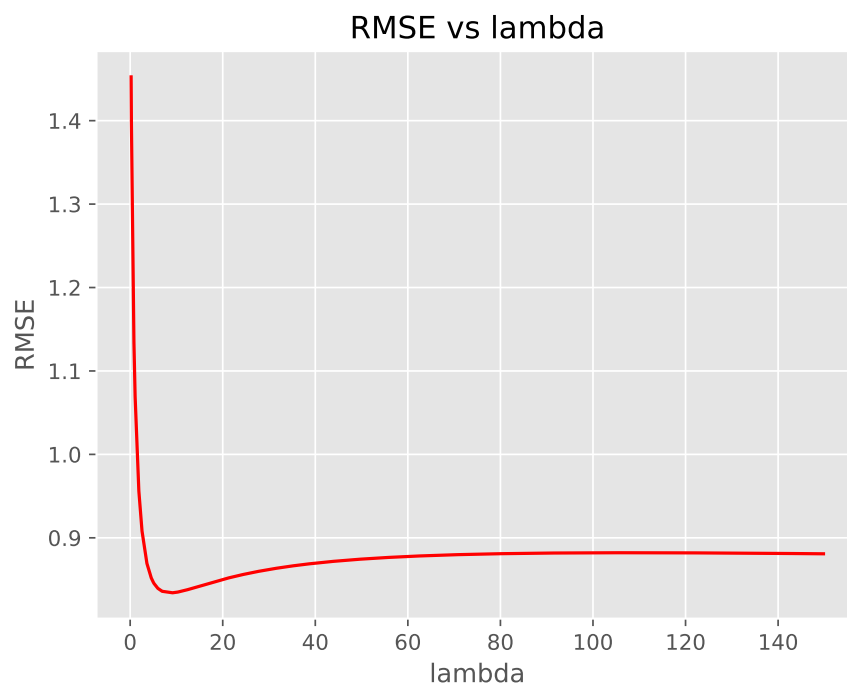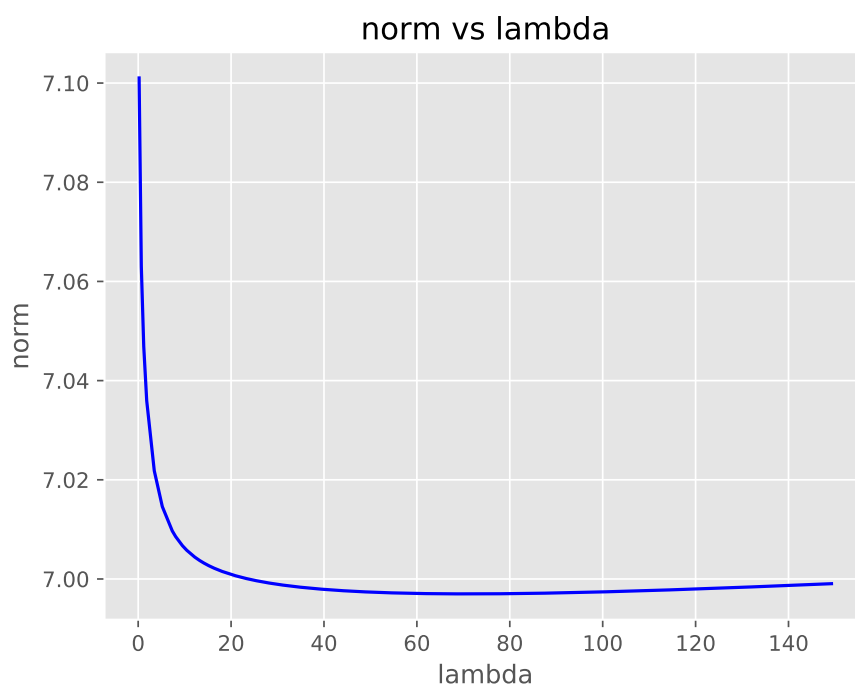$$( A^\top A + \Gamma^\top \Gamma) \mathbf{x} = A^\top \mathbf{b}$$

suppose $( A^\top A + \Gamma^\top \Gamma)$ is invertible. Then

$$\boxed{ \mathbf{x} = \big( A^\top A + \Gamma^\top \Gamma \big)^{-1} A^\top \mathbf{b} }$$

having looked at the coding portion and become confused, I checked the answer key, and saw that we were supposed to state that in the special case where $\Gamma = \sqrt{\lambda} I$, we obtain an expression of the form

$$\boxed{ \mathbf{x} = \big( A^\top A + \lambda I \big)^{-1} A^\top \mathbf{b} }$$

(c) We find $\lambda^\star = 9.0314$, with RMSE on the test set being .8628. See plots on next page.

## RMSE vs lambda



Figure 1: RMSE vs. $\lambda$

## norm vs lambda



Figure 2: $\|\boldsymbol{\theta}^{\star}\|_2^2$ vs. $\lambda$

(d) I'm still a little unclear about why we can't recycle our answer from (b). Anyways, we want to minimize

$$f(\mathbf{x}) = (A\mathbf{x} + b\mathbf{1} - \mathbf{y})^\top (A\mathbf{x} + b\mathbf{1} - \mathbf{y}) + \mathbf{x}^\top \Gamma^\top \Gamma \mathbf{x}$$
$$= \mathbf{x}^\top A^\top A\mathbf{x} + \mathbf{x}^\top A^\top b\mathbf{1} - \mathbf{x}^\top A^\top \mathbf{y} + b\mathbf{1}^\top A\mathbf{x} + b^2 \mathbf{1}^\top \mathbf{1} - b\mathbf{1}^\top \mathbf{y} - \mathbf{y}^\top A\mathbf{x} - \mathbf{y}^\top b\mathbf{1} + \mathbf{y}^\top \mathbf{y} + \mathbf{x}^\top \Gamma^\top \Gamma \mathbf{x}$$

hence

$$\nabla_\mathbf{x} f = 2(A^\top A + \Gamma^\top \Gamma)\mathbf{x} + A^\top (2b\mathbf{1} - 2\mathbf{y})$$
$$= 0$$

and

$$\nabla_b f = 2\mathbf{x}^\top A^\top \mathbf{1} + 2bn - 2\mathbf{1}^\top \mathbf{y}$$
$$= 0$$

thus, rearranging terms, transposing the $\mathbf{x}$ term and dividing by $2n$, we have

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

plugging back into the $\nabla_\mathbf{x}$ equation, we have

$$0 = \left(A^\top A + \Gamma^\top \Gamma\right)\mathbf{x} + A^\top \left(\frac{\mathbf{1}\mathbf{1}^\top (\mathbf{y} - A\mathbf{x})}{n} - \mathbf{y}\right)$$
$$= (A^\top A + \Gamma^\top \Gamma)\mathbf{x} + \frac{1}{n}A^\top \mathbf{1}\mathbf{1}^\top \mathbf{y} - \frac{1}{n}A^\top \mathbf{1}\mathbf{1}^\top A\mathbf{x} - A^\top \mathbf{y}$$
$$= \left(A^\top A + \Gamma^\top \Gamma - \frac{1}{n}A^\top \mathbf{1}\mathbf{1}^\top A\right)\mathbf{x} + \frac{1}{n}A^\top \mathbf{1}\mathbf{1}^\top \mathbf{y} - A^\top \mathbf{y}$$

hence

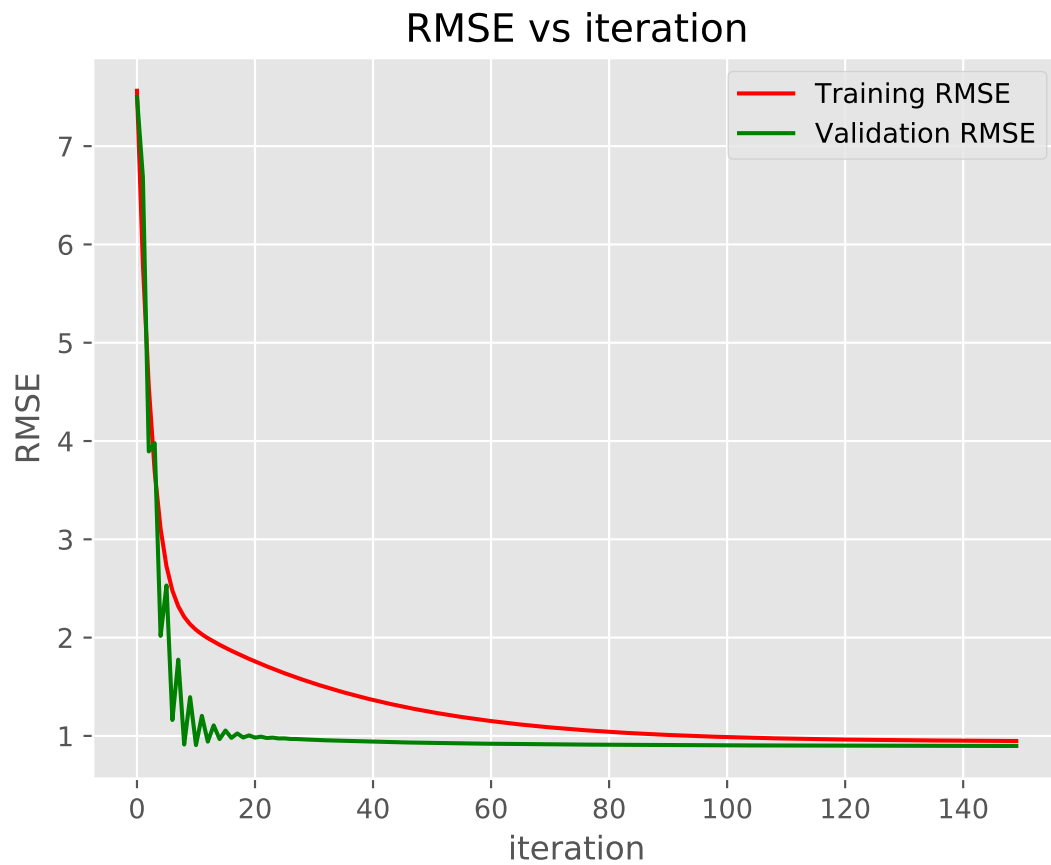$$\left(A^\top A + \Gamma^\top \Gamma - \frac{1}{n}A^\top \mathbf{1}\mathbf{1}^\top A\right)\mathbf{x} = A^\top \mathbf{y} - \frac{1}{n}A^\top \mathbf{1}\mathbf{1}^\top \mathbf{y}$$
$$= \left(A^\top - \frac{1}{n}A^\top \mathbf{1}\mathbf{1}^\top\right)\mathbf{y}$$

and so

$$\mathbf{x}^\star = \left(A^\top A + \Gamma^\top \Gamma - \frac{1}{n}A^\top \mathbf{1}\mathbf{1}^\top A\right)^{-1}\left(A^\top - \frac{1}{n}A^\top \mathbf{1}\mathbf{1}^\top\right)\mathbf{y}$$
$$= \boxed{\left(A^\top \left[I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top\right]A + \Gamma^\top \Gamma\right)^{-1}A^\top \left(I - \frac{1}{n}\mathbf{1}\mathbf{1}^\top\right)\mathbf{y}}$$

this yields a difference in bias of just $3.3939 \cdot 10^{-11}$, with the difference in weights being similar. Success!

(e) See plot

## RMSE vs iteration



Figure 3: $\ell_2$ norm vs. iterations