

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) Starting with the definition of $\sigma'(x)$:

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

as desired.

(b) Let the negative log likelihood function by $f(\theta)$. Note that the function is:

$$f(\theta) = - \sum_i y_i \log \sigma(\theta^T x_i) + (1 - y_i) \log(1 - \sigma(\theta^T x_i))$$

Now, taking the gradient:

$$\begin{aligned} \nabla_{\theta} f(\theta) &= - \sum_i y_i \frac{1}{\sigma(\theta^T x_i)} \sigma'(\theta^T x_i) + (1 - y_i) \frac{1}{(1 - \sigma(\theta^T x_i))} \sigma'(\theta^T x_i) \\ &= - \sum_i (1 - \sigma(\theta^T x_i)) x_i - (1 - y_i) \sigma(\theta^T x_i) x_i \\ &= \sum_i (\sigma(\theta^T x_i) - y_i) x_i \\ &= X^T (\sigma(\theta^T x) - y) \end{aligned}$$

(c) Using the formula of the Hessian we get:

$$\begin{aligned} H_{\theta} &= \nabla_{\theta} (\nabla_{\theta} f(\theta))^T \\ &= \nabla_{\theta} (X^T (\sigma(\theta^T X^T) - y))^T \\ &= \nabla_{\theta} (\sigma(\theta^T X^T)^T X - y^T X) \\ &= \nabla_{\theta} (\sigma(X\theta)^T) X \\ &= X^T S X \end{aligned}$$

In order to be positive semidefinite, the eigenvalues of S must be nonnegative. However S is a diagonal matrix, so its eigenvalues are just the entries in the following form:

$$\sigma(\theta^T x_i)(1 - \sigma(\theta^T x_i))$$

Note that the range of $\sigma(x)$ is strictly between 0 and 1. Thus $1 - \sigma(x) \geq 0$, and therefore $\sigma(x)(1 - \sigma(x)) \geq 0$ and S is positive semidefinite as desired. ■

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 \mathbb{P} to be a valid density function, it's integral must be 1. Therefore we have:

$$\begin{aligned} 1 &= \int \frac{1}{Z} \exp\left(-\frac{x^2}{2\sigma^2}\right) dx \\ &= \frac{1}{Z} \int \exp\left(-\frac{x^2}{2\sigma^2}\right) dx \\ Z &= \int \exp\left(-\frac{x^2}{2\sigma^2}\right) dx \end{aligned}$$

Consider 2d coordinates we have:

$$\begin{aligned} Z^2 &= \int \exp\left(-\frac{x^2}{2\sigma^2}\right) dx \int \exp\left(-\frac{y^2}{2\sigma^2}\right) dy \\ &= \int \int \exp\left(-\frac{x^2 + y^2}{2\sigma^2}\right) dx dy \end{aligned}$$

Using polar coordinates:

$$\begin{aligned} Z^2 &= \int_0^\infty \int_0^{2\pi} \exp\left(-\frac{r^2}{2\sigma^2}\right) r d\theta dr \\ &= 2\pi \int_0^\infty \exp\left(-\frac{r^2}{2\sigma^2}\right) r dr \\ &= 2\pi(-\sigma)^2 \int_0^\infty \exp\left(-\frac{r^2}{2\sigma^2}\right) \left(-\frac{r}{\sigma^2}\right) dr \quad \text{Take out a factor of } (-\sigma)^2 \\ &= -2\pi\sigma^2 \left[\exp\left(-\frac{r^2}{2\sigma^2}\right) \right]_0^\infty \\ &= -2\pi\sigma^2(0 - 1) \\ &= 2\pi\sigma^2 \end{aligned}$$

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:

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

- (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 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 $\|\boldsymbol{\theta}^*\|_2$ where $\boldsymbol{\theta}$ 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 $\hat{\mathbf{y}} = \boldsymbol{\theta}^\top \mathbf{x}$ 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 ℓ_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) 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:

$$\arg \max_{\mathbf{w}} \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(y_i - \omega_0 - \mathbf{w}^T \mathbf{x}_i)^2}{2\sigma^2}\right) + \sum_{j=1}^D \log \frac{1}{\sqrt{2\pi\tau}} \exp\left(-\frac{\omega_j^2}{2\tau^2}\right)$$

Which is equal to:

$$\arg \max_{\mathbf{w}} \sum_{i=1}^N \left(-\frac{(y_i - \omega_0 - \mathbf{w}^T \mathbf{x}_i)^2}{2\sigma^2} - \log \sqrt{2\pi\sigma} \right) + \sum_{j=1}^D \left(-\frac{\omega_j^2}{2\tau^2} - \log \sqrt{2\pi\tau} \right)$$

$$\arg \max_{\mathbf{w}} - \left((N + D) \log \sqrt{2\pi\sigma} + \sum_{i=1}^N \left(\frac{(y_i - \omega_0 - \mathbf{w}^T \mathbf{x}_i)^2}{2\sigma^2} \right) + \sum_{j=1}^D \left(\frac{\omega_j^2}{2\tau^2} \right) \right)$$

The constant $(N + D) \log \sqrt{2\pi\sigma}$ will not change the optimal value, so it can be ignored. Also, finding the maximum of a negative value is the same as finding the minimum of that value, giving us:

$$\arg \min_{\mathbf{w}} \sum_{i=1}^N (y_i - \omega_0 - \mathbf{w}^T \mathbf{x}_i)^2 + \frac{\sigma^2}{\tau^2} \sum_{j=1}^D \omega_j^2$$

With $\lambda = \frac{\sigma^2}{\tau^2}$ we have:

$$\arg \min_{\mathbf{w}} \sum_{i=1}^N (y_i - \omega_0 - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \sum_{j=1}^D \omega_j^2$$

$$\arg \min_{\mathbf{w}} \sum_{i=1}^N (y_i - \omega_0 - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_2^2$$

as desired

- (b) To find the closed form solution \mathbf{x}^* , we find the gradient of f with respect to \mathbf{x} and set it to 0.

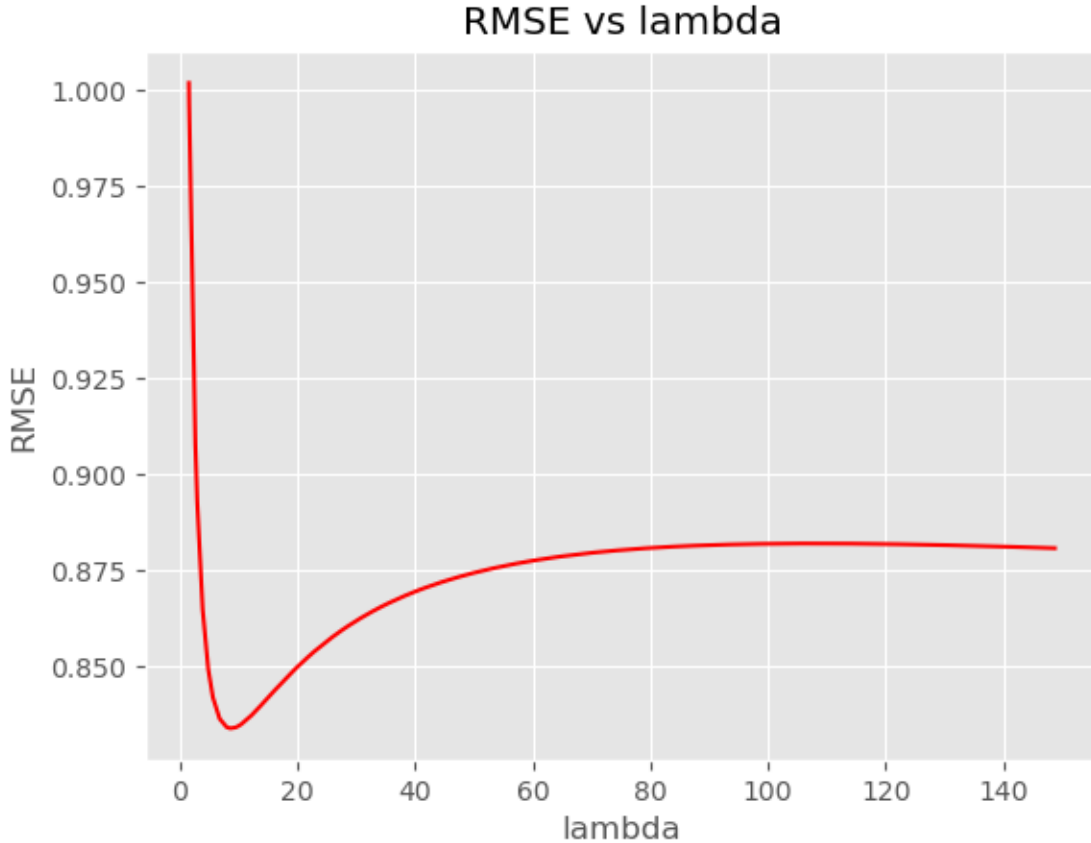
$$\begin{aligned} \nabla_{\mathbf{x}} f &= \nabla_{\mathbf{x}} \left[(\mathbf{A}\mathbf{x} - \mathbf{b})^T (\mathbf{A}\mathbf{x} - \mathbf{b}) + (\mathbf{\Gamma}\mathbf{x})^T (\mathbf{\Gamma}\mathbf{x}) \right] \\ &= \nabla_{\mathbf{x}} \left[(\mathbf{x}^T \mathbf{A}^T - \mathbf{b}^T) (\mathbf{A}\mathbf{x}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

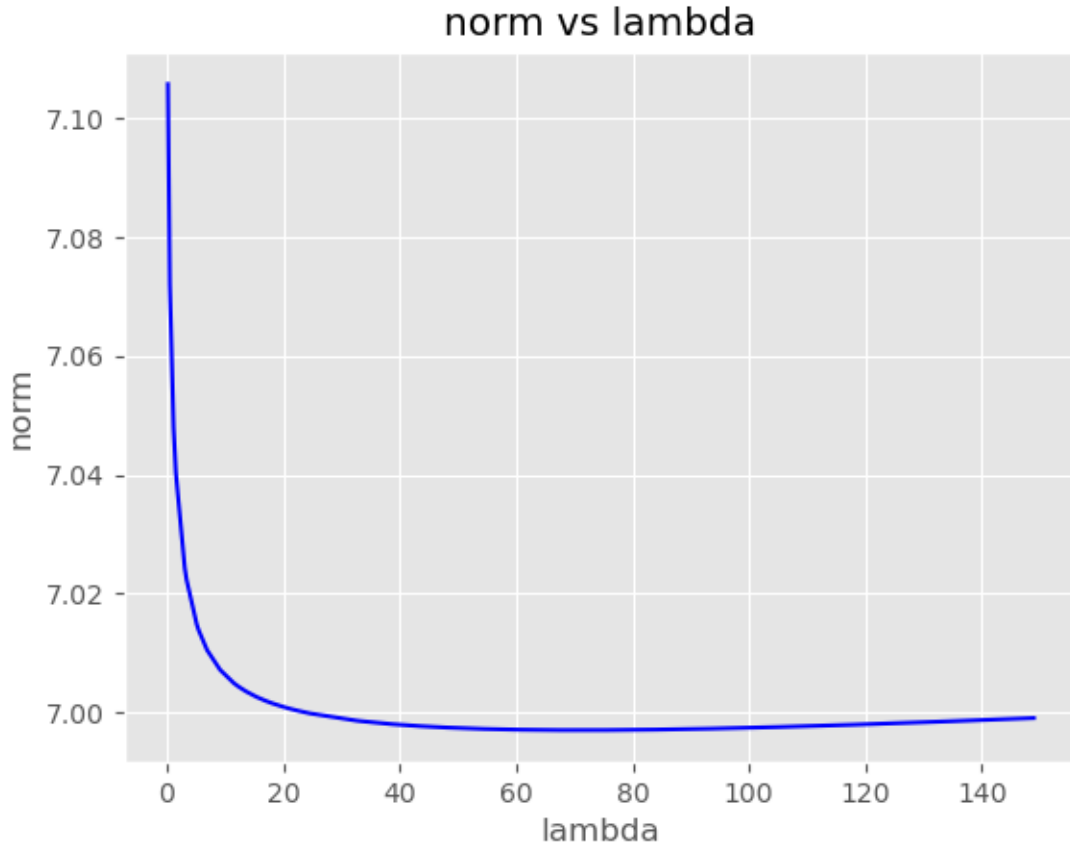
When the gradient is equal to 0 we have:

$$\mathbf{A}^T \mathbf{b} = \mathbf{x}$$

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

- (c) The optimal regulation parameter is 8.4990. Validation RMSE is 0.8340 and test RMSE is 0.8628. Plots below:





(d) First expanding the term we have:

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

Find the gradient of f in terms of b and set it to 0 to solve for b^* :

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

$$2bn = 2\mathbf{1}^T \mathbf{y} - 2\mathbf{1}^T A \mathbf{x}$$

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

Find the gradient of f in terms of \mathbf{x} and set it to 0 to solve for \mathbf{x}^* :

$$\nabla_{\mathbf{x}} f = 0 = 2A^T A \mathbf{x} + 2bA^T \mathbf{1} - 2A^T \mathbf{y} + 2\Gamma^T \Gamma \mathbf{x}$$

Plugging in the value of b^*

$$0 = 2A^T A \mathbf{x} + \frac{2}{n}(\mathbf{1}^T \mathbf{y} - \mathbf{1}^T A \mathbf{x}) A^T \mathbf{1} - 2A^T \mathbf{y} + 2\Gamma^T \Gamma \mathbf{x}$$

$$0 = 2A^T A \mathbf{x} + \frac{2}{n}A^T \mathbf{1} \mathbf{1}^T \mathbf{y} - \frac{2}{n}A^T \mathbf{1} \mathbf{1}^T A \mathbf{x} - 2A^T \mathbf{y} + 2\Gamma^T \Gamma \mathbf{x}$$

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

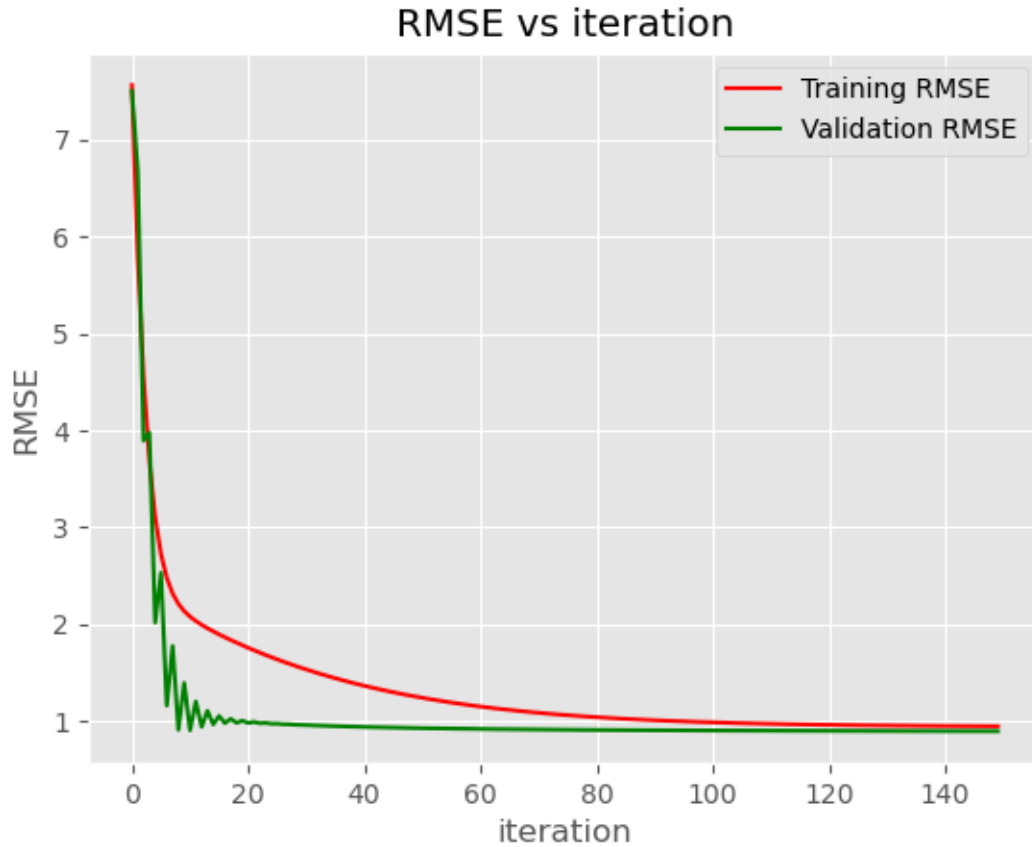
$$(A^T A - \frac{1}{n}A^T \mathbf{1} \mathbf{1}^T A + \Gamma^T \Gamma) \mathbf{x} = A^T \mathbf{y} - \frac{1}{n}A^T \mathbf{1} \mathbf{1}^T \mathbf{y}$$

$$\mathbf{x}^* = (A^T A - \frac{1}{n}A^T \mathbf{1} \mathbf{1}^T A + \Gamma^T \Gamma)^{-1} (A^T \mathbf{y} - \frac{1}{n}A^T \mathbf{1} \mathbf{1}^T \mathbf{y})$$

Difference in bias : $1.9807E - 10$

Difference in weights : $2.5702E - 10$

(e) Plot below:



Difference in bias : $1.5387E - 01$

Difference in weights : $8.0140E - 01$

■