

HWO

Name: Mengying Lin.

SID: 3038737132.

2. Course Policies

Go to the course website and read the course policies carefully. Leave a followup in the Homework 0, Question 2 thread on Ed if you have any questions. Are the following situations violations of course policy? Write "Yes" or "No", and a short explanation for each.

- (a) Alice and Bob work on a problem in a study group. They write up a solution together and submit it, noting on their submissions that they wrote up their homework answers together.
- (b) Carol goes to a homework party and listens to Dan describe his approach to a problem on the board, taking notes in the process. She writes up her homework submission from her notes, crediting Dan.
- (c) Erin gets frustrated by the fact that a homework problem given seems to have nothing in the lecture, notes, or discussion that is parallel to it. So, she starts searching for the problem online. She finds a solution to the homework problem on a website. She reads it and then, after she has understood it, writes her own solution using the same approach. She submits the homework with a citation to the website.
- (d) Frank is having trouble with his homework and asks Grace for help. Grace lets Frank look at her written solution. Frank copies it onto his notebook and uses the copy to write and submit his homework, crediting Grace.
- (e) Heidi has completed her homework. Her friend Irene has been working on a homework problem for hours, and asks Heidi for help. Heidi sends Irene her photos of her solution through an instant messaging service, and Irene uses it to write her own solution with a citation to Heidi.

a> Yes. Each students should write up their solutions individually.

b> No. She has credited external resources.

c> ~~No.~~ Yes. She finally made the information clear to herself and wrote up the homework on her own.

d> Yes. He shouldn't directly copy others' solutions.

e> ~~No.~~ Yes. It is OK to work in group but Heidi shouldn't send her homework to others directly.

Yes. Erin shouldn't look at a soln to hw on a website directly

3. Gradient Descent Doesn't Go Nuts with Ill-Conditioning

Consider a linear regression problem with n training points and d features. When $n = d$, the feature matrix $F \in \mathbb{R}^{n \times n}$ has some maximum singular value α and an extremely tiny minimum singular value. We have noisy observations $\mathbf{y} = F\mathbf{w}^* + \boldsymbol{\epsilon}$. If we compute $\hat{\mathbf{w}}_{inv} = F^{-1}\mathbf{y}$, then due to the tiny singular value of F and the presence of noise we observe that $\|\hat{\mathbf{w}}_{inv} - \mathbf{w}^*\|_2 = 10^{10}$.

Suppose instead of inverting the matrix we decide to use gradient descent instead. We run k iterations of gradient descent to minimize the loss $\ell(w) = \frac{1}{2}\|\mathbf{y} - F\mathbf{w}\|^2$ starting from $\mathbf{w}_0 = \mathbf{0}$. We use a learning rate η which is *small enough* that gradient descent cannot possibly diverge for the given problem. (**This is important. You will need to use this.**)

The gradient-descent update for $t > 0$ is:

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta \left(F^\top (F\mathbf{w}_{t-1} - \mathbf{y}) \right).$$

We are interested in the error $\|\mathbf{w}_k - \mathbf{w}^*\|_2^2$. We want to show that in the worst case, this error can grow at most linearly with iterations k and in particular $\|\mathbf{w}_k - \mathbf{w}^*\|_2 \leq k\eta\alpha\|\mathbf{y}\|_2 + \|\mathbf{w}^*\|_2$.

i.e. The error cannot go “nuts,” at least not very fast.

For the purposes of the homework, you only have to prove the key idea, since the rest follows by applying induction and the triangle inequality.

Show that for $t > 0$, $\|\mathbf{w}_t\|_2 \leq \|\mathbf{w}_{t-1}\|_2 + \eta\alpha\|\mathbf{y}\|_2$.

(HINT: What do you know about $(I - \eta F^\top F)$ if gradient descent cannot diverge? What are its eigenvalues like? Use this fact.)

Proof: $\|\mathbf{w}_t\|_2 = \|\mathbf{w}_{t-1} - \eta(F^\top(F\mathbf{w}_{t-1} - \mathbf{y}))\|_2 = \|\mathbf{w}_{t-1} - \eta F^\top F \mathbf{w}_{t-1} + \eta F^\top \mathbf{y}\|_2$
 $\leq \|(I - \eta F^\top F)\mathbf{w}_{t-1}\|_2 + \eta \|F^\top \mathbf{y}\|_2. \quad \textcircled{1}$

Now that α is the max singular value for F ,

$$\text{we have } \|F^\top \mathbf{y}\|_2 \leq \alpha \|\mathbf{y}\|_2. \quad \textcircled{2}$$

In addition, $\|(I - \eta F^\top F)\mathbf{w}_{t-1}\|_2 \leq \|(I - \eta F^\top F)\|_2 \|\mathbf{w}_{t-1}\|_2$
 $= |\sigma_{\max}(I - \eta F^\top F)| \|\mathbf{w}_{t-1}\|_2. \quad \textcircled{3}$

Assume η is small enough, $\rightarrow \eta$ is chosen so that gradient descent does not diverge,

$$\sigma_{\max}(I - \eta F^\top F) = 1 - \bar{\eta} \sigma_{\max}(F^\top F) > 0$$

$$\text{so } |\lambda_{\max}(I - F^\top F)| < 1.$$

$$\text{Thus } \textcircled{3} \leq \|\mathbf{w}_{t-1}\|_2. \quad \textcircled{4}$$

According to $\textcircled{1}, \textcircled{2}, \textcircled{4}$, we have $\|\mathbf{w}_t\|_2 \leq \|\mathbf{w}_{t-1}\|_2 + \eta\alpha\|\mathbf{y}\|_2$.

4. Regularization from the Augmentation Perspective

Assume \mathbf{w} is a d -dimensional Gaussian random vector $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma)$ and Σ is symmetric positive-definite. Our model for how the $\{y_i\}$ training data is generated is

$$y = \mathbf{w}^\top \mathbf{x} + Z, \quad Z \sim \mathcal{N}(0, 1), \quad (1)$$

where the noise variables Z are independent of \mathbf{w} and iid across training samples. Notice that all the training $\{y_i\}$ and the parameters \mathbf{w} are jointly normal/Gaussian random variables conditioned on the training inputs $\{\mathbf{x}_i\}$. Let us define the standard data matrix and measurement vector:

$$X = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}.$$

In this model, the MAP estimate of \mathbf{w} is given by the Tikhonov regularization counterpart of ridge regression:

$$\hat{\mathbf{w}} = (X^\top X + \Sigma^{-1})^{-1} X^\top \mathbf{y}, \quad (2)$$

In this question, we explore Tikhonov regularization from the data augmentation perspective.

Define the matrix Γ as a $d \times d$ matrix that satisfies $\Gamma^\top \Gamma = \Sigma^{-1}$. Consider the following augmented design matrix (data) \hat{X} and augmented measurement vector $\hat{\mathbf{y}}$:

$$\hat{X} = \begin{bmatrix} X \\ \Gamma \end{bmatrix} \in \mathbb{R}^{(n+d) \times d}, \quad \text{and} \quad \hat{\mathbf{y}} = \begin{bmatrix} \mathbf{y} \\ \mathbf{0}_d \end{bmatrix} \in \mathbb{R}^{n+d},$$

where $\mathbf{0}_d$ is the zero vector in \mathbb{R}^d . Show that the ordinary least squares problem

$$\underset{\mathbf{w}}{\operatorname{argmin}} \|\hat{\mathbf{y}} - \hat{X}\mathbf{w}\|_2^2$$

has the same solution as (2).

(HINT: Feel free to just use the formula you know for the OLS solution. You don't have to rederive that. This problem is not intended to be hard or time consuming.)

$$\begin{aligned}
 \text{Define } f(\mathbf{w}) &= \|\hat{\mathbf{y}} - \hat{X}\mathbf{w}\|_2^2 \\
 &= \sum_{i=1}^{n+d} (\hat{y}_i - \hat{x}_i^\top \mathbf{w})^2 = \sum_{i=1}^n (\hat{y}_i - \hat{x}_i^\top \mathbf{w})^2 + \sum_{i=n+1}^{n+d} (\hat{y}_i - \hat{x}_i^\top \mathbf{w})^2 \\
 &= (\mathbf{x}\mathbf{w} - \mathbf{y})^\top (\mathbf{x}\mathbf{w} - \mathbf{y}) + (\mathbf{I}\mathbf{w})^\top (\mathbf{I}\mathbf{w}). \\
 &= \mathbf{w}^\top \mathbf{x}^\top \mathbf{x}\mathbf{w} - \mathbf{y}^\top \mathbf{x}\mathbf{w} - \mathbf{w}^\top \mathbf{x}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y} + \mathbf{w}^\top \mathbf{I}^\top \mathbf{I}\mathbf{w}. \\
 \frac{\partial f}{\partial \mathbf{w}} &= 2\mathbf{x}^\top \mathbf{x}\mathbf{w} - 2\mathbf{x}^\top \mathbf{y} - 2\mathbf{x}^\top \mathbf{y} + 2\mathbf{I}^\top \mathbf{I}\mathbf{w} \\
 \text{let } \frac{\partial f}{\partial \mathbf{w}} &\approx 0. \quad \mathbf{w} = (\mathbf{x}^\top \mathbf{x} + \mathbf{I}^\top \mathbf{I})^{-1} \mathbf{x}^\top \mathbf{y} \\
 &= (\mathbf{x}^\top \mathbf{x} + \Sigma)^{-1} \mathbf{x}^\top \mathbf{y}.
 \end{aligned}$$

So the solution is same as (2).

5. Vector Calculus Review

Let $\mathbf{x}, \mathbf{c} \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times n}$. For the following parts, before taking any derivatives, identify what the derivative looks like (is it a scalar, vector, or matrix?) and how we calculate each term in the derivative. Then carefully solve for an arbitrary entry of the derivative, then stack/arrange all of them to get the final result. Note that the convention we will use going forward is that vector derivatives of a scalar (with respect to a column vector) are expressed as a row vector, i.e. $\frac{\partial f}{\partial \mathbf{x}} = [\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n}]$ since a row acting on a column gives a scalar. You may have seen alternative conventions before, but the important thing is that you need to understand the types of objects and how they map to the shapes of the multidimensional arrays we use to represent those types.

- (a) Show $\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{c}) = \mathbf{c}^T$
- (b) Show $\frac{\partial}{\partial \mathbf{x}} \|\mathbf{x}\|_2^2 = 2\mathbf{x}^T$
- (c) Show $\frac{\partial}{\partial \mathbf{x}} (\mathbf{A}\mathbf{x}) = \mathbf{A}$
- (d) Show $\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{A}\mathbf{x}) = \mathbf{x}^T (\mathbf{A} + \mathbf{A}^T)$
- (e) Under what condition is the previous derivative equal to $2\mathbf{x}^T \mathbf{A}$?

$$\text{a)} \quad \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{c}) = \left(\frac{\partial(\mathbf{x}^T \mathbf{c})}{\partial x_1}, \dots, \frac{\partial(\mathbf{x}^T \mathbf{c})}{\partial x_n} \right)$$

$$\text{Note } \mathbf{x}^T \mathbf{c} = \sum_{i=1}^n x_i c_i,$$

$$\text{Therefore } \frac{\partial(\mathbf{x}^T \mathbf{c})}{\partial x_i} = c_i. \quad \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{c}) = (c_1, \dots, c_n) \\ = \mathbf{c}^T.$$

$$\text{b)} \quad \frac{\partial \|\mathbf{x}\|_2^2}{\partial \mathbf{x}} = \frac{\partial (\sum x_i^2)}{\partial \mathbf{x}} = (2x_1, 2x_2, \dots, 2x_n) = 2\mathbf{x}^T.$$

$$\text{c)} \quad \frac{\partial}{\partial \mathbf{x}} (\mathbf{A}\mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} \left(\begin{array}{c} \sum a_{1i} x_i \\ \vdots \\ \sum a_{ni} x_i \end{array} \right) = \left(\begin{array}{c} \frac{\partial}{\partial x_1} \sum a_{1i} x_i \\ \vdots \\ \frac{\partial}{\partial x_n} \sum a_{ni} x_i \end{array} \right) = \begin{pmatrix} a_{11}, & \dots, & a_{1n} \\ \vdots & & \vdots \\ a_{n1}, & \dots, & a_{nn} \end{pmatrix}$$

$$\text{d)} \quad \frac{\partial(\mathbf{x}^T \mathbf{A}\mathbf{x})}{\partial \mathbf{x}} = \frac{\partial(\mathbf{x}^T (\mathbf{A}\mathbf{x}))}{\partial \mathbf{x}} = (\mathbf{A}\mathbf{x})^T + \mathbf{x}^T \frac{\partial(\mathbf{A}\mathbf{x})}{\partial \mathbf{x}} \\ = \mathbf{x}^T (\mathbf{A}^T + \mathbf{A}).$$

$$\text{e)} \quad \text{Let } \mathbf{x}^T (\mathbf{A}^T + \mathbf{A}) = 2\mathbf{x}^T \mathbf{A}.$$

$$\mathbf{x}^T \mathbf{A} + \mathbf{x}^T \mathbf{A}^T = 2\mathbf{x}^T \mathbf{A}.$$

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

$$\mathbf{x}^T (\mathbf{A}^T - \mathbf{A}) = \mathbf{0}_n. \quad \text{for all } \mathbf{x} \in \mathbb{R}^n.$$

Therefore when $\mathbf{A}^T = \mathbf{A}$, i.e. \mathbf{A} is symmetric,

the previous derivative equals to $2\mathbf{x}^T \mathbf{A}$.

6. ReLU Elbow Update under SGD

In this question we will explore the behavior of the ReLU nonlinearity with Stochastic Gradient Descent (SGD) updates. The hope is that this problem should help you build a more intuitive understanding for how SGD works and how it iteratively adjusts the learned function.

We want to model a 1D function $y = f(x)$ using a 1-hidden layer network with ReLU activations and no biases in the linear output layer. Mathematically, our network is

$$\hat{f}(x) = \mathbf{W}^{(2)} \Phi(\mathbf{W}^{(1)} x + \mathbf{b})$$

where $x, y \in \mathbb{R}$, $\mathbf{b} \in \mathbb{R}^d$, $\mathbf{W}^{(1)} \in \mathbb{R}^{d \times 1}$, and $\mathbf{W}^{(2)} \in \mathbb{R}^{1 \times d}$. We define our loss function to be the squared error,

$$\ell(x, y, \mathbf{W}^{(1)}, \mathbf{b}, \mathbf{W}^{(2)}) = \frac{1}{2} \|\hat{f}(x) - y\|_2^2.$$

For the purposes of this problem, we define the gradient of a ReLU at 0 to be 0.

(a) Let's start by examining the behavior of a single ReLU with a linear function of x as the input,

$$\phi(x) = \begin{cases} wx + b, & wx + b > 0 \\ 0, & \text{else} \end{cases}.$$

Notice that the slope of $\phi(x)$ is w in the non-zero domain.

We define a loss function $\ell(x, y, \phi) = \frac{1}{2} \|\phi(x) - y\|_2^2$. Find the following:

(i) The location of the 'elbow' e of the function, where it transitions from 0 to something else.

(ii) The derivative of the loss w.r.t. $\phi(x)$, namely $\frac{d\ell}{d\phi}$

(iii) The partial derivative of the loss w.r.t. w , namely $\frac{\partial \ell}{\partial w}$

(iv) The partial derivative of the loss w.r.t. b , namely $\frac{\partial \ell}{\partial b}$

a) i). let $\phi(x) = 0$.

$x = -\frac{b}{w}$. The elbow is $(-\frac{b}{w}, 0)$.

ii).

$$\frac{\partial \ell}{\partial \phi} = \frac{\partial [\frac{1}{2} \|\phi - y\|_2^2]}{\partial \phi} = \frac{1}{2} \frac{\partial (\phi^2 - 2\phi y + y^2)}{\partial \phi} = \phi - y.$$

$$\text{iii)} \quad l = \begin{cases} \frac{1}{2} \|wx + b - y\|_2^2, & wx + b > 0 \\ 0, & \text{else} \end{cases}$$

$$\frac{\partial l}{\partial w} = \frac{\partial l}{\partial \phi} \cdot \frac{\partial \phi}{\partial w} = \begin{cases} (wx + b - y) \cdot x, & wx + b > 0 \\ 0 & \text{else} \end{cases}$$

$$\text{iv)} \quad \frac{\partial l}{\partial b} = \frac{\partial l}{\partial \phi} \cdot \frac{\partial \phi}{\partial b} = \begin{cases} (wx + b - y), & b > -wx \\ 0 & \text{else} \end{cases}$$



- (b) Now suppose we have some training point (x, y) such that $\phi(x) - y = 1$. In other words, the prediction $\phi(x)$ is 1 unit above the target y — we are too high and are trying to pull the function downward.

Describe what happens to the slope and elbow of $\phi(x)$ when we perform gradient descent in the following cases:

- (i) $\phi(x) = 0$.
- (ii) $w > 0, x > 0$, and $\phi(x) > 0$. It is fine to check the behavior of the elbow numerically in this case.
- (iii) $w > 0, x < 0$, and $\phi(x) > 0$.
- (iv) $w < 0, x > 0$, and $\phi(x) > 0$. It is fine to check the behavior of the elbow numerically in this case.

Additionally, draw and label $\phi(x)$, the elbow, and the qualitative changes to the slope and elbow after a gradient update to w and b . You should label the elbow location and a candidate (x, y) pair. Remember that the update for some parameter vector p and loss ℓ under SGD is

$$p' = p - \lambda \nabla_p(\ell), \lambda > 0.$$

b) Gradient descent :

$$\begin{cases} w \leftarrow w - \eta \frac{\partial L}{\partial w} \\ b \leftarrow b - \eta \frac{\partial L}{\partial b}. \end{cases}$$

Denote the updated w and b as w', b' :

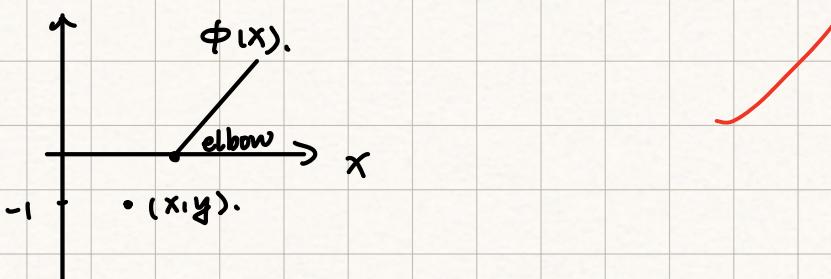
$$\begin{cases} w' = w - \eta (wx + b - y) x = w - \eta x \\ b' = b - \eta (wx + b - y) = b - \eta \end{cases} \quad \textcircled{1}$$

$$\begin{cases} w' = w \\ b' = b \end{cases} \quad \text{so.} \quad \textcircled{2}$$

i. We have $wx + b = 0$.

Note in this case $\frac{\partial L}{\partial b}, \frac{\partial L}{\partial w} = 0$. so w and b remain unchanged.

The slope and the elbow do not change.

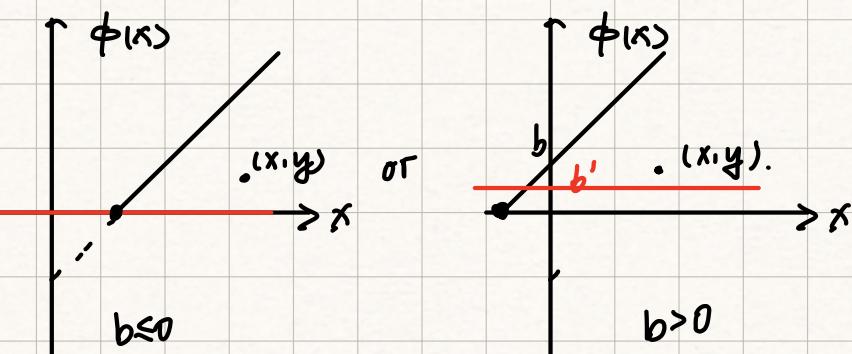


* Red line indicates the changed $\phi(x)$.

ii). As indicated in ①, $w' = w - \eta x$ and $x > 0$, so the slope is dropping.

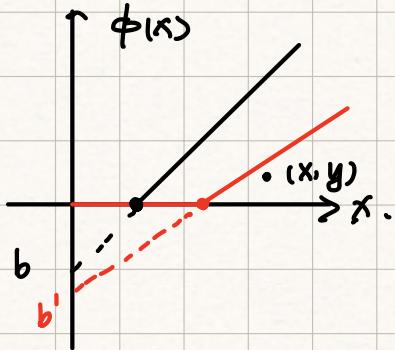
① If $w' = 0$, $\phi = \begin{cases} b' & (b' > 0) \\ 0 & (\text{else}) \end{cases}$. the elbow disappears.

Corresponding graphs:



② if $w' \neq 0$. Let $\phi = w'x + b' = 0$, $x = -\frac{b'}{w'} = -\frac{b - \eta}{w - \eta x}$.

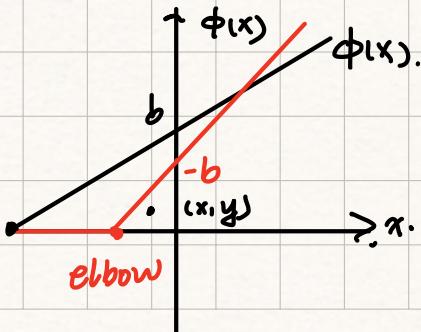
The elbow is moving towards right.



iii). $w' = w - \eta x > w > 0$, so the slope is increasing.

$$b' = b - \eta.$$

Let $\phi = w'x + b' = 0$, $x = -\frac{b - \eta}{w - \eta x}$.



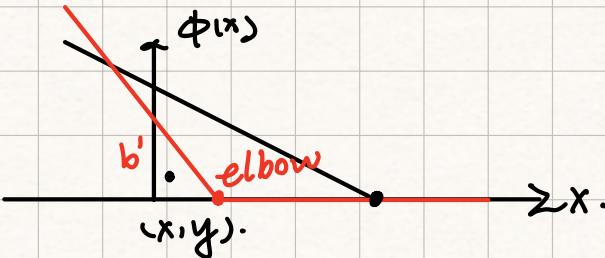
The elbow is moving towards right.

iv). $w' = w - \eta x < w$. so the slope is decreasing numerically.
 $= \text{gets steeper } (w < 0)$.

$$b' = b - \eta$$

$$\text{Let } \phi = w'x + b' = 0. \quad x = -\frac{b - \eta}{w - \eta x}.$$

The elbow is moving towards left.



(c) Now we return to the full network function $\hat{f}(x)$. Derive the location e_i of the elbow of the i 'th elementwise ReLU activation.

(d) Derive the new elbow location e'_i of the i 'th elementwise ReLU activation after one stochastic gradient update with learning rate λ .

c). Denote $W^{(i)} = \begin{pmatrix} w_1^{(i)} \\ \vdots \\ w_d^{(i)} \end{pmatrix}$. $b = \begin{pmatrix} b_1 \\ \vdots \\ b_d \end{pmatrix}$.

$$\text{Let } w_i^{(i)} e_i + b_i = 0.$$

$$\text{If } w_i^{(i)} \neq 0, e_i = -\frac{b_i}{w_i^{(i)}}$$

$$\ell(x, y, W^{(1)}, b, W^{(2)}) = \frac{1}{2} \|\hat{f}(x) - y\|_2^2.$$

$$\hat{f}(x) = W^{(2)} \Phi(W^{(1)}x + b)$$

d). $\frac{\partial \ell}{\partial w_i^{(i)}} = \frac{\partial \ell}{\partial \hat{f}} \cdot \frac{\partial \hat{f}}{\partial \phi} \cdot \frac{\partial \phi}{\partial w_i} \cdot \frac{\partial W^{(i)}}{\partial w_i^{(i)}}$

$$\begin{cases} (\hat{f}(x) - y) W_i^{(2)} x \quad (w_i^{(i)} x + b_i > 0) \\ 0 \quad (\text{else}) \end{cases}$$

$$\frac{\partial \ell}{\partial b} = \frac{\partial \ell}{\partial \hat{f}} \cdot \frac{\partial \hat{f}}{\partial \phi} \cdot \frac{\partial \phi}{\partial b}$$

$$\begin{cases} (\hat{f}(x) - y) W_i^{(2)} \quad (w_i^{(i)} x + b_i > 0) \\ 0 \quad (\text{else}) \end{cases}$$

Hence if $w_i^{(i)} x + b_i > 0$

$$\begin{cases} W_i^{(i)} \leftarrow W_i^{(i)} - \lambda (\hat{f}(x) - y) W_i^{(2)} x \\ b_i \leftarrow b_i - \lambda (\hat{f}(x) - y) W_i^{(2)} \end{cases}$$

$$\text{let } w_i^{(i)} x + b_i = 0.$$

$$x = -\frac{b_i - \lambda (\hat{f}(x) - y) W_i^{(2)}}{W_i^{(i)} - \lambda (\hat{f}(x) - y) W_i^{(2)}} x$$

Else if $w_i^{(i)} x + b_i \leq 0$.

$$\begin{cases} W_i^{(i)} \leftarrow W_i^{(i)} \\ b_i \leftarrow b_i \end{cases}$$

$$x = -\frac{b_i}{W_i^{(i)}}$$

7. Using PyTorch to Learn the Color Organ

a) The value I found: 200

b). The value I found: 200.

c) Yes

R-init (\$)	Iterations	Final R (\$).
0	Fail to train	
100	20321	197
200	0	200
300	27235	200
400	43670	200
1000	77121	200

② It does not always converge to the same value.

④ If lr is too small, (e.g. lr = 20, R-init = 1000), it will fail to converge within 1×10^5 iterations. If lr is too large (e.g. lr = 2×10^7), it diverges.

⑥ Generally, within a moderate interval (e.g. lr $\in [200, 2 \times 10^5]$). the larger the learning rate is, the faster it converges.

I test on lr = 2×10^5 , and it converges in 87 iterations.

a) learned value: R = 339 \$

e) Yes, The loss I find is:

$$\ell(\hat{y}, y) = (1-y) \times \max(\hat{y} - \text{cutoff-mag}, 0) + y \times \max(\text{cutoff-mag}, 0)$$

This loss has 2 term.

① Let us first look at $y \times \max(\text{cutoff-mag}, 0)$

For positive sample ($y=1$), we only need the predicted value to be greater than threshold. If it does, we force the loss to equal to zero.

② $(1-y) \times \max(\hat{y} - \text{cutoff-mag}, 0)$ aims at negative sample,

and the key idea is similar to ①.

+). The learned value is $32\ \Omega$.

② The learned value:

$$\begin{cases} R_{\text{low}} = 40\ \Omega \\ R_{\text{high}} = 160\ \Omega \end{cases}$$

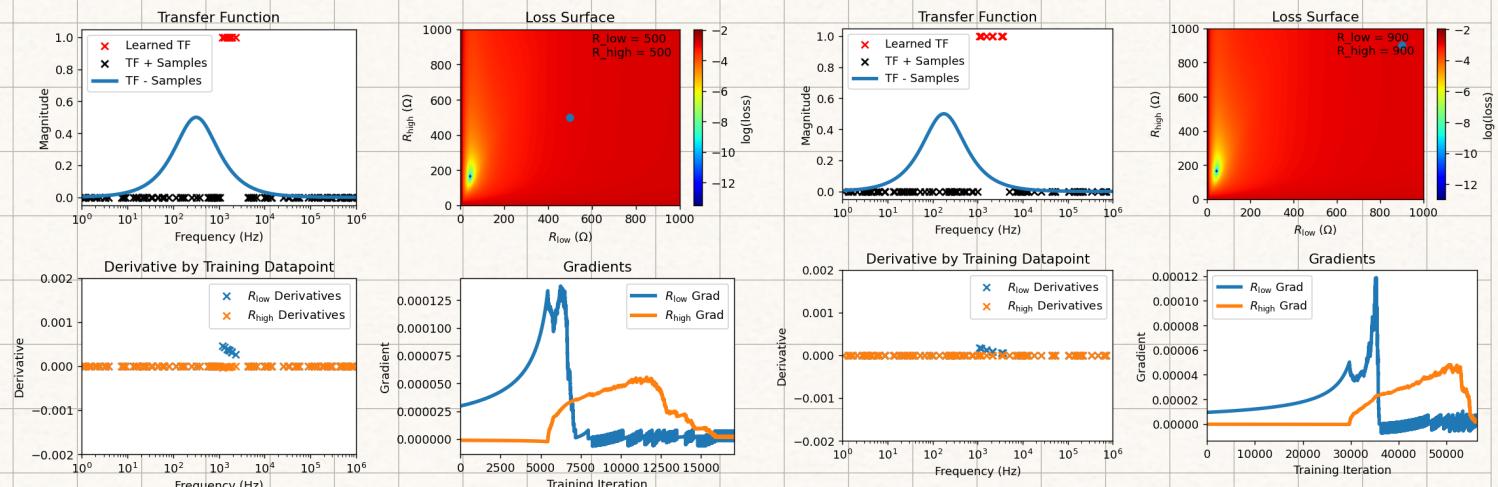
③ If the initial values are too far from the solution, the circuit may not converge.

When the initial values are both set to $900\ \Omega$, the circuit converges

with $R_{\text{low}} = 40\ \Omega$, $R_{\text{high}} = 160\ \Omega$.

④ When initialized to $500\ \Omega$, it only takes 17090 iterations, and 59873 iterations when set to $900\ \Omega$.

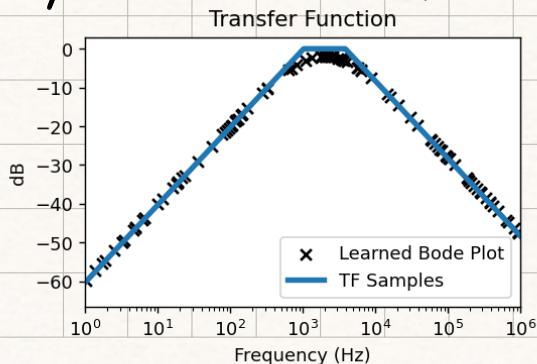
⑤ In this case loss is relatively high and the start point is far from low-loss region as shown in loss surface



Initialized with $500\ \Omega$.

Initialized with $900\ \Omega$.

⑥ ① The bode plot match the expected curve.



② The learned uncutoff low pass frequency is 3864 Hz, and the high pass frequency is 1029 Hz, which are slightly different from the given 4000 Hz and 1000 Hz.

i) It takes similar amount of iterations because the parameters will all be updated in a single iteration, but it evaluated with computing time. it'll definitely takes longer.

8. Homework Process and Study Group

Citing sources and collaborators are an important part of life, including being a student!

We also want to understand what resources you find helpful and how much time homework is taking, so we can change things in the future if possible.

- (a) **What sources (if any) did you use as you worked through the homework?**
- (b) **If you worked with someone on this homework, who did you work with?**
List names and student ID's. (In case of homework party, you can also just describe the group.)
- (c) **Roughly how many total hours did you work on this homework? Write it down here where you'll need to remember it for the self-grade form.**

a)

<https://eecs16b.org/notes/fa21/note6.pdf>

b)

Names	Xiang Fei
SID	3038733024.

c) 15 h

