

**Due Date: March 17th 23:00, 2020**

### Instructions

- For all questions, show your work!
- Starred questions are **hard** questions, not **bonus** questions.
- Please use a document preparation system such as LaTeX, unless noted otherwise.
- Unless noted that questions are related, assume that notation and definitions for each question are self-contained and independent
- All norms denote Euclidean norms unless otherwise specified.
- Submit your answers electronically via Gradescope.
- TAs for this assignment are **Jessica Thompson, Jonathan Cornford and Lluis Castrejon**.

**Question 1** (4-4-4). In this question you will demonstrate that an estimate of the first moment of the gradient using an (exponential) running average is equivalent to using momentum, and is biased by a scaling factor. The goal of this question is for you to consider the relationship between different optimization schemes, and to practice noting and quantifying the effect (particularly in terms of bias/variance) of *estimating* a quantity.

Let  $\mathbf{g}_t$  be an unbiased sample of gradient at time step  $t$  and  $\Delta\boldsymbol{\theta}_t$  be the update to be made. Initialize  $\mathbf{v}_0$  to be a vector of zeros.

1. For  $t \geq 1$ , consider the following update rules:

- SGD with momentum:

$$\mathbf{v}_t = \alpha \mathbf{v}_{t-1} + \epsilon \mathbf{g}_t \quad \Delta\boldsymbol{\theta}_t = -\mathbf{v}_t$$

where  $\epsilon > 0$  and  $\alpha \in (0, 1)$ .

- SGD with running average of  $\mathbf{g}_t$ :

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) \mathbf{g}_t \quad \Delta\boldsymbol{\theta}_t = -\delta \mathbf{v}_t$$

where  $\beta \in (0, 1)$  and  $\delta > 0$ .

Express the two update rules recursively ( $\Delta\boldsymbol{\theta}_t$  as a function of  $\Delta\boldsymbol{\theta}_{t-1}$ ). Show that these two update rules are equivalent; i.e. express  $(\alpha, \epsilon)$  as a function of  $(\beta, \delta)$ .

- Unroll the running average update rule, i.e. express  $\mathbf{v}_t$  as a linear combination of  $\mathbf{g}_i$ 's ( $1 \leq i \leq t$ ).
- Assume  $\mathbf{g}_t$  has a stationary distribution independent of  $t$ . Show that the running average is biased, i.e.  $\mathbb{E}[\mathbf{v}_t] \neq \mathbb{E}[\mathbf{g}_t]$ . Propose a way to eliminate such a bias by rescaling  $\mathbf{v}_t$ .

**Answer 1.**

Q1

Momentum

$$-\Delta\boldsymbol{\theta}_t = \alpha(-\Delta\boldsymbol{\theta}_{t-1}) + \epsilon \mathbf{g}_t$$

$$\Delta\boldsymbol{\theta}_t = \alpha \Delta\boldsymbol{\theta}_{t-1} - \epsilon \mathbf{g}_t$$

Running Average

$$-\frac{1}{\delta}\Delta\boldsymbol{\theta}_t = \beta(-\frac{1}{\delta}\Delta\boldsymbol{\theta}_{t-1}) + (1-\beta)\mathbf{g}_t$$

$$\Delta\boldsymbol{\theta}_t = \beta\Delta\boldsymbol{\theta}_{t-1} - \delta(1-\beta)\mathbf{g}_t$$

Hence we could express them in similar fashion, whereby,

$$\alpha = \beta$$

$$\epsilon = \delta(1-\beta)$$

Q2

$$\mathbf{v}_1 = \beta\mathbf{v}_0 + (1-\beta)\mathbf{g}_1 = (1-\beta)\mathbf{g}_1$$

$$\mathbf{v}_2 = \beta\mathbf{v}_1 + (1-\beta)\mathbf{g}_2 = \beta(1-\beta)\mathbf{g}_1 + (1-\beta)\mathbf{g}_2$$

$$\mathbf{v}_3 = \beta\mathbf{v}_2 + (1-\beta)\mathbf{g}_3 = \beta(\beta(1-\beta)\mathbf{g}_1 + (1-\beta)\mathbf{g}_2) + (1-\beta)\mathbf{g}_3 = \beta^2(1-\beta)\mathbf{g}_1 + \beta(1-\beta)\mathbf{g}_2 + (1-\beta)\mathbf{g}_3$$

We already see a pattern here, so for  $\mathbf{g}_i$ 's ( $1 \leq i \leq t$ ),

$$\mathbf{v}_t = (1-\beta)(\beta^{t-1}\mathbf{g}_1 + \beta^{t-2}\mathbf{g}_2 + \dots + \beta^0\mathbf{g}_t)$$

Q3 We assume  $\mathbf{g}_t$  to be stationary, which means it is constant throughout all values. Recall from last question

$$E[v_t] = (1-\beta)(\beta^{t-1}\mathbf{g}_1 + \beta^{t-2}\mathbf{g}_2 + \dots + \beta^0\mathbf{g}_t)$$

throughout all values. Recall from last question

$$E[v_t] = (1-\beta)(\beta^{t-1}\mathbf{g}_t + \beta^{t-2}\mathbf{g}_t + \dots + \beta^0\mathbf{g}_t)$$

$$E[v_t] = [(1-\beta)(\beta^{t-1} + \beta^{t-2} + \dots + \beta^0)]\mathbf{g}_t$$

And since  $E[\mathbf{g}_t] = \mathbf{g}_t$ , they obviously don't equal. One way to fix that is to scale  $E[v_t]$  by manipulating the constant value in the formula.

**Question 2** (7-5-5-3). The point of this question is to understand and compare the effects of different regularizers (specifically dropout and weight decay) on the weights of a network. Consider a linear regression problem with input data  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , weights  $\mathbf{w} \in \mathbb{R}^{d \times 1}$  and targets  $\mathbf{y} \in \mathbb{R}^{n \times 1}$ . Suppose that dropout is applied to the input (with probability  $1-p$  of dropping the unit i.e. setting it to 0). Let  $\mathbf{R} \in \mathbb{R}^{n \times d}$  be the dropout mask such that  $\mathbf{R}_{ij} \sim \text{Bern}(p)$  is sampled i.i.d. from the Bernoulli distribution.

For a squared error loss function with dropout, we then have:

$$L(\mathbf{w}) = \|\mathbf{y} - (\mathbf{X} \odot \mathbf{R})\mathbf{w}\|^2$$

1. Let  $\Gamma$  be a diagonal matrix with  $\Gamma_{ii} = (\mathbf{X}^\top \mathbf{X})_{ii}^{1/2}$ . Show that the *expectation (over  $\mathbf{R}$ )* of the loss function can be rewritten as  $\mathbb{E}[L(\mathbf{w})] = \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1-p)\|\Gamma\mathbf{w}\|^2$ . *Hint: Note we are trying to find the expectation over a squared term and use  $\text{Var}(Z) = \mathbb{E}[Z^2] - \mathbb{E}[Z]^2$ .*

2. Show that the solution  $\mathbf{w}^{\text{dropout}}$  that minimizes the expected loss from question 2.1 satisfies

$$p\mathbf{w}^{\text{dropout}} = (\mathbf{X}^\top \mathbf{X} + \lambda^{\text{dropout}} \Gamma^2)^{-1} \mathbf{X}^\top \mathbf{y}$$

where  $\lambda^{\text{dropout}}$  is a regularization coefficient depending on  $p$ . How does the value of  $p$  affect the regularization coefficient,  $\lambda^{\text{dropout}}$ ?

3. Express the loss function for a linear regression problem without dropout and with  $L^2$  regularization, with regularization coefficient  $\lambda^{L^2}$ . Derive its closed form solution  $\mathbf{w}^{L^2}$ .
4. Compare the results of 2.2 and 2.3: identify specific differences in the equations you arrive at, and discuss qualitatively what the equations tell you about the similarities and differences in the effects of weight decay and dropout (1-3 sentences).

**Answer 2.**

Q1. Say  $Z = \|\mathbf{y} - (\mathbf{X} \odot \mathbf{R})\mathbf{w}\|$

$$\begin{aligned} \mathbb{E}[Z^2] &= \text{Var}(Z) + \mathbb{E}[Z]^2 \\ &= \text{Var}(\|\mathbf{y} - (\mathbf{X} \odot \mathbf{R})\mathbf{w}\|) + \mathbb{E}[\|\mathbf{y} - (\mathbf{X} \odot \mathbf{R})\mathbf{w}\|]^2 \\ &= \|\text{Var}(\mathbf{y}) - \text{Var}(\mathbf{y} - (\mathbf{X} \odot \mathbf{R})\mathbf{w})\| + \|\mathbb{E}(\mathbf{y}) - \mathbb{E}(\mathbf{y} - (\mathbf{X} \odot \mathbf{R})\mathbf{w})\| \\ &= \|0 - \omega^2 \text{Var}(\mathbf{X} \odot \mathbf{R})\| + \|\mathbf{y} - X\omega \mathbb{E}(\mathbf{R})\| \\ &= \|0 - X^2 \omega^2 \text{Var}(\mathbf{R})\| + \|\mathbf{y} - X\omega \mathbb{E}(\mathbf{R})\| \\ &= \|0 - X^2 \omega^2 p(1-p)\| + \|\mathbf{y} - pX\omega\| \\ &= \|X^T X \omega^2 p(1-p)\| + \|\mathbf{y} - pX\omega\| \\ &= p(1-p)\|\Gamma^2 \omega^2 p\| + \|\mathbf{y} - pX\omega\| \\ &= \|\mathbf{y} - pX\omega\| + p(1-p)\|\Gamma^2 \omega^2 p\| \end{aligned}$$

Q2.

$$\begin{aligned} \frac{\partial}{\partial \omega^{\text{dropout}}} \mathbb{E}[L(\omega)] &= 0 \\ \frac{\partial}{\partial \omega^{\text{dropout}}} (\|\mathbf{y} - pX\omega\|^2 + p(1-p)\|\Gamma\omega\|^2) &= 0 \\ 2p(1-p)\Gamma(\Gamma\omega^{\text{dropout}} - 2pX^T(\mathbf{y} - pX\omega^{\text{dropout}})) &= 0 \\ (1-p)\Gamma(\Gamma\omega^{\text{dropout}} - X^T(\mathbf{y} - pX\omega^{\text{dropout}})) &= 0 \\ (1-p)\Gamma^2\omega^{\text{dropout}} + pX^T X\omega^{\text{dropout}} &= X^T \mathbf{y} \\ (X^T X + \frac{1-p}{p}\Gamma^2)p\omega^{\text{dropout}} &= X^T \mathbf{y} \\ p\omega^{\text{dropout}} &= (X^T X + \frac{1-p}{p}\Gamma^2)^{-1} X^T \mathbf{y} \end{aligned}$$

And  $\lambda^{\text{dropout}} = \frac{1-p}{p}$ , which means as  $p$  increases in range  $(0, 1]$ , regularization coefficient decreases in range  $(-\infty, 0]$

Q3.

$$\begin{aligned}
& (\|y - X\omega\|^2 + \lambda^{L_2}\|\omega\|^2) = 0 \\
& -2X^T(y - X\omega^{L_2}) + 2\lambda^{L_2}\omega^{L_2} = 0 \\
& X^T(y - X\omega^{L_2}) = \lambda^{L_2}\omega^{L_2} \\
& \lambda^{L_2}\omega^{L_2} + X^T X\omega^{L_2} = X^T y \\
& (\lambda^{L_2} + X^T X)\omega^{L_2} = X^T y \\
& \omega^{L_2} = (\lambda^{L_2}I + X^T X)^{-1}X^T y
\end{aligned}$$

Q4. We can see that the two have similar expression, in that they both comprise of a regularization term which helps with invertibility (stability).

However, where weight decay applies a linear penalty, dropout can cause the penalty to grow exponentially and to infinity. This property of dropout can lead to catastrophic failures.

**Question 3** (6-10-2). The goal of this question is for you to understand the reasoning behind different parameter initializations for deep networks, particularly to think about the ways that the initialization affects the activations (and therefore the gradients) of the network. Consider the following equation for the  $t$ -th layer of a deep network:

$$\mathbf{h}^{(t)} = g(\mathbf{a}^{(t)}) \quad \mathbf{a}^{(t)} = \mathbf{W}^{(t)}\mathbf{h}^{(t-1)} + \mathbf{b}^{(t)}$$

where  $\mathbf{a}^{(t)}$  are the pre-activations and  $\mathbf{h}^{(t)}$  are the activations for layer  $t$ ,  $g$  is an activation function,  $\mathbf{W}^{(t)}$  is a  $d^{(t)} \times d^{(t-1)}$  matrix, and  $\mathbf{b}^{(t)}$  is a  $d^{(t)} \times 1$  bias vector. The bias is initialized as a constant vector  $\mathbf{b}^{(t)} = [c, \dots, c]^T$  for some  $c \in \mathbb{R}$ , and the entries of the weight matrix are initialized by sampling i.i.d. from a Gaussian distribution  $W_{ij}^{(t)} \sim \mathcal{N}(\mu, \sigma^2)$ .

Your task is to design an initialization scheme that would achieve a vector of **pre-activations** at layer  $t$  whose elements are zero-mean and unit variance (i.e.:  $\mathbb{E}[a_i^{(t)}] = 0$  and  $\text{Var}(a_i^{(t)}) = 1$ ,  $1 \leq i \leq d^{(t)}$ ) for the assumptions about either the activations or pre-activations of layer  $t-1$  listed below. Note we are not asking for a general formula; you just need to provide one setting that meets these criteria (there are many possibilities).

- First assume that the activations of the previous layer satisfy  $\mathbb{E}[h_i^{(t-1)}] = 0$  and  $\text{Var}(h_i^{(t-1)}) = 1$  for  $1 \leq i \leq d^{(t-1)}$ . Also, assume entries of  $\mathbf{h}^{(t-1)}$  are uncorrelated (the answer should not depend on  $g$ ).
  - Show  $\text{Var}(XY) = \text{Var}(X)\text{Var}(Y) + \text{Var}(X)\mathbb{E}[Y]^2 + \text{Var}(Y)\mathbb{E}[X]^2$  when  $X \perp Y$
  - Write  $\mathbb{E}[a_i^{(t)}]$  and  $\text{Var}(a_i^{(t)})$  in terms of  $c, \mu, \sigma^2, \text{Var}(h_i^{(t-1)}), \mathbb{E}[h_i^{(t-1)}]$ .
  - Give values for  $c, \mu$ , and  $\sigma^2$  as a function of  $d^{(t-1)}$  such that  $\mathbb{E}[a_i^{(t)}] = 0$  and  $\text{Var}(a_i^{(t)}) = 1$  for  $1 \leq i \leq d^{(t)}$ .
- Now assume that the pre-activations of the previous layer satisfy  $\mathbb{E}[a_i^{(t-1)}] = 0$ ,  $\text{Var}(a_i^{(t-1)}) = 1$  and  $a_i^{(t-1)}$  has a symmetric distribution for  $1 \leq i \leq d^{(t-1)}$ . Assume entries of  $\mathbf{a}^{(t-1)}$  are uncorrelated. Consider the case of ReLU activation:  $g(x) = \max\{0, x\}$ .
  - Derive  $\mathbb{E}[(h_i^{(t-1)})^2]$

- (b) Using the result from (a), give values for  $c$ ,  $\mu$ , and  $\sigma^2$  as a function of  $d^{(t-1)}$  such that  $\mathbb{E}[a_i^{(t)}] = 0$  and  $\text{Var}(a_i^{(t)}) = 1$  for  $1 \leq i \leq d^{(t)}$ .
- (c) What popular initialization scheme has this form?
- (d) Why do you think this initialization would work well in practice? Answer in 1-2 sentences.
3. For both assumptions (1,2) give values  $\alpha, \beta$  for  $W_{ij}^{(t)} \sim \text{Uniform}(\alpha, \beta)$  such that  $\mathbb{E}[a_i^{(t)}] = 0$  and  $\text{Var}(a_i^{(t)}) = 1$ .

**Answer 3.**

Q1.

a)

$$\text{Var}(X) = E(X^2) - (E(X))^2$$

$$\text{Var}(Y) = E(Y^2) - (E(Y))^2$$

We have then 1):

$$\text{Var}(X)\text{Var}(Y) = E(X^2)E(Y^2) - (E(X))^2E(Y^2) - (E(Y))^2E(X^2) + (E(X)E(Y))^2$$

2):

$$\text{Var}(X)(E(Y))^2 = E(X^2)(E(Y))^2 - (E(X))^2(E(Y))^2$$

3):

$$\text{Var}(Y)(E(X))^2 = E(Y^2)(E(X))^2 - (E(Y))^2(E(X))^2$$

Adding these three,

$$\begin{aligned} \text{Var}(X)\text{Var}(Y) + \text{Var}(X)E(Y^2) + \text{Var}(Y)E(X^2) &= E(X^2)E(Y^2) - (E(X))^2E(Y^2) - (E(Y))^2E(X^2) + (E(X)E(Y))^2 \\ &= E(X^2)E(Y^2) + (E(X)E(Y))^2 - 2(E(Y))^2(E(X))^2 \\ &= E(X^2)E(Y^2) + (E(X)E(Y))^2 - 2(E(Y)E(X))^2 \\ &= E(X^2)E(Y^2) - (E(X)E(Y))^2 \\ &= E(X^2Y^2) - (E(XY))^2 = E((XY)^2) - (E(XY))^2 = \text{Var}(XY) \end{aligned}$$

b) Expected Value

$$\begin{aligned} E[a_i^{(t)}] &= E\left[\sum_j (W_{ij}^{(t-1)} h_j^{(t-1)} + b_i^{(t)})\right] \\ &= \sum_j E[W_{ij}^{(t-1)} h_j^{(t-1)}] + E[b_i^{(t)}] \\ &= d^{(t-1)} E[W_{ij}^{(t-1)}] E[h_j^{(t-1)}] + c \\ E[a_i^{(t)}] &= c \end{aligned}$$

Variance

$$\text{Var}(a_i^{(t)}) = \text{Var}\left(\sum_j (W_{ij}^{(t-1)} h_j^{(t-1)} + b_i^{(t)})\right) = \sum_j \text{Var}(W_{ij}^{(t-1)} h_j^{(t-1)}) + \text{Var}(b_i^{(t)})$$

$$= d^{(t-1)} \text{Var}(W_{ij}^{(t-1)} h_j^{(t-1)}) + \text{Var}(b_i^{(t)})$$

Use what we proved from last question for  $\text{Var}(W_{ij}^{(t-1)} h_j^{(t-1)})$ ,

$$\text{Var}(a_i^{(t)}) = d^{(t-1)}(\sigma^2 * 1 + \sigma^2 * 0 + 1 * \mu) + 0 = d^{(t-1)}(\sigma^2 + \mu^2)$$

c) We want 'Zero mean'. That means  $c = 0$ . We also want 'unit-variance', that means:

$$\begin{aligned} d^{(t-1)}(\sigma^2 + \mu^2) &= d^{(t-1)}(\sigma^2) = 1 \\ \sigma &= (d^{(t-1)})^{-1} \end{aligned}$$

Q2.

a)

$$E[(h_i^{(t-1)})^2] = \int_{-\infty}^{\infty} \max(0, a_i^{(t-1)})^2 p(a_i^{(t-1)}) da_i^{(t-1)}$$

Now, in -ve region, the integral is 0. We can just evaluate the integral as follows:

$$\begin{aligned} &= \int_0^{\infty} (a_i^{(t-1)})^2 p(a_i^{(t-1)}) da_i^{(t-1)} \\ &= \frac{1}{2} \int_{-\infty}^{\infty} (a_i^{(t-1)})^2 p(a_i^{(t-1)}) da_i^{(t-1)} \\ &= \frac{1}{2} \int_{-\infty}^{\infty} (a_i^{(t-1)} - E[a_i^{(t-1)}])^2 p(a_i^{(t-1)}) da_i^{(t-1)} \\ &= \frac{1}{2} E[(a_i^{(t-1)} - E[a_i^{(t-1)}])^2] \\ E[(h_i^{(t-1)})^2] &= \frac{1}{2} \text{Var}(a_i^{(t-1)}) = \frac{1}{2} \end{aligned}$$

b)

$$E[(h_i^{(t-1)})] = \int_{-\infty}^{\infty} \max(0, a_i^{(t-1)}) p(a_i^{(t-1)}) da_i^{(t-1)}$$

Again, similar to last question,

$$\begin{aligned} &= \frac{1}{2} \int_{-\infty}^{\infty} (a_i^{(t-1)}) p(a_i^{(t-1)}) da_i^{(t-1)} \\ E[(h_i^{(t-1)})] &= \frac{1}{2} E[(a_i^{(t-1)})] = 0 \end{aligned}$$

Use the variance-expectation formula, we get,

$$\text{Var}[(h_i^{(t-1)})] = \frac{1}{2} - 0 = \frac{1}{2}$$

Now if we re-visit Part 1b) where we got the expected value,

$$E[a_i^{(t)}] = d^{(t-1)} E[W_{ij}^{(t-1)}] E[h_j^{(t-1)}] + c$$

$$c = 0$$

Same for the variance,

$$\text{Var}(a_i^{(t)}) = \frac{1}{2}d^{(t-1)}(\sigma^2) = 1$$

$$\sigma^2 = \frac{2}{d^{(t-1)}}$$

c) He Normal. Note the activation is ReLU

d) It defines everything very simply. We find a good variance for the distribution from which the initial parameters are drawn. This variance is adapted to the activation function used and is derived without explicitly considering the type of the distribution.

Q3. Case 1 from part a)

$$E[a_i^{(t)}] = d^{(t-1)} E[W_{ij}^{(t-1)}] E[h_j^{(t-1)}] + c$$

$$0 = d^{(t-1)} * \frac{\alpha + \beta}{2} * 0 + c$$

$$c = 0$$

From Variance formula,

$$1 = d^{(t-1)} \left( \frac{(\beta - \alpha)^2}{12} + \left( \frac{\alpha + \beta}{2} \right)^2 \right)$$

Solve,

$$\alpha^2 + \beta^2 + \alpha\beta = \frac{3}{d^{(t-1)}}$$

Now if  $\beta = 0$ ,

$$\alpha = \pm \sqrt{\left( \frac{3}{d^{(t-1)}} \right)}$$

Case 2 from part b)

$$E[a_i^{(t)}] = d^{(t-1)} E[W_{ij}^{(t-1)}] E[h_j^{(t-1)}] + c$$

$$0 = d^{(t-1)} * \frac{\alpha + \beta}{2} * 0 + c$$

$$c = 0$$

From Variance formula,

$$1 = d^{(t-1)} \left( \frac{(\beta - \alpha)^2}{12} * \frac{1}{2} + \frac{1}{2} * \left( \frac{\alpha + \beta}{2} \right)^2 \right)$$

Again, we solve it like before and consider  $\beta = 0$ ,

$$\alpha = \pm \sqrt{\left( \frac{6}{d^{(t-1)}} \right)}$$

**Question 4 (4-6-6).** This question is about normalization techniques.

1. Batch normalization, layer normalization and instance normalization all involve calculating the mean  $\boldsymbol{\mu}$  and variance  $\boldsymbol{\sigma}^2$  with respect to different subsets of the tensor dimensions. Given the following 3D tensor, calculate the corresponding mean and variance tensors for each normalization technique:  $\boldsymbol{\mu}_{batch}$ ,  $\boldsymbol{\mu}_{layer}$ ,  $\boldsymbol{\mu}_{instance}$ ,  $\boldsymbol{\sigma}_{batch}^2$ ,  $\boldsymbol{\sigma}_{layer}^2$ , and  $\boldsymbol{\sigma}_{instance}^2$ .

$$\begin{bmatrix} \begin{bmatrix} 1, 3, 2 \\ 1, 2, 3 \end{bmatrix}, \begin{bmatrix} 3, 3, 2 \\ 2, 4, 4 \end{bmatrix}, \begin{bmatrix} 4, 2, 2 \\ 1, 2, 4 \end{bmatrix}, \begin{bmatrix} 3, 3, 2 \\ 3, 3, 2 \end{bmatrix} \end{bmatrix}$$

The size of this tensor is 4 x 2 x 3 which corresponds to the batch size, number of channels, and number of features respectively.

2. For the next two subquestions, we consider the following parameterization of a weight vector  $\boldsymbol{w}$ :

$$\boldsymbol{w} := \gamma \frac{\boldsymbol{u}}{\|\boldsymbol{u}\|}$$

where  $\gamma$  is scalar parameter controlling the magnitude and  $\boldsymbol{u}$  is a vector controlling the direction of  $\boldsymbol{w}$ .

Consider one layer of a neural network, and omit the bias parameter. To carry out batch normalization, one normally standardizes the preactivation and performs elementwise scale and shift  $\hat{y} = \gamma \cdot \frac{y - \mu_y}{\sigma_y} + \beta$  where  $y = \boldsymbol{u}^\top \boldsymbol{x}$ . Assume the data  $\boldsymbol{x}$  (a random vector) is whitened ( $\text{Var}(\boldsymbol{x}) = \boldsymbol{I}$ ) and centered at 0 ( $\mathbb{E}[\boldsymbol{x}] = \mathbf{0}$ ). Show that  $\hat{y} = \boldsymbol{w}^\top \boldsymbol{x} + \beta$ .

3. Show that the gradient of a loss function  $L(\boldsymbol{u}, \gamma, \beta)$  with respect to  $\boldsymbol{u}$  can be written in the form  $\nabla_{\boldsymbol{u}} L = s \boldsymbol{W}^\perp \nabla_{\boldsymbol{w}} L$  for some  $s$ , where  $\boldsymbol{W}^\perp = \left( \boldsymbol{I} - \frac{\boldsymbol{u} \boldsymbol{u}^\top}{\|\boldsymbol{u}\|^2} \right)$ . Note that <sup>1</sup>  $\boldsymbol{W}^\perp \boldsymbol{u} = \mathbf{0}$ .

**Answer 4.**

Q1.

$$\boldsymbol{\mu}_{batch} = \begin{bmatrix} 2.5 \\ 2.583 \end{bmatrix}$$

$$\boldsymbol{\sigma}_{batch} = \begin{bmatrix} 0.5830.1875, 0 \\ 0.6875, 0.6875, 0.6875 \end{bmatrix}$$

$$\boldsymbol{\mu}_{layer} = \begin{bmatrix} 2 \\ 3 \\ 2.5 \\ 2.6667 \end{bmatrix}$$

$$\boldsymbol{\sigma}_{layer} = \begin{bmatrix} 0.6667 \\ 0.6667 \\ 1.25 \\ 0.2222 \end{bmatrix}$$

1. As a side note:  $\boldsymbol{W}^\perp$  is an orthogonal complement that projects the gradient away from the direction of  $\boldsymbol{w}$ , which is usually (empirically) close to a dominant eigenvector of the covariance of the gradient. This helps to condition the landscape of the objective that we want to optimize.



$$\boldsymbol{\mu}_{instance} = \begin{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \begin{bmatrix} 2.6667 \\ 3.3333 \end{bmatrix}, \begin{bmatrix} 2.6667 \\ 2.3333 \end{bmatrix}, \begin{bmatrix} 2.6667 \\ 2.6667 \end{bmatrix} \end{bmatrix}$$

$$\boldsymbol{\sigma}_{instance} = \begin{bmatrix} \begin{bmatrix} 0.6667 \\ 0.6667 \end{bmatrix}, \begin{bmatrix} 0.2222 \\ 0.8889 \end{bmatrix}, \begin{bmatrix} 0.8889 \\ 1.5556 \end{bmatrix}, \begin{bmatrix} 0.2222 \\ 0.2222 \end{bmatrix} \end{bmatrix}$$

Q2.

$$Var(y) = Var(\mathbf{u}\mathbf{x}) = \mathbf{u}^T \mathbf{u} Var(\mathbf{x}) = \mathbf{u}^T \mathbf{u} = \|\mathbf{u}\|^2 = \boldsymbol{\sigma}^2$$

$$\boldsymbol{\sigma} = \|\mathbf{u}\|$$

$$\mu_y = \mathbb{E}[\mathbf{y}] = \mathbb{E}[\mathbf{u}^T \mathbf{x}] = \mathbb{E}[\mathbf{u}^T] \mathbb{E}[\mathbf{x}] = \mathbf{0}$$

$$\hat{y} = \gamma \frac{\boldsymbol{\mu}^T \mathbf{x}}{\|\mathbf{u}\|} + \beta$$

Given  $\mathbf{w} := \gamma \frac{\mathbf{u}}{\|\mathbf{u}\|}$ ,

$$\mathbf{w}^T := \gamma \frac{\mathbf{u}^T}{\|\mathbf{u}\|}$$

$$\hat{y} = \mathbf{w}^T \mathbf{x} + \beta$$

Q3. From chain rule,

$$\frac{\partial L}{\partial u} = \frac{\partial L}{\partial \omega} \frac{\partial \omega}{\partial u}$$

All we need to find is the  $\frac{\partial \omega}{\partial u}$ ,

$$\begin{aligned} \frac{\partial \omega}{\partial u} &= \\ \frac{\partial}{\partial u} \left[ \gamma \frac{U}{\|U\|} \right] &= \gamma \frac{\partial}{\partial u} \left[ \frac{U}{\|U\|} \right] \\ &= \gamma * \frac{\|U\| * I - U * \frac{U^T}{\|U\|}}{\|U\|^2} \\ &= (\gamma) * \frac{\|U\| * I - U * \frac{U^T}{\|U\|}}{\|U\|^2} \\ &= (\gamma) * \frac{I - U^T U}{\|U\|^2} \end{aligned}$$

This is the format we are expecting to see. Hence proven.

**Question 5** (4-6-4). This question is about activation functions and vanishing/exploding gradients in recurrent neural networks (RNNs). Let  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  be an activation function. When the argument is a vector, we apply  $\sigma$  element-wise. Consider the following recurrent unit:

$$\mathbf{h}_t = \mathbf{W} \sigma(\mathbf{h}_{t-1}) + \mathbf{U} \mathbf{x}_t + \mathbf{b}$$

1. Show that applying the activation function in this way is equivalent to the conventional way of applying the activation function:  $\mathbf{g}_t = \sigma(\mathbf{W}\mathbf{g}_{t-1} + \mathbf{U}\mathbf{x}_t + \mathbf{b})$  (i.e. express  $\mathbf{g}_t$  in terms of  $\mathbf{h}_t$ ). More formally, you need to prove it using mathematical induction. You only need to prove the induction step in this question, assuming your expression holds for time step  $t - 1$ .
- \*2. Let  $\|\mathbf{A}\|$  denote the  $L_2$  operator norm<sup>2</sup> of matrix  $\mathbf{A}$  ( $\|\mathbf{A}\| := \max_{\mathbf{x}: \|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|$ ). Assume  $\sigma(x)$  has bounded derivative, i.e.  $|\sigma'| \leq \gamma$  for some  $\gamma > 0$  and for all  $x$ . We denote as  $\lambda_1(\cdot)$  the largest eigenvalue of a symmetric matrix. Show that if the largest eigenvalue of the weights is bounded by  $\frac{\delta^2}{\gamma^2}$  for some  $0 \leq \delta < 1$ , gradients of the hidden state will vanish over time, i.e.

$$\lambda_1(\mathbf{W}^\top \mathbf{W}) \leq \frac{\delta^2}{\gamma^2} \implies \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| \rightarrow 0 \text{ as } T \rightarrow \infty$$

Use the following properties of the  $L_2$  operator norm

$$\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\| \quad \text{and} \quad \|\mathbf{A}\| = \sqrt{\lambda_1(\mathbf{A}^\top \mathbf{A})}$$

3. What do you think will happen to the gradients of the hidden state if the condition in the previous question is reversed, i.e. if the largest eigenvalue of the weights is larger than  $\frac{\delta^2}{\gamma^2}$ ? Is this condition *necessary* or *sufficient* for the gradient to explode? (Answer in 1-2 sentences).

**Answer 5.**

Q1. Base Case:  $t = 1$

$$\begin{aligned} h_1 &= W\sigma(h_0) + Ux_0 + b \\ h_2 &= W\sigma(h_1) + Ux_t + b = W\sigma(W\sigma(h_0) + Ux_t + b) + Ux_t + b \\ g_2 &= \sigma(W\sigma(g_0) + Ux_t + b) \end{aligned}$$

Induction step (Assume we know it holds for  $k = t - 1$ , we need to look at  $k = t$ ). The new activation function is:

$$h_{t+1} = W\sigma(h_t) + Ux_0 + b$$

And that,

$$h_{t+2} = W\sigma(h_{t+1}) + Ux_t + b = W\sigma(W\sigma(h_t) + Ux_t + b) + Ux_t + b$$

Similarly, we also have at  $t$  for conventional activation function:

$$g_{t+1} = \sigma(W\sigma(g_t) + Ux_t + b)$$

From above, we can see that in  $h_{t+1}$ , the activation part has exactly the same formulation in terms of  $h_t$  and in terms of  $g_t$ . Hence proved.

Q2. We can see that

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial W\sigma(h_{t-1})}{\partial h_{t-1}} = W \frac{\partial \sigma(h_{t-1})}{\partial h_{t-1}} = W\sigma'$$

Again from chain rule

$$\frac{\partial h_T}{\partial h_0} = \frac{\partial h_T}{\partial h_{T-1}} \frac{\partial h_{T-1}}{\partial h_{T-2}} \cdots \frac{\partial h_1}{\partial h_0}$$

---

2. The  $L_2$  operator norm of a matrix  $\mathbf{A}$  is an *induced norm* corresponding to the  $L_2$  norm of vectors. You can try to prove the given properties as an exercise.

Also,

$$\lambda_1(W^T W) = \|W\|^2 \leq \frac{\delta^2}{\gamma^2}$$

$$\|W\| \leq \frac{\delta}{\gamma}$$

Then from given property I,

$$\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| = \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \right\| * \left\| \frac{\partial \mathbf{h}_{T-1}}{\partial \mathbf{h}_{T-2}} \right\| \dots = (\|W\| * \|\sigma'\|) * (\|W\| * \|\sigma'\|) * \dots$$

Also given the bound for  $\sigma'$

$$= (\|W\| * \|\sigma'\|)^T \leq \left(\frac{\delta}{\gamma}\gamma\right)^T = (\delta)^T$$

$$\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| \leq (\delta)^T$$

Now, since  $\delta < 1$ , if  $T$  approaches  $\infty$ ,  $\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\|$  approaches 0 thanks to properties of exponents of a number between 0 and 1.

Q3. The gradient is no longer upper-bounded by a constant, which exposes it to the risk of exploding. However, it is only a necessary condition as it means that the norm will not be bounded by a constant. It will become a sufficient condition whereby  $\gamma$  is infinitely close to 0 making the expression's limit to be evaluated to infinity.

**Question 6** (4-8-8). Consider the following Bidirectional RNN:

$$\begin{aligned} \mathbf{h}_t^{(f)} &= \sigma(\mathbf{W}^{(f)} \mathbf{x}_t + \mathbf{U}^{(f)} \mathbf{h}_{t-1}^{(f)}) \\ \mathbf{h}_t^{(b)} &= \sigma(\mathbf{W}^{(b)} \mathbf{x}_t + \mathbf{U}^{(b)} \mathbf{h}_{t+1}^{(b)}) \\ \mathbf{y}_t &= \mathbf{V}^{(f)} \mathbf{h}_t^{(f)} + \mathbf{V}^{(b)} \mathbf{h}_t^{(b)} \end{aligned}$$

where the superscripts  $f$  and  $b$  correspond to the forward and backward RNNs respectively and  $\sigma$  denotes the logistic sigmoid function. Let  $\mathbf{z}_t$  be the true target of the prediction  $\mathbf{y}_t$  and consider the sum of squared loss  $L = \sum_t L_t$  where  $L_t = \|\mathbf{z}_t - \mathbf{y}_t\|_2^2$ .

In this question our goal is to obtain an expression for the gradients  $\nabla_{\mathbf{W}^{(f)}} L$  and  $\nabla_{\mathbf{U}^{(b)}} L$ .

1. First, complete the following computational graph for this RNN, unrolled for 3 time steps (from  $t = 1$  to  $t = 3$ ). Label each node with the corresponding hidden unit and each edge with the corresponding weight. Note that it includes the initial hidden states for both the forward and backward RNNs.
2. Using total derivatives we can express the gradients  $\nabla_{\mathbf{h}_t^{(f)}} L$  and  $\nabla_{\mathbf{h}_t^{(b)}} L$  recursively in terms of  $\nabla_{\mathbf{h}_{t+1}^{(f)}} L$  and  $\nabla_{\mathbf{h}_{t-1}^{(b)}} L$  as follows:

$$\nabla_{\mathbf{h}_t^{(f)}} L = \nabla_{\mathbf{h}_t^{(f)}} L_t + \left( \frac{\partial \mathbf{h}_{t+1}^{(f)}}{\partial \mathbf{h}_t^{(f)}} \right)^\top \nabla_{\mathbf{h}_{t+1}^{(f)}} L$$

$$\nabla_{\mathbf{h}_t^{(b)}} L = \nabla_{\mathbf{h}_t^{(b)}} L_t + \left( \frac{\partial \mathbf{h}_{t-1}^{(b)}}{\partial \mathbf{h}_t^{(b)}} \right)^\top \nabla_{\mathbf{h}_{t-1}^{(b)}} L$$

Derive an expression for  $\nabla_{\mathbf{h}_t^{(f)}} L_t$ ,  $\nabla_{\mathbf{h}_t^{(b)}} L_t$ ,  $\frac{\partial \mathbf{h}_{t+1}^{(f)}}{\partial \mathbf{h}_t^{(f)}}$  and  $\frac{\partial \mathbf{h}_{t-1}^{(b)}}{\partial \mathbf{h}_t^{(b)}}$ .

3. Now derive  $\nabla_{\mathbf{W}^{(f)}} L$  and  $\nabla_{\mathbf{U}^{(b)}} L$  as functions of  $\nabla_{\mathbf{h}_t^{(f)}} L$  and  $\nabla_{\mathbf{h}_t^{(b)}} L$ , respectively.

*Hint: It might be useful to consider the contribution of the weight matrices when computing the recurrent hidden unit at a particular time  $t$  and how those contributions might be aggregated.*

**Answer 6.**

Q1 Please see the Figure below

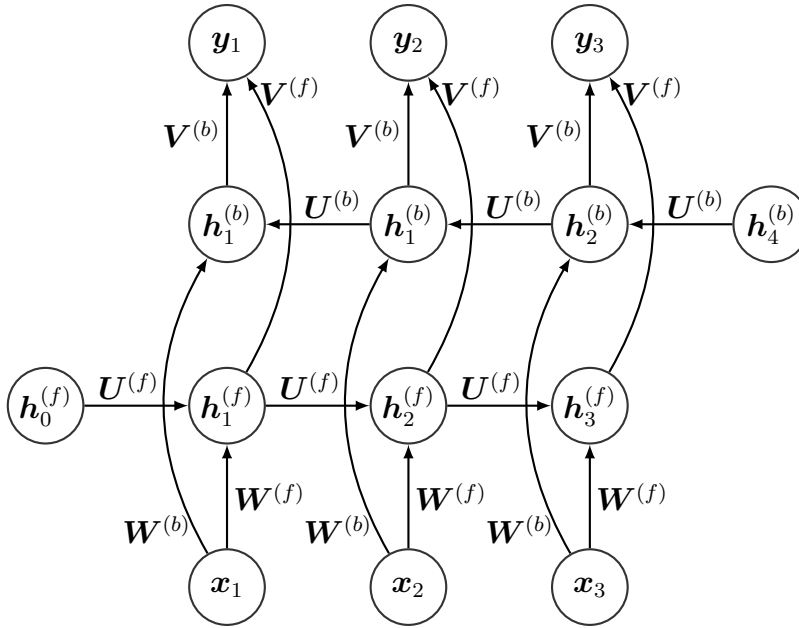


FIGURE 1 – Computational graph of the bidirectional RNN unrolled for three timesteps.

Q2 We can use chain rule:

$$\frac{\partial L_t}{\partial h_t^f} = \frac{\partial L_t}{\partial y_t} \frac{\partial y_t}{\partial h_t^f} = -2(Z_t - y_t)V^f$$

$$\frac{\partial L_t}{\partial h_t^b} = \frac{\partial L_t}{\partial y_t} \frac{\partial y_t}{\partial h_t^b} = -2(Z_t - y_t)V^b$$

$$\frac{\partial h_t^{f+1}}{\partial h_t^f} = \frac{\partial h_t^{f+1}}{\partial \sigma} \frac{\partial \sigma}{\partial h_t^f} = \sigma(\mathbf{W}^{(f)}\mathbf{x}_t + \mathbf{U}^{(f)}\mathbf{h}_{t-1}^{(f)})(1 - \sigma(\mathbf{W}^{(f)}\mathbf{x}_t + \mathbf{U}^{(f)}\mathbf{h}_{t-1}^{(f)}))\mathbf{U}^{(f)}$$

$$\frac{\partial h_t^{b-1}}{\partial h_t^b} = \frac{\partial h_t^{b-1}}{\partial \sigma} \frac{\partial \sigma}{\partial h_t^b} = \sigma(\mathbf{W}^{(b)}\mathbf{x}_t + \mathbf{U}^{(b)}\mathbf{h}_{t+1}^{(b)})(1 - \sigma(\mathbf{W}^{(b)}\mathbf{x}_t + \mathbf{U}^{(b)}\mathbf{h}_{t+1}^{(b)}))\mathbf{U}^{(b)}$$

Q3 We again use chain rule to find their relations.