

**Due Date: March 10th 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
- Submit your answers electronically via Gradescope.
- TAs for this assignment are **Jessica Thompson, Jonathan Cornford and Lluís 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.**

- For SGD with momentum, we have;

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

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

For SGD with running average,

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

Given

$$\alpha = \beta, \quad \epsilon = \delta(1 - \beta)$$

the two updates are equivalents.

2. By unrolling the average update rule, we have ;

$$\begin{aligned}
 \mathbf{v}_t &= \beta \mathbf{v}_{t-1} + (1 - \beta) \mathbf{g}_t \\
 &= \beta(\beta \mathbf{v}_{t-2} + (1 - \beta) \mathbf{g}_{t-1}) + (1 - \beta) \mathbf{g}_t \\
 &= \beta^2 \mathbf{v}_{t-2} + \beta(1 - \beta) \mathbf{g}_{t-1} + (1 - \beta) \mathbf{g}_t \\
 &= \beta^2(\beta \mathbf{v}_{t-3} + (1 - \beta) \mathbf{g}_{t-2}) + \beta(1 - \beta) \mathbf{g}_{t-1} + (1 - \beta) \mathbf{g}_t \\
 &= \beta^3 \mathbf{v}_{t-3} + \beta^2(1 - \beta) \mathbf{g}_{t-2} + \beta(1 - \beta) \mathbf{g}_{t-1} + (1 - \beta) \mathbf{g}_t \\
 &= \dots \\
 &= \beta^t \mathbf{v}_0 + \sum_{i=1}^t \beta^{t-i} (1 - \beta) \mathbf{g}_i \\
 &= (1 - \beta) \sum_{i=1}^t \beta^{t-i} \mathbf{g}_i \quad (\text{since } \mathbf{v}_0 = 0)
 \end{aligned}$$

3. Given  $\mathbf{g}_t$  has a stationary distribution independent of  $t$ , we know from the previous subquestion that ;

$$\begin{aligned}
 \mathbf{v}_t &= (1 - \beta) \sum_{i=1}^t \beta^{t-i} \mathbf{g}_i \\
 \mathbb{E}[\mathbf{v}_t] &= \mathbb{E}[(1 - \beta) \sum_{i=1}^t \beta^{t-i} \mathbf{g}_i] \\
 &= (1 - \beta) \sum_{i=1}^t \beta^{t-i} \mathbb{E}[\mathbf{g}_i] \quad (\text{linearity of expectation}) \\
 &= \mathbb{E}[\mathbf{g}_t] (1 - \beta) \sum_{i=1}^t \beta^{t-i} \quad (\text{stationarity}) \\
 &= \mathbb{E}[\mathbf{g}_t] (1 - \beta) \frac{1 - \beta^t}{1 - \beta} \quad (\text{sum of finite geometric series}) \\
 &= \mathbb{E}[\mathbf{g}_t] (1 - \beta^t)
 \end{aligned}$$

Thus, we can observe that the running average is biased. If we want an unbiased estimate, we can rescale  $\mathbf{v}_t$  by this bias, which would give us  $\tilde{\mathbf{v}}_t = \mathbf{v}_t(1 - \beta^t)$ ,  $\Delta \boldsymbol{\theta}_t = -\delta \mathbf{v}_t$

**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 2.2 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 2.3 and 2.4: 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.**

1. Given the hint,  $\text{Var}(Z) = \mathbb{E}[Z^2] - \mathbb{E}[Z]^2$ , we can rearrange it as follows;  $\mathbb{E}[Z^2] = \mathbb{E}[Z]^2 + \text{Var}(Z)$

$$\begin{aligned} \mathbb{E}_R[L(\mathbf{w})] &= \mathbb{E}_R[\|\mathbf{y} - (\mathbf{X} \odot \mathbf{R})\mathbf{w}\|^2] \\ &= \sum_{i=1}^n \mathbb{E}_R[(\mathbf{y}_i - (\mathbf{X}_i \odot \mathbf{R}_i)\mathbf{w})^2] \\ &= \sum_{i=1}^n (\mathbb{E}_R[\mathbf{y}_i - (\mathbf{X}_i \odot \mathbf{R}_i)\mathbf{w}]^2 + \text{Var}_R(\mathbf{y}_i - (\mathbf{X}_i \odot \mathbf{R}_i)\mathbf{w})) \\ &= \sum_{i=1}^n (\mathbf{y}_i - \mathbb{E}_R[\mathbf{R}_i] \odot \mathbf{X}_i)^2 + \text{Var}_R(\mathbf{y}_i - (\mathbf{X}_i \odot \mathbf{R}_i)\mathbf{w}) \end{aligned}$$

Deriving the first block and keeping in mind that the expected value of a Bernoulli random variable  $R$  is  $\mathbb{E}(R) = \Pr(R=1) \cdot 1 + \Pr(R=0) \cdot 0 = p \cdot 1 + q \cdot 0 = p$  we get;

$$\begin{aligned} &\sum_{i=1}^n (\mathbf{y}_i - \mathbb{E}_R[\mathbf{R}_i] \odot \mathbf{X}_i)^2 \\ &= \sum_{i=1}^n (\mathbf{y}_i - p\mathbf{X}_i)^2 \\ &= \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 \end{aligned}$$

Now doing the same thing for the second block with  $\mathbb{E}[R^2] = \Pr(R = 1) \cdot 1^2 + \Pr(R = 0) \cdot 0^2 = p \cdot 1^2 + q \cdot 0^2 = p$  we get ;

$$\begin{aligned}
& \text{Var}_R(\mathbf{y}_i - (\mathbf{X}_i \odot \mathbf{R}_i)\mathbf{w}) \\
&= \sum_{i=1}^n \mathbb{E}_R[(\mathbf{y}_i - (\mathbf{X}_i \odot \mathbf{R}_i)\mathbf{w} - \mathbb{E}_R[\mathbf{y}_i - (\mathbf{X}_i \odot \mathbf{R}_i)\mathbf{w}])^2] \\
&= \sum_{i=1}^n \mathbb{E}_R[(\mathbf{y}_i - (\mathbf{X}_i \odot \mathbf{R}_i)\mathbf{w} - \mathbf{y}_i + p\mathbf{X}_i\mathbf{w})^2] \\
&= \sum_{i=1}^n \mathbb{E}_R[((p\mathbf{X}_i - \mathbf{R}_i \odot \mathbf{X}_i)\mathbf{w})^2] \\
&= \sum_{i=1}^n \mathbb{E}_R[\mathbf{w}^T(p\mathbf{X}_i^T - \mathbf{R}_i^T \odot \mathbf{X}_i^T)(p\mathbf{X}_i - \mathbf{R}_i \odot \mathbf{X}_i)\mathbf{w}] \\
&= \sum_{i=1}^n (\mathbf{w}^T(p^2\mathbf{X}_i^T\mathbf{X}_i - p\mathbb{E}_R[\mathbf{R}_i]\mathbf{X}_i^T\mathbf{X}_i - p\mathbb{E}_R[\mathbf{R}_i^T]\mathbf{X}_i^T\mathbf{X}_i + \mathbb{E}[\mathbf{R}^T\mathbf{R}](\mathbf{X}_i^T\mathbf{X}_i))\mathbf{w}) \\
&= \sum_{i=1}^n (\mathbf{w}^T(p^2\mathbf{X}_i^T\mathbf{X}_i - p^2\mathbf{X}_i^T\mathbf{X}_i - p^2\mathbf{X}_i^T\mathbf{X}_i + p(\mathbf{X}_i^T\mathbf{X}_i))\mathbf{w}) \\
&= \sum_{i=1}^n (\mathbf{w}^T p(1-p)(\mathbf{X}_i^T\mathbf{X}_i)\mathbf{w}) \\
&= p(1-p)(\mathbf{w}^T \text{diag}(\mathbf{X}^T\mathbf{X})\mathbf{w}) \\
&= p(1-p)\|\Gamma\mathbf{w}\|^2
\end{aligned}$$

Combining the two together, we get ;

$$\mathbb{E}[L(\mathbf{w})] = \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1-p)\|\Gamma\mathbf{w}\|^2$$

2.

$$\mathbb{E}[L(\mathbf{w})] = \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1-p)\|\Gamma\mathbf{w}\|^2$$

Setting the gradient to 0, we have ;

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

Using the following formulas ;

$$\begin{aligned}
\frac{\partial \mathbf{X}\mathbf{w}}{\partial \mathbf{w}} &= \mathbf{X}^T, \quad \frac{\partial \mathbf{w}^T\mathbf{X}\mathbf{w}}{\partial \mathbf{w}} = (\mathbf{X} + \mathbf{X}^T)\mathbf{w}, \quad \frac{\partial \mathbf{w}^T\mathbf{X}}{\partial \mathbf{w}} = \mathbf{X} \\
&= 0 - p(\mathbf{y}^T\mathbf{X})^T - p(\mathbf{X}^T\mathbf{y}) + p^2(\mathbf{X}^T\mathbf{X} + \mathbf{X}^T\mathbf{X})\mathbf{w} + p(1-p)((\Gamma^2 + (\Gamma^2)^T)\mathbf{w}) \\
&= -2p(\mathbf{X}^T\mathbf{y}) + 2(p^2(\mathbf{X}^T\mathbf{X}) + p(1-p)\Gamma^2)\mathbf{w}
\end{aligned}$$

Now equaling this to 0;

$$\begin{aligned}
 p^2 \mathbf{X}^T \mathbf{X} \mathbf{w} + p(1-p)\Gamma^2 \mathbf{w} &= p \mathbf{X}^T \mathbf{y} \\
 (p \mathbf{X}^T \mathbf{X} + (1-p)\Gamma^2) p \mathbf{w} &= p \mathbf{X}^T \mathbf{y} \\
 (\mathbf{X}^T \mathbf{X} + (\frac{1-p}{p})\Gamma^2) p \mathbf{w} &= \mathbf{X}^T \mathbf{y} \\
 p \mathbf{w} &= (\mathbf{X}^T \mathbf{X} + (\frac{1-p}{p})\Gamma^2)^{-1} \mathbf{X}^T \mathbf{y}
 \end{aligned}$$

Letting  $\lambda^{\text{dropout}} = \frac{1-p}{p}$  we get;

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

When  $p = 1$ ,  $\lambda^{\text{dropout}} = 0$  so there is no contribution from the regularization term in  $L(\mathbf{w})$ . When  $p$  is close to 0,  $\lambda^{\text{dropout}} \rightarrow \infty$ , thus the contribution from the regularization increases.

3.

With  $L_2$  regularization we can write the loss as follows;

$$L(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda^{L_2} \|\mathbf{w}\|^2$$

$$\begin{aligned}
 \frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}) &= \frac{\partial}{\partial \mathbf{w}} (\|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda^{L_2} \|\mathbf{w}\|^2) \\
 &= \frac{\partial}{\partial \mathbf{w}} ((\mathbf{y}^T - \mathbf{w}^T \mathbf{X}^T)(\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda^{L_2} \mathbf{w}^T \mathbf{w}) \\
 &= \frac{\partial}{\partial \mathbf{w}} (\mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} + \lambda^{L_2} \mathbf{w}^T \mathbf{w}) \\
 &= 0 - (\mathbf{y}^T \mathbf{X})^T - \mathbf{X}^T \mathbf{y} + (\mathbf{X}^T \mathbf{X} + (\mathbf{X}^T \mathbf{X})^T) \mathbf{w} + 2\lambda^{L_2} \mathbf{w} \\
 &= -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \mathbf{w} + 2\lambda^{L_2} \mathbf{w}
 \end{aligned}$$

Equaling this to 0, we find  $\mathbf{w}^{L_2}$  to a minimal  $L(\mathbf{w})$ ;

$$\begin{aligned}
 -\mathbf{X}^T \mathbf{y} + \mathbf{X}^T \mathbf{X} \mathbf{w} + \lambda^{L_2} \mathbf{w} &= 0 \\
 \mathbf{X}^T \mathbf{X} \mathbf{w} + \lambda^{L_2} \mathbf{w} &= \mathbf{X}^T \mathbf{y} \\
 \mathbf{w}(\mathbf{X}^T \mathbf{X} + \lambda^{L_2}) &= \mathbf{X}^T \mathbf{y} \\
 \mathbf{w}^{L_2} &= (\mathbf{X}^T \mathbf{X} + \lambda^{L_2})^{-1} \mathbf{X}^T \mathbf{y}
 \end{aligned}$$

4. We can observe that dropout with linear regression is somewhat equivalent to a scaled version of ridge regression. In fact, in dropout, the penalty on the weight  $\mathbf{w}_i$  is scaled by the standard deviation of the  $i$ 'th feature. As for  $L^2$ , the weight is penalized uniformly. Also, with dropout, when the probability of dropping units is higher, the regularization coefficient grows, but when the probability decreases, we can see it as an ordinary least square problem.

**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  $\mathbf{b}^{(t)} = [c, \dots, c]^\top$  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  $\mathbb{E}[(h_i^{(t-1)})^2]$
  - 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  $1 \leq i \leq d^{(t)}$ .
  - What popular initialization scheme has this form?
  - Why do you think this initialization would work well in practice? Answer in 1-2 sentences.
- 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.**

- a) Given that  $X \perp Y$ , we can use some useful properties in our calculations.

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

1. b) Keeping in mind that  $\mathbb{E}[h_i^{(t-1)}] = 0$ ,  $\text{Var}(h_i^{(t-1)}) = 1$  for  $1 \leq i \leq d^{(t-1)}$  and that  $X \perp Y$  we have ;

$$\begin{aligned}
 \mathbb{E}[a_i^{(t)}] &= \mathbb{E}[\mathbf{W}_i^{(t)} \mathbf{h}^{(t-1)} + \mathbf{b}_i^{(t)}] \\
 &= \mathbb{E}\left[\sum_{j=1}^{d^{(t-1)}} \mathbf{W}_{ij}^{(t)} \mathbf{h}_j^{(t-1)}\right] + \mathbb{E}[\mathbf{b}_i^{(t)}] \\
 &= \sum_{j=1}^{d^{(t-1)}} \mathbb{E}[\mathbf{W}_{ij}^{(t)} \mathbf{h}_j^{(t-1)}] + \mathbb{E}[\mathbf{b}_i^{(t)}] \\
 &= \sum_{j=1}^{d^{(t-1)}} \mathbb{E}[\mathbf{W}_{ij}^{(t)}] \mathbb{E}[\mathbf{h}_j^{(t-1)}] + \mathbb{E}[\mathbf{b}_i^{(t)}] \\
 &= \mathbb{E}[\mathbf{b}_i^{(t)}] \\
 &= c
 \end{aligned}$$

As for  $\text{Var}(\mathbf{a}_i^{(t)})$  we have ;

$$\begin{aligned}
 \text{Var}(\mathbf{a}_i^{(t)}) &= \text{Var}(\mathbf{W}_i^{(t)} \mathbf{h}^{(t-1)} + \mathbf{b}_i^{(t)}) \\
 &= \text{Var}\left(\sum_{j=1}^{d^{(t-1)}} \mathbf{W}_{ij}^{(t)} \mathbf{h}_j^{(t-1)}\right) + \text{Var}(\mathbf{b}_i^{(t)}) \\
 &= \sum_{j=1}^{d^{(t-1)}} \text{Var}(\mathbf{W}_{ij}^{(t)} \mathbf{h}_j^{(t-1)}) + \text{Var}(\mathbf{b}_i^{(t)}) \\
 &= \sum_{j=1}^{d^{(t-1)}} (\text{Var}(\mathbf{W}_{ij}^{(t)}) \text{Var}(\mathbf{h}_j^{(t-1)}) + \text{Var}(\mathbf{W}_{ij}^{(t)}) \mathbb{E}[\mathbf{h}_j^{(t-1)}]^2 + \text{Var}(\mathbf{h}_j^{(t-1)}) \mathbb{E}[\mathbf{W}_{ij}^{(t)}]^2) + \text{Var}(\mathbf{b}_i^{(t)}) \\
 &= d^{(t-1)} (\text{Var}(\mathbf{W}_{ij}^{(t)}) + \mathbb{E}[\mathbf{W}_{ij}^{(t)}]^2)
 \end{aligned}$$

1. c) We could initialize the parameter  $c = 0$ , such that  $\mathbb{E}[\mathbf{a}_i^{(t)}] = 0$  and  $\text{Var}(\mathbf{b}_i^{(t)}) = 0$  and let  $\mathbb{E}[\mathbf{W}_{ij}^{(t)}] = 0$  and  $\text{Var}(\mathbf{W}_{ij}^{(t)}) = \frac{1}{d^{(t-1)}}$  so that  $\text{Var}(\mathbf{a}_i^{(t)}) = 1$ . With the Gaussian weight distribution  $\mathbf{W}_{ij}^{(t)} \sim \mathcal{N}(\mu, \sigma^2)$ ,

$$\mathbb{E}[\mathbf{W}_{ij}^{(t)}] = \mu \quad \text{Var}(\mathbf{W}_{ij}^{(t)}) = \sigma^2$$

Thus,

$$c = 0, \quad \mu = 0, \quad \sigma^2 = \frac{1}{d^{(t-1)}}$$

2. a)

$$\begin{aligned}
\mathbb{E}[(\mathbf{h}_i^{(t-1)})^2] &= \int_{-\infty}^{\infty} \max\{0, \mathbf{a}_i^{(t-1)}\}^2 p(\mathbf{a}_i^{(t-1)}) d\mathbf{a}_i^{(t-1)} \\
&= \int_0^{\infty} (\mathbf{a}_i^{(t-1)})^2 p(\mathbf{a}_i^{(t-1)}) d\mathbf{a}_i^{(t-1)} \\
&= \frac{1}{2} \int_{-\infty}^{\infty} (\mathbf{a}_i^{(t-1)})^2 p(\mathbf{a}_i^{(t-1)}) d\mathbf{a}_i^{(t-1)} \quad \text{since } \mathbf{a}_i^{(t-1)} \text{ is symmetric} \\
&= \frac{1}{2} \mathbb{E}[(\mathbf{a}_i^{(t-1)})^2] \\
&= \frac{1}{2} (\mathbb{E}[(\mathbf{a}_i^{(t-1)})^2] - \mathbb{E}[(\mathbf{a}_i^{(t-1)})]^2) \\
&= \frac{1}{2} \text{Var}(\mathbf{a}_i^{(t-1)}) \\
&= \frac{1}{2} \quad \text{since } \text{Var}(\mathbf{a}_i^{(t-1)}) = 1
\end{aligned}$$

2. b) Based on the previous questions we know ;

$$\mathbb{E}[\mathbf{a}_i^{(t)}] = \sum_{j=1}^{d^{(t-1)}} \mathbb{E}[\mathbf{W}_{ij}^{(t)}] \mathbb{E}[\mathbf{h}_j^{(t-1)}] + \mathbb{E}[\mathbf{b}_i^{(t)}]$$

$$\begin{aligned}
\text{Var}(\mathbf{a}_i^{(t)}) &= \sum_{j=1}^{d^{(t-1)}} (\text{Var}(\mathbf{W}_{ij}^{(t)}) \text{Var}(\mathbf{h}_j^{(t-1)}) + \text{Var}(\mathbf{W}_{ij}^{(t)}) \mathbb{E}[\mathbf{h}_j^{(t-1)}]^2 + \text{Var}(\mathbf{h}_j^{(t-1)}) \mathbb{E}[\mathbf{W}_{ij}^{(t)}]^2 + \text{Var}(\mathbf{b}_i^{(t)})) \\
&= \sum_{j=1}^{d^{(t-1)}} (\text{Var}(\mathbf{W}_{ij}^{(t)}) (\mathbb{E}[(\mathbf{h}_j^{(t-1)})^2] - \mathbb{E}[\mathbf{h}_j^{(t-1)}]^2) + \text{Var}(\mathbf{W}_{ij}^{(t)}) \mathbb{E}[\mathbf{h}_j^{(t-1)}]^2 + \\
&\quad \text{Var}(\mathbf{h}_j^{(t-1)}) \mathbb{E}[\mathbf{W}_{ij}^{(t)}]^2 + \text{Var}(\mathbf{b}_i^{(t)})) \\
&= \sum_{j=1}^{d^{(t-1)}} (\text{Var}(\mathbf{W}_{ij}^{(t)}) \mathbb{E}[(\mathbf{h}_j^{(t-1)})^2] + \text{Var}(\mathbf{h}_j^{(t-1)}) \mathbb{E}[\mathbf{W}_{ij}^{(t)}]^2 + \text{Var}(\mathbf{b}_i^{(t)})) \\
&= d^{(t-1)} (\text{Var}(\mathbf{W}_{ij}^{(t)}) \mathbb{E}[(\mathbf{h}_j^{(t-1)})^2] + \text{Var}(\mathbf{h}_j^{(t-1)}) \mathbb{E}[\mathbf{W}_{ij}^{(t)}]^2 + \text{Var}(\mathbf{b}_i^{(t)}))
\end{aligned}$$

We could initialize the parameters to  $c = 0$  and  $\mathbb{E}[\mathbf{W}_{ij}^{(t)}] = 0$ , so that  $\mathbb{E}[\mathbf{a}_i^{(t)}] = 0$ ,  $\text{Var}(\mathbf{b}_i^{(t)}) = 0$  and simplify  $\text{Var}(\mathbf{a}_i^{(t)})$ .

$$\text{Var}(\mathbf{a}_i^{(t)}) = \frac{d^{(t-1)}}{2} \text{Var}(\mathbf{W}_{ij}^{(t)}) = 1$$

Thus,

$$\text{Var}(\mathbf{W}_{ij}^{(t)}) = \frac{2}{d^{(t-1)}}$$

With the Gaussian weight distribution  $\mathbf{W}_{ij}^{(t)} \sim \mathcal{N}(\mu, \sigma^2)$ ,

$$\mathbb{E}[\mathbf{W}_{ij}^{(t)}] = \mu \quad \text{Var}(\mathbf{W}_{ij}^{(t)}) = \sigma^2$$



Thus,

$$c = 0 \quad \mu = 0 \quad \sigma^2 = \frac{2}{d^{(t-1)}}$$

2. c) The He Initialization has this form (ReLU version of the Glorot initialization).

2. d) This initialization scheme would work well in practice since with the He Initialization, the scale of the activations is controlled, which decreases the probability of the gradient exploding or vanishing, thus enabling the model to train better.

3. For the assumptions in (1) and given  $\alpha, \beta$  for  $W_{ij}^{(t)} \sim \text{Uniform}(\alpha, \beta)$  such that  $\mathbb{E}[a_i^{(t)}] = 0$  and  $\text{Var}(a_i^{(t)}) = 1$ ,

$$\mathbb{E}[\mathbf{W}_{ij}^{(t)}] = \frac{\alpha + \beta}{2}, \quad \text{Var}(\mathbf{W}_{ij}^{(t)}) = \frac{(\beta - \alpha)^2}{12}$$

Therefore,

$$c = 0, \quad \beta = -\alpha > 0, \quad \beta = \sqrt{\frac{3}{d^{(t-1)}}}$$

For the assumptions in (2) and given  $\alpha, \beta$  for  $W_{ij}^{(t)} \sim \text{Uniform}(\alpha, \beta)$  such that  $\mathbb{E}[a_i^{(t)}] = 0$  and  $\text{Var}(a_i^{(t)}) = 1$ ,

$$\mathbb{E}[\mathbf{W}_{ij}^{(t)}] = \frac{\alpha + \beta}{2}, \quad \text{Var}(\mathbf{W}_{ij}^{(t)}) = \frac{(\beta - \alpha)^2}{12}$$

Thus,

$$c = 0, \quad \beta = -\alpha > 0, \quad \beta = \sqrt{\frac{6}{d^{(t-1)}}}$$

**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:  $\boldsymbol{\mu}_{\text{batch}}$ ,  $\boldsymbol{\mu}_{\text{layer}}$ ,  $\boldsymbol{\mu}_{\text{instance}}$ ,  $\boldsymbol{\sigma}_{\text{batch}}^2$ ,  $\boldsymbol{\sigma}_{\text{layer}}^2$ , and  $\boldsymbol{\sigma}_{\text{instance}}^2$ .

$$\left[ \begin{bmatrix} 1, 3, 2 \\ 1, 2, 3 \end{bmatrix}, \begin{bmatrix} 3, 3, 2 \\ 2, 4, 4 \end{bmatrix}, \begin{bmatrix} 4, 2, 2 \\ 1, 2, 4 \end{bmatrix}, \begin{bmatrix} 3, 3, 2 \\ 3, 3, 2 \end{bmatrix} \right]$$

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  $\mathbf{w}$ :

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

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

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  $y = \mathbf{u}^\top \mathbf{x}$ . Assume the data  $\mathbf{x}$  (a random vector) is whitened ( $\text{Var}(\mathbf{x}) = \mathbf{I}$ ) and centered at 0 ( $\mathbb{E}[\mathbf{x}] = \mathbf{0}$ ). Show that  $\hat{y} = \mathbf{w}^\top \mathbf{x} + \beta$ .

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

**Answer 4.**

1. Batch normalization in a the case of a convolutional neural network is different than batch normalization for usual networks since filter weights are shared across the input image. In that case, it is reasonable to normalize the output in the same way so that each output value takes the mean and variance of batch \* features, at different location.

With that in mind,

$$\mu_{\text{batch}} = \begin{bmatrix} (1+3+2+3+3+2+4+2+2+3+3+2)/12 \\ (1+2+3+2+4+4+1+2+4+3+3+2)/12 \end{bmatrix}$$

$$\mu_{\text{batch}} = \begin{bmatrix} 2.5 \\ 2.5833 \end{bmatrix}$$

$$\begin{aligned} \sigma_{\text{batch}}^2 &= \begin{bmatrix} \frac{(1-2.5)^2 + (3-2.5)^2 + (2-2.5)^2 + (3-2.5)^2 + \dots + (2-2.5)^2}{12} \\ \frac{(1-2.5833)^2 + (2-2.5833)^2 + (3-2.5833)^2 + \dots + (2-2.5833)^2}{12} \end{bmatrix} \\ &= \begin{bmatrix} 0.5833 \\ 1.0764 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \mu_{\text{layer}} &= \begin{bmatrix} [(1+3+2)/3], [(3+3+2)/3], [(4+2+2)/3], [(3+3+2)/3] \\ [(1+2+3)/3], [(2+4+4)/3], [(1+2+4)/3], [(3+3+2)/3] \end{bmatrix} \\ &= \begin{bmatrix} [2], [2.6667], [2.6667], [2.6667] \\ [2], [3.3333], [2.3333], [2.6667] \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \sigma_{\text{layer}}^2 &= \begin{bmatrix} \left[ \frac{(1-2)^2 + (3-2)^2 + (2-2)^2}{3} \right], \begin{bmatrix} \dots \\ \dots \end{bmatrix}, \begin{bmatrix} \dots \\ \dots \end{bmatrix}, \left[ \frac{(3-2.6667)^2 + (3-2.6667)^2 + (2-2.6667)^2}{3} \right] \\ \left[ \frac{(1-2)^2 + (2-2)^2 + (3-2)^2}{3} \right], \begin{bmatrix} \dots \\ \dots \end{bmatrix}, \begin{bmatrix} \dots \\ \dots \end{bmatrix}, \left[ \frac{(3-2.6667)^2 + (3-2.6667)^2 + (2-2.6667)^2}{3} \right] \end{bmatrix} \\ &= \begin{bmatrix} [0.6666], [0.2222], [0.8889], [0.2222] \\ [0.6666], [1.8889], [1.5556], [0.3334] \end{bmatrix} \end{aligned}$$

$$\mu_{\text{instance}} = \begin{bmatrix} [2], [2.6667], [2.6667], [2.6667] \\ [2], [3.3333], [2.3333], [2.6667] \end{bmatrix}$$

$$\sigma_{\text{instance}}^2 = \begin{bmatrix} [0.6666], [0.2222], [0.8889], [0.2222] \\ [0.6666], [1.8889], [1.5556], [0.3334] \end{bmatrix}$$

1. As a side note:  $\mathbf{W}^\perp$  is an orthogonal complement that projects the gradient away from the direction of  $\mathbf{w}$ , 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.

2. Given the following definitions and properties ;  $\mathbf{w} = \gamma \frac{\mathbf{u}}{\|\mathbf{u}\|}$ ,  $\mathbf{y} = \mathbf{u}^T \mathbf{x}$ ,  $\hat{\mathbf{y}} = \gamma \frac{\mathbf{y} - \mu_y}{\sigma_y} + \beta$ ,  $\text{Var}(x) = \mathbf{I}$ , and  $\mathbb{E}[x] = 0$ .

We can derive the following ;

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

$$\begin{aligned}\text{Var}(y) &= \text{Var}(\mathbf{u}^T \mathbf{x}) \\ &= \mathbf{u}^T \text{Var}(x) \mathbf{u} \\ &= \mathbf{u}^T \mathbf{u} \\ &= \|\mathbf{u}\|^2 \\ &= \sigma_y \\ &= \|\mathbf{u}\|\end{aligned}$$

Thus,

$$\begin{aligned}\hat{\mathbf{y}} &= \gamma \frac{\mathbf{y} - \mu_y}{\sigma_y} + \beta \\ &= \gamma \frac{\mathbf{u}^T \mathbf{x} - 0}{\|\mathbf{u}\|} + \beta \\ &= (\gamma \frac{\mathbf{u}^T}{\|\mathbf{u}\|}) \mathbf{x} + \beta \\ &= \mathbf{w}^T \mathbf{x} + \beta\end{aligned}$$

3. First off,  $\nabla_{\mathbf{u}} L(\mathbf{u}) = \nabla_{\mathbf{u}} \mathbf{w} \nabla_{\mathbf{w}} L = \nabla_{\mathbf{u}} (\gamma \frac{\mathbf{u}}{\|\mathbf{u}\|}) \nabla_{\mathbf{w}} L = \gamma (\nabla_{\mathbf{u}} \frac{\mathbf{u}}{\|\mathbf{u}\|}) \nabla_{\mathbf{w}} L$

Using the quotient rule and the previous definitions, we get:

$$\begin{aligned}\nabla_{\mathbf{u}} \frac{\mathbf{u}}{\|\mathbf{u}\|} &= \frac{\|\mathbf{u}\| \frac{d\mathbf{u}}{du} - \mathbf{u} \frac{d\|\mathbf{u}\|}{du}}{\|\mathbf{u}\|^2} \\ &= \frac{\|\mathbf{u}\| \mathbf{I} - \mathbf{u} \frac{d}{du} ((\mathbf{u}^T \mathbf{u})^{\frac{1}{2}})}{\|\mathbf{u}\|^2} \\ &= \frac{1}{\|\mathbf{u}\|^2} (\|\mathbf{u}\| \mathbf{I} - \mathbf{u} \frac{1}{2} (\mathbf{u}^T \mathbf{u})^{-\frac{1}{2}} \frac{d(\mathbf{u}^T \mathbf{u})}{du}) \\ &= \frac{1}{\|\mathbf{u}\|^2} (\|\mathbf{u}\| \mathbf{I} - \mathbf{u} \frac{1}{2\|\mathbf{u}\|} 2\mathbf{u}^T) \\ &= \frac{1}{\|\mathbf{u}\|^2} (\|\mathbf{u}\| \mathbf{I} - \frac{\mathbf{u} \mathbf{u}^T}{\|\mathbf{u}\|}) \\ &= \frac{1}{\|\mathbf{u}\|} (\mathbf{I} - \frac{\mathbf{u} \mathbf{u}^T}{\|\mathbf{u}\|^2}) \\ &= \frac{1}{\|\mathbf{u}\|} \mathbf{W}^\perp\end{aligned}$$

With that in mind,

$$\nabla_{\mathbf{u}} L(\mathbf{u}) = \frac{\partial \mathbf{w}}{\partial \mathbf{u}} \nabla_{\mathbf{w}} L = \frac{\gamma}{\|\mathbf{u}\|} (\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^T}{\|\mathbf{u}\|^2}) \nabla_{\mathbf{w}} L = s \mathbf{W}^\perp \nabla_{\mathbf{w}} L$$

where  $s = \frac{\gamma}{\|\mathbf{u}\|}$  and  $\mathbf{W}^\perp = (\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^T}{\|\mathbf{u}\|^2})$ . Also, note that  $\mathbf{W}^\perp \mathbf{u} = \mathbf{u} - \mathbf{u} \frac{\mathbf{u}\mathbf{u}^T}{\|\mathbf{u}\|^2} = \mathbf{u} - \mathbf{u} = 0$ .

**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.**

1. By observing the two functions, we can see that;

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

This lets us see that  $\mathbf{g}_t = \sigma(\mathbf{h}_t)$ , for each  $t \geq 0$ . Assuming that it is also true that  $\mathbf{g}_{t-1} = \sigma(\mathbf{h}_{t-1})$  and by fixing  $t \geq 1$  then we can prove that  $\mathbf{g}_t = \sigma(\mathbf{h}_t)$ .

We have;

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

<sup>2</sup> 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.

2.

$$\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| = \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \right\| \cdot \left\| \frac{\partial \mathbf{h}_{T-1}}{\partial \mathbf{h}_{T-2}} \right\| \cdots \left\| \frac{\partial \mathbf{h}_1}{\partial \mathbf{h}_0} \right\|$$

With the two properties given,

$$\begin{aligned} \left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\| &= \left\| \mathbf{W} \frac{\partial \mathbf{h}_{t-1}}{\partial \mathbf{h}_{t-1}} \right\| = \|\mathbf{W} \sigma'\| \\ &\leq \|\mathbf{W}\| \|\sigma'\| \quad (\text{Since } \|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|) \\ &\leq \|\mathbf{W}\| \cdot \gamma \quad (\text{Since } |\sigma'| \leq \gamma) \\ &= \sqrt{\lambda_1(\mathbf{W}^\top \mathbf{W})} \cdot \gamma \quad (\text{Since } \|\mathbf{A}\| = \sqrt{\lambda_1(\mathbf{A}^\top \mathbf{A})}) \\ &\leq \sqrt{\frac{\delta^2}{\gamma^2}} \cdot \gamma \quad (\text{Since } \lambda_1(\mathbf{W}^\top \mathbf{W}) \leq \frac{\delta^2}{\gamma^2}) \\ &= \left| \frac{\delta}{\gamma} \right| \cdot \gamma \\ &= \frac{\delta}{\gamma} \cdot \gamma \quad (\text{Since } \delta > 0, \quad \gamma > 0) \\ &= \delta \end{aligned}$$

Thus,

$$\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\| \leq \delta \implies \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| \leq \delta^T \implies \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| \rightarrow 0 \quad \text{as } T \rightarrow \infty$$

3. The largest eigenvalue of the weights being larger than  $\frac{\delta^2}{\gamma^2}$  is a necessary condition because otherwise the gradient of the hidden state would not become arbitrarily large. However, it is not sufficient because the product of the norms can be greater than the norm of the product.

**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.

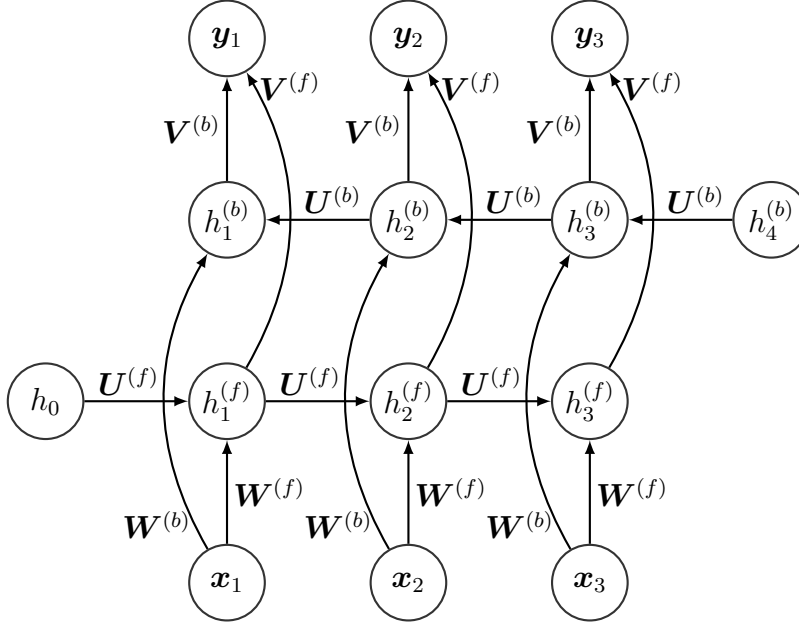


FIGURE 1 – Computational graph of the bidirectional RNN unrolled for three timesteps.

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.**

2. The gradient of  $L_t$  with respect to  $\mathbf{h}_t^{(f)}$  is

$$\nabla_{\mathbf{h}_t^{(f)}} L_t = \left( \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t^{(f)}} \right)^\top \nabla_{\mathbf{y}_t} L_t$$

$$\begin{aligned}
\nabla_{\mathbf{y}_t} L_t &= \frac{\partial \|\mathbf{z}_t - \mathbf{y}_t\|_2^2}{\partial \mathbf{y}_t} \\
&= \frac{\partial (\mathbf{z}_t - \mathbf{y}_t)^T (\mathbf{z}_t - \mathbf{y}_t)}{\partial \mathbf{y}_t} \\
&= \frac{\partial (\mathbf{z}_t^T - \mathbf{y}_t^T) (\mathbf{z}_t - \mathbf{y}_t)}{\partial \mathbf{y}_t} \\
&= \frac{\partial}{\partial \mathbf{y}_t} (\mathbf{z}_t^T \mathbf{z}_t - \mathbf{z}_t^T \mathbf{y}_t - \mathbf{y}_t^T \mathbf{z}_t + \mathbf{y}_t^T \mathbf{y}_t) \\
&= -\mathbf{z}_t^T - \mathbf{z}_t + 2\mathbf{y}_t \\
&= -2(\mathbf{z}_t - \mathbf{y}_t)
\end{aligned}$$

$$\begin{aligned}
\left( \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t^{(f)}} \right)^T &= \left( \partial \mathbf{h}_t^{(f)} (\mathbf{V}^{(f)} \mathbf{h}_t^{(f)} + \mathbf{V}^{(b)} \mathbf{h}_t^{(b)}) \right)^T \\
&= \mathbf{V}^{(f)\top}
\end{aligned}$$

Combining the two together, we get ;

$$\nabla_{\mathbf{h}_t^{(f)}} L_t = -2\mathbf{V}^{(f)\top} (\mathbf{z}_t - \mathbf{y}_t)$$

Now for the gradient of  $L_t$  with respect to  $\mathbf{h}_t^{(b)}$  is  $\nabla_{\mathbf{h}_t^{(b)}} L_t = \left( \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t^{(b)}} \right)^T \nabla_{\mathbf{y}_t} L_t$ .

$$\begin{aligned}
\left( \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t^{(b)}} \right)^T &= \left( \partial \mathbf{h}_t^{(b)} (\mathbf{V}^{(f)} \mathbf{h}_t^{(f)} + \mathbf{V}^{(b)} \mathbf{h}_t^{(b)}) \right)^T \\
&= \mathbf{V}^{(b)\top}
\end{aligned}$$

Reusing what we already computed and combining, we get ;

$$\nabla_{\mathbf{h}_t^{(b)}} L_t = -2\mathbf{V}^{(b)\top} (\mathbf{z}_t - \mathbf{y}_t)$$

Now for  $\frac{\partial \mathbf{h}_{t+1}^{(f)}}{\partial \mathbf{h}_t^{(f)}}$  we have ;

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

$$\begin{aligned}
\frac{\partial \mathbf{h}_{t+1}^{(f)}}{\partial \mathbf{h}_t^{(f)}} &= \frac{\partial (\sigma(\mathbf{W}^{(f)} \mathbf{x}_{t+1} + \mathbf{U}^{(f)} \mathbf{h}_t^{(f)}))}{\partial \mathbf{h}_t^{(f)}} \\
&= \sigma(\mathbf{W}^{(f)} \mathbf{x}_{t+1} + \mathbf{U}^{(f)} \mathbf{h}_t^{(f)}) (1 - \sigma(\mathbf{W}^{(f)} \mathbf{x}_{t+1} + \mathbf{U}^{(f)} \mathbf{h}_t^{(f)})) \cdot \mathbf{U}^{(f)} \\
&\quad \text{Given that } \nabla_a \sigma(a) = \sigma(a)(1 - \sigma(a)) \text{ and then applying the chain rule.} \\
&= \text{diag}(\mathbf{h}_{t+1}^{(f)} (1 - \mathbf{h}_{t+1}^{(f)})) \mathbf{U}^{(f)}
\end{aligned}$$

Similarly, for  $\frac{\partial \mathbf{h}_{t-1}^{(b)}}{\partial \mathbf{h}_t^{(b)}}$  we get ;

$$\begin{aligned} \frac{\partial \mathbf{h}_{t-1}^{(b)}}{\partial \mathbf{h}_t^{(b)}} &= \frac{\partial(\sigma(\mathbf{W}^{(b)} \mathbf{x}_{t-1} + \mathbf{U}^{(b)} \mathbf{h}_t^{(b)}))}{\partial \mathbf{h}_t^{(b)}} \\ &= \sigma(\mathbf{W}^{(b)} \mathbf{x}_{t-1} + \mathbf{U}^{(b)} \mathbf{h}_t^{(b)}) (1 - \sigma(\mathbf{W}^{(b)} \mathbf{x}_{t-1} + \mathbf{U}^{(b)} \mathbf{h}_t^{(b)})) \cdot \mathbf{U}^{(b)} \\ &\quad \text{Given that } \nabla_a \sigma(a) = \sigma(a)(1 - \sigma(a)) \text{ and then applying the chain rule.} \\ &= \text{diag}(\mathbf{h}_{t-1}^{(b)}(1 - \mathbf{h}_{t-1}^{(b)})) \mathbf{U}^{(b)} \end{aligned}$$

3. As suggested in the question, let  $t$  denote the contribution of weight matrices when computing the recurrent hidden unit, i.e,  $\mathbf{W}_t^{(f)}$  and  $\mathbf{U}_t^{(f)}$  we have ;

$$\begin{aligned} \nabla_{\mathbf{W}^{(f)}} L &= \sum_t \nabla_{\mathbf{W}_t^{(f)}} L \\ &= \sum_t \left( \frac{\partial \mathbf{h}_t^{(f)}}{\partial \mathbf{W}_t^{(f)}} \right)^T \nabla_{\mathbf{h}_t^{(f)}} L \\ &= \sum_t \text{diag}(\mathbf{h}_t^{(f)}(1 - \mathbf{h}_t^{(f)})) (\nabla_{\mathbf{h}_t^{(f)}} L) \mathbf{x}_t^T \end{aligned}$$

Similarly, we have ;

$$\begin{aligned} \nabla_{\mathbf{U}^{(b)}} L &= \sum_t \nabla_{\mathbf{U}_t^{(b)}} L \\ &= \sum_t \left( \frac{\partial \mathbf{h}_t^{(b)}}{\partial \mathbf{U}_t^{(b)}} \right)^T \nabla_{\mathbf{h}_t^{(b)}} L \\ &= \sum_t \text{diag}(\mathbf{h}_t^{(b)}(1 - \mathbf{h}_t^{(b)})) (\nabla_{\mathbf{h}_t^{(b)}} L) \mathbf{h}_{t+1}^T \end{aligned}$$