

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 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 (α, ϵ) as a function of (β, δ) .

- 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. 1. Express the two update rules recursively:

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

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

→ Both update rules are equivalent, with $(\alpha, \epsilon) = (\beta, \Delta(1-\beta))$ and $(\beta, \Delta) = (\alpha, \frac{\epsilon}{1-\alpha})$

2. SGD with running average:

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

3. We know that \mathbf{g}_t has a stationary distribution independent of t and we need to show that the running average is biased. Lets start with calculating $\mathbb{E}[\mathbf{v}_t]$.

$$\begin{aligned}
\mathbb{E}[\mathbf{v}_t] &= \mathbb{E}[(1 - \beta) \sum_{i=1}^t \beta^{t-i} \mathbf{g}_i] \\
&= (1 - \beta) \sum_{i=1}^t \beta^{t-i} \mathbb{E}[\mathbf{g}_i] \xrightarrow{g-is-stationary} \\
&= (1 - \beta) \sum_{i=1}^t \beta^{t-i} \mathbb{E}[\mathbf{g}_t] \\
&= (1 - \beta) \mathbb{E}[\mathbf{g}_t] \sum_{i=1}^t \beta^{t-i} \\
&= (1 - \beta) \mathbb{E}[\mathbf{g}_t] \sum_{i=0}^t \beta^{t-1} \xrightarrow{Telescoping-Sum} \\
&= (1 - \beta) \mathbb{E}[\mathbf{g}_t] \frac{1 - \beta^t}{1 - \beta} \xrightarrow{(1-\beta)cancelled-out} \\
&= \mathbb{E}[\mathbf{g}_t] (1 - \beta^t)
\end{aligned} \tag{2}$$

We showed that the running average is biased. If we rescale \mathbf{v}_t by $\frac{1}{1-\beta^t}$ then the bias will be removed and we will have an unbiased estimation of running average.

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 Γ 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 λ^{dropout} is a regularization coefficient depending on p . How does the value of p affect the regularization coefficient, λ^{dropout} ?

3. Express the loss function for a linear regression problem without dropout and with L^2 regularization, with regularization coefficient λ^{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. 1. Since we are trying to find the expectation over a squared term we can re-write this formula $\text{Var}(Z) = \mathbb{E}[Z^2] - \mathbb{E}[Z]^2$ in the following form $\mathbb{E}[Z^2] = \mathbb{E}[Z]^2 + \text{Var}(Z)$. $L(\mathbf{w})$ have a squared form so let $Z = L(\mathbf{w}) = \|\mathbf{y} - (\mathbf{X} \odot \mathbf{R})\mathbf{w}\|^2$.

\mathbf{R}_{ij} is a Bernoulli variable so by definition we know $\mathbb{E}[r_{ij}] = p$, and $\text{Var}(r_{ij}) = (1-p)(0-p)^2 + p(1-p)^2 = p^2(1-p) + p(1-p-p+p^2) = p - p^2 = p(1-p)$.

$$\begin{aligned} \mathbb{E}[L(\mathbf{w})] &= \mathbb{E}[\|\mathbf{y}_i - (\mathbf{x}_i \odot \mathbf{r}_i)\mathbf{w}\|^2] = \sum_{i=1}^n \mathbb{E}[(\mathbf{y}_i - (\mathbf{x}_i \odot \mathbf{r}_i)\mathbf{w})^2] \\ &= \sum_{i=1}^n \mathbb{E}[(\mathbf{y}_i - (\mathbf{x}_i \odot \mathbf{r}_i)\mathbf{w})^2] + \text{Var}(\mathbf{y}_i - (\mathbf{x}_i \odot \mathbf{r}_i)\mathbf{w}) \\ &= \sum_{i=1}^n (\mathbf{y}_i - (\mathbf{x}_i \odot \mathbb{E}[\mathbf{r}_i])\mathbf{w})^2 + \text{Var}(\mathbf{y}_i - (\mathbf{x}_i \odot \mathbf{r}_i)\mathbf{w}) \xrightarrow{\text{Bernoulli-Variable}} \\ &= \sum_{i=1}^n (\mathbf{y}_i - p\mathbf{x}_i\mathbf{w})^2 + \sum_{i=1}^n \mathbf{w}^T p(1-p)(\mathbf{X}_i^T \mathbf{X}_i) \mathbf{w} \\ &= \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1-p)\text{diag}(\mathbf{X}^T \mathbf{X})\mathbf{w} \\ &= \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1-p)\|\Gamma\mathbf{w}\|^2 \end{aligned} \quad (3)$$

2. How does p affects the regularization coefficient, λ^{dropout} ? In order to find the $\mathbf{w}^{\text{dropout}}$ which minimize expected loss $\mathbb{E}[L(\mathbf{w})] = \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1-p)\|\Gamma\mathbf{w}\|^2$ we need to set the gradient of $\mathbb{E}[L(\mathbf{w})]$ to zero:

$$\begin{aligned}
\frac{d\mathbb{E}[L(\mathbf{w})]}{d\mathbf{w}} &= -2(-p\mathbf{X}^T)(\mathbf{y} - p\mathbf{X}\mathbf{w}^{dropout}) + 2p(1-p)\Gamma^T(\Gamma\mathbf{w}^{dropout}) \\
&= -2p\mathbf{X}^T\mathbf{y} + 2p^2\mathbf{X}^T\mathbf{X}\mathbf{w}^{dropout} + 2p(1-p)(\Gamma^T\Gamma)\mathbf{w}^{dropout} = 0 \xrightarrow{\text{Note-that:}} \\
&(\Gamma^T\Gamma = (\text{diag}(\mathbf{X}^T\mathbf{X}))^T(\text{diag}(\mathbf{X}^T\mathbf{X})^{0.5}) = \text{diag}(\mathbf{X}^T\mathbf{X})^{0.5}\text{diag}(\mathbf{X}^T\mathbf{X})^{0.5} = \Gamma^2) \\
&= -p\mathbf{X}^T\mathbf{y} + p^2\mathbf{X}^T\mathbf{X}\mathbf{w}^{dropout} + p(1-p)(\Gamma^T\Gamma)\mathbf{w}^{dropout} \\
&= p^2\mathbf{X}^T\mathbf{X}\mathbf{w}^{dropout} + p(1-p)\Gamma^2\mathbf{w}^{dropout} = p\mathbf{X}^T\mathbf{y} \xrightarrow{\text{divide-by-p}} \\
&(\mathbf{X}^T\mathbf{X} + (\frac{1}{p} - 1)\Gamma^2)p\mathbf{w}^{dropout} = \mathbf{X}^T\mathbf{y} \longrightarrow \\
&p\mathbf{w}^{dropout} = (\mathbf{X}^T\mathbf{X} + (\frac{1}{p} - 1)\Gamma^2)^{-1}\mathbf{X}^T\mathbf{y}
\end{aligned} \tag{4}$$

Until now we showed that the solution $\mathbf{w}^{dropout}$ that minimizes expected loss satisfies the form mentioned in the problem. $\lambda^{dropout}$ is equal to $(\frac{1}{p} - 1)$.

How does the value of p affect the regularization coefficient, $\lambda^{dropout}$? If we set p to 1 then the value of $\lambda^{dropout}$ will be zero which means $L(\mathbf{w})$ have no regularization(it turned regularization off). If p is small and goes to zero then the value of $(\frac{1}{p} - 1)$ goes to infinity and the contribution from dropout increases.

3. The loss function is $L(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda^{L_2}\|\mathbf{w}\|^2$. In order to derive its closed form solution we need to take the derivative with respect to \mathbf{w} and set it to zero:

$$\begin{aligned}
\frac{dL(\mathbf{w})}{d\mathbf{w}} &= (-2\mathbf{X}^T)(\mathbf{y} - \mathbf{X}\mathbf{w}^{L_2}) + 2\lambda^{L_2}\mathbf{w}^{L_2} = -\mathbf{X}^T\mathbf{y} + \mathbf{X}^T\mathbf{X}\mathbf{w}^{L_2} + \lambda^{L_2}\mathbf{w}^{L_2} = 0 \xrightarrow{\text{factore-}\mathbf{w}^{L_2}} \\
&\mathbf{w}^{L_2}(\mathbf{X}^T\mathbf{X} + \lambda^{L_2}) = \mathbf{X}^T\mathbf{y} \longrightarrow \mathbf{w}^{L_2} = (\mathbf{X}^T\mathbf{X} + \lambda^{L_2})^{-1}\mathbf{X}^T\mathbf{y}
\end{aligned} \tag{5}$$

4. In part 2.2, the regularization term $\lambda^{dropout}$ is scaled by Γ^2 which is the standard deviation of each feature i . If p is equal to 1, the dropout is off and we will only have least square and when p is small then the regularization term increases. In part 2.3. λ^{L_2} is scaled by I and are penalized uniformly. One's might think of dropout as scaled L^2 .

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, μ , and σ^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, μ , and σ^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 α, β 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. 1. (a) X and Y are independent then $X|Y = X, Y|X = Y$.

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

(b) Write $\mathbb{E}[a_i^{(t)}]$ and $\text{Var}(a_i^{(t)})$ in terms of $c, \mu, \sigma^2, \text{Var}(h_i^{(t-1)}), \mathbb{E}[h_i^{(t-1)}]$.

$$\begin{aligned}
\mathbb{E}[a_i^{(t)}] &= \mathbb{E}[W_i^{(t)} h^{(t-1)} + b_i^{(t)}] \xrightarrow{E[X+Y]=E[X]+E[Y]} \\
&= \mathbb{E}[W_i^{(t)} h^{(t-1)}] + \mathbb{E}[b_i^{(t)}] \xrightarrow{\text{if } (X \perp Y) \text{ then } (E[XY]=E[X]E[Y])} \\
&= \mathbb{E}[W_i^{(t)}] \mathbb{E}[h^{(t-1)}] + \mathbb{E}[b_i^{(t)}] \xrightarrow{E[h^{(t-1)}]=0, (1 \leq i \leq d^{(t-1)})} \\
&= \mathbb{E}[W_i^{(t)}] \mathbb{E}[h^{(t-1)}] + \mathbb{E}[b_i^{(t)}] \xrightarrow{(b_i^{(t)}=c) \text{ so: } (E[b_i^{(t)}]=c)} \\
&= \mathbb{E}[b_i^{(t)}] \xrightarrow{E[a_i^{(t)}]=0} \\
&= c = 0
\end{aligned} \tag{7}$$

$$\begin{aligned}
\text{Var}[a_i^{(t)}] &= \text{Var}(W_i^{(t)} h^{(t-1)} + b_i^{(t)}) \xrightarrow{\text{Var}(X+Y)=\text{Var}(x)+\text{Var}(Y)} \\
&= \text{Var}(W_i^{(t)} h^{(t-1)}) + \text{Var}(b_i^{(t)}) \xrightarrow{\text{part-a}} \\
&= \text{Var}(W_i^{(t)}) \text{Var}(h^{(t-1)}) + \text{Var}(W_i^{(t)}) \mathbb{E}[h^{(t-1)}]^2 + \mathbb{E}[W_i^{(t)}]^2 \text{Var}(h^{(t-1)}) + \text{Var}(b_i^{(t)}) \xrightarrow{\text{problem def.}} \\
&= \text{Var}(W_i^{(t)}) \text{Var}(h^{(t-1)}) + \text{Var}(W_i^{(t)}) \mathbb{E}[h^{(t-1)}]^2 + \mathbb{E}[W_i^{(t)}]^2 \text{Var}(h^{(t-1)}) + \text{Var}(b_i^{(t)}) \xrightarrow{0} \\
&= \text{Var}(W_i^{(t)}) 1_{d^{(t-1)}} + \mathbb{E}[W_i^{(t)}]^2 1_{d^{(t-1)}} \\
&= 1
\end{aligned} \tag{8}$$

We need to write the above mentioned formula in terms of $c, \mu, \sigma^2, \text{Var}(h_i^{(t-1)}), \mathbb{E}[h_i^{(t-1)}]$. We already find out that $\mathbb{E}[a_i^{(t)}] = c = 0$ and we also have:

$$\text{Var}[a_i^{(t)}] = \text{Var}(W_i^{(t)}) 1_{d^{(t-1)}} + \mathbb{E}[W_i^{(t)}]^2 1_{d^{(t-1)}}.$$

From problem definition we can derive $\mathbb{E}[W_i^{(t)}] = \mu 1_{d^{(t-1)}}^T$, $\text{Var}(W_i^{(t)}) = \sigma^2 1_{d^{(t-1)}}^T$, and from $\mathbb{E}[a_i^{(t)}] = 0$ we have $\mathbb{E}[W_i^{(t)}] = 0$ so $\mu = 0$

Now, we need to rewrite $\text{Var}[a_i^{(t)}]$ formula.

$$\begin{aligned}
\text{Var}[a_i^{(t)}] &= \text{Var}(W_i^{(t)}) 1_{d^{(t-1)}} + \mathbb{E}[W_i^{(t)}]^2 1_{d^{(t-1)}} \\
&= \sigma^2 1_{d^{(t-1)}}^T 1_{d^{(t-1)}} + \mu^2 1_{d^{(t-1)}}^T 1_{d^{(t-1)}} \\
&= (\sigma^2 + \mu^2) d^{(t-1)} \xrightarrow{\text{problem def.}} \\
&= \sigma^2 d^{(t-1)} = 1
\end{aligned} \tag{9}$$

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

(c) I take these properties into account in part b after deriving the formula.

Here is all I find out:

$$c = \mu = 0$$

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

2. (a) Derive $\mathbb{E}[(h_i^{(t-1)})^2]$:

$$\begin{aligned}
 \mathbb{E}[(h_i^{(t-1)})^2] &= \mathbb{E}[g(a_i^{(t-1)})^2] \\
 &= \int_{-\infty}^{+\infty} (\max(0, a_i^{(t-1)}))^2 f(a_i^{(t-1)}) da_i^{(t-1)} = \int_0^{+\infty} (a_i^{(t-1)})^2 f(a_i^{(t-1)}) da_i^{(t-1)} \xrightarrow{\text{symmetric}} \\
 &= \frac{1}{2} \int_{-\infty}^{+\infty} (a_i^{(t-1)})^2 f(a_i^{(t-1)}) da_i^{(t-1)} = \frac{1}{2} \mathbb{E}[(a_i^{(t-1)})^2] \\
 &= \frac{1}{2} (\text{Var}(a_i^{(t-1)}) + (E[a_i^{(t-1)}])^2) \xrightarrow{E[a_i^{(t-1)}]=0} \\
 &= \frac{1}{2} (\text{Var}(a_i^{(t-1)}) + \cancel{(E[a_i^{(t-1)}])^2}^0) = \frac{1}{2} \text{Var}(a_i^{(t-1)}) \xrightarrow{\text{Var}(a_i^{(t-1)})=1} \\
 &= \frac{1}{2} \cancel{\text{Var}(a_i^{(t-1)})}^1 = \frac{1}{2}
 \end{aligned} \tag{10}$$

(b) Now we need to find out the values for c , μ , and σ^2 as a function of $d^{(t-1)}$.

$$\begin{aligned}
 \mathbb{E}[a_i^{(t)}] &= \mathbb{E}[W_i^{(t)} h_i^{(t-1)} + b_i^{(t)}] \xrightarrow{E[X+Y]=E[X]+E[Y]} \\
 &= \mathbb{E}[W_i^{(t)} h_i^{(t-1)}] + \mathbb{E}[b_i^{(t)}] \xrightarrow{\text{if } (X \perp Y) \text{ then } (E[XY]=E[X]E[Y])} \\
 &= \mathbb{E}[W_i^{(t)}] \mathbb{E}[h_i^{(t-1)}] + \mathbb{E}[b_i^{(t)}] \xrightarrow{(b_i^{(t)}=c) \text{ so: } (E[b_i^{(t)}]=c)} \\
 &= \mathbb{E}[W_i^{(t)}] \mathbb{E}[h_i^{(t-1)}] + c = 0 \quad (*) \\
 \text{Var}(h_i^{(t-1)}) &= \frac{1}{2} - \mathbb{E}[h_i^{(t-1)}]^2 \longrightarrow \text{Var}(h^{(t-1)}) = \frac{1}{2} 1_{d^{(t-1)}} - \mathbb{E}[h^{(t-1)}]^2
 \end{aligned} \tag{11}$$

$$\begin{aligned}
 \text{Var}[a_i^{(t)}] &= \text{Var}(W_i^{(t)} h_i^{(t-1)} + b_i^{(t)}) \xrightarrow{\text{Var}(X+Y)=\text{Var}(x)+\text{Var}(Y)} \\
 &= \text{Var}(W_i^{(t)} h_i^{(t-1)}) + \text{Var}(b_i^{(t)}) \xrightarrow{\text{part-a}} \\
 &= \text{Var}(W_i^{(t)}) (\frac{1}{2} 1_{d^{(t-1)}} - \mathbb{E}[h^{(t-1)}]^2) + \text{Var}(W_i^{(t)}) \mathbb{E}[h_i^{(t-1)}]^2 \\
 &\quad + \mathbb{E}[W_i^{(t)}]^2 (\frac{1}{2} 1_{d^{(t-1)}} - \mathbb{E}[h^{(t-1)}]^2) + \text{Var}(b_i^{(t)}) \xrightarrow{\text{problem def.}} \\
 &= \frac{1}{2} \text{Var}(W_i^{(t)}) 1_{d^{(t-1)}} - \text{Var}(W_i^{(t)}) \mathbb{E}[h^{(t-1)}]^2 + \text{Var}(W_i^{(t)}) \mathbb{E}[h_i^{(t-1)}]^2 \\
 &\quad + \frac{1}{2} \mathbb{E}[W_i^{(t)}]^2 1_{d^{(t-1)}} - \mathbb{E}[W_i^{(t)}]^2 \mathbb{E}[h^{(t-1)}]^2 + \cancel{\text{Var}(b_i^{(t)})}^0 \\
 &= \frac{1}{2} \text{Var}(W_i^{(t)}) 1_{d^{(t-1)}} + \frac{1}{2} \mathbb{E}[W_i^{(t)}]^2 1_{d^{(t-1)}} - \mathbb{E}[W_i^{(t)}]^2 \mathbb{E}[h^{(t-1)}]^2 \xrightarrow{\text{from } (*)} \\
 &= \frac{1}{2} \text{Var}(W_i^{(t)}) 1_{d^{(t-1)}} + \frac{1}{2} \mathbb{E}[W_i^{(t)}]^2 1_{d^{(t-1)}} - \cancel{\mathbb{E}[W_i^{(t)}]^2} (\frac{-c}{\cancel{\mathbb{E}[W_i^{(t)}]} d^{(t-1)}})^2 d^{(t-1)} \\
 &= \frac{1}{2} \text{Var}(W_i^{(t)}) 1_{d^{(t-1)}} + \frac{1}{2} \mathbb{E}[W_i^{(t)}]^2 1_{d^{(t-1)}} - \frac{-c^2}{d^{(t-1)}} \\
 &= 1
 \end{aligned} \tag{12}$$

Lets assume that $\text{Var}(W_i^{(t)}) = \sigma^2$ and $E[W_i^{(t)}] = \mu$ since W_{ij} sampled form Gaussian

distribution.

$$\begin{aligned}
 Var[a_i^{(t)}] &= \frac{1}{2} Var(W_i^{(t)}) 1_{d^{(t-1)}} + \frac{1}{2} \mathbb{E}[W_i^{(t)}]^2 1_{d^{(t-1)}} - \frac{-c^2}{d^{(t-1)}} \\
 &= \frac{1}{2} \sigma^2 1_{d^{(t-1)}} 1_{d^{(t-1)}} + \frac{1}{2} \mu^2 1_{d^{(t-1)}} 1_{d^{(t-1)}} - \frac{-c^2}{d^{(t-1)}} \\
 &= \frac{1}{2} ((\sigma^2 + \mu^2)(d^{(t-1)})) - \frac{-c^2}{d^{(t-1)}} = 1
 \end{aligned} \tag{13}$$

One solution is when we set $c = 0$, and $\mu = 0$ which results $\sigma^2 = \frac{2}{d^{(t-1)}}$.

(c) The He initialization also have this form.

(d) Using this normalization method, the scale of the output activation is controlled and it is great because it will give us unit variance pre-activations. The ReLU turns approximately half of the neurons (the negative ones) into zeros and it help to control the scale of the activations.

3. For part 1 we know $Var(W_i^{(t)}) = \frac{(\beta-\alpha)^2}{12}$ and $E[W_i^{(t)}] = \frac{\alpha+\beta}{2}$ since W_{ij} is sampled from Uniform distribution. We also know that $\mathbf{W}_i^{(t)}$'s are independent and identically distributed random variables, $Var(\mathbf{W}_i^{(t)}) = \frac{1}{d^{(t-1)}}$.

We can write:

$$\begin{aligned}
 \mathbb{E}[a_i^t] &= c = 0 \\
 Var[a_i^{(t)}] &= Var(\mathbf{W}_i^{(t)}) + \mathbb{E}[\mathbf{W}_i^{(t)}]^2 \\
 &= \frac{(\beta - \alpha)^2}{12} d^{(t-1)} + \left(\frac{\alpha + \beta}{2}\right)^2 d^{(t-1)} \\
 &= \frac{4(\alpha^2 + \beta^2 - \alpha\beta)}{12} = \frac{(\alpha - \beta)(\alpha + \beta)}{3} d^{(t-1)} = 1
 \end{aligned} \tag{14}$$

$\alpha = -\sqrt{\frac{3}{d^{(t-1)}}}$ $\beta = \sqrt{\frac{3}{d^{(t-1)}}}$ is a possible solution for the equation above having the constraints in mind. This seems to be Glorot Uniform initialization.

For part 2 we know $Var(W_i^{(t)}) = \frac{(\beta-\alpha)^2}{12}$ and $E[W_i^{(t)}] = \frac{\alpha+\beta}{2}$ since W_{ij} is sampled from Uniform distribution.

We can write:

$$\begin{aligned}
 \mathbb{E}[a_i^t] &= c = 0 \\
 Var[a_i^{(t)}] &= \frac{1}{2} Var(\mathbf{W}_i^{(t)}) + \frac{1}{2} \mathbb{E}[\mathbf{W}_i^{(t)}]^2 \\
 &= \frac{1}{2} \frac{(\beta - \alpha)^2}{12} d^{(t-1)} + \frac{1}{2} \left(\frac{\alpha + \beta}{2}\right)^2 d^{(t-1)} \\
 &= \frac{4(\alpha^2 + \beta^2 - \alpha\beta)}{24} = \frac{(\alpha - \beta)(\alpha + \beta)}{6} d^{(t-1)} = 1
 \end{aligned} \tag{15}$$

$\alpha = -\sqrt{\frac{6}{d^{(t-1)}}}$ $\beta = \sqrt{\frac{6}{d^{(t-1)}}}$ is a possible solution for the equation above having the constraints in mind. This is He Uniform initialization.

As far as I know one can solve this problem differently, I just wrote one possible answer for the questions in hand.

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 γ 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 ¹ $\boldsymbol{W}^\perp \boldsymbol{u} = \mathbf{0}$.

Answer 4. 1. (a) **Batch Normalization:** Batch normalization normalizes activations in a network across the mini-batch. The mean and variance for this technique will be of size *number of channels* \times *number of features*. Using the formula $\mu_{batch} = \frac{1}{batch-size} \sum_i^{batch-size} \boldsymbol{x}_i$ we have:

$$\mu_{batch} = \begin{bmatrix} 11/4, 11/4, 8/4 \\ 7/4, 11/4, 13/4 \end{bmatrix} = \begin{bmatrix} 2.75, 2.75, 2 \\ 1.75, 2.75, 3.25 \end{bmatrix}$$

We can calculate the variance using:

$$\sigma_{batch}^2 = \frac{1}{batch-size} \sum_i^{batch-size} (\boldsymbol{x}_i - \mu_{batch})^2 = \begin{bmatrix} 1.1875, 1.1875, 0 \\ 0.6875, 0.6875, 0.6875 \end{bmatrix}$$

- (b) **Layer Normalization:** Layer normalization normalizes input across the features. Each of our training examples have 3 features so the mean and variance for this technique will be of size *batch size* \times *number of features*. Using the formula $\mu_{layer} = \frac{1}{feature-size} \sum_j^{feature-size} \boldsymbol{x}_{ij}$ we have:

$$\mu_{layer} = \begin{bmatrix} 1, 2.5, 2.5 \\ 2.5, 3.5, 3 \\ 2.5, 2, 3 \\ 3, 3, 2 \end{bmatrix}$$

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.

$$\sigma_{layer}^2 = \begin{bmatrix} 0, 0.25, 0.25 \\ 0.25, 0.25, 1 \\ 2.25, 0, 1 \\ 0, 0, 0 \end{bmatrix}$$

- (c) **Instance Normalization:** Instance normalization normalizes across each channel in each training example in a minibatch instead of normalizing across input features in a training example. Each of our training examples have 2 channels so the mean and variance for this technique will be of shape *batch size* \times *number of channels*. Using the formula in the slides we have:

$$\mu_{instance} = \frac{1}{HW} \sum_{l=1}^W \sum_{m=1}^H x_{tilm} = \begin{bmatrix} 2, 2 \\ 8/3, 10/3 \\ 8/3, 8/3 \\ 8/3, 8/3 \end{bmatrix} = \begin{bmatrix} 2.66, 3.33 \\ 2.66, 2.66 \\ 2.66, 2.66 \end{bmatrix}$$

We can calculate the variance using:

$$\sigma_{instance}^2 = \frac{1}{HW} \sum_{l=1}^W \sum_{m=1}^H (x_{tilm} - \mu_{ti})^2 = \begin{bmatrix} 0.66, 0.66 \\ 0.22, 0.88 \\ 0.88, 1.55 \\ 0.22, 0.22 \end{bmatrix}$$

2. We know $\mathbf{w} := \gamma \frac{\mathbf{u}}{\|\mathbf{u}\|}$ where γ is scalar parameter controlling the magnitude and \mathbf{u} is a vector controlling the direction of \mathbf{w} . We also know that $\hat{y} = \gamma \cdot \frac{y - \mu_y}{\sigma_y} + \beta$ where $y = \mathbf{u}^T \mathbf{x}$, and for vector \mathbf{x} we know that $(\text{Var}(\mathbf{x}) = \mathbf{I})$, and $(\mathbb{E}[\mathbf{x}] = \mathbf{0})$. Since \mathbf{x} is centered at 0 ($\mathbb{E}[\mathbf{x}] = \mathbf{0}$), $\mu_y = \mathbb{E}[\mathbf{y}] = \mathbb{E}[\mathbf{u}^T \mathbf{x}] = \mathbf{u}^T \mathbb{E}[\mathbf{x}] = 0$, and since \mathbf{x} is whitened ($\text{Var}(\mathbf{x}) = \mathbf{I}$) we have $\sigma_y = \text{Var}(y)^{1/2} = \text{Var}(\mathbf{u}^T \mathbf{x})^{1/2} = (\mathbf{u}^T \text{Var}(\mathbf{x}) \mathbf{u})^{1/2} = (\mathbf{u}^T \mathbf{I} \mathbf{u})^{1/2} = (\mathbf{u}^T \mathbf{u})^{1/2} = \|\mathbf{u}\|$. So we can write $\hat{y} = \gamma \cdot \frac{y - \mu_y}{\sigma_y} + \beta = \gamma \cdot \frac{\mathbf{u}^T \mathbf{x} - 0}{\|\mathbf{u}\|} + \beta = \gamma \cdot \frac{\mathbf{u}^T \mathbf{x}}{\|\mathbf{u}\|} + \beta = \mathbf{w}^T \mathbf{x} + \beta$
3. Our objective is to show that: $\nabla_{\mathbf{u}} L = s \mathbf{W}^\perp \nabla_{\mathbf{w}} L$ for some s .

$$\begin{aligned} \nabla_{\mathbf{u}} L &= \nabla_{\mathbf{u}} \mathbf{w} \cdot \nabla_{\mathbf{w}} L = \nabla_{\mathbf{u}} \left(\gamma \frac{\mathbf{u}}{\|\mathbf{u}\|} \right) \cdot \nabla_{\mathbf{w}} L \xrightarrow{\text{quotient-rule}} \\ \nabla_{\mathbf{u}} \frac{\mathbf{u}}{\|\mathbf{u}\|} &= \frac{\|\mathbf{u}\| \mathbf{I} - \mathbf{u} \frac{d}{d\mathbf{u}} (\mathbf{u}^T \mathbf{u})^{0.5}}{\|\mathbf{u}\|^2} = \frac{\|\mathbf{u}\| \mathbf{I} - \frac{1}{2} \mathbf{u} \cdot (\mathbf{u}^T \mathbf{u})^{-0.5} \frac{d(\mathbf{u}^T \mathbf{u})}{d\mathbf{u}}}{\|\mathbf{u}\|^2} \\ &= \frac{\|\mathbf{u}\| \mathbf{I} - \mathbf{u} \frac{1}{2\|\mathbf{u}\|} (2\mathbf{u}^T)}{\|\mathbf{u}\|^2} = \frac{\|\mathbf{u}\| \mathbf{I} - \frac{\mathbf{u} \mathbf{u}^T}{\|\mathbf{u}\|}}{\|\mathbf{u}\|^2} = \frac{\mathbf{I} - \frac{\mathbf{u} \mathbf{u}^T}{\|\mathbf{u}\|^2}}{\|\mathbf{u}\|} = \frac{1}{\|\mathbf{u}\|} \mathbf{W}^\perp \\ \nabla_{\mathbf{u}} L &= \frac{\gamma}{\|\mathbf{u}\|} \left(\mathbf{I} - \frac{\mathbf{u} \mathbf{u}^T}{\|\mathbf{u}\|^2} \right) \nabla_{\mathbf{w}} L = \frac{\gamma}{\|\mathbf{u}\|} \mathbf{W}^\perp \nabla_{\mathbf{w}} L = s \mathbf{W}^\perp \nabla_{\mathbf{w}} L, \quad \text{where } s = \frac{\gamma}{\|\mathbf{u}\|} \end{aligned} \tag{16}$$

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 σ 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 norm² 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. $\mathbf{g}_t = \sigma(\mathbf{W}\mathbf{g}_{t-1} + \mathbf{U}\mathbf{x}_t + \mathbf{b})$ and $\mathbf{h}_t = \mathbf{W}\sigma(\mathbf{h}_{t-1}) + \mathbf{U}\mathbf{x}_t + \mathbf{b}$ are given.

Initial point: $\mathbf{g}_0 = \sigma(\mathbf{h}_0)$ satisfies, assuming that it is also true for $\mathbf{g}_{t-1} = \sigma(\mathbf{h}_{t-1})$ then we should prove that $\mathbf{g}_t = \sigma(\mathbf{h}_t)$.

We can write:

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

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

2.

$$\frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} = \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \frac{\partial \mathbf{h}_{T-1}}{\partial \mathbf{h}_{T-2}} \dots \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \dots \frac{\partial \mathbf{h}_1}{\partial \mathbf{h}_0}$$

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} = \frac{\partial(\mathbf{W}\sigma(\mathbf{h}_{t-1}) + \mathbf{U}\mathbf{x}_t + \mathbf{b})}{\partial \mathbf{h}_{t-1}} = \mathbf{W} \frac{\partial \sigma(\mathbf{h}_{t-1})}{\partial \mathbf{h}_{t-1}} = \mathbf{W}\sigma'$$

Recall the properties of the L_2 operator norm:

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

From the equations above we have:

$$\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\| = \left\| \mathbf{W}\sigma' \right\| = \left\| \mathbf{W} \right\| \gamma = \sqrt{\lambda_1(\mathbf{W}^\top \mathbf{W})} \gamma \leq \sqrt{\frac{\delta^2}{\gamma^2}} \gamma = \left| \frac{\delta}{\gamma} \right| \cdot \gamma = \delta$$

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.

Until now, we showed that 2-norm is bounded by δ . Now we need to calculate it for the

$$\left\| \frac{\partial h_T}{\partial h_0} \right\| \leq \prod_{t=1}^T \left\| \frac{\partial h_t}{\partial h_{t-1}} \right\| \leq \delta^T$$

We know $0 \leq \delta < 1$ so if T goes to ∞ the value of δ^T will go to zero and the gradients will vanish.

3. We know that the largest eigenvalue of weights is not larger than $\frac{\delta^2}{\gamma^2}$. And because of that the gradient of hidden states cannot explode. So it is a necessary condition for gradient explosion. But it cannot always prevent explosion. One possible case where it is not sufficient is when the largest eigenvalue is orthogonal to the hidden state.

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 σ 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.

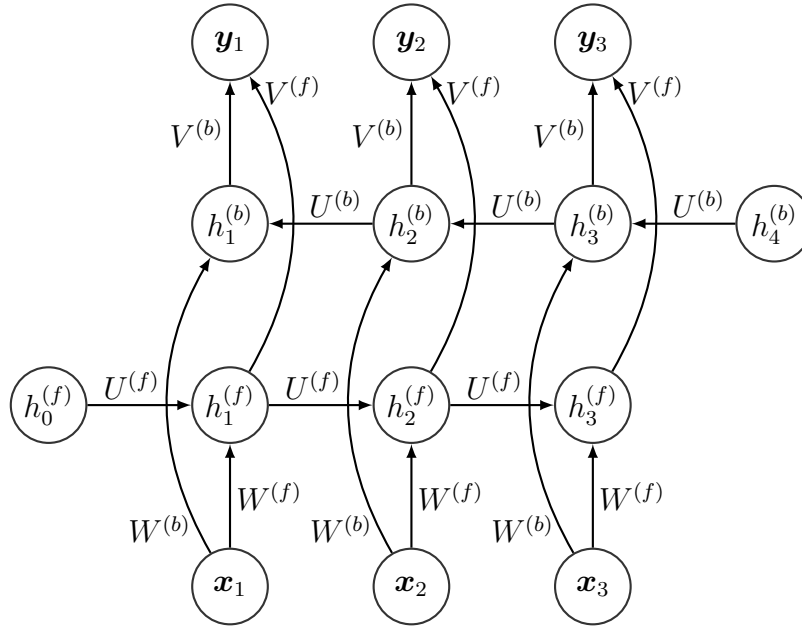


FIGURE 1 – Computational graph of the bidirectional RNN unrolled for three timesteps.

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:

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

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. 1. The graph is completed.

2. Given $L = \sum_t L_t$ where $L_t = \|\mathbf{z}_t - \mathbf{y}_t\|_2^2$ and

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

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

we need to find out the values 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 $\frac{\partial \mathbf{h}_{t-1}^{(b)}}{\partial \mathbf{h}_t^{(b)}}$. Lets start with $\nabla_{\mathbf{h}_t^{(f)}} L_t$:

- $\nabla_{\mathbf{h}_t^{(f)}} L_t = \nabla_{\mathbf{h}_t^{(f)}} \|\mathbf{z}_t - \mathbf{y}_t\|_2^2 = 2\|\mathbf{z}_t - \mathbf{y}_t\|_2 \nabla_{\mathbf{h}_t^{(f)}} (\mathbf{z}_t - (\mathbf{V}^{(f)} \mathbf{h}_t^{(f)} + \mathbf{V}^{(b)} \mathbf{h}_t^{(b)}))$
 $= 2\|\mathbf{z}_t - \mathbf{y}_t\|_2 (-\mathbf{V}^{(f)}) = -2\mathbf{V}^{(f)T} \|\mathbf{z}_t - \mathbf{y}_t\|_2 = -2\mathbf{V}^{(f)T} \cdot (\mathbf{z}_t - \mathbf{y}_t)$
- $\nabla_{\mathbf{h}_t^{(b)}} L_t = \nabla_{\mathbf{h}_t^{(b)}} \|\mathbf{z}_t - \mathbf{y}_t\|_2^2 = 2\|\mathbf{z}_t - \mathbf{y}_t\|_2 \nabla_{\mathbf{h}_t^{(b)}} (\mathbf{z}_t - (\mathbf{V}^{(f)} \mathbf{h}_t^{(f)} + \mathbf{V}^{(b)} \mathbf{h}_t^{(b)}))$
 $= 2\|\mathbf{z}_t - \mathbf{y}_t\|_2 (-\mathbf{V}^{(b)}) = -2\mathbf{V}^{(b)T} \|\mathbf{z}_t - \mathbf{y}_t\|_2 = -2\mathbf{V}^{(b)T} \cdot (\mathbf{z}_t - \mathbf{y}_t)$
- $\frac{\partial \mathbf{h}_{t+1}^{(f)}}{\partial \mathbf{h}_t^{(f)}} = \nabla_{\mathbf{h}_t^{(f)}} (\sigma(\mathbf{W}^{(f)} \mathbf{x}_{t+1} + \mathbf{U}^{(f)} \mathbf{h}_t^{(f)}))$
 $= \sigma(\mathbf{W}^{(f)} \mathbf{x}_{t+1} + \mathbf{U}^{(f)} \mathbf{h}_t^{(f)}) (1 - \sigma(\mathbf{W}^{(f)} \mathbf{x}_{t+1} + \mathbf{U}^{(f)} \mathbf{h}_t^{(f)})) \mathbf{U}^{(f)}$
 $= \mathbf{h}_{t+1}^{(f)} (1 - \mathbf{h}_{t+1}^{(f)}) \mathbf{U}^{(f)} = \text{diag}(\mathbf{h}_{t+1}^{(f)} (1 - \mathbf{h}_{t+1}^{(f)})) \mathbf{U}^{(f)}$
- $\frac{\partial \mathbf{h}_{t-1}^{(b)}}{\partial \mathbf{h}_t^{(b)}} = \nabla_{\mathbf{h}_t^{(b)}} (\sigma(\mathbf{W}^{(b)} \mathbf{x}_{t-1} + \mathbf{U}^{(b)} \mathbf{h}_t^{(b)}))$
 $= \sigma(\mathbf{W}^{(b)} \mathbf{x}_{t-1} + \mathbf{U}^{(b)} \mathbf{h}_t^{(b)}) (1 - \sigma(\mathbf{W}^{(b)} \mathbf{x}_{t-1} + \mathbf{U}^{(b)} \mathbf{h}_t^{(b)})) \mathbf{U}^{(b)}$
 $= \mathbf{h}_{t-1}^{(b)} (1 - \mathbf{h}_{t-1}^{(b)}) \mathbf{U}^{(b)} = \text{diag}(\mathbf{h}_{t-1}^{(b)} (1 - \mathbf{h}_{t-1}^{(b)})) \mathbf{U}^{(b)}$

3. Derive $\nabla_{\mathbf{W}^{(f)}} L$ and $\nabla_{\mathbf{U}^{(b)}} L$ as a functions of $\nabla_{\mathbf{h}_t^{(f)}} L$ and $\nabla_{\mathbf{h}_t^{(b)}} L$:

•

$$\nabla_{\mathbf{W}^{(f)}} L = \sum_{t=1}^{\dim} \frac{\partial \mathbf{h}_t^{(f)}}{\partial \mathbf{W}_t^{(f)}} \nabla_{\mathbf{h}_t^{(f)}} L$$

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

$$\nabla_{\mathbf{W}^{(f)}} L = \sum_{t=1}^{\dim} \nabla_{\mathbf{w}_t^{(f)}} L = \mathbf{h}_t^{(f)} (1 - \mathbf{h}_t^{(f)}) \mathbf{x}_t \nabla_{\mathbf{h}_t^{(f)}} L$$

$$\nabla_{\mathbf{W}^{(f)}} L = \sum_{t=1}^{\dim} \text{diag}(\mathbf{h}_t^{(f)} (1 - \mathbf{h}_t^{(f)})) \cdot \nabla_{\mathbf{h}_t^{(f)}} L \cdot \mathbf{x}_t^T$$

•

$$\nabla_{\mathbf{U}^{(b)}} L = \sum_{t=1}^{dim} \frac{\partial h_t^{(b)}}{\partial \mathbf{U}_t^{(b)}} \nabla_{h_t^{(f)}} L$$

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

$$\nabla_{\mathbf{U}^{(b)}} L = \sum_{t=1}^{dim} \nabla_{\mathbf{U}_t^{(b)}} L = \mathbf{h}_t^{(b)} (1 - \mathbf{h}_t^{(b)}) h_{t+1}^{(T)} \nabla_{h_t^{(b)}} L$$

$$\nabla_{\mathbf{U}^{(b)}} L = \sum_{t=1}^{dim} \text{diag}(\mathbf{h}_t^{(b)} (1 - \mathbf{h}_t^{(b)})) \cdot \nabla_{h_t^{(b)}} L \cdot h_{t+1}^T$$