

Due Date: March 22nd 23:59, 2019

Name: Aditya Joshi

Matricule:

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 David Krueger, Tegan Maharaj, and Chin-Wei Huang.**

Question 1 (6-10). 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 preactivations 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 either (a) a Gaussian distribution $\mathbf{W}_{ij}^{(t)} \sim \mathcal{N}(\mu, \sigma^2)$, or (b) a Uniform distribution $\mathbf{W}_{ij}^{(t)} \sim U(\alpha, \beta)$.

For both of the assumptions (1 and 2) about the distribution of the inputs to layer t listed below, and for both (a) Gaussian, and (b) Uniform sampling, design an initialization scheme that would achieve preactivations with zero-mean and unit variance at layer t , i.e.: $\mathbb{E}[\mathbf{a}_i^{(t)}] = 0$ and $\text{Var}(\mathbf{a}_i^{(t)}) = 1$, for $1 \leq i \leq d^{(t)}$.

(Hint: if $X \perp Y$, $\text{Var}(XY) = \text{Var}(X)\text{Var}(Y) + \text{Var}(X)\mathbb{E}[Y]^2 + \text{Var}(Y)\mathbb{E}[X]^2$)

1. Assume $\mathbb{E}[\mathbf{h}_i^{(t-1)}] = 0$ and $\text{Var}(\mathbf{h}_i^{(t-1)}) = 1$ for $1 \leq i \leq d^{(t-1)}$. Assume entries of $\mathbf{h}^{(t-1)}$ are uncorrelated (the answer should not depend on g).

(a) Gaussian: give the values for c , μ , and σ^2 as a function of $d^{(t-1)}$.

Consider a single entry in the preactivation $\mathbf{a}_i^{(t)}$. This equals

$$\begin{aligned}\mathbf{a}_i^{(t)} &= \sum_{j=1}^{d^{(t-1)}} \mathbf{W}_{ij}^{(t)} \mathbf{h}_j^{(t-1)} + \mathbf{b}_i^{(t)} \\ \mathbb{E}[\mathbf{a}_i^{(t)}] &= \mathbb{E}\left[\sum_{j=1}^{d^{(t-1)}} \mathbf{W}_{ij}^{(t)} \mathbf{h}_j^{(t-1)} + \mathbf{b}_i^{(t)}\right] \\ \mathbb{E}[\mathbf{a}_i^{(t)}] &= \sum_{j=1}^{d^{(t-1)}} \mathbb{E}[\mathbf{W}_{ij}^{(t)} \mathbf{h}_j^{(t-1)}] + \mathbb{E}[\mathbf{b}_i^{(t)}] \\ 0 &= \sum_{j=1}^{d^{(t-1)}} \mathbb{E}[\mathbf{W}_{ij}^{(t)}] \cdot \mathbb{E}[\mathbf{h}_j^{(t-1)}] + \mathbb{E}[\mathbf{b}_i^{(t)}] \\ 0 &= \mu \cdot 0 + c \\ c &= 0\end{aligned}$$

Therefore we have $c = 0$ and value of μ does not matter to achieve $\mathbb{E}[\mathbf{a}_i^{(t)}] = 0$.

$$\begin{aligned}\mathbf{a}_i^{(t)} &= \sum_{j=1}^{d^{(t-1)}} \mathbf{W}_{ij}^{(t)} \mathbf{h}_j^{(t-1)} + \mathbf{b}_i^{(t)} \\ \text{Var}[\mathbf{a}_i^{(t)}] &= \text{Var}\left[\sum_{j=1}^{d^{(t-1)}} \mathbf{W}_{ij}^{(t)} \mathbf{h}_j^{(t-1)} + \mathbf{b}_i^{(t)}\right] \\ 1 &= \text{Var}\left[\sum_{j=1}^{d^{(t-1)}} \mathbf{W}_{ij}^{(t)} \mathbf{h}_j^{(t-1)}\right] + \text{Var}[\mathbf{b}_i^{(t)}] \\ &= \sum_{j=1}^{d^{(t-1)}} \text{Var}[\mathbf{W}_{ij}^{(t)} \mathbf{h}_j^{(t-1)}] + \text{Var}[\mathbf{b}_i^{(t)}] \\ &= \sum_{j=1}^{d^{(t-1)}} \text{Var}[\mathbf{W}_{ij}^{(t)}] \cdot \text{Var}[\mathbf{h}_j^{(t-1)}] + \text{Var}[\mathbf{W}_{ij}^{(t)}] \cdot (\mathbb{E}[\mathbf{h}_j^{(t-1)}])^2 + \text{Var}[\mathbf{h}_j^{(t-1)}] \cdot (\mathbb{E}[\mathbf{W}_{ij}^{(t)}])^2 + 0 \\ &= \sum_{j=1}^{d^{(t-1)}} \sigma^2 \cdot 1 + \sigma^2 \cdot 0^2 + 1 \cdot \mu^2 \\ &= \sum_{j=1}^{d^{(t-1)}} \sigma^2 + \mu^2 \\ 1 &= d^{(t-1)}(\sigma^2 + \mu^2)\end{aligned}$$

Thus, if we set $\mu = 0$ we find that the variance must equal $\sigma^2 = \frac{1}{d^{(t-1)}}$. The initialization

scheme is then $\boxed{c = 0, \mu = 0, \sigma^2 = \frac{1}{d^{(t-1)}}}$

(b) Uniform: give the values for c , α , and β as a function of $d^{(t-1)}$.

We consider the same equation as part (a).

$$\begin{aligned}
 \mathbf{a}_i^{(t)} &= \sum_{j=1}^{d^{(t-1)}} \mathbf{W}_{ij}^{(t)} \mathbf{h}_j^{(t-1)} + \mathbf{b}_i^{(t)} \\
 0 &= \sum_{j=1}^{d^{(t-1)}} \mathbb{E}[\mathbf{W}_{ij}^{(t)}] \cdot \mathbb{E}[\mathbf{h}_j^{(t-1)}] + \mathbb{E}[\mathbf{b}_i^{(t)}] \\
 &= \sum_{j=1}^{d^{(t-1)}} \left(\frac{\alpha + \beta}{2} \right) \cdot 0 + c \\
 c &= 0
 \end{aligned}$$

For the variance, we have

$$\text{Var}[\mathbf{a}_i^{(t)}] = \sum_{j=1}^{d^{(t-1)}} \text{Var}[\mathbf{W}_{ij}^{(t)}] \cdot \text{Var}[\mathbf{h}_j^{(t-1)}] + \text{Var}[\mathbf{W}_{ij}^{(t)}] \cdot (\mathbb{E}[\mathbf{h}_j^{(t-1)}])^2 + \text{Var}[\mathbf{h}_j^{(t-1)}] \cdot (\mathbb{E}[\mathbf{W}_{ij}^{(t)}])^2 + 0 \quad (1)$$

$$= \sum_{j=1}^{d^{(t-1)}} \left(\frac{(\beta - \alpha)^2}{12} \cdot 1 + 0 + \frac{\alpha + \beta}{2} \right) \quad (2)$$

$$1 = d^{(t-1)} \cdot \frac{(\beta - \alpha)^2}{12} + \frac{\alpha + \beta}{2} \quad (3)$$

If we have $\beta = \text{constant}$ and $\beta = -\text{constant}$, then $\frac{\alpha + \beta}{2}$ will equal 0. Then we just have to find the appropriate value of the constant. With some experimentation, we realize that if we set $\beta = \sqrt{\frac{3}{d^{(t-1)}}}$ and $\alpha = -\sqrt{\frac{3}{d^{(t-1)}}}$, we have:

$$1 = d^{(t-1)} \cdot \frac{\left(\frac{\sqrt{3}}{\sqrt{d^{(t-1)}}} + \frac{\sqrt{3}}{\sqrt{d^{(t-1)}}} \right)^2}{12} \quad (4)$$

$$(5)$$

$$1 = d^{(t-1)} \cdot \frac{\left(\frac{2\sqrt{3}}{\sqrt{d^{(t-1)}}} \right)^2}{12} \quad (6)$$

$$(7)$$

$$1 = \cancel{d^{(t-1)}} \cdot \frac{\cancel{12}}{\cancel{d^{(t-1)}} \cdot \cancel{12}} \quad (8)$$

Thus, we can choose the following values $c = 0, \alpha = -\sqrt{\frac{3}{d^{(t-1)}}}, \beta = \sqrt{\frac{3}{d^{(t-1)}}}$

2. Assume that the preactivations of the previous layer satisfy $\mathbb{E}[\mathbf{a}_i^{(t-1)}] = 0$, $\text{Var}(\mathbf{a}_i^{(t-1)}) = 1$ and $\mathbf{a}_i^{(t-1)}$ 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: $g(x) = \max\{0, x\}$.

(a) Gaussian: give the values for c , μ , and σ^2 as a function of $d^{(t-1)}$.

$$\mathbf{a}_i^{(t)} = \sum_{j=1}^{d^{(t-1)}} \mathbf{W}_{ij}^{(t)} g(\mathbf{a}_j^{(t-1)}) + \mathbf{b}_i^{(t)} \quad (9)$$

$$\mathbb{E}[\mathbf{a}_i^{(t)}] = \sum_{j=1}^{d^{(t-1)}} \mathbb{E}[\mathbf{W}_{ij}^{(t)}] \mathbb{E}[g(\mathbf{a}_j^{(t-1)})] + \mathbb{E}[\mathbf{b}_i^{(t)}] \quad (10)$$

We want to ensure that $\mathbb{E}[\mathbf{a}_i^{(t)}] = 0$. Note that $\mathbb{E}[g(\mathbf{a}_j^{(t-1)})]$ may not equal zero. If we sample $\mathbf{W}_{ij}^{(t)}$ from a gaussian with zero mean, the first product will be 0 as $\mathbb{E}[\mathbf{W}_{ij}^{(t)}] = 0$. We can also initialize $\mathbf{b}_i^{(t)} = c = 0$.

Let us now compute the variance. (Note that $\text{Var}[XY] = \mathbb{E}[X^2Y^2] - (\mathbb{E}[XY])^2$).

In the following steps we also assume that X and Y are independent, and thus, $\mathbb{E}[XY] = \mathbb{E}[X] \cdot \mathbb{E}[Y]$.

$$\text{Var}[\mathbf{a}_i^{(t)}] = \sum_{j=1}^{d^{(t-1)}} \text{Var}[\mathbf{W}_{ij}^{(t)} g(\mathbf{a}_j^{(t-1)})] + \text{Var}[\mathbf{b}_i^{(t)}] \quad (11)$$

$$1 = \sum_{j=1}^{d^{(t-1)}} \frac{\text{Var}[\mathbf{W}_{ij}^{(t)}] + \overbrace{(\mathbb{E}[(\mathbf{W}_{ij}^{(t)})])^2}^0}{\mathbb{E}[(\mathbf{W}_{ij}^{(t)})^2]} \mathbb{E}[g(\mathbf{a}_j^{(t-1)})^2] - \underbrace{(\mathbb{E}[\mathbf{W}_{ij}^{(t)}] \cdot \mathbb{E}[g(\mathbf{a}_j^{(t-1)})])^2}_0 \quad (12)$$

$$= d^{(t-1)} \cdot \underbrace{\text{Var}[\mathbf{W}_{ij}^{(t)}]}_{\sigma^2} \cdot \mathbb{E}[g(\mathbf{a}_j^{(t-1)})^2] \quad (13)$$

If we compute the value of $\mathbb{E}[g(\mathbf{a}_j^{(t-1)})^2]$, we can find an appropriate value for σ^2 .

We can compute $\mathbb{E}[g(\mathbf{a}_j^{(t-1)})^2]$ as follows:

$$\mathbb{E}[g(\mathbf{a}_j^{(t-1)})^2] = \int_{-\infty}^{\infty} g(\mathbf{a}_j^{(t-1)})^2 f(\mathbf{a}_j^{(t-1)}) d\mathbf{a} \quad (14)$$

$$= \int_{-\infty}^0 g(\mathbf{a}_j^{(t-1)})^2 f(\mathbf{a}_j^{(t-1)}) d\mathbf{a} + \int_0^{\infty} g(\mathbf{a}_j^{(t-1)})^2 f(\mathbf{a}_j^{(t-1)}) d\mathbf{a} \quad (15)$$

$$= \int_0^{\infty} (\mathbf{a}_j^{(t-1)})^2 f(\mathbf{a}_j^{(t-1)}) d\mathbf{a} = \frac{1}{2} \underbrace{\text{Var}[\mathbf{a}_j^{(t-1)}]}_1 \quad (16)$$

Therefore we have:

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

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

Thus we can choose the values $c = 0, \mu = 0, \sigma^2 = \frac{2}{d^{(t-1)}}$

- (b) Uniform: give the values for c , α , and β as a function of $d^{(t-1)}$.

Using the same reasoning as in 10, we can have $c = 0$. The values of α and β should be $-\text{constant}$ and constant so that their mean is 0. We find the exact values using the same reasoning we had in 13.

$$\text{Var}[\mathbf{a}_i^{(t)}] = d^{(t-1)} \cdot \underbrace{\text{Var}[\mathbf{W}_{ij}^{(t)}]}_{\sigma^2} \cdot \underbrace{\mathbb{E}[g(\mathbf{a}_j^{(t-1)})^2]}_{\frac{1}{2} \text{Var}[\mathbf{a}_j^{(t-1)}]} \quad (17)$$

$$1 = d^{(t-1)} \cdot \frac{(\beta - \alpha)^2}{12} \cdot \frac{1}{2} \quad (18)$$

After some experimentation like in equation 3, we can set the value of $\alpha = -\sqrt{\frac{6}{d^{(t-1)}}}$ and

$$\beta = \sqrt{\frac{6}{d^{(t-1)}}}.$$

Thus we have $\boxed{c = 0, \alpha = -\sqrt{\frac{6}{d^{(t-1)}}}, \beta = \sqrt{\frac{6}{d^{(t-1)}}}}$

- (c) What popular initialization scheme has this form?

The **Glorot** [1] and **Kaiming He** [2] initialization schemes have this form.

- (d) Why do you think this initialization would work well in practice? Answer in 1-2 sentences. In the case of using sigmoid or tanh activation function (as was originally analyzed in *Understanding the difficulty of training deep feedforward neural networks* [1]), the gradients of the above mentioned activation functions saturate very quickly, that is, the gradients are very close to zero beyond a certain range. This prevents any gradients from flowing backwards thus preventing any learning. Maintaining the activations around 0 with unit variance helps the flow of a gradient.

In general, we can say that maintaining the values around a certain range (close to 0) helps to keep the values (and their gradients) from exploding or vanishing.

Question 2 (4-6-4-4-3). 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.

1. For squared error loss, express the loss function $L(\mathbf{w})$ in matrix form (in terms of \mathbf{X} , \mathbf{y} , \mathbf{w} , and \mathbf{R}).

The dropout mask \mathbf{R} is sampled from a Bernoulli distribution. It will have an entry of 1 if the element is to be kept (with probability p) and 0, if we drop it (with probability $(1 - p)$). Thus, in matrix form, the resultant input data after applying the dropout is simply the element wise (*hadamard*) product of the two matrices ($\tilde{\mathbf{X}} = \mathbf{X} \odot \mathbf{R}$). Thus, the final equation for the loss is as follows.

$$L(\mathbf{w}) = \|\mathbf{y} - \tilde{\mathbf{X}}\mathbf{w}\|^2 \quad (19)$$

$$= (\mathbf{y} - \tilde{\mathbf{X}}\mathbf{w})^T (\mathbf{y} - \tilde{\mathbf{X}}\mathbf{w}) \quad (20)$$

2. Let Γ be a diagonal matrix with $\Gamma_{ii} = (\mathbf{X}^T \mathbf{X})_{ii}^{1/2}$. Show that the *expectation (over \mathbf{R})* of the loss function can be rewritten as $L(\mathbf{w}) = \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1 - p)\|\Gamma\mathbf{w}\|^2$.

Consider the loss function from equation 20,

$$\begin{aligned} L(\mathbf{w}) &= (\mathbf{y} - \tilde{\mathbf{X}}\mathbf{w})^T (\mathbf{y} - \tilde{\mathbf{X}}\mathbf{w}) \\ &= (\mathbf{y}^T - \mathbf{w}^T \tilde{\mathbf{X}}^T) (\mathbf{y} - \tilde{\mathbf{X}}\mathbf{w}) \\ &= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \tilde{\mathbf{X}}\mathbf{w} - \mathbf{w}^T \tilde{\mathbf{X}}^T \mathbf{y} + \mathbf{w}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}}\mathbf{w} \end{aligned}$$

Note that $\mathbf{y}^T \tilde{\mathbf{X}}\mathbf{w}$ is a scalar, and for a scalar $a^T = a$. Thus we have $(\mathbf{y}^T \tilde{\mathbf{X}}\mathbf{w})^T = \mathbf{w}^T \tilde{\mathbf{X}}^T \mathbf{y}$. Substituting this in the equation above we get (using properties of linearity of expectation):

$$L(\mathbf{w}) = \mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \tilde{\mathbf{X}}^T \mathbf{y} + \mathbf{w}^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}}\mathbf{w} \quad (21)$$

$$\mathbb{E}_{\mathbf{R}}[L(\mathbf{w})] = \mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbb{E}[\tilde{\mathbf{X}}^T] \mathbf{y} + \mathbf{w}^T \mathbb{E}[\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}] \mathbf{w} \quad (22)$$

$$= \mathbf{y}^T \mathbf{y} - 2p\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbb{E}[\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}] \mathbf{w} \quad (23)$$

Let us calculate the value of $\mathbb{E}[\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}]$. $\tilde{\mathbf{X}}^T$ and $\tilde{\mathbf{X}}$ are not independent. Thus, we can use the property that $\mathbb{E}[X \cdot Y] = \mathbb{E}[X] \cdot \mathbb{E}[Y] + \text{Cov}(X, Y)$.

$$\mathbb{E}[\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}] = \mathbb{E}[\tilde{\mathbf{X}}^T] \mathbb{E}[\tilde{\mathbf{X}}] + \text{Cov}(\tilde{\mathbf{X}}^T, \tilde{\mathbf{X}}) \quad (24)$$

$$= p^2 \mathbf{X}^T \mathbf{X} + p(1 - p)\Gamma \quad (25)$$

Substituting 25 in 23,

$$= \mathbf{y}^T \mathbf{y} - 2p\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T (p^2 \mathbf{X}^T \mathbf{X} + p(1 - p)\Gamma) \mathbf{w} \quad (26)$$

$$= \mathbf{y}^T \mathbf{y} - 2p\mathbf{w}^T \mathbf{X}^T \mathbf{y} + p^2 \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} + \mathbf{w}^T p(1 - p)\Gamma \mathbf{w} \quad (27)$$

$$= \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1 - p)\|\Gamma\mathbf{w}\|^2 \quad (28)$$

3. Show that the solution $\mathbf{w}^{\text{dropout}}$ that minimizes the expected loss from question 2.2 satisfies

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

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

Consider the equation in 27. Let us differentiate that with respect to \mathbf{w}^T .

$$\begin{aligned} \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}^T} &= -2p\mathbf{X}^T \mathbf{y} + 2p^2 \mathbf{X}^T \mathbf{X} \mathbf{w} + 2 \cdot \frac{p^2(1-p)}{p} \Gamma^2 \mathbf{w} \\ 0 &= -\mathbf{X}^T \mathbf{y} + p\mathbf{X}^T \mathbf{X} \mathbf{w} + \frac{p \cdot (1-p)}{p} \Gamma^2 \mathbf{w} \\ \mathbf{X}^T \mathbf{y} &= p\mathbf{w} \left(\mathbf{X}^T \mathbf{X} + \frac{(1-p)}{p} \Gamma^2 \right) \\ p\mathbf{w} &= \left(\mathbf{X}^T \mathbf{X} + \frac{(1-p)}{p} \Gamma^2 \right)^{-1} \mathbf{X}^T \mathbf{y} \\ p\mathbf{w}^{\text{dropout}} &= \left(\mathbf{X}^T \mathbf{X} + \lambda^{\text{dropout}} \Gamma^2 \right)^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

Effect of p : Here λ^{dropout} is a regularization coefficient that depends on p . For values of p close to 1, the Γ^2 term disappears, and there is no regularization term. That is, as $p \rightarrow 1$, $\lambda^{\text{dropout}} \rightarrow 0$. As $p \rightarrow 0$, λ^{dropout} gets arbitrarily large ($\rightarrow \infty$) and the regularization increases.

4. Express the solution \mathbf{w}^{L_2} for a linear regression problem without dropout and with L^2 regularization, with regularization coefficient λ^{L_2} in closed form.

Linear regression with L_2 regularization is often called *ridge regression*. In this form of regression, we add an additional term $\|\mathbf{w}\|^2$ to the cost function to penalize the weights from growing too large. The closed form solution for this is as follows:

$$\begin{aligned} L(\mathbf{w}) &= (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda^{L_2} \|\mathbf{w}\|^2 \\ &= \mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} + \lambda^{L_2} \|\mathbf{w}\|^2 \\ \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} &= -2\mathbf{X}^T \mathbf{y} + \mathbf{X}^T \mathbf{X} \mathbf{w} + 2\lambda^{L_2} \mathbf{w} \\ \mathbf{X}^T \mathbf{y} &= \mathbf{X}^T \mathbf{X} \mathbf{w} + \lambda^{L_2} \mathbf{w} \\ \mathbf{w}^{L_2} &= (\mathbf{X}^T \mathbf{X} + \lambda^{L_2} \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

5. 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).

We see that the two equations are identical except that in the dropout case, we have $\lambda^{\text{dropout}} \Gamma$ and for L_2 regularization we have $\lambda^{L_2} \mathbf{I}$. One difference that we observe is that the L_2 regularization is independent of the input features \mathbf{X} . In the dropout case, it is a function of the input through the existence of Γ^2 .

Question 3 (5-5-5). 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) .

- SGD with momentum:

$$\Delta\boldsymbol{\theta}_t = -\mathbf{v}_t \quad \mathbf{v}_t = -\Delta\boldsymbol{\theta}_t \quad (29)$$

$$\Delta\boldsymbol{\theta}_{t-1} = -\mathbf{v}_{t-1} \quad \mathbf{v}_{t-1} = -\Delta\boldsymbol{\theta}_{t-1} \quad (30)$$

Now let us write \mathbf{v}_t as a function of $\Delta\boldsymbol{\theta}$.

$$\mathbf{v}_t = \alpha\mathbf{v}_{t-1} + \epsilon\mathbf{g}_t \quad (31)$$

$$-\Delta\boldsymbol{\theta}_t = -\alpha\Delta\boldsymbol{\theta}_{t-1} + \epsilon\mathbf{g}_t \quad (32)$$

$$\Delta\boldsymbol{\theta}_t = \alpha\Delta\boldsymbol{\theta}_{t-1} - \epsilon\mathbf{g}_t \quad (33)$$

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

$$\Delta\boldsymbol{\theta}_t = -\delta\mathbf{v}_t \quad \mathbf{v}_t = \frac{-\Delta\boldsymbol{\theta}_t}{\delta} \quad (34)$$

$$\Delta\boldsymbol{\theta}_{t-1} = -\delta\mathbf{v}_{t-1} \quad \mathbf{v}_{t-1} = \frac{-\Delta\boldsymbol{\theta}_{t-1}}{\delta} \quad (35)$$

Now let us write \mathbf{v}_t as a function of $\Delta\boldsymbol{\theta}$.

$$\mathbf{v}_t = \beta\mathbf{v}_{t-1} + (1 - \beta)\mathbf{g}_t \quad (36)$$

$$\frac{-\Delta\boldsymbol{\theta}_t}{\delta} = \beta\frac{-\Delta\boldsymbol{\theta}_{t-1}}{\delta} + (1 - \beta)\mathbf{g}_t \quad (37)$$

$$\Delta\boldsymbol{\theta}_t = \beta\Delta\boldsymbol{\theta}_{t-1} - \delta \cdot (1 - \beta)\mathbf{g}_t \quad (38)$$

We can see the equivalence in equations 33 and 38 with the relations between hyperparameters as follows:

$$\begin{aligned} \alpha &= \beta \\ \epsilon &= \delta \cdot (1 - \beta) \end{aligned}$$

2. Unroll the running average update rule, i.e. express \mathbf{v}_t as a linear combination of \mathbf{g}_i 's ($1 \leq i \leq t$).

The running average update is given as:

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) \mathbf{g}_t \quad (39)$$

$$\mathbf{v}_{t-1} = \beta \mathbf{v}_{t-2} + (1 - \beta) \mathbf{g}_{t-1} \quad (40)$$

$$\vdots \quad (41)$$

$$\mathbf{v}_1 = (1 - \beta) \mathbf{g}_1 \quad (42)$$

We can unroll the steps as follows:

$$\mathbf{v}_t = \beta^2 \mathbf{v}_{t-2} + \beta \cdot (1 - \beta) \mathbf{g}_{t-1} + (1 - \beta) \mathbf{g}_t$$

$$\mathbf{v}_t = \beta^2 \left(\beta \mathbf{v}_{t-3} + (1 - \beta) \mathbf{g}_{t-2} \right) + \beta \cdot (1 - \beta) \mathbf{g}_{t-1} + (1 - \beta) \mathbf{g}_t$$

$$\mathbf{v}_t = \beta^3 \mathbf{v}_{t-3} + \beta^2 \cdot (1 - \beta) \mathbf{g}_{t-2} + \beta \cdot (1 - \beta) \mathbf{g}_{t-1} + (1 - \beta) \mathbf{g}_t$$

$$\vdots$$

$$\mathbf{v}_t = (1 - \beta) \sum_{i=0}^{t-1} \beta^i \mathbf{g}_{t-i}$$

3. 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 .

From the above equation we have that:

$$\begin{aligned} \mathbf{v}_t &= (1 - \beta) \sum_{i=0}^{t-1} \beta^i \mathbf{g}_{t-i} \\ \mathbb{E}[\mathbf{v}_t] &= (1 - \beta) \sum_{i=0}^{t-1} \underbrace{\mathbb{E}[\mathbf{g}_{t-i}]}_{\text{stationary}} \beta^i \\ &= (1 - \beta) \cdot \mathbb{E}[\mathbf{g}_t] \cdot \underbrace{\sum_{i=0}^{t-1} \beta^i}_{\text{geometric series}} \\ &= (1 - \beta) \cdot \mathbb{E}[\mathbf{g}_t] \cdot \frac{1 - \beta^t}{1 - \beta} \\ &= \mathbb{E}[\mathbf{g}_t] \cdot (1 - \beta^t) \end{aligned}$$

Thus $\mathbb{E}[\mathbf{v}_t] \neq \mathbb{E}[\mathbf{g}_t]$ and we can eliminate this bias by rescaling \mathbf{v}_t by a value of $(1 - \beta^t)$. That is,

$$\mathbf{v}_t^{\text{new}} = \frac{\mathbf{v}_t^{\text{old}}}{(1 - \beta^t)}$$

Question 4 (5-5-5). This question is about weight normalization. We consider the following parameterization of a weight vector \mathbf{w} :

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

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

1. 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 $\hat{y} = \mathbf{w}^\top \mathbf{x} + \beta$.

The equation $\hat{y} = \gamma \cdot \frac{y - \mu_y}{\sigma_y} + \beta$ can be written as follows:

$$\begin{aligned} \hat{y} &= \gamma \cdot \frac{y - \mathbb{E}[y]}{\sqrt{\text{Var}[y]}} + \beta \\ &= \gamma \cdot \frac{y - \mathbb{E}[\mathbf{u}^\top \mathbf{x}]}{\sqrt{\text{Var}[\mathbf{u}^\top \mathbf{x}]}} + \beta \\ &= \gamma \cdot \frac{y - \mathbf{u}^\top \overbrace{\mathbb{E}[\mathbf{x}]^0}}{\sqrt{\|\mathbf{u}\|^2 \underbrace{\text{Var}[\mathbf{x}]_I}}} + \beta \\ &= \gamma \cdot \frac{\mathbf{u}^\top \mathbf{x}}{\sqrt{\|\mathbf{u}\|^2}} + \beta \\ &= \gamma \cdot \underbrace{\frac{\mathbf{u}}{\|\mathbf{u}\|}}_{\mathbf{w}}^\top \mathbf{x} + \beta \\ \hat{y} &= \mathbf{w}^\top \mathbf{x} + \beta \end{aligned}$$

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

Using the derivative chain rule, we can express the gradient $\nabla_{\mathbf{u}} L$ as follows:

$$\nabla_{\mathbf{u}} L = \nabla_{\mathbf{u}} \mathbf{w} \nabla_{\mathbf{w}} L \quad (43)$$

Let us find the value of $\nabla_{\mathbf{u}} \mathbf{w}$ by differentiating each component of the vector valued function.

$$\underline{i = j} \quad \underline{i \neq j} \quad (44)$$

$$(45)$$

$$\frac{\partial \mathbf{w}_i}{\partial \mathbf{u}_i} = \gamma \cdot \frac{\partial}{\partial \mathbf{u}_i} \left(\frac{\mathbf{u}_i}{\|\mathbf{u}\|^2} \right) \quad \frac{\partial \mathbf{w}_i}{\partial \mathbf{u}_j} = \gamma \cdot \frac{\partial}{\partial \mathbf{u}_j} \left(\frac{\mathbf{u}_j}{\|\mathbf{u}\|^2} \right) \quad (46)$$

$$= \frac{\gamma}{\|\mathbf{u}\|} \cdot \frac{1 - \mathbf{u}_i^2}{\|\mathbf{u}\|^2} \quad = \frac{\gamma}{\|\mathbf{u}\|} \cdot \frac{0 - \mathbf{u}_i \mathbf{u}_j}{\|\mathbf{u}\|^2} \quad (47)$$

$$(48)$$

$$\nabla_{\mathbf{u}} L = \overbrace{\frac{\gamma}{\|\mathbf{u}\|}}^{\text{scalar}} \cdot \left(\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2} \right) \quad (49)$$

$$= s \mathbf{W}^\perp \quad (50)$$

Substituting 50 in 43 we get:

$$\nabla_{\mathbf{u}} L = s \mathbf{W}^\perp \nabla_{\mathbf{w}} L$$

$$\text{where } s = \frac{\gamma}{\|\mathbf{u}\|} \text{ and } \mathbf{W}^\perp = \left(\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2} \right)$$

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.

3. Figure 1 shows the norm of \mathbf{u} as a function of number of updates made to a two-layer MLP using gradient descent. Different curves correspond to models trained with different log-learning rate. Explain why (1) the norm is increasing, and (2) why larger learning rate corresponds to faster growth.

(Hint: Use the Pythagorean theorem and the fact that $\mathbf{W}^\perp \mathbf{u} = 0$ from question 4.2).

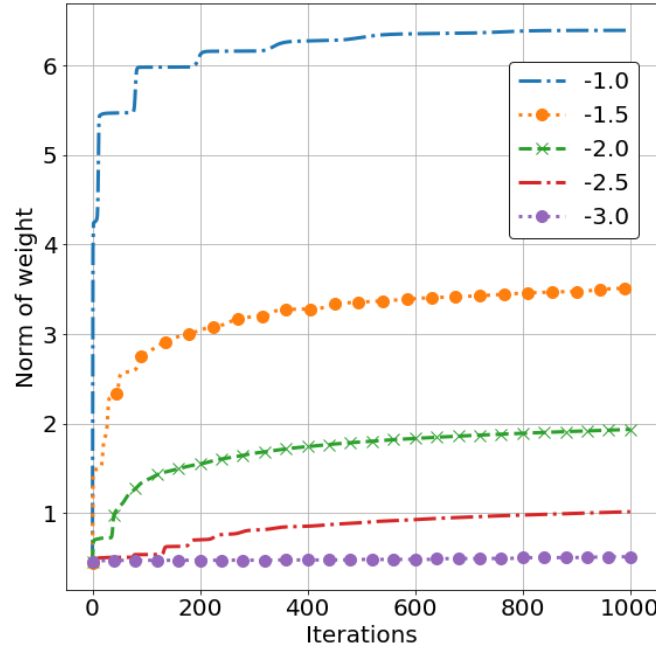


FIGURE 1 – Norm of parameters with different learning rate.

For a gradient descent update we have that:

$$\begin{aligned}\mathbf{u}_{t+1} &= \mathbf{u}_t - \eta \cdot \nabla_{\mathbf{u}} L \\ &= \mathbf{u}_t - \eta \cdot s \mathbf{W}^\perp \nabla_{\mathbf{w}} L\end{aligned}$$

The matrix \mathbf{W}^\perp projects the gradient of \mathbf{w} and thus gradient of \mathbf{u} in a direction orthogonal to \mathbf{u} . Thus, by pythagoras theorem, we have that:

$$\|\mathbf{u}_{i+1}\| > \|\mathbf{u}_t - \eta \cdot s \mathbf{W}^\perp \nabla_{\mathbf{w}} L\|$$

Thus, the updated vector will have a norm that is greater than the initial value. This explains why the norm increases.

Secondly, the norm of the vector \mathbf{u}_{t+1} grows proportional to the learning rate η . This implies that for larger values of η , for the inequality to hold, the norm of the updated vector grows faster. This explains why the larger learning rate corresponds to faster growth.

Question 5 (5-5-5). 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).

Let us assume that $\sigma(\mathbf{h}_0)$ and \mathbf{g}_0 are initialized to the $\mathbf{0}$ vector. Thus we have that:

$$\begin{aligned}\mathbf{g}_t &= \sigma(\mathbf{W}\mathbf{g}_{t-1} + \mathbf{U}\mathbf{x}_t + \mathbf{b}) \\ \mathbf{g}_1 &= \sigma(\mathbf{W}\mathbf{g}_0 + \underbrace{\mathbf{U}\mathbf{x}_1 + \mathbf{b}}_{\mathbf{h}_1 - \mathbf{W}\sigma(\mathbf{h}_0)}) \\ \mathbf{g}_1 &= \sigma(\underbrace{\mathbf{W}\mathbf{g}_0}_0 + \mathbf{h}_1 - \underbrace{\mathbf{W}\sigma(\mathbf{h}_0)}_0) \\ \mathbf{g}_1 &= \sigma(\mathbf{h}_1)\end{aligned}$$

For the general induction case, let us assume that $\mathbf{g}_t = \sigma(\mathbf{h}_t)$. Now we have:

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

Thus, we have the base case that $\mathbf{g}_0 = \sigma(\mathbf{h}_0)$ and have proved the induction step above. Therefore in general we can say that $\boxed{\mathbf{g}_t = \sigma(\mathbf{h}_t)}$

- *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. $|\sigma'(x)| \leq \gamma$ 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 upper-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

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

By the multivariate chain rule, we have that:

$$\frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} = \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \cdot \frac{\partial \mathbf{h}_{T-1}}{\partial \mathbf{h}_{T-2}} \cdot \frac{\partial \mathbf{h}_{T-2}}{\partial \mathbf{h}_{T-3}} \cdot \dots \cdot \frac{\partial \mathbf{h}_1}{\partial \mathbf{h}_0} = \prod_{i=0}^{T-1} \frac{\partial \mathbf{h}_{i+1}}{\partial \mathbf{h}_i} \quad (51)$$

Let us consider one term in the product $\frac{\partial \mathbf{h}_{i+1}}{\partial \mathbf{h}_i}$. This is equal to:

$$\frac{\partial \mathbf{h}_{i+1}}{\partial \mathbf{h}_i} = \frac{\partial \mathbf{h}_{i+1}}{\partial \mathbf{a}_{i+1}} \cdot \frac{\partial \mathbf{a}_{i+1}}{\partial \mathbf{h}_i}$$

Now we have that:

$$\begin{aligned} \mathbf{h}_{i+1} &= \sigma(\mathbf{a}_{i+1}) & \therefore \frac{\partial \mathbf{h}_{i+1}}{\partial \mathbf{a}_{i+1}} &= \sigma'(\mathbf{a}_{i+1}) \\ \mathbf{a}_{i+1} &= \mathbf{W}\mathbf{h}_i + \mathbf{U}\mathbf{x}_{i+1} + \mathbf{b} & \therefore \frac{\partial \mathbf{a}_{i+1}}{\partial \mathbf{h}_i} &= \mathbf{W} \end{aligned}$$

Thus we have:

$$\frac{\partial \mathbf{h}_{i+1}}{\partial \mathbf{h}_i} = \sigma'(\mathbf{a}_{i+1}) \mathbf{W} \quad (52)$$

$$\left\| \frac{\partial \mathbf{h}_{i+1}}{\partial \mathbf{h}_i} \right\| = \left\| \sigma'(\mathbf{a}_{i+1}) \mathbf{W} \right\| \quad (53)$$

$$\leq \left\| \sigma'(\mathbf{a}_{i+1}) \right\| \cdot \left\| \mathbf{W} \right\| \quad (54)$$

$$\leq c \cdot \gamma \cdot \sqrt{\lambda_1 \mathbf{W}^\top \mathbf{W}} \quad (55)$$

$$\leq c \cdot \gamma \cdot \frac{\delta}{\gamma} = c \cdot \delta \quad (56)$$

Substituting 56 in 51 we have:

$$\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| = \prod_{i=0}^{T-1} c \cdot \delta = c \cdot \delta^{T-1}$$

Thus we see that as $T \rightarrow \infty$, the term δ^{T-1} where $0 \leq \delta < 1$ will go towards zero, and as a consequence, the gradients of the hidden state will vanish over time.

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.

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).

The inequalities above would not hold if $\sqrt{\lambda_1 \mathbf{W}^\top \mathbf{W}} > \frac{\delta^2}{\gamma^2}$. The gradient *could* explode. It is a *necessary* but not *sufficient* condition for the gradient to explode. It is a *necessary* condition because we have already proved that if $\sqrt{\lambda_1 \mathbf{W}^\top \mathbf{W}} < \frac{\delta^2}{\lambda^2}$, then the gradients vanish.

Question 6 (6-12). Denote by σ the logistic sigmoid function. 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.

1. Draw the computational graph for this RNN, unrolled for 3 time steps (from $t = 1$ to $t = 3$). Include and label the initial hidden states for both the forward and backward RNNs, $\mathbf{h}_0^{(f)}$ and $\mathbf{h}_4^{(b)}$ respectively. You may draw this by hand; you may also use a computer rendering package such as TikZ, but you are not required to do so. Label each node and edge with the corresponding hidden unit or weight.

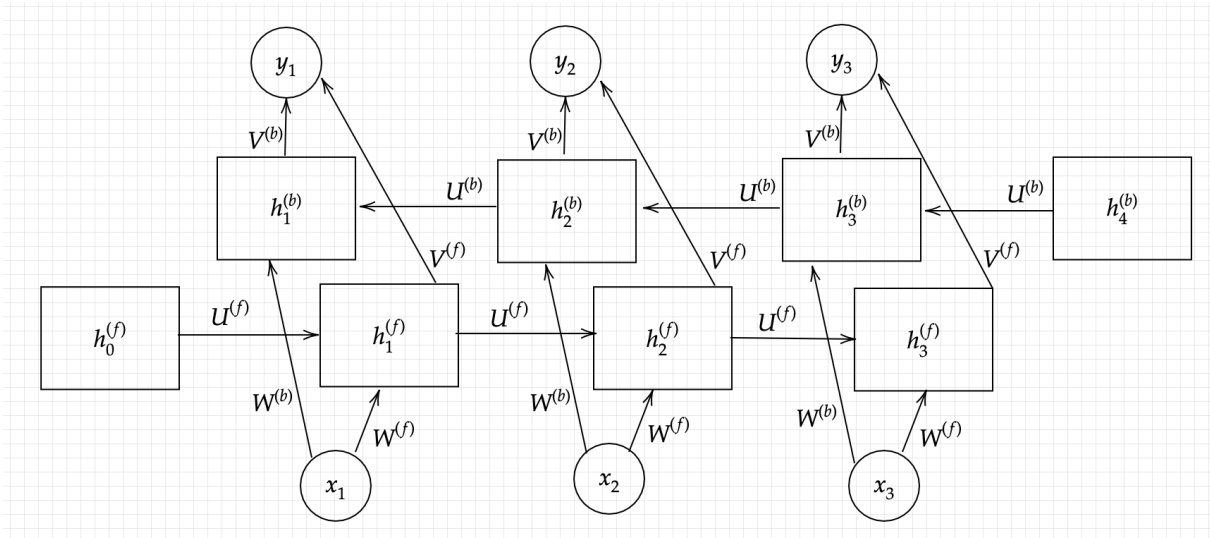


FIGURE 2 – Bidirectional RNN

- *2. 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 $L_t = \|\mathbf{z}_t - \mathbf{y}_t\|_2^2$. 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 $\nabla_{\mathbf{h}_{t-1}^{(b)}} L$ respectively). Then derive $\nabla_{\mathbf{W}^{(f)}} L$ and $\nabla_{\mathbf{U}^{(b)}} L$.

Let us first find the gradient $\nabla_{\mathbf{h}_t^{(f)}} L$. In general $\nabla_{\mathbf{h}_t} L$ does not depend on the hidden states of previous time steps. By the multivariate chain rule, we have the following.

$$\nabla_{\mathbf{h}_3^{(f)}} L = \nabla_{\mathbf{h}_3^{(f)}} L_3$$

$$\nabla_{\mathbf{h}_3^{(f)}} L = \frac{\partial L_3}{\partial \mathbf{y}_3} \cdot \frac{\partial \mathbf{y}_3}{\partial \mathbf{h}_3^{(f)}}$$

$$\nabla_{\mathbf{h}_2^{(f)}} L = \nabla_{\mathbf{h}_2^{(f)}} L_2 + \nabla_{\mathbf{h}_2^{(f)}} L_3$$

$$\nabla_{\mathbf{h}_2^{(f)}} L = \frac{\partial L_2}{\partial \mathbf{y}_2} \cdot \frac{\partial \mathbf{y}_2}{\partial \mathbf{h}_2^{(f)}} + \underbrace{\frac{\partial L_3}{\partial \mathbf{y}_3} \cdot \frac{\partial \mathbf{y}_3}{\partial \mathbf{h}_3^{(f)}}}_{\nabla_{\mathbf{h}_3^{(f)}} L} \cdot \frac{\mathbf{h}_3^{(f)}}{\mathbf{h}_2^{(f)}}$$

$$\nabla_{\mathbf{h}_1^{(f)}} L = \nabla_{\mathbf{h}_1^{(f)}} L_1 + \nabla_{\mathbf{h}_1^{(f)}} L_2 + \nabla_{\mathbf{h}_1^{(f)}} L_3$$

$$\begin{aligned} \nabla_{\mathbf{h}_1^{(f)}} L &= \frac{\partial L_1}{\partial \mathbf{y}_1} \cdot \frac{\partial \mathbf{y}_1}{\partial \mathbf{h}_1^{(f)}} + \frac{\partial L_2}{\partial \mathbf{y}_2} \cdot \frac{\partial \mathbf{y}_2}{\partial \mathbf{h}_2^{(f)}} \cdot \frac{\partial \mathbf{h}_2^{(f)}}{\partial \mathbf{h}_1^{(f)}} + \frac{\partial L_3}{\partial \mathbf{y}_3} \cdot \frac{\partial \mathbf{y}_3}{\partial \mathbf{h}_3^{(f)}} \cdot \frac{\partial \mathbf{h}_3^{(f)}}{\partial \mathbf{h}_2^{(f)}} \cdot \frac{\partial \mathbf{h}_2^{(f)}}{\partial \mathbf{h}_1^{(f)}} \\ &= \frac{\partial L_1}{\partial \mathbf{y}_1} \cdot \frac{\partial \mathbf{y}_1}{\partial \mathbf{h}_1^{(f)}} + \frac{\partial \mathbf{h}_2^{(f)}}{\partial \mathbf{h}_1^{(f)}} \left(\underbrace{\frac{\partial L_2}{\partial \mathbf{y}_2} \cdot \frac{\partial \mathbf{y}_2}{\partial \mathbf{h}_2^{(f)}} + \frac{\partial L_3}{\partial \mathbf{y}_3} \cdot \frac{\partial \mathbf{y}_3}{\partial \mathbf{h}_3^{(f)}} \cdot \frac{\partial \mathbf{h}_3^{(f)}}{\partial \mathbf{h}_2^{(f)}}}_{\nabla_{\mathbf{h}_2^{(f)}} L} \right) \end{aligned}$$

From this pattern, we can infer the following equation:

$$\nabla_{\mathbf{h}_t^{(f)}} L = \nabla_{\mathbf{h}_t^{(f)}} L_t + \nabla_{\mathbf{h}_t^{(f)}} \mathbf{h}_{t+1}^{(f)} \cdot \nabla_{\mathbf{h}_{t+1}^{(f)}} L \quad (57)$$

Now let us calculate the actual gradients. The first term $\nabla_{\mathbf{h}_t^{(f)}} L_t$ can be decomposed as follows:

$$\nabla_{\mathbf{h}_t^{(f)}} L_t = \nabla_{\mathbf{y}_t} L \cdot \nabla_{\mathbf{h}_t^{(f)}} \mathbf{y}_t \quad (58)$$

$$= -2(\mathbf{z}_t - \mathbf{y}_t) \cdot \mathbf{V}^{(f)} \quad (59)$$

We can work out the first part of the second term $\nabla_{\mathbf{h}_t^{(f)}} \mathbf{h}_{t+1}^{(f)}$ as follows (using the fact that $\sigma'(x) = \sigma(x) \cdot (1 - \sigma(x))$)

$$\nabla_{\mathbf{h}_t^{(f)}} \mathbf{h}_{t+1}^{(f)} = \mathbf{h}_{t+1}^{(f)}(1 - \mathbf{h}_{t+1}^{(f)}) \cdot \mathbf{U}^{(f)} \quad (60)$$

Finally combining 59 and 60 in 57 we have:

$$\boxed{\nabla_{\mathbf{h}_t^{(f)}} L = -2(\mathbf{z}_t - \mathbf{y}_t) \cdot \mathbf{V}^{(f)} + \mathbf{h}_{t+1}^{(f)}(1 - \mathbf{h}_{t+1}^{(f)}) \cdot \mathbf{U}^{(f)} \cdot \nabla_{\mathbf{h}_{t+1}^{(f)}} L} \quad (61)$$

Similarly, we can derive the values for $\nabla_{\mathbf{h}_t^{(b)}} L$:

$$\nabla_{\mathbf{h}_t^{(b)}} L = \nabla_{\mathbf{h}_t^{(b)}} L_t + \nabla_{\mathbf{h}_t^{(b)}} \mathbf{h}_{t-1}^{(b)} \cdot \nabla_{\mathbf{h}_{t-1}^{(b)}} L \quad (62)$$

$$= \boxed{-2(\mathbf{z}_t - \mathbf{y}_t) \cdot \mathbf{V}^{(b)} + \mathbf{h}_t^{(b)}(1 - \mathbf{h}_t^{(b)}) \cdot \mathbf{U}^{(b)} \cdot \nabla_{\mathbf{h}_{t-1}^{(b)}} L} \quad (63)$$

Now we find the gradient with respect to the parameter $\nabla_{\mathbf{W}^{(f)}} L$. This equals:

$$\nabla_{\mathbf{W}^{(f)}} L = \sum_t \underbrace{\nabla_{\mathbf{h}_t^{(f)}} L}_{\text{calculated in 61}} \cdot \nabla_{\mathbf{W}^{(f)}} \mathbf{h}_t^{(f)}$$

Note that we have calculated the quantity $\nabla_{\mathbf{h}_t^{(f)}} L$ in 61. The second quantity $\nabla_{\mathbf{W}^{(f)}} \mathbf{h}_t^{(f)}$ is the derivative of a vector with respect to a matrix which is a tensor. In the following derivation we ease the notation for simplicity.

$$\nabla_{\mathbf{W}^{(f)}} L = \boxed{\sum_t \left(\nabla_{\mathbf{h}_t^{(f)}} L = -2(\mathbf{z}_t - \mathbf{y}_t) \cdot \mathbf{V}^{(f)} + \mathbf{h}_{t+1}^{(f)}(1 - \mathbf{h}_{t+1}^{(f)}) \cdot \mathbf{U}^{(f)} \cdot \nabla_{\mathbf{h}_{t+1}^{(f)}} L \right) \cdot \left(\mathbf{h}_t^{(f)}(1 - \mathbf{h}_t^{(f)}) \cdot \mathbf{x}_t \right)}$$

Similary we find the gradient with respect to the parameter $\nabla_{\mathbf{U}^{(b)}} L$. This equals:

$$\nabla_{\mathbf{U}^{(b)}} L = \sum_t \nabla_{\mathbf{h}_t^{(b)}} L \cdot \nabla_{\mathbf{U}^{(b)}} \mathbf{h}_t^{(b)}$$

We have calculated the first quantity $\nabla_{\mathbf{h}_t^{(b)}} L$ in 63. The second quantity works out to be $\mathbf{h}_t^{(b)}(1 - \mathbf{h}_t^{(b)}) \cdot \mathbf{h}_{t+1}^{(b)}$ using the fact that $\sigma'(x) = \sigma(x) \cdot (1 - \sigma(x))$. Thus we have:

$$\nabla_{\mathbf{U}^{(b)}} L = \boxed{\sum_t \left(-2(\mathbf{z}_t - \mathbf{y}_t) \cdot \mathbf{V}^{(b)} + \mathbf{h}_t^{(b)}(1 - \mathbf{h}_t^{(b)}) \cdot \mathbf{U}^{(b)} \cdot \nabla_{\mathbf{h}_{t-1}^{(b)}} L \right) \cdot \left(\mathbf{h}_t^{(b)}(1 - \mathbf{h}_t^{(b)}) \cdot \mathbf{h}_{t+1}^{(b)} \right)}$$

Références

- [1] X. GLOROT AND Y. BENGIO, *Understanding the difficulty of training deep feedforward neural networks*, in Proceedings of the thirteenth international conference on artificial intelligence and statistics, 2010, pp. 249–256.
- [2] K. HE, X. ZHANG, S. REN, AND J. SUN, *Delving deep into rectifiers: Surpassing human-level performance on imagenet classification*, CoRR, abs/1502.01852 (2015).