

**Due Date: March 22nd 23:59, 2019**

### 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  $c$ ,  $\mu$ , and  $\sigma^2$  as a function of  $d^{(t-1)}$ .
  - (b) Uniform: give the values for  $c$ ,  $\alpha$ , and  $\beta$  as a function of  $d^{(t-1)}$ .
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  $c$ ,  $\mu$ , and  $\sigma^2$  as a function of  $d^{(t-1)}$ .
  - (b) Uniform: give the values for  $c$ ,  $\alpha$ , and  $\beta$  as a function of  $d^{(t-1)}$ .
  - (c) What popular initialization scheme has this form?
  - (d) Why do you think this initialization would work well in practice? Answer in 1-2 sentences.

**Answer 1.**

1. For  $1 \leq i \leq d^{(t)}, 1 \leq j \leq d^{(t-1)}$ ,  $\mathbf{W}_{ij}^{(t)}$  from Gaussian or Uniform distributions means  $\mathbf{W}_{ij}^{(t)}$  and  $\mathbf{h}_j^{(t-1)}$  are independent:  $\mathbf{W}_{ij}^{(t)} \perp \mathbf{h}_j^{(t-1)}$ .

Given  $\mathbb{E}[\mathbf{h}_i^{(t-1)}] = 0$  and  $\text{Var}(\mathbf{h}_i^{(t-1)}) = 1$ , for both distributions:

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

and,

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

One possible parameters initialize scheme is to let  $c = 0$ , such that  $\mathbb{E}[\mathbf{a}_i^{(t)}] = 0$  and  $\text{Var}(\mathbf{b}_i^{(t)}) = 0$ , and to let  $\mathbb{E}[\mathbf{W}_{ij}^{(t)}] = 0$  and  $\text{Var}(\mathbf{W}_{ij}^{(t)}) = \frac{1}{d^{(t-1)}}$  such that  $\text{Var}(\mathbf{a}_i^{(t)}) = 1$

- (a) For a Gaussian distribution:  $\mathbf{W}_{ij}^{(t)} \sim \mathcal{N}(\mu, \sigma^2)$ ,

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

Therefore, the scheme is equivalent to:

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

- (b) For a Uniform distribution:  $\mathbf{W}_{ij}^{(t)} \sim U(\alpha, \beta)$ ,

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

Therefore, the scheme is equivalent to:

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

2. Given  $\mathbb{E}[\mathbf{a}_i^{(t-1)}] = 0$ ,  $\text{Var}(\mathbf{a}_i^{(t-1)}) = 1$  and  $\mathbf{a}_i^{(t-1)}$  has a symmetric distribution for  $1 \leq i \leq d^{(t-1)}$ , we have:

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

From previous question, we know:

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

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

One possible parameters initialize scheme is to let  $c = 0$  and  $\mathbb{E}[\mathbf{W}_{ij}^{(t)}] = 0$ , such that  $\mathbb{E}[\mathbf{a}_i^{(t)}] = 0$ ,  $\text{Var}(\mathbf{b}_i^{(t)}) = 0$  and simply  $\text{Var}(\mathbf{a}_i^{(t)})$  to:

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

Therefore, we need

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

- (a) For a Gaussian distribution:  $\mathbf{W}_{ij}^{(t)} \sim \mathcal{N}(\mu, \sigma^2)$ ,

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

Therefore, the scheme is equivalent to:

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

(b) For a Uniform distribution:  $\mathbf{W}_{ij}^{(t)} \sim U(\alpha, \beta)$ ,

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

Therefore, the scheme is equivalent to:

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

(c) He Normal (He-et-al) Initialization scheme has this form.

(d) By using this initialization scheme, the expectation and variance of the output of each layer will be very close to 0 and 1, respectively; which maintains a certain level of gradient for each layer, decreases the possibility of the gradient exploding or vanishing, and therefore gets a better trained model.

**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  $\mathbf{R}$ ).
2. 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  $\mathbb{E}[L(\mathbf{w})] = \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1-p)\|\Gamma\mathbf{w}\|^2$ .
3. Show that the solution  $\mathbf{w}^{\text{dropout}}$  that minimizes the expected loss from question 2.2 satisfies

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

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

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

**Answer 2.**

1.  $L(\mathbf{w}) = \|\mathbf{y} - (\mathbf{R} \odot \mathbf{X})\mathbf{w}\|^2$ , where  $\odot$  is an element-wise multiplication.
- 2.

$$\begin{aligned} \mathbb{E}[L(\mathbf{w})] &= \mathbb{E}[\|\mathbf{y} - (\mathbf{R} \odot \mathbf{X})\mathbf{w}\|^2] \\ &= \mathbb{E}[(\mathbf{y} - (\mathbf{R} \odot \mathbf{X})\mathbf{w})^\top (\mathbf{y} - (\mathbf{R} \odot \mathbf{X})\mathbf{w})] \\ &= \mathbb{E}[\mathbf{y}^\top \mathbf{y} - 2\mathbf{w}^\top (\mathbf{R} \odot \mathbf{X})^\top \mathbf{y} + \mathbf{w}^\top (\mathbf{R} \odot \mathbf{X})^\top (\mathbf{R} \odot \mathbf{X})\mathbf{w}] \\ &= \mathbf{y}^\top \mathbf{y} - 2\mathbf{w}^\top \mathbb{E}[(\mathbf{R} \odot \mathbf{X})]^\top \mathbf{y} + \mathbf{w}^\top \mathbb{E}[(\mathbf{R} \odot \mathbf{X})^\top (\mathbf{R} \odot \mathbf{X})]\mathbf{w} \end{aligned}$$

For the second item,

$$(\mathbb{E}[\mathbf{R} \odot \mathbf{X}])_{ij} = \mathbb{E}[(\mathbf{R} \odot \mathbf{X})_{ij}] = p\mathbf{X}_{ij}$$

For the last item,

$$\begin{aligned} & \mathbb{E}[(\mathbf{R} \odot \mathbf{X})^\top (\mathbf{R} \odot \mathbf{X})]_{ij} \\ &= \mathbb{E}\left[\sum_{k=1}^N (\mathbf{R} \odot \mathbf{X})_{ki} (\mathbf{R} \odot \mathbf{X})_{kj}\right] \\ &= \sum_{k=1}^d \mathbb{E}[\mathbf{R}_{ki} \mathbf{R}_{kj}] \mathbf{X}_{ki} \mathbf{X}_{kj} \\ &= \begin{cases} p^2 (\mathbf{X}^\top \mathbf{X})_{ij} & \text{if } i \neq j \\ \sum_{k=1}^d \mathbb{E}[\mathbf{R}_{ki}^2] \mathbf{X}_{ki}^2 = p(\mathbf{X}^\top \mathbf{X})_{ii} & \text{if } i = j \end{cases} \end{aligned}$$

With the two items re-written, we have:

$$\begin{aligned} \mathbb{E}[L(\mathbf{w})] &= \mathbf{y}^\top \mathbf{y} - 2\mathbf{w}^\top \mathbb{E}[(\mathbf{R} \odot \mathbf{X})]^\top \mathbf{y} + \mathbf{w}^\top \mathbb{E}[(\mathbf{R} \odot \mathbf{X})^\top (\mathbf{R} \odot \mathbf{X})] \mathbf{w} \\ &= \mathbf{y}^\top \mathbf{y} - 2p\mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \mathbf{w}^\top \mathbb{E}[(\mathbf{R} \odot \mathbf{X})^\top (\mathbf{R} \odot \mathbf{X})] \mathbf{w} \\ &= \mathbf{y}^\top \mathbf{y} - 2p\mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + p^2 \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - p^2 \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} + \mathbf{w}^\top \mathbb{E}[(\mathbf{R} \odot \mathbf{X})^\top (\mathbf{R} \odot \mathbf{X})] \mathbf{w} \\ &= \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + \mathbf{w}^\top \left( \mathbb{E}[(\mathbf{R} \odot \mathbf{X})^\top (\mathbf{R} \odot \mathbf{X})] - p^2 \mathbf{X}^\top \mathbf{X} \right) \mathbf{w} \\ &= \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + \mathbf{w}^\top (p(1-p) \text{diag}[(\mathbf{X}^\top \mathbf{X})_{11}, (\mathbf{X}^\top \mathbf{X})_{22}, \dots, (\mathbf{X}^\top \mathbf{X})_{dd}]^\top) \mathbf{w} \\ &= \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + (p(1-p) \mathbf{w}^\top \Gamma^\top \Gamma) \mathbf{w} \\ &= \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + (p(1-p) (\Gamma \mathbf{w})^\top \Gamma \mathbf{w}) \\ &= \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1-p) \|\Gamma \mathbf{w}\|^2 \end{aligned}$$

3. For a matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  and a vector  $\mathbf{w} \in \mathbb{R}^{d \times 1}$ , formulas of matrix derivatives (denominator layout):

$$\frac{\partial \mathbf{A} \mathbf{w}}{\partial \mathbf{w}} = \mathbf{A}^\top, \quad \frac{\partial \mathbf{w}^\top \mathbf{A}}{\partial \mathbf{w}} = \mathbf{A}, \quad \frac{\partial \mathbf{w}^\top \mathbf{A} \mathbf{w}}{\partial \mathbf{w}} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{w}$$

are used for the following calculation:

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

Let  $\frac{\partial(\mathbb{E}[L(\mathbf{w})])}{\partial \mathbf{w}} = 0$ , we can find  $\mathbf{w}^{\text{dropout}}$  to a minimal  $L(\mathbf{w})$ :

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

Obviously,  $\lambda^{\text{dropout}} = \frac{1-p}{p}$ . When  $p$  is close to 0,  $\lambda^{\text{dropout}}$  grows to  $\infty$ ; when  $p$  is close to 1,  $\lambda^{\text{dropout}}$  decreases to 0.

4. When use  $L_2$  regularization, the loss function can be written as:

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

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

Let  $\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = 0$ , we find  $\mathbf{w}^{L_2}$  to a minimal  $L(\mathbf{w})$ :

$$\mathbf{w}^{L_2} = (\mathbf{X}^\top \mathbf{X} + \lambda^{L_2})^{-1} \mathbf{X}^\top \mathbf{y}$$

5. Dropout with linear regression is equivalent to ridge regression with a particular form for  $\Gamma$ . This form essentially scales the weight cost for weight  $\mathbf{w}_i$  by the standard deviation of the  $i_{\text{th}}$  dimension of the data. Both dropout and weight decay can be considered a kind of regularization method in training a network, and the theoretical solution of the optimal weights from two methods have similar structures, indicating regularization method can be considered a special case of dropout.

**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  $\Delta\theta_t$  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:

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

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

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

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

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

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

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

**Answer 3.**

1. • for SGD with momentum,

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

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

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

Given:

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

the two update rules are equivalent.

2.

$$\begin{aligned}\Delta\theta_t &= \beta\Delta\theta_{t-1} - \delta(1-\beta)\mathbf{g}_t \\ &= \beta(\beta\Delta\theta_{t-2} - \delta(1-\beta)\mathbf{g}_{t-1}) - \delta(1-\beta)\mathbf{g}_t \\ &= \beta^2\Delta\theta_{t-2} - \delta\beta(1-\beta)\mathbf{g}_{t-1} - \delta(1-\beta)\mathbf{g}_t \\ &\dots \\ &= \beta^t\Delta\theta_0 - \delta(1-\beta)(\beta^{t-1}\mathbf{g}_1 + \beta^{t-2}\mathbf{g}_2 + \dots + \beta\mathbf{g}_{t-1} + \beta^0\mathbf{g}_t) \\ &= \beta^t(-\delta\mathbf{v}_0) - \delta(1-\beta)\sum_{i=1}^t\beta^{t-i}\mathbf{g}_i \\ &= -\delta(1-\beta)\sum_{i=1}^t\beta^{t-i}\mathbf{g}_i \\ \mathbf{v}_t &= -\frac{\Delta\theta_t}{\delta} = (1-\beta)\sum_{i=1}^t\beta^{t-i}\mathbf{g}_i\end{aligned}$$

3. Given  $\mathbf{g}_t$  has a stationary distribution independent of  $t$ ,

$$\begin{aligned}\mathbb{E}[\mathbf{v}_t] &= \mathbb{E}\left[\frac{\Delta\boldsymbol{\theta}_t}{-\delta}\right] \\ &= \mathbb{E}\left[(1-\beta) \sum_{i=1}^t \beta^{t-i} \mathbf{g}_i\right] \\ &= \mathbb{E}[\mathbf{g}_t](1-\beta) \sum_{i=1}^t \beta^{t-i} \\ &= (1-\beta) \frac{1-\beta^t}{1-\beta} \mathbb{E}[\mathbf{g}_t] \\ &= (1-\beta^t) \mathbb{E}[\mathbf{g}_t]\end{aligned}$$

By importing an exponential weighted mean factor  $\frac{1}{1-\beta^t}$ ; that is, let  $\tilde{\mathbf{v}}_t = \frac{1}{1-\beta^t} \mathbf{v}_t$ , then:

$$\mathbb{E}[\tilde{\mathbf{v}}_t] = \mathbb{E}[\mathbf{g}_t]$$

**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  $\gamma$  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  $\hat{y} = \mathbf{w}^\top \mathbf{x} + \beta$ .
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  $\mathbf{W}^\perp \mathbf{u} = \mathbf{0}$ .
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  $\mathbf{W}^\perp \mathbf{u} = \mathbf{0}$  from question 4.2).

**Answer 4.**

1. Given  $\mathbf{w} := \gamma \frac{\mathbf{u}}{\|\mathbf{u}\|}$ ,  $y = \mathbf{u}^\top \mathbf{x}$ ,  $\text{Var}(\mathbf{x}) = \mathbf{I}$ , and  $\mathbb{E}[\mathbf{x}] = \mathbf{0}$ , we have:

$$\mu_y = \mathbb{E}[y] = \mathbb{E}[\mathbf{u}^\top \mathbf{x}] = \mathbb{E}\left[\sum_i \mathbf{u}_i \mathbf{x}_i\right] = \sum_i \mathbb{E}[\mathbf{u}_i \mathbf{x}_i] = \sum_i \mathbf{u}_i \mathbb{E}[\mathbf{x}_i] = 0$$

---

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.



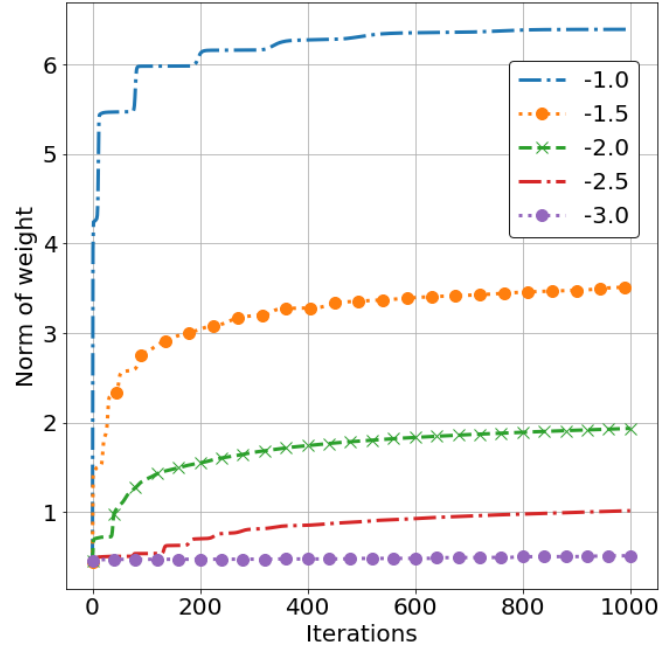


FIGURE 1 – Norm of parameters with different learning rate.

$$\begin{aligned}
 \gamma \cdot \frac{y - \mu_y}{\sigma_y} &= \gamma \frac{\mathbf{u}^\top \mathbf{x} - 0}{\sqrt{\text{Var}(\mathbf{u}^\top \mathbf{x})}} \\
 &= \gamma \frac{\mathbf{u}^\top \mathbf{x}}{\sqrt{\text{Var}(\sum_i \mathbf{u}_i \mathbf{x}_i)}} \\
 &= \gamma \frac{\mathbf{u}^\top \mathbf{x}}{\sqrt{\sum_i \mathbf{u}_i^2 \text{Var}(\mathbf{x}_i)}} \quad (\text{Var}(\mathbf{x}) = \mathbf{I}, \quad \text{Cov}(\mathbf{x}_i, \mathbf{x}_j) = 0 \text{ if } i \neq j) \\
 &= \gamma \frac{\mathbf{u}^\top \mathbf{x}}{\sqrt{\sum_i \mathbf{u}_i^2}} \\
 &= \gamma \frac{\mathbf{u}^\top \mathbf{x}}{\|\mathbf{u}\|} \\
 &= \mathbf{w}^\top \mathbf{x}
 \end{aligned}$$

Therefore,  $\hat{y} = \mathbf{w}^\top \mathbf{x} + \beta$

2. First,

$$\begin{aligned}
 \frac{\partial \|\mathbf{u}\|}{\partial \mathbf{u}} &= \frac{\partial (\mathbf{u}^\top \mathbf{u})^{1/2}}{\partial \mathbf{u}} \\
 &= \frac{1}{2} (\mathbf{u}^\top \mathbf{u})^{-1/2} \frac{\partial \mathbf{u}^\top \mathbf{u}}{\partial \mathbf{u}} \\
 &= \frac{1}{2} \frac{2 \mathbf{u}^\top}{\|\mathbf{u}\|} \\
 &= \frac{\mathbf{u}^\top}{\|\mathbf{u}\|}
 \end{aligned}$$

Given  $\mathbf{W}^\perp = \left(\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2}\right)$ , let  $s = \frac{\gamma}{\|\mathbf{u}\|}$ , then:

$$\begin{aligned}
\nabla_{\mathbf{u}} L &= \frac{\partial \mathbf{w}}{\partial \mathbf{u}} \nabla_{\mathbf{w}} L \\
&= \gamma \frac{\partial}{\partial \mathbf{u}} \left( \frac{\mathbf{u}}{\|\mathbf{u}\|} \right) \nabla_{\mathbf{w}} L \\
&= \gamma \frac{\|\mathbf{u}\| \mathbf{I} - \mathbf{u} \partial \|\mathbf{u}\| / \partial \mathbf{u}}{\|\mathbf{u}\|^2} \nabla_{\mathbf{w}} L \\
&= \gamma \left( \frac{\mathbf{I}}{\|\mathbf{u}\|} - \frac{\mathbf{u}}{\|\mathbf{u}\|^2} \frac{\partial \|\mathbf{u}\|}{\partial \mathbf{u}} \right) \nabla_{\mathbf{w}} L \\
&= \gamma \left( \frac{\mathbf{I}}{\|\mathbf{u}\|} - \frac{\mathbf{u}}{\|\mathbf{u}\|^2} \frac{\mathbf{u}^\top}{\|\mathbf{u}\|} \right) \nabla_{\mathbf{w}} L \\
&= \frac{\gamma}{\|\mathbf{u}\|} \left( \mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2} \right) \nabla_{\mathbf{w}} L \\
&= s \mathbf{W}^\perp \nabla_{\mathbf{w}} L
\end{aligned}$$

3. First,

$$\begin{aligned}
\mathbf{u}^\top \nabla_{\mathbf{u}} L &= \mathbf{u}^\top s \mathbf{W}^\perp \nabla_{\mathbf{w}} L \\
&= s \mathbf{u}^\top \left( \mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2} \right) \nabla_{\mathbf{w}} L \\
&= s \left( \mathbf{u}^\top - \frac{\mathbf{u}^\top \mathbf{u} \mathbf{u}^\top}{\|\mathbf{u}\|^2} \right) \nabla_{\mathbf{w}} L \\
&= s (\mathbf{u}^\top - \mathbf{u}^\top) \nabla_{\mathbf{w}} L \\
&= 0
\end{aligned}$$

which also means

$$(\nabla_{\mathbf{u}} L)^\top \mathbf{u} = (\mathbf{u}^\top \nabla_{\mathbf{u}} L)^\top = 0$$

Let  $\mathbf{u}'$  be the one time step updated  $\mathbf{u}$ , that is:

$$\mathbf{u}' = \mathbf{u} - \alpha \nabla_{\mathbf{u}} L$$

then,

$$\begin{aligned}
\|\mathbf{u}'\|^2 - \|\mathbf{u}\|^2 &= (\mathbf{u} - \alpha \nabla_{\mathbf{u}} L)^\top (\mathbf{u} - \alpha \nabla_{\mathbf{u}} L) - \mathbf{u}^\top \mathbf{u} \\
&= \mathbf{u}^\top \mathbf{u} - \mathbf{u}^\top \alpha \nabla_{\mathbf{u}} L - \alpha (\nabla_{\mathbf{u}} L)^\top \mathbf{u} + \alpha^2 (\nabla_{\mathbf{u}} L)^\top \nabla_{\mathbf{u}} L - \mathbf{u}^\top \mathbf{u} \\
&= \alpha^2 (\nabla_{\mathbf{u}} L)^\top \nabla_{\mathbf{u}} L - \alpha (\mathbf{u}^\top \nabla_{\mathbf{u}} L + (\nabla_{\mathbf{u}} L)^\top \mathbf{u}) \\
&= \alpha^2 \|\nabla_{\mathbf{u}} L\|^2 + \alpha(0 + 0) \\
&= \|\alpha \nabla_{\mathbf{u}} L\|^2
\end{aligned}$$

This indicates that  $\Delta \mathbf{u} = -\alpha \nabla_{\mathbf{u}} L$  ( $\alpha$  is learning rate) is always orthonormal to  $\mathbf{u}$ .  $\|\mathbf{u}'\|$ ,  $\|\mathbf{u}\|$ , and  $\|\Delta \mathbf{u}\|$  can form a right triangle. According to the Pythagorean theorem,  $\|\mathbf{u}'\|$  is the longest. That is, if  $\|\Delta \mathbf{u}\| > 0$ , the norm of updated weight  $\mathbf{u}'$  satisfies:

$$\|\mathbf{u}'\| = \|\mathbf{u} + \Delta \mathbf{u}\| > \|\mathbf{u}\|$$

This explains that the norm of weight is always increasing.

Larger learning rate corresponds to larger  $\|\Delta \mathbf{u}\|$  in early iteration, and therefore leading to a larger  $\|\mathbf{u}'\|$ . This explains larger learning rate corresponds to faster growth of  $\|\mathbf{u}\|$  in early several iterations.

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

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

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. I will first make a guess, then prove it by mathematical induction.

Assume existing a function  $\sigma^{-1}$  such that:

$$\sigma^{-1}(\sigma(x)) = x$$

from

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

we can derive:

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

take  $\mathbf{U}\mathbf{x}_t + \mathbf{b}$  into  $\mathbf{h}_t$ , we have:

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

---

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.

This implies:

$$\mathbf{g}_t = \sigma(\mathbf{h}_t)$$

Prove: At time step  $t = 0$ ,  $\mathbf{g}_0$  and  $\mathbf{h}_0$  are initialized to 0.

At time step  $t = 1$ , based on the formulas give for  $\mathbf{g}_t$  and  $\mathbf{h}_t$ , we have:

$$\begin{aligned}\mathbf{h}_1 &= \mathbf{W}\sigma(\mathbf{h}_0) + \mathbf{U}\mathbf{x}_1 + \mathbf{b} = \mathbf{U}\mathbf{x}_1 + \mathbf{b} \\ \mathbf{g}_1 &= \sigma(\mathbf{W}\mathbf{g}_0 + \mathbf{U}\mathbf{x}_1 + \mathbf{b}) = \sigma(\mathbf{U}\mathbf{x}_1 + \mathbf{b}) \\ &= \sigma(\mathbf{h}_1)\end{aligned}$$

which means  $\mathbf{g}_t = \sigma(\mathbf{h}_t)$  holds for  $t = 1$ .

Assume for the time  $t - 1$ ,  $\mathbf{g}_{t-1} = \sigma(\mathbf{h}_{t-1})$  holds, then for the time step  $t$ :

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

which means  $\mathbf{g}_t = \sigma(\mathbf{h}_t)$  holds as well for time step  $t$ .

2.

$$\begin{aligned}\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| &= \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \frac{\partial \mathbf{h}_{T-1}}{\partial \mathbf{h}_{T-2}} \cdots \frac{\partial \mathbf{h}_1}{\partial \mathbf{h}_0} \right\| \\ &= \left\| \mathbf{W}^\top \sigma'(\mathbf{h}_{T-1}) \mathbf{W}^\top \sigma'(\mathbf{h}_{T-1}) \cdots \mathbf{W}^\top \sigma'(\mathbf{h}_0) \right\| \\ &= \left\| \prod_{i=0}^{T-1} \mathbf{W}^\top \sigma'(\mathbf{h}_i) \right\| \\ &\leq \prod_{i=0}^{T-1} \left\| \mathbf{W}^\top \sigma'(\mathbf{h}_i) \right\| \\ &\leq \left\| \mathbf{W}^\top \right\|^T \prod_{i=0}^{T-1} \left\| \sigma'(\mathbf{h}_i) \right\| \\ &= \left\| \mathbf{W} \right\|^T \prod_{i=0}^{T-1} \left\| \sigma'(\mathbf{h}_i) \right\| \\ &\leq \left( \sqrt{\lambda_1(\mathbf{W}^\top \mathbf{W})} \right)^T \gamma^T \\ &\leq \left( \frac{\delta}{\gamma} \right)^T \gamma^T \\ &= \delta^T\end{aligned}$$

Since  $0 \leq \delta < 1$ , as  $T \rightarrow \infty$ , the norm of the gradient of the hidden states will be 0, meaning the gradient of the hidden states will vanish.

3. The condition of the largest eigenvalue of the weights being larger than  $\frac{\delta^2}{\gamma^2}$  is a necessary but not sufficient for the gradient to explode. There are some other factors which affect the gradient, including the relation of the directions of matrices  $\mathbf{W}$  and  $\sigma'(\mathbf{h}_i)$  (for  $0 \leq i \leq T - 1$ ).

**Question 6** (6-12). Denote by  $\sigma$  the logistic sigmoid function. Consider the following Bidirectional RNN:

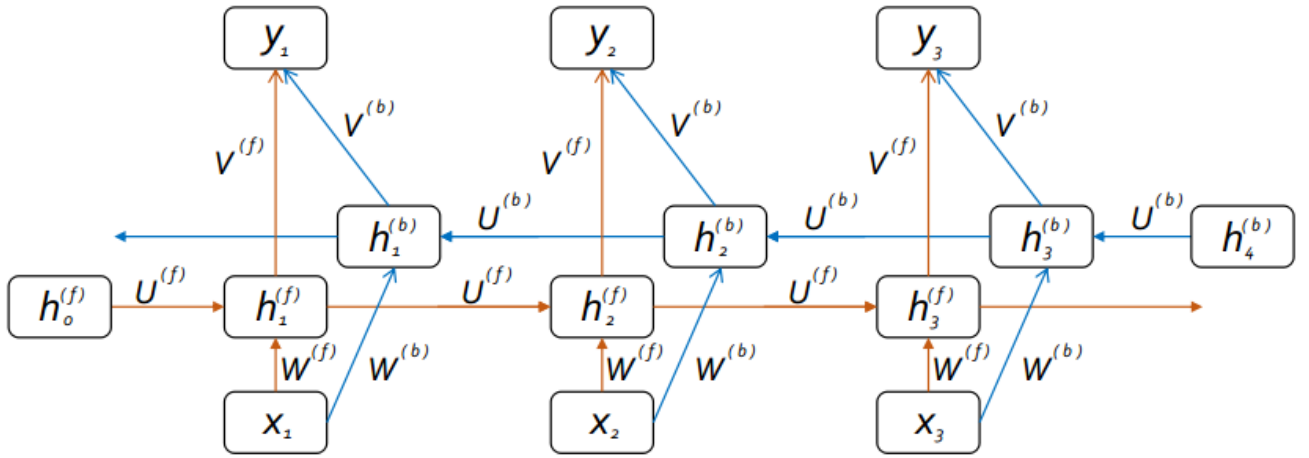
$$\begin{aligned} \mathbf{h}_t^{(f)} &= \sigma(\mathbf{W}^{(f)}\mathbf{x}_t + \mathbf{U}^{(f)}\mathbf{h}_{t-1}^{(f)}) \\ \mathbf{h}_t^{(b)} &= \sigma(\mathbf{W}^{(b)}\mathbf{x}_t + \mathbf{U}^{(b)}\mathbf{h}_{t+1}^{(b)}) \\ \mathbf{y}_t &= \mathbf{V}^{(f)}\mathbf{h}_t^{(f)} + \mathbf{V}^{(b)}\mathbf{h}_t^{(b)}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

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,  $h_0^{(f)}$  and  $h_4^{(b)}$  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.
- \*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  $\nabla_{\mathbf{U}^{(b)}} L$ .

**Answer 6.**

1. The computational graph is as follow:



2. Note: denominator display of matrix derivative is used for the following derivation.

First we have:

$$\frac{\partial L}{\partial L_t} = 1$$

then,

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

At the final time step  $\tau$ ,  $\mathbf{h}_\tau^{(f)}$  only has  $\mathbf{y}_\tau$  as a descendent, its gradient is:

$$\nabla_{\mathbf{h}_\tau^{(f)}} L = \mathbf{V}^{(f)\top} \nabla_{\mathbf{y}_\tau} L = 2\mathbf{V}^{(f)\top} (\mathbf{y}_\tau - \mathbf{z}_\tau)$$

Similarly, at the first time step ( $t = 1$ ),  $\mathbf{h}_1^{(b)}$  only has  $\mathbf{y}_1$  as a descendent, its gradient is:

$$\nabla_{\mathbf{h}_1^{(b)}} L = \mathbf{V}^{(b)\top} \nabla_{\mathbf{y}_1} L = 2\mathbf{V}^{(b)\top} (\mathbf{y}_1 - \mathbf{z}_1)$$

For the time step  $t$  ( $1 < t < \tau$ ), the gradients  $\nabla_{\mathbf{h}_t^{(f)}} L$  and  $\nabla_{\mathbf{h}_t^{(b)}} L$  are composed of two parts: one from the next or previous time step, and the other from the output  $\mathbf{y}_t$ .

So, we have:

$$\begin{aligned}
\nabla_{\mathbf{h}_t^{(f)}} L &= \frac{\partial L}{\partial \mathbf{h}_t^{(f)}} \\
&= \frac{\partial \mathbf{h}_{t+1}^{(f)}}{\partial \mathbf{h}_t^{(f)}} \frac{\partial L}{\partial \mathbf{h}_{t+1}^{(f)}} + \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t^{(f)}} \frac{\partial L_t}{\partial \mathbf{y}_t} \frac{\partial L}{\partial L_t} \\
&= (\mathbf{U}^{(f)})^\top \text{diag}(\mathbf{h}_{t+1}^{(f)} - (\mathbf{h}_{t+1}^{(f)})^2) \nabla_{\mathbf{h}_{t+1}^{(f)}} L + 2(\mathbf{V}^{(f)})^\top (\mathbf{y}_t - \mathbf{z}_t)
\end{aligned}$$

$$\begin{aligned}
\nabla_{\mathbf{h}_t^{(b)}} L &= \frac{\partial L}{\partial \mathbf{h}_t^{(b)}} \\
&= \frac{\partial \mathbf{h}_{t-1}^{(b)}}{\partial \mathbf{h}_t^{(b)}} \frac{\partial L}{\partial \mathbf{h}_{t-1}^{(b)}} + \frac{\partial \mathbf{y}_t}{\partial \mathbf{h}_t^{(b)}} \frac{\partial L_t}{\partial \mathbf{y}_t} \frac{\partial L}{\partial L_t} \\
&= (\mathbf{U}^{(b)})^\top \text{diag}(\mathbf{h}_{t-1}^{(b)} - (\mathbf{h}_{t-1}^{(b)})^2) \nabla_{\mathbf{h}_{t-1}^{(b)}} L + 2(\mathbf{V}^{(b)})^\top (\mathbf{y}_t - \mathbf{z}_t)
\end{aligned}$$

Introduce dummy variables  $\mathbf{W}_t^{(f)}$  and  $\mathbf{U}_t^{(b)}$  that are defined to be copies of  $\mathbf{W}^{(f)}$  and  $\mathbf{U}^{(b)}$  respectively, but with each  $\mathbf{W}_t^{(f)}$  and  $\mathbf{U}_t^{(b)}$  used only at time step  $t$ .

$$\begin{aligned}
\nabla_{\mathbf{W}^{(f)}} L &= \sum_t \nabla_{\mathbf{W}_t^{(f)}} L \\
&= \sum_t \left( \frac{\partial \mathbf{h}_t^{(f)}}{\partial \mathbf{W}_t^{(f)}} \right)^\top (\nabla_{\mathbf{h}_t^{(f)}} L) \\
&= \sum_t \text{diag}(\mathbf{h}_t^{(f)} - (\mathbf{h}_t^{(f)})^2) (\nabla_{\mathbf{h}_t^{(f)}} L) \mathbf{x}_t^\top
\end{aligned}$$

$$\begin{aligned}
\nabla_{\mathbf{U}^{(b)}} L &= \sum_t \nabla_{\mathbf{U}_t^{(b)}} L \\
&= \sum_t \left( \frac{\partial \mathbf{h}_t^{(b)}}{\partial \mathbf{U}_t^{(b)}} \right)^\top (\nabla_{\mathbf{h}_t^{(b)}} L) \\
&= \sum_t \text{diag} \left( \mathbf{h}_t^{(b)} - (\mathbf{h}_t^{(b)})^2 \right) (\nabla_{\mathbf{h}_t^{(b)}} L) \mathbf{h}_{t+1}^\top
\end{aligned}$$