

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 $d^{(t-1)}$.
 - (b) Uniform: give the values for c , α , and β 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 , μ , and σ^2 as a function of $d^{(t-1)}$.
 - (b) Uniform: give the values for c , α , and β 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. Assume $\mathbf{W}_i^{(t)}$ is the i th row of $\mathbf{W}^{(t)}$ for $1 \leq i \leq d^{(t)}$.

$$\begin{aligned} \mathbb{E}[\mathbf{a}_i^{(t)}] &= \mathbb{E}[\mathbf{W}_i^{(t)} \mathbf{h}^{(t-1)} + \mathbf{b}_i^{(t)}] \\ &= \mathbb{E}[\mathbf{W}_i^{(t)} \mathbf{h}^{(t-1)}] + \mathbb{E}[\mathbf{b}_i^{(t)}] \quad (\because \mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]) \\ &= \mathbb{E}[\mathbf{W}_i^{(t)}] \mathbb{E}[\mathbf{h}^{(t-1)}] + \mathbb{E}[\mathbf{b}_i^{(t)}] \quad (\because \mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y] \text{ if } X \perp Y) \end{aligned}$$

$$\begin{aligned} \text{Var}(\mathbf{a}_i^{(t)}) &= \text{Var}(\mathbf{W}_i^{(t)} \mathbf{h}^{(t-1)} + \mathbf{b}_i^{(t)}) \\ &= \text{Var}(\mathbf{W}_i^{(t)} \mathbf{h}^{(t-1)}) + \text{Var}(\mathbf{b}_i^{(t)}) \quad (\because \text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)) \\ &= \text{Var}(\mathbf{W}_i^{(t)}) \text{Var}(\mathbf{h}^{(t-1)}) + \text{Var}(\mathbf{W}_i^{(t)}) \mathbb{E}[\mathbf{h}^{(t-1)}]^2 + \mathbb{E}[\mathbf{W}_i^{(t)}]^2 \text{Var}(\mathbf{h}^{(t-1)}) + \text{Var}(\mathbf{b}_i^{(t)}) \\ &\quad (\because \text{Var}(XY) = \text{Var}(X)\text{Var}(Y) + \text{Var}(X)\mathbb{E}[Y]^2 + \mathbb{E}[X]^2\text{Var}(Y) \text{ if } X \perp Y) \end{aligned}$$

1.

$$\begin{aligned}
 \mathbb{E}[\mathbf{h}_i^{(t-1)}] &= 0, \text{Var}(\mathbf{h}_i^{(t-1)}) = 1 \text{ for } 1 \leq i \leq d^{(t-1)} \\
 \implies \mathbb{E}[\mathbf{a}_i^{(t)}] &= \mathbb{E}[\mathbf{W}_i^{(t)}] \mathbb{E}[\mathbf{h}^{(t-1)}] + \mathbb{E}[\mathbf{b}_i^{(t)}] = c, \\
 (\mathbf{b}_i^{(t)} = c &\implies \mathbb{E}[\mathbf{b}_i^{(t)}] = c) \\
 \text{Var}(\mathbf{a}_i^{(t)}) &= \text{Var}(\mathbf{W}_i^{(t)}) \mathbb{E}[\mathbf{h}^{(t-1)}] \mathbb{E}[\mathbf{h}^{(t-1)}]^\top + \text{Var}(\mathbf{W}_i^{(t)}) \mathbb{E}[\mathbf{h}^{(t-1)}] \mathbb{E}[\mathbf{h}^{(t-1)}]^\top + \mathbb{E}[\mathbf{W}_i^{(t)}]^2 \text{Var}(\mathbf{h}^{(t-1)}) \\
 &\quad + \text{Var}(\mathbf{b}_i^{(t)}) \quad (\mathbf{b}_i^{(t)} = c \implies \text{Var}(\mathbf{b}_i^{(t)}) = 0) \\
 &= \text{Var}(\mathbf{W}_i^{(t)}) \mathbf{1}_{d^{(t-1)}} + \mathbb{E}[\mathbf{W}_i^{(t)}]^2 \mathbf{1}_{d^{(t-1)}}
 \end{aligned}$$

(a)

$$\begin{aligned}
 \mathbf{W}_{ij}^{(t)} &\sim \mathcal{N}(\mu, \sigma^2) \implies \mathbb{E}[\mathbf{W}_i^{(t)}] = \mu \mathbf{1}_{d^{(t-1)}}^\top, \text{Var}(\mathbf{W}_i^{(t)}) = \sigma^2 \mathbf{1}_{d^{(t-1)}}^\top \\
 \implies \mathbb{E}[\mathbf{a}_i^{(t)}] &= c, \\
 \text{Var}(\mathbf{a}_i^{(t)}) &= \text{Var}(\mathbf{W}_i^{(t)}) \mathbf{1}_{d^{(t-1)}} + \mathbb{E}[\mathbf{W}_i^{(t)}]^2 \mathbf{1}_{d^{(t-1)}} = \sigma^2 \mathbf{1}_{d^{(t-1)}}^\top \mathbf{1}_{d^{(t-1)}} + \mu^2 \mathbf{1}_{d^{(t-1)}}^\top \mathbf{1}_{d^{(t-1)}} \\
 &= (\sigma^2 + \mu^2) d^{(t-1)} \quad (\because \mathbf{1}_{d^{(t-1)}}^\top \mathbf{1}_{d^{(t-1)}} = d^{(t-1)})
 \end{aligned}$$

$$\mathbb{E}[\mathbf{a}_i^{(t)}] = 0, \text{Var}(\mathbf{a}_i^{(t)}) = 1 \implies c = 0, \sigma^2 + \mu^2 = 1/d^{(t-1)}$$

One possible combination of μ and σ that satisfies this is $c = 0, \mu = 0, \sigma = 1/\sqrt{d^{(t-1)}}$. This is Glorot Normal initialization, considering only fan_in.

(b)

$$\begin{aligned}
 \mathbf{W}_{ij}^{(t)} &\sim U(\alpha, \beta) \implies \mathbb{E}[\mathbf{W}_i^{(t)}] = \frac{\alpha + \beta}{2} \mathbf{1}^\top, \text{Var}(\mathbf{W}_i^{(t)}) = \frac{(\beta - \alpha)^2}{12} \mathbf{1}^\top \\
 \implies \mathbb{E}[\mathbf{a}_i^{(t)}] &= c, \\
 \text{Var}(\mathbf{a}_i^{(t)}) &= \text{Var}(\mathbf{W}_i^{(t)}) \mathbf{1} + \mathbb{E}[\mathbf{W}_i^{(t)}]^2 \mathbf{1} = \frac{(\beta - \alpha)^2}{12} \mathbf{1}^\top \mathbf{1} + \left(\frac{\alpha + \beta}{2}\right)^2 \mathbf{1}^\top \mathbf{1} \\
 &= \left(\frac{\alpha^2 + \beta^2 - 2\alpha\beta}{12}\right) d^{(t-1)} + \left(\frac{\alpha^2 + \beta^2 + 2\alpha\beta}{4}\right) d^{(t-1)} = \left(\frac{\alpha^2 + \beta^2 - 2\alpha\beta + 3\alpha^2 + 3\beta^2 + 6\alpha\beta}{12}\right) d^{(t-1)} \\
 &= \left(\frac{4\alpha^2 + 4\beta^2 + 4\alpha\beta}{12}\right) d^{(t-1)} = (\alpha^2 + \beta^2 + \alpha\beta) d^{(t-1)} / 3
 \end{aligned}$$

$$\mathbb{E}[\mathbf{a}_i^{(t)}] = 0, \text{Var}(\mathbf{a}_i^{(t)}) = 1 \implies c = 0, (\alpha^2 + \beta^2 + \alpha\beta) d^{(t-1)} / 3 = 1$$

One possible combination of α and β that satisfies this is $c = 0, \alpha = -\sqrt{\frac{3}{d^{(t-1)}}}, \beta = \sqrt{\frac{3}{d^{(t-1)}}}$. This is Glorot Uniform initialization, considering only fan_in.

2. $\mathbb{E}[\mathbf{a}_i^{(t-1)}] = 0, \text{Var}(\mathbf{a}_i^{(t-1)}) = 1$ 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. ReLU activation: $g(x) = \max\{0, x\}$.

We know that $\mathbb{E}[\mathbf{b}_i^{(t)}] = c$ and $\text{Var}(\mathbf{b}_i^{(t)}) = 0 \because \mathbf{b}_i = c \forall 1 < i < d^{(t)}$.

$$\text{Var}(\mathbf{h}_i^{(t)}) = \mathbb{E}[(\mathbf{h}_i^{(t)})^2] - \mathbb{E}[\mathbf{h}_i^{(t)}]^2$$

Here,

$$\begin{aligned}
 \mathbb{E}[(\mathbf{h}_i^{(t)})^2] &= \int_{-\infty}^{+\infty} \max\{0, \mathbf{a}_i^{(t)}\}^2 p(\mathbf{a}_i^{(t)}) d\mathbf{a}_i^{(t)} = \int_0^{+\infty} (\mathbf{a}_i^{(t)})^2 p(\mathbf{a}_i^{(t)}) d\mathbf{a}_i^{(t)} \\
 &= \frac{1}{2} \int_{-\infty}^{+\infty} (\mathbf{a}_i^{(t)})^2 p(\mathbf{a}_i^{(t)}) d\mathbf{a}_i^{(t)} \quad [\because \mathbf{a}_i^{(t)} \text{ is symmetric}] \\
 &= \frac{1}{2} \mathbb{E}[(\mathbf{a}_i^{(t)})^2] = \frac{1}{2} \left(\mathbb{E}[(\mathbf{a}_i^{(t)})^2] - \mathbb{E}[\mathbf{a}_i^{(t)}]^2 \right) \quad [\because \mathbb{E}[\mathbf{a}_i^{(t)}] = 0 \text{ acc. to the question}] \\
 &= \frac{1}{2} \text{Var}(\mathbf{a}_i^{(t)}) = \frac{1}{2} \quad [\because \text{Var}(\mathbf{a}_i^{(t)}) = 1 \text{ acc. to the question}]
 \end{aligned}$$

$$\therefore \text{Var}(\mathbf{h}_i^{(t)}) = \mathbb{E}[(\mathbf{h}_i^{(t)})^2] - \mathbb{E}[\mathbf{h}_i^{(t)}]^2 = \frac{1}{2} - \mathbb{E}[\mathbf{h}_i^{(t)}]^2 \implies \text{Var}(\mathbf{h}^{(t)}) = \frac{1}{2} \mathbf{1}_{d^{(t)}} - \mathbb{E}[\mathbf{h}^{(t)}]^2$$

From the equations derived at the beginning,

$$\begin{aligned}
 \mathbb{E}[\mathbf{a}_i^{(t)}] &= \mathbb{E}[\mathbf{W}_i^{(t)}] \mathbb{E}[\mathbf{h}^{(t-1)}] + \mathbb{E}[\mathbf{b}_i^{(t)}] = \mathbb{E}[\mathbf{W}_i^{(t)}] \mathbb{E}[\mathbf{h}^{(t-1)}] + c \quad [\because \mathbb{E}[\mathbf{b}_i^{(t)}] = c] \\
 \implies 0 &= \mathbb{E}[\mathbf{W}_i^{(t)}] \mathbb{E}[\mathbf{h}^{(t-1)}] + c \\
 \implies \mathbb{E}[\mathbf{W}_i^{(t)}] \mathbb{E}[\mathbf{h}^{(t-1)}] &= -c \\
 \implies (\mathbb{E}[\mathbf{W}_{ij}^{(t)}] \mathbf{1}_{d^{(t-1)}}^\top) (\mathbb{E}[\mathbf{h}_i^{(t-1)}] \mathbf{1}_{d^{(t-1)}}) &= -c \\
 \implies \mathbb{E}[\mathbf{W}_{ij}^{(t)}] \mathbb{E}[\mathbf{h}_i^{(t-1)}] \mathbf{1}_{d^{(t-1)}}^\top \mathbf{1}_{d^{(t-1)}} &= -c \\
 \implies \mathbb{E}[\mathbf{W}_{ij}^{(t)}] \mathbb{E}[\mathbf{h}_i^{(t-1)}] d^{(t-1)} &= -c \\
 \implies \mathbb{E}[\mathbf{h}_i^{(t-1)}] &= -\frac{c}{\mathbb{E}[\mathbf{W}_{ij}^{(t)}] d^{(t-1)}}
 \end{aligned}$$

$$\begin{aligned}
 \text{Var}(\mathbf{a}_i^{(t)}) &= \text{Var}(\mathbf{W}_i^{(t)}) \text{Var}(\mathbf{h}^{(t-1)}) + \text{Var}(\mathbf{W}_i^{(t)}) \mathbb{E}[\mathbf{h}^{(t-1)}]^2 + \mathbb{E}[\mathbf{W}_i^{(t)}]^2 \text{Var}(\mathbf{h}^{(t-1)}) + \text{Var}(\mathbf{b}_i^{(t)}) \\
 \implies 1 &= \text{Var}(\mathbf{W}_i^{(t)}) \left(\frac{1}{2} \mathbf{1}_{d^{(t-1)}} - \mathbb{E}[\mathbf{h}^{(t-1)}]^2 \right) + \text{Var}(\mathbf{W}_i^{(t)}) \mathbb{E}[\mathbf{h}^{(t-1)}]^2 \\
 &\quad + \mathbb{E}[\mathbf{W}_i^{(t)}]^2 \left(\frac{1}{2} \mathbf{1}_{d^{(t-1)}} - \mathbb{E}[\mathbf{h}^{(t-1)}]^2 \right) + 0 \\
 \implies 1 &= \frac{1}{2} \text{Var}(\mathbf{W}_i^{(t)}) \mathbf{1}_{d^{(t-1)}} - \cancel{\text{Var}(\mathbf{W}_i^{(t)}) \mathbb{E}[\mathbf{h}^{(t-1)}]^2} + \cancel{\text{Var}(\mathbf{W}_i^{(t)}) \mathbb{E}[\mathbf{h}^{(t-1)}]^2} \\
 &\quad + \frac{1}{2} \mathbb{E}[\mathbf{W}_i^{(t)}]^2 \mathbf{1}_{d^{(t-1)}} - \mathbb{E}[\mathbf{W}_i^{(t)}]^2 \mathbb{E}[\mathbf{h}^{(t-1)}]^2 \\
 \implies 1 &= \frac{1}{2} \text{Var}(\mathbf{W}_i^{(t)}) \mathbf{1}_{d^{(t-1)}} + \frac{1}{2} \mathbb{E}[\mathbf{W}_i^{(t)}]^2 \mathbf{1}_{d^{(t-1)}} - \mathbb{E}[\mathbf{W}_i^{(t)}]^2 \mathbb{E}[\mathbf{h}^{(t-1)}]^2 \\
 \implies 1 &= \frac{1}{2} \text{Var}(\mathbf{W}_{ij}^{(t)}) \mathbf{1}_{d^{(t-1)}}^\top \mathbf{1}_{d^{(t-1)}} + \frac{1}{2} \mathbb{E}[\mathbf{W}_{ij}^{(t)}]^2 \mathbf{1}_{d^{(t-1)}}^\top \mathbf{1}_{d^{(t-1)}} \\
 &\quad - (\mathbb{E}[\mathbf{W}_{ij}^{(t)}]^2 \mathbf{1}_{d^{(t-1)}}^\top) (\mathbb{E}[\mathbf{h}_i^{(t-1)}]^2 \mathbf{1}_{d^{(t-1)}}) \\
 \implies 1 &= \frac{1}{2} \text{Var}(\mathbf{W}_{ij}^{(t)}) d^{(t-1)} + \frac{1}{2} \mathbb{E}[\mathbf{W}_{ij}^{(t)}]^2 d^{(t-1)} - \mathbb{E}[\mathbf{W}_{ij}^{(t)}]^2 \mathbb{E}[\mathbf{h}_i^{(t-1)}]^2 d^{(t-1)} \\
 \implies 1 &= \frac{1}{2} \text{Var}(\mathbf{W}_{ij}^{(t)}) d^{(t-1)} + \frac{1}{2} \mathbb{E}[\mathbf{W}_{ij}^{(t)}]^2 d^{(t-1)} - \cancel{\mathbb{E}[\mathbf{W}_{ij}^{(t)}]^2} \left(-\frac{c}{\mathbb{E}[\mathbf{W}_{ij}^{(t)}] d^{(t-1)}} \right)^2 d^{(t-1)} \\
 \implies 1 &= \frac{1}{2} \text{Var}(\mathbf{W}_{ij}^{(t)}) d^{(t-1)} + \frac{1}{2} \mathbb{E}[\mathbf{W}_{ij}^{(t)}]^2 d^{(t-1)} - \frac{c^2}{d^{(t-1)}}
 \end{aligned}$$

(a)

$$\begin{aligned} \mathbf{W}_{ij}^{(t)} \sim \mathcal{N}(\mu, \sigma^2) &\implies \mathbb{E}[\mathbf{W}_{ij}^{(t)}] = \mu, \text{Var}(\mathbf{W}_{ij}^{(t)}) = \sigma^2 \\ &\implies 1 = \frac{1}{2} \sigma^2 d^{(t-1)} + \frac{1}{2} \mu^2 d^{(t-1)} - \frac{c^2}{d^{(t-1)}} \\ &\implies c^2 = \frac{1}{2} (\mu^2 + \sigma^2) (d^{(t-1)})^2 - d^{(t-1)} \end{aligned}$$

One possible combination of c , μ and σ that satisfies this is $c = 0$, $\mu = 0$, $\sigma = \sqrt{\frac{2}{d^{(t-1)}}}$. This is He Normal initialization.

(b)

$$\begin{aligned} \mathbf{W}_{ij}^{(t)} \sim U(\alpha, \beta) &\implies \mathbb{E}[\mathbf{W}_{ij}^{(t)}] = \frac{\alpha + \beta}{2}, \text{Var}(\mathbf{W}_{ij}^{(t)}) = \frac{(\beta - \alpha)^2}{12} \\ &\implies 1 = \frac{1}{2} \frac{(\beta - \alpha)^2}{12} d^{(t-1)} + \frac{1}{2} \left(\frac{\alpha + \beta}{2} \right)^2 d^{(t-1)} - \frac{c^2}{d^{(t-1)}} \\ &\implies c^2 = \frac{1}{2} \frac{\alpha^2 + \beta^2 - 2\alpha\beta}{12} (d^{(t-1)})^2 + \frac{1}{2} \frac{\alpha^2 + \beta^2 + 2\alpha\beta}{4} (d^{(t-1)})^2 - d^{(t-1)} \\ &\implies c^2 = \frac{1}{2} \frac{\alpha^2 + \beta^2 - 2\alpha\beta + 3\alpha^2 + 3\beta^2 + 6\alpha\beta}{12} (d^{(t-1)})^2 - d^{(t-1)} \\ &\implies c^2 = \frac{1}{2} \frac{4\alpha^2 + 4\beta^2 + 4\alpha\beta}{12} (d^{(t-1)})^2 - d^{(t-1)} \\ &\implies c^2 = \frac{\alpha^2 + \beta^2 + \alpha\beta}{6} (d^{(t-1)})^2 - d^{(t-1)} \end{aligned}$$

One possible combination of c , μ and σ that satisfies this is $c = 0$, $\alpha = -\sqrt{\frac{6}{d^{(t-1)}}}$, $\beta = \sqrt{\frac{6}{d^{(t-1)}}}$. This is He Uniform initialization.

(c) The values proposed above are He initialization.

(d) This would work well in practice since ReLU has proven to be a very popular non-linearity because of its convenient gradient, and it is desirable to have unit variance of pre-activations of non-linear activations. This was missing in the previous (Glorot) initialization scheme.

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 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 $L(\mathbf{w}) = \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1 - p)\|\Gamma\mathbf{w}\|^2$.
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 λ^{dropout} is a regularization coefficient depending on p . How does the value of p affect the regularization coefficient, λ^{dropout} ?

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

$$2. \mathbb{E}_{\mathbf{R}}[L(\mathbf{w})] = \mathbb{E}_{\mathbf{R}}[\|\mathbf{y} - (\mathbf{R} * \mathbf{X}).\mathbf{w}\|^2] = \sum_{i=1}^n \mathbb{E}_{\mathbf{R}}[(\mathbf{y}_i - (\mathbf{R}_i * \mathbf{X}_i).\mathbf{w})^2]$$

We know (from Assignment 0) that:

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

$$\sum_{i=1}^n (\mathbf{y}_i - (\mathbb{E}_{\mathbf{R}}[\mathbf{R}_i] * \mathbf{X}_i).\mathbf{w})^2 = \sum_{i=1}^n (\mathbf{y}_i - p\mathbf{X}_i.\mathbf{w})^2 = \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2$$

$$\begin{aligned} \sum_{i=1}^n \text{Var}_{\mathbf{R}}(\mathbf{y}_i - (\mathbf{R}_i * \mathbf{X}_i).\mathbf{w}) &= \sum_{i=1}^n \mathbb{E}_{\mathbf{R}}[(\mathbf{y}_i - (\mathbf{R}_i * \mathbf{X}_i).\mathbf{w} - \mathbb{E}_{\mathbf{R}}[\mathbf{y}_i - (\mathbf{R}_i * \mathbf{X}_i).\mathbf{w}])^2] \\ &= \sum_{i=1}^n \mathbb{E}_{\mathbf{R}}[(\mathbf{y}_i - (\mathbf{R}_i * \mathbf{X}_i).\mathbf{w} - \mathbf{y}_i + p\mathbf{X}_i.\mathbf{w})^2] = \sum_{i=1}^n \mathbb{E}_{\mathbf{R}}[(p\mathbf{X}_i - \mathbf{R}_i * \mathbf{X}_i).\mathbf{w}]^2 \\ &= \sum_{i=1}^n \mathbb{E}_{\mathbf{R}}[\mathbf{w}^\top \cdot (p\mathbf{X}_i^\top - \mathbf{R}_i^\top * \mathbf{X}_i^\top) \cdot (p\mathbf{X}_i - \mathbf{R}_i * \mathbf{X}_i).\mathbf{w}] \\ &= \sum_{i=1}^n (\mathbf{w}^\top \cdot (p^2 \mathbf{X}_i^\top \mathbf{X}_i - p\mathbb{E}_{\mathbf{R}}[\mathbf{R}_i] \mathbf{X}_i^\top \mathbf{X}_i - p\mathbb{E}_{\mathbf{R}}[\mathbf{R}_i] \mathbf{X}_i^\top \mathbf{X}_i + \mathbb{E}[\mathbf{R}^\top \mathbf{R}](\mathbf{X}_i^\top \mathbf{X}_i)).\mathbf{w}) \\ &= \sum_{i=1}^n (\mathbf{w}^\top \cdot (p^2 \mathbf{X}_i^\top \mathbf{X}_i - p^2 \mathbf{X}_i^\top \mathbf{X}_i - p^2 \mathbf{X}_i^\top \mathbf{X}_i + p(\mathbf{X}_i^\top \mathbf{X}_i)).\mathbf{w}) \\ &[\because \text{expectation of square of Bernoulli random variable of mean } p \text{ is also } p] \\ &= \sum_{i=1}^n (\mathbf{w}^\top \cdot p(1 - p)(\mathbf{X}_i^\top \mathbf{X}_i).\mathbf{w}) \\ &= p(1 - p)(\mathbf{w}^\top \cdot \text{diag}(\mathbf{X}^\top \mathbf{X}).\mathbf{w}) \\ &= p(1 - p)\|\Gamma\mathbf{w}\|^2 \end{aligned}$$

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

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

$$\begin{aligned} \frac{d\mathbb{E}_{\mathbf{R}}[L(\mathbf{w})]}{d\mathbf{w}} = 0 &\implies 2(-p\mathbf{X}^\top)(\mathbf{y} - p\mathbf{X}\mathbf{w}^{\text{dropout}}) + p(1-p)2\Gamma^\top(\Gamma\mathbf{w}^{\text{dropout}}) = 0 \\ &\implies -2p\mathbf{X}^\top\mathbf{y} + 2p^2\mathbf{X}^\top\mathbf{X}\mathbf{w}^{\text{dropout}} + 2p(1-p)(\Gamma^\top\Gamma)\mathbf{w}^{\text{dropout}} = 0 \end{aligned}$$

$$\begin{aligned} \Gamma^\top\Gamma &= (\text{diag}(\mathbf{X}^\top\mathbf{X})^{1/2})^\top (\text{diag}(\mathbf{X}^\top\mathbf{X})^{1/2}) = \text{diag}(\mathbf{X}^\top\mathbf{X})^{1/2} \text{diag}(\mathbf{X}^\top\mathbf{X})^{1/2} \\ &= (\text{diag}(\mathbf{X}^\top\mathbf{X})^{1/2})^2 = \Gamma^2 \end{aligned}$$

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

$$\therefore p\mathbf{w}^{\text{dropout}} = (\mathbf{X}^\top\mathbf{X} + \lambda^{\text{dropout}}\Gamma^2)^{-1}\mathbf{X}^\top\mathbf{y}, \text{ with } \lambda^{\text{dropout}} = \left(\frac{1}{p} - 1\right).$$

When $p = 1$, $\lambda^{\text{dropout}} = 0 \implies$ no contribution from the regularization term in $L(\mathbf{w})$ involving Γ .

As $p \rightarrow 0$, $\lambda^{\text{dropout}} \rightarrow \infty$, so the contribution from the regularization increases.

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

$$\begin{aligned} \frac{dL(\mathbf{w})}{d\mathbf{w}} = 0 &\implies 2(-\mathbf{X}^\top)(\mathbf{y} - \mathbf{X}\mathbf{w}^{L_2}) + 2\lambda^{L_2}\mathbf{w}^{L_2} = 0 \\ &\implies -\mathbf{X}^\top\mathbf{y} + \mathbf{X}^\top\mathbf{X}\mathbf{w}^{L_2} + \lambda^{L_2}\mathbf{w}^{L_2} = 0 \\ &\implies \mathbf{w}^{L_2} = (\mathbf{X}^\top\mathbf{X} + \lambda^{L_2})^{-1}\mathbf{X}^\top\mathbf{y} \end{aligned}$$

5. From the equation in 2.4, it is clear that the regularization offered by L_2 merely reduces the value of the optimal closed-form weight — by increasing the denominator value by λ^{L_2} , the overall values of the weights are decreased.

However, from the equation in 2.3, it is clear that the regularization offered by dropout is data-dependent. Γ qualitatively represents the standard deviation of the data in each of its dimensions. Hence, in the dimensions where the data varies more, the values of the weights are reduced more.

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) .

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 \mathbf{v}_t to eliminate such a bias by rescaling.

Answer 3. 1. Expressing both update rules in recursive form:

— SGD with momentum: $\mathbf{v}_t = \alpha\mathbf{v}_{t-1} + \epsilon\mathbf{g}_t$, $\Delta\boldsymbol{\theta}_t = -\mathbf{v}_t$, where $\epsilon > 0$ and $\alpha \in (0, 1)$.

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

— SGD with running average of \mathbf{g}_t : $\mathbf{v}_t = \beta\mathbf{v}_{t-1} + (1 - \beta)\mathbf{g}_t$, $\Delta\boldsymbol{\theta}_t = -\delta\mathbf{v}_t$, where $\beta \in (0, 1)$ and $\delta > 0$.

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

Thus, we can see that both update rules are equivalent, with $(\alpha, \epsilon) = (\beta, \delta(1 - \beta))$.

2. Running average update rule has:

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

3. If each sample of the gradient \mathbf{g}_t is independent of time t , we know from 3.2 that:

$$\begin{aligned}\mathbf{v}_t &= (1 - \beta) \sum_{i=1}^t \beta^{t-i} \mathbf{g}_i \\ \mathbb{E}[\mathbf{v}_t] &= \mathbb{E}\left[(1 - \beta) \sum_{i=1}^t \beta^{t-i} \mathbf{g}_i\right] \\ &= (1 - \beta) \sum_{i=1}^t \beta^{t-i} \mathbb{E}[\mathbf{g}_t] \quad [\because \mathbb{E} \text{ of sum is sum of } \mathbb{E}, \text{ and } \beta \text{ is a constant}] \\ &= (1 - \beta) \mathbb{E}[\mathbf{g}_t] \sum_{i=1}^t \beta^{t-i} = (1 - \beta) \mathbb{E}[\mathbf{g}_t] \sum_{i=0}^{t-1} \beta^i = (1 - \beta) \mathbb{E}[\mathbf{g}_t] \frac{(1 - \beta^t)}{(1 - \beta)} \\ &= \mathbb{E}[\mathbf{g}_t] (1 - \beta^t)\end{aligned}$$

Thus, we can see that the running average is biased.

We can see that this bias can be removed if we rescale \mathbf{v}_t by this bias (as is done in Adam for each of the moments calculated):

$$\tilde{\mathbf{v}}_t = \mathbf{v}_t / (1 - \beta^t) \quad \Delta \boldsymbol{\theta}_t = -\delta \tilde{\mathbf{v}}_t$$

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 $\hat{y} = \mathbf{w}^\top \mathbf{x} + \beta$.
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 $^1 \mathbf{W}^\perp \mathbf{u} = \mathbf{0}$.
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).

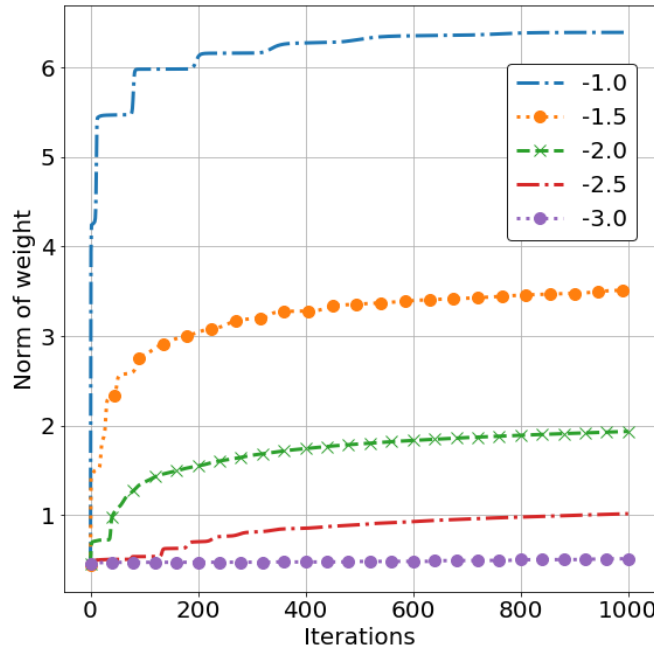


FIGURE 1 – Norm of parameters with different learning rate.

Answer 4. 1. Given: $\mathbf{w} := \gamma \frac{\mathbf{u}}{\|\mathbf{u}\|}$, $\hat{y} = \gamma \cdot \frac{y - \mu_y}{\sigma_y} + \beta$, $y = \mathbf{u}^\top \mathbf{x}$, $\mathbb{E}[\mathbf{x}] = \mathbf{0}$, $\text{Var}(\mathbf{x}) = \mathbf{I}$.

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

$$\text{Var}(y) = \text{Var}(\mathbf{u}^\top \mathbf{x}) = \mathbf{u}^\top \text{Var}(\mathbf{x}) \mathbf{u} = \mathbf{u}^\top \mathbf{I} \mathbf{u} = \mathbf{u}^\top \mathbf{u} = \|\mathbf{u}\|^2 \implies \sigma_y = \|\mathbf{u}\|$$

$$\therefore \hat{y} = \gamma \cdot \frac{y - \mu_y}{\sigma_y} + \beta = \gamma \cdot \frac{\mathbf{u}^\top \mathbf{x} - 0}{\|\mathbf{u}\|} + \beta = \left(\gamma \frac{\mathbf{u}^\top}{\|\mathbf{u}\|} \right) \mathbf{x} + \beta = \mathbf{w}^\top \mathbf{x} + \beta$$

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

2. To show that $\nabla_{\mathbf{u}} L = s \mathbf{W}^\perp \nabla_{\mathbf{w}} L$, given $\mathbf{W}^\perp = \left(\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2} \right)$.

$$\nabla_{\mathbf{u}} L = \nabla_{\mathbf{u}} \mathbf{w} \cdot \nabla_{\mathbf{w}} L = \nabla_{\mathbf{u}} \left(\gamma \frac{\mathbf{u}}{\|\mathbf{u}\|} \right) \cdot \nabla_{\mathbf{w}} L = \gamma \left(\nabla_{\mathbf{u}} \frac{\mathbf{u}}{\|\mathbf{u}\|} \right) \cdot \nabla_{\mathbf{w}} L$$

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

$$\therefore \nabla_{\mathbf{u}} L = \gamma \left(\frac{1}{\|\mathbf{u}\|} \mathbf{W}^\perp \right) \cdot \nabla_{\mathbf{w}} L = \frac{\gamma}{\|\mathbf{u}\|} \mathbf{W}^\perp \nabla_{\mathbf{w}} L = s \mathbf{W}^\perp \nabla_{\mathbf{w}} L, \text{ where } s = \frac{\gamma}{\|\mathbf{u}\|}.$$

3. Weight update using gradient descent:

$$\mathbf{u} \leftarrow \mathbf{u} - \eta \nabla_{\mathbf{u}} L = \mathbf{u} - \eta s \mathbf{W}^\perp \nabla_{\mathbf{w}} L = \mathbf{u} - \frac{\eta \gamma}{\|\mathbf{u}\|} \mathbf{W}^\perp \nabla_{\mathbf{w}} L$$

We know that \mathbf{u} and \mathbf{W}^\perp are orthogonal to each other. Hence, the updated \mathbf{u} is a vector addition of two vectors \mathbf{u} and $-\frac{\eta \gamma}{\|\mathbf{u}\|} \mathbf{W}^\perp \nabla_{\mathbf{w}} L$ that are orthogonal to each other. We know by Pythagoras theorem that the length of the diagonal, is always greater than both the sides of the right-angled triangle. In this case, there is a right angle formed by the vectors \mathbf{u} and $-\frac{\eta \gamma}{\|\mathbf{u}\|} \mathbf{W}^\perp \nabla_{\mathbf{w}} L$, hence their vector sum, the diagonal of said right-angled triangle, is always greater than the previous weight, one of the sides of the triangle. Hence, the norm of the weights must increase with more gradient updates.

However, we also observe that the second vector has $\|\mathbf{u}\|$ in its denominator. This means that as the norm of the weights increases, this denominator also increases, and so the second vector's magnitude decreases. Hence, the magnitude of the vector sum would be changing lesser and lesser from the magnitude of the first vector, \mathbf{u} . Hence, we also observe that the rate of increase of the weight norm decreases with updates.

Finally, we observe that the learning rate is in the numerator of the second vector. Hence, at every update, the magnitude of the second vector is greater for a greater rate. By Pythagoras' theorem, this implies that the magnitude of the sum of this vector with a vector orthogonal to it would also be greater. Hence, the rate of increase of the norm of the weights is greater for greater learning rate.

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

If we assume that $\mathbf{g}_{t-1} = \sigma(\mathbf{h}_{t-1})$, then $\mathbf{g}_t = \sigma(\mathbf{h}_t + \mathbf{W}(\mathbf{0})) = \sigma(\mathbf{h}_t)$

$\therefore \mathbf{g}_t = \sigma(\mathbf{h}_t)$ (assuming/necessitating $\mathbf{g}_0 = \sigma(\mathbf{h}_0)$)

2.

$$\begin{aligned} \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| &= \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \right\| \cdot \left\| \frac{\partial \mathbf{h}_{T-1}}{\partial \mathbf{h}_{T-2}} \right\| \cdots \left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\| \cdots \left\| \frac{\partial \mathbf{h}_1}{\partial \mathbf{h}_0} \right\| \\ \left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\| &= \left\| \frac{\partial (\mathbf{W}\sigma(\mathbf{h}_{t-1}) + \mathbf{U}\mathbf{x}_t + \mathbf{b})}{\partial \mathbf{h}_{t-1}} \right\| = \left\| \mathbf{W} \frac{\partial \sigma(\mathbf{h}_t)}{\partial \mathbf{h}_{t-1}} \right\| = \|\mathbf{W} \sigma'\| \\ &\leq \|\mathbf{W}\| \|\sigma'\| \quad [\because \|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|] \\ &\leq \|\mathbf{W}\| \cdot \gamma \quad [\because |\sigma'| \leq \gamma] \\ &= \sqrt{\lambda_1(\mathbf{W}^\top \mathbf{W})} \cdot \gamma \quad [\because \|\mathbf{A}\| = \sqrt{\lambda_1(\mathbf{A}^\top \mathbf{A})}] \\ &\leq \sqrt{\frac{\delta^2}{\gamma^2}} \cdot \gamma \quad [\because \lambda_1(\mathbf{W}^\top \mathbf{W}) \leq \frac{\delta^2}{\gamma^2}] \\ &= \left| \frac{\delta}{\gamma} \right| \cdot \gamma = \frac{\delta}{\gamma} \cdot \gamma \quad [\because \delta > 0, \gamma > 0] \\ &= \delta \\ \therefore \left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\| \leq \delta &\implies \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| \leq \delta^T \implies \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| \rightarrow 0 \text{ as } T \rightarrow \infty \because 0 \leq \delta < 1 \end{aligned}$$

3. $\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\| \leq \sqrt{\lambda_1(\mathbf{W}^\top \mathbf{W})} \cdot \gamma$; if $\lambda_1(\mathbf{W}^\top \mathbf{W}) > \frac{\delta^2}{\gamma^2}$, this does not say anything about the relation between $\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\|$ and $\frac{\delta^2}{\gamma^2}$. Hence, this offers neither a necessary nor a sufficient condition for gradients to vanish or explode.

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.

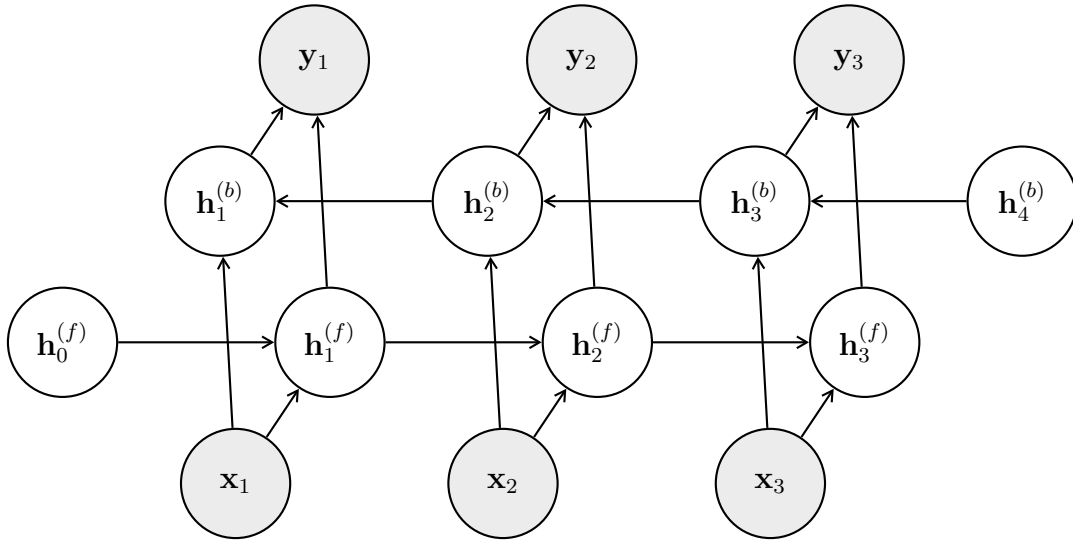
Question 6 (6-12). Denote by σ the logistic sigmoid function. Consider the following Bidirectional RNN:

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

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.
- *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 unfolded bidirectional recurrent neural network for 3 time steps is:



2.

$$\begin{aligned}
 \nabla_{\mathbf{h}_t^{(f)}} L &= \nabla_{\mathbf{h}_t^{(f)}} \sum_{i=t}^T L_i \quad [\because \text{there is no path from } \mathbf{h}_t^{(f)} \text{ to any } L_{i < t}] \\
 &= \nabla_{\mathbf{h}_t^{(f)}} L_t + \nabla_{\mathbf{h}_t^{(f)}} \sum_{i=t+1}^T L_i \\
 &= \nabla_{\mathbf{h}_t^{(f)}} L_t + \nabla_{\mathbf{h}_t^{(f)}} \mathbf{h}_{t+1}^{(f)} \cdot \nabla_{\mathbf{h}_{t+1}^{(f)}} \sum_{i=t+1}^T L_i \\
 &= \nabla_{\mathbf{h}_t^{(f)}} L_t + \nabla_{\mathbf{h}_t^{(f)}} \mathbf{h}_{t+1}^{(f)} \cdot \nabla_{\mathbf{h}_{t+1}^{(f)}} L \quad [\because \nabla_{\mathbf{h}_{t+1}^{(f)}} \sum_{i=t+1}^T L_i = \nabla_{\mathbf{h}_{t+1}^{(f)}} L]
 \end{aligned}$$

$$\begin{aligned}
 \nabla_{\mathbf{h}_t^{(f)}} L_t &= \nabla_{\mathbf{h}_t^{(f)}} \|\mathbf{z}_t - \mathbf{y}_t\|_2^2 = 2\|\mathbf{z}_t - \mathbf{y}_t\|_2 \nabla_{\mathbf{h}_t^{(f)}} (\mathbf{z}_t - (\mathbf{V}^{(f)} \mathbf{h}_t^{(f)} + \mathbf{V}^{(b)} \mathbf{h}_t^{(b)})) \\
 &= 2\|\mathbf{z}_t - \mathbf{y}_t\|_2 (-\mathbf{V}^{(f)}) = -2\mathbf{V}^{(f)} \|\mathbf{z}_t - \mathbf{y}_t\|_2
 \end{aligned}$$

$$\begin{aligned}
 \nabla_{\mathbf{h}_t^{(f)}} \mathbf{h}_{t+1}^{(f)} &= \nabla_{\mathbf{h}_t^{(f)}} (\sigma(\mathbf{W}^{(f)} \mathbf{x}_{t+1} + \mathbf{U}^{(f)} \mathbf{h}_t^{(f)})) \\
 &= \sigma(\mathbf{W}^{(f)} \mathbf{x}_{t+1} + \mathbf{U}^{(f)} \mathbf{h}_t^{(f)}) \cdot (1 - \sigma(\mathbf{W}^{(f)} \mathbf{x}_{t+1} + \mathbf{U}^{(f)} \mathbf{h}_t^{(f)})) \cdot \mathbf{U}^{(f)} \\
 &\quad [\because \nabla_a \sigma(a) = \sigma(a)(1 - \sigma(a)), \text{ and then applying chain rule }] \\
 &= \mathbf{h}_{t+1}^{(f)} (1 - \mathbf{h}_{t+1}^{(f)}) \mathbf{U}^{(f)}
 \end{aligned}$$

$$\therefore \nabla_{\mathbf{h}_t^{(f)}} L = -2\mathbf{V}^{(f)} \|\mathbf{z}_t - \mathbf{y}_t\|_2 + \mathbf{h}_{t+1}^{(f)} (1 - \mathbf{h}_{t+1}^{(f)}) \mathbf{U}^{(f)} \cdot \nabla_{\mathbf{h}_{t+1}^{(f)}} L$$

Similarly,

$$\begin{aligned}
 \nabla_{\mathbf{h}_t^{(b)}} L &= \nabla_{\mathbf{h}_t^{(b)}} \sum_{i=0}^t L_i \quad [\because \text{there is no path from } \mathbf{h}_t^{(b)} \text{ to any } L_{i > t}] \\
 &= \nabla_{\mathbf{h}_t^{(b)}} L_t + \nabla_{\mathbf{h}_t^{(b)}} \sum_{i=0}^{t-1} L_i \\
 &= \nabla_{\mathbf{h}_t^{(b)}} L_t + \nabla_{\mathbf{h}_t^{(b)}} \mathbf{h}_{t-1}^{(b)} \cdot \nabla_{\mathbf{h}_{t-1}^{(b)}} \sum_{i=0}^{t-1} L_i \\
 &= \nabla_{\mathbf{h}_t^{(b)}} L_t + \nabla_{\mathbf{h}_t^{(b)}} \mathbf{h}_{t-1}^{(b)} \cdot \nabla_{\mathbf{h}_{t-1}^{(b)}} L \quad [\because \nabla_{\mathbf{h}_{t-1}^{(b)}} \sum_{i=0}^{t-1} L_i = \nabla_{\mathbf{h}_{t-1}^{(b)}} L]
 \end{aligned}$$

$$\begin{aligned}
 \nabla_{\mathbf{h}_t^{(b)}} L_t &= \nabla_{\mathbf{h}_t^{(b)}} \|\mathbf{z}_t - \mathbf{y}_t\|_2^2 = 2\|\mathbf{z}_t - \mathbf{y}_t\|_2 \nabla_{\mathbf{h}_t^{(b)}} (\mathbf{z}_t - (\mathbf{V}^{(f)} \mathbf{h}_t^{(f)} + \mathbf{V}^{(b)} \mathbf{h}_t^{(b)})) \\
 &= 2\|\mathbf{z}_t - \mathbf{y}_t\|_2 (-\mathbf{V}^{(b)}) = -2\mathbf{V}^{(b)} \|\mathbf{z}_t - \mathbf{y}_t\|_2
 \end{aligned}$$

$$\begin{aligned}
 \nabla_{\mathbf{h}_t^{(b)}} \mathbf{h}_{t-1}^{(b)} &= \nabla_{\mathbf{h}_t^{(b)}} (\sigma(\mathbf{W}^{(b)} \mathbf{x}_{t-1} + \mathbf{U}^{(b)} \mathbf{h}_t^{(b)})) \\
 &= \sigma(\mathbf{W}^{(b)} \mathbf{x}_{t-1} + \mathbf{U}^{(b)} \mathbf{h}_t^{(b)}) \cdot (1 - \sigma(\mathbf{W}^{(b)} \mathbf{x}_{t-1} + \mathbf{U}^{(b)} \mathbf{h}_t^{(b)})) \cdot \mathbf{U}^{(b)} \\
 &\quad [\because \nabla_a \sigma(a) = \sigma(a)(1 - \sigma(a)), \text{ and then applying chain rule }] \\
 &= \mathbf{h}_{t-1}^{(b)} (1 - \mathbf{h}_{t-1}^{(b)}) \mathbf{U}^{(b)}
 \end{aligned}$$

$$\therefore \nabla_{\mathbf{h}_t^{(b)}} L = -2\mathbf{V}^{(b)} \|\mathbf{z}_t - \mathbf{y}_t\|_2 + \mathbf{h}_{t-1}^{(b)} (1 - \mathbf{h}_{t-1}^{(b)}) \mathbf{U}^{(b)} \cdot \nabla_{\mathbf{h}_{t-1}^{(b)}} L$$

$$\nabla_{\mathbf{W}^{(f)}} L = \nabla_{\mathbf{W}^{(f)}} \mathbf{h}_t^{(f)} \cdot \nabla_{\mathbf{h}_t^{(f)}} L \quad [\because \text{the path from } L \text{ to } \mathbf{W}^{(f)} \text{ only involves } \mathbf{h}_t^{(f)} \text{ (and not } \mathbf{h}_t^{(b)} \text{)}]$$

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

$$\implies \nabla_{\mathbf{W}^{(f)}} L = \mathbf{h}_t^{(f)} (1 - \mathbf{h}_t^{(f)}) \mathbf{x}_t (- 2\mathbf{V}^{(f)} \|\mathbf{z}_t - \mathbf{y}_t\|_2 + \mathbf{h}_{t+1}^{(f)} (1 - \mathbf{h}_{t+1}^{(f)}) \mathbf{U}^{(f)} \cdot \nabla_{\mathbf{h}_{t+1}^{(f)}} L)$$

Similarly,

$$\nabla_{\mathbf{U}^{(b)}} L = \nabla_{\mathbf{U}^{(b)}} \mathbf{h}_t^{(b)} \cdot \nabla_{\mathbf{h}_t^{(b)}} L \quad [\because \text{the path from } L \text{ to } \mathbf{U}^{(b)} \text{ only involves } \mathbf{h}_t^{(b)} \text{ (and not } \mathbf{h}_t^{(f)} \text{)}]$$

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

$$\implies \nabla_{\mathbf{U}^{(b)}} L = \mathbf{h}_t^{(b)} (1 - \mathbf{h}_t^{(b)}) \mathbf{h}_{t-1}^{(b)} (- 2\mathbf{V}^{(b)} \|\mathbf{z}_t - \mathbf{y}_t\|_2 + \mathbf{h}_{t-1}^{(b)} (1 - \mathbf{h}_{t-1}^{(b)}) \mathbf{U}^{(b)} \cdot \nabla_{\mathbf{h}_{t-1}^{(b)}} L)$$