

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 preacti

vations 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 $\mathbf{b}^{(t)} = [c, \dots, c]^T$ 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 $1 \leq i \leq d^{(t)}$. Note we are not asking for a general formula; you just need to provide **one setting** that meets these criteria (there are many possibilities).

(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 a value for c , μ , and σ^2 as a function of $d^{(t-1)}$.
 - (b) Uniform: give a value 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 a value for c , μ , and σ^2 as a function of $d^{(t-1)}$.
 - (b) Uniform: give a value 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. In both questions let $c = 0$, $\mu = 0$ and $\alpha = -\beta$ where $\beta > 0$ such that $\mathbb{E}[\mathbf{W}_{ij}^{(t)}] = 0$ and thus the first moment condition $\mathbb{E}[\mathbf{a}_i^{(t)}] = 0$ is satisfied. This also means for $j \neq k$, if we write $a_{ij} = \mathbf{W}_{ij}^{(t)} \mathbf{h}_j^{(t-1)}$,

$$\text{Cov}(\mathbf{W}_{ij}^{(t)} \mathbf{h}_j^{(t-1)}, \mathbf{W}_{ik}^{(t)} \mathbf{h}_k^{(t-1)}) = \mathbb{E}[(a_{ij} - \mathbb{E}[a_{ij}])(a_{ik} - \mathbb{E}[a_{ik}])] = \mathbb{E}[a_{ij} a_{ik}] = 0$$

where we use the independence and zero mean of $\mathbf{W}_{ij}^{(t)}$ and $\mathbf{W}_{ik}^{(t)}$:

$$\mathbb{E}[a_{ij}] = \mathbb{E}[\mathbf{W}_{ij}^{(t)}] \mathbb{E}[\mathbf{h}_j^{(t-1)}] = 0 \quad \mathbb{E}[a_{ij} a_{ik}] = \mathbb{E}[\mathbf{W}_{ij}^{(t)}] \mathbb{E}[\mathbf{W}_{ik}^{(t)}] \mathbb{E}[\mathbf{h}_j^{(t-1)} \mathbf{h}_k^{(t-1)}] = 0$$

Note that the variance of $U(-b, b)$ is $\frac{(b-(-b))^2}{12} = \frac{b^2}{3}$.

1. Now the second moment condition:

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

where the second equality is because $\text{Cov}(\mathbf{W}_{ij}^{(t)} \mathbf{h}_j^{(t-1)}, \mathbf{W}_{ik}^{(t)} \mathbf{h}_k^{(t-1)}) = 0$ for $j \neq k$ (uncorrelated), and the third equality is due to $\mathbf{W}_{ij}^{(t)} \perp \mathbf{h}_j^{(t-1)}$. Since $\mathbf{W}_{ij}^{(t)}$'s are i.i.d., $\text{Var}(\mathbf{W}_{ij}^{(t)}) = \frac{1}{d^{(t-1)}}$. Hence, to satisfy the second-moment constraint, we would need $\sigma^2 = \frac{1}{d^{(t-1)}}$ or $\beta = \sqrt{\frac{3}{d^{(t-1)}}}$.

2. Let $A = \{\mathbf{a}_i^{(t-1)} \geq 0\}$, and $\delta_A \in \{0, 1\}$ be the indicator function that is equal to 1 iff A is satisfied. First, let's observe that

$$\mathbb{E}[(\mathbf{h}_i^{(t-1)})^2] = \mathbb{E}[g(\mathbf{a}_i^{(t-1)})^2] = \mathbb{E}[g(\mathbf{a}_i^{(t-1)})^2 \delta_A] = \mathbb{E}[\mathbf{a}_i^{(t-1)2} \delta_A] = \frac{1}{2} \mathbb{E}[\mathbf{a}_i^{(t-1)2}] = \frac{1}{2}$$

The second last equality is due to the symmetric distribution of $\mathbf{a}_i^{(t-1)}$ and the symmetric nature of the function $x \mapsto x^2$.

Now the second moment condition:

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

Likewise, to have $\text{Var}(\mathbf{W}_{ij}^{(t)}) = \frac{2}{d(t-1)}$, we would need $\sigma^2 = \frac{2}{d(t-1)}$ or $\beta = \sqrt{\frac{6}{d(t-1)}}$.

This is known as the He Initialization (which is the ReLU version of Glorot init). If the parameters are initialized this way, roughly half of the neurons are activated, and the scale of the activations is controlled (i.e. it will not grow as the dimensionality of the input increases).

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 $\mathbb{E}[L(\mathbf{w})] = \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1-p)\|\Gamma\mathbf{w}\|^2$.
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 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. When a particular dropout mask \mathbf{R} is applied, we are dealing with the matrix $\mathbf{X}' = \mathbf{X} \odot \mathbf{R}$. The loss function is

$$L(\mathbf{w}) = \|\mathbf{y} - (\mathbf{X} \odot \mathbf{R})\mathbf{w}\|^2 = \|\mathbf{y} - \mathbf{X}'\mathbf{w}\|^2$$

2. Let \mathbf{x}_i and \mathbf{r}_i be the i 'th rows of \mathbf{X} and \mathbf{R} . Note that $\mathbb{E}[\mathbf{r}_{ij}] = p$ and $\text{Var}(\mathbf{r}_{ij}) = p(1-p)$. Since the quadratic term is a sum of squared error, using the relation $\mathbb{E}[Z^2] = \text{Var}(Z) + \mathbb{E}[Z]^2$ and the fact that \mathbf{r}_{ij} is i.i.d., we have

$$\begin{aligned} \mathbb{E}[L(\mathbf{w})] &= \sum_{i=1}^n \mathbb{E}[(\mathbf{y}_i - \mathbf{w}^\top (\mathbf{x}_i \odot \mathbf{r}_i))^2] \\ &= \sum_{i=1}^n \mathbb{E}[\mathbf{y}_i - \mathbf{w}^\top (\mathbf{x}_i \odot \mathbf{r}_i)]^2 + \text{Var}(\mathbf{y}_i - \mathbf{w}^\top (\mathbf{x}_i \odot \mathbf{r}_i)) \\ &= \sum_{i=1}^n (\mathbf{y}_i - p\mathbf{w}^\top \mathbf{x}_i)^2 + p(1-p) \sum_{j=1}^d (\mathbf{w}_j \mathbf{x}_{ij})^2 \\ &= \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1-p)\|\Gamma\mathbf{w}\|^2 \end{aligned}$$

3. Setting the gradient of the above equation to zero, we have

$$\nabla_{\mathbf{w}} \mathbb{E}[L(\mathbf{w})] = -p \mathbf{X}^\top (\mathbf{y} - p \mathbf{X} \mathbf{w}) + p(1-p) \Gamma^2 \mathbf{w} = 0$$

After rearrangement,

$$p \left(\mathbf{X}^\top \mathbf{X} + \frac{1-p}{p} \Gamma^2 \right) \mathbf{w} = \mathbf{X}^\top \mathbf{y}$$

Letting $\lambda^{\text{dropout}} = \frac{1-p}{p}$ and inverting the matrix on the left yield

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

4. The solution to L^2 regularized linear regression problem is

$$\mathbf{w}^{L^2} = (\mathbf{X}^\top \mathbf{X} + \lambda^{L^2} \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

5. In L^2 the weight is penalized uniformly, whereas in the case of dropout, the penalty on weight parameter \mathbf{w}_i is scaled by the standard deviation of the i 'th feature. When the probability of dropping out the units is higher (lower p), the regularization coefficient grows; when less dropout is applied (larger p), it gets back to the ordinary least square problem, i.e. we can view dropout as scaled L^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) .

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,

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

For SGD with running average,

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

To equalize the two, set $(\alpha, \epsilon) = (\beta, \delta(1 - \beta))$, and reversely $(\beta, \delta) = (\alpha, \frac{\epsilon}{1 - \alpha})$.

2. For SGD with running average we have:

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

3. Taking the expectation of the running average yields

$$\begin{aligned} \mathbb{E}[\mathbf{v}_t] &= \mathbb{E} \left[(1 - \beta) \sum_{i=1}^t \beta^{t-i} \mathbf{g}_i \right] \\ &= (1 - \beta) \sum_{i=1}^t \beta^{t-i} \mathbb{E}[\mathbf{g}_i] && \text{(Linearity)} \\ &= \mathbb{E}[\mathbf{g}_t] (1 - \beta) \sum_{i=1}^t \beta^{t-i} && \text{(Stationarity)} \\ &= \mathbb{E}[\mathbf{g}_t] (1 - \beta^t) && \text{(Telescoping sum)} \end{aligned}$$

One can use $\frac{\mathbf{v}_t}{1 - \beta^t}$ instead as an unbiased estimate.

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

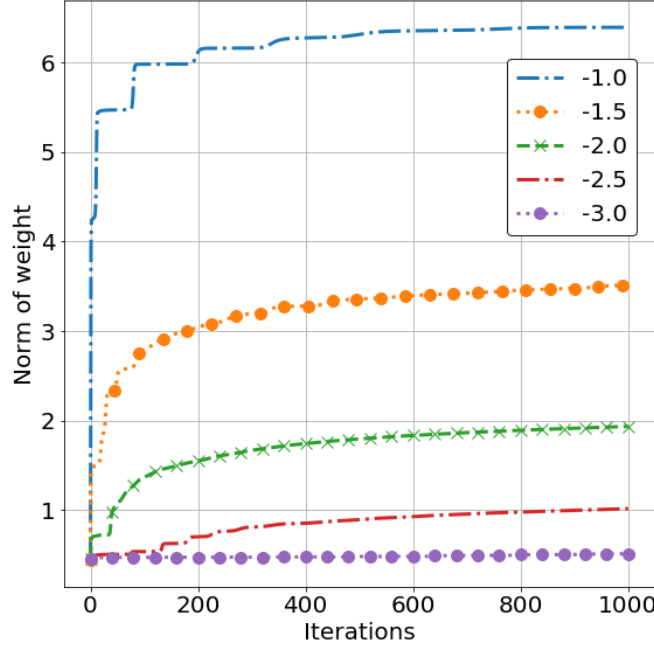


FIGURE 1 – Norm of parameters with different learning rate.

Answer 4.

1. Since the data is centered and whitened, $\mu_y = 0$ and $\sigma_y = \text{Var}(\mathbf{u}^\top \mathbf{x})^{0.5} = \mathbb{E}[\mathbf{u}^\top \mathbf{x} \mathbf{x}^\top \mathbf{u}]^{0.5} = (\mathbf{u}^\top \mathbf{u})^{0.5} = \|\mathbf{u}\|$. Thus,

$$\hat{y} = \gamma \cdot \frac{\mathbf{u}^\top \mathbf{x}}{\|\mathbf{u}\|} + \beta = \mathbf{w}^\top \mathbf{x} + \beta$$

2. The gradient of the norm is the unit vector itself $\nabla_{\mathbf{u}} \|\mathbf{u}\| = \frac{\mathbf{u}}{\|\mathbf{u}\|}$. By the quotient rule, $\frac{\partial}{\partial \mathbf{u}} \frac{\mathbf{u}}{\|\mathbf{u}\|} = \frac{\|\mathbf{u}\| \mathbf{I} - \mathbf{u} \mathbf{u}^\top / \|\mathbf{u}\|}{\|\mathbf{u}\|^2} = \frac{1}{\|\mathbf{u}\|} \left(\mathbf{I} - \frac{\mathbf{u} \mathbf{u}^\top}{\|\mathbf{u}\|^2} \right)$. Now

$$\nabla_{\mathbf{u}} L(\mathbf{u}) = \frac{\partial \mathbf{w}}{\partial \mathbf{u}} \nabla_{\mathbf{w}} L = \frac{\gamma}{\|\mathbf{u}\|} \left(\mathbf{I} - \frac{\mathbf{u} \mathbf{u}^\top}{\|\mathbf{u}\|^2} \right) \nabla_{\mathbf{w}} L = s \mathbf{W}^* \nabla_{\mathbf{w}} L$$

where $s = \frac{\gamma}{\|\mathbf{u}\|}$ and $\mathbf{W}^\perp = \left(\mathbf{I} - \frac{\mathbf{u} \mathbf{u}^\top}{\|\mathbf{u}\|^2} \right)$. Also $\mathbf{W}^\perp \mathbf{u} = \mathbf{u} - \mathbf{u} \frac{\mathbf{u}^\top \mathbf{u}}{\|\mathbf{u}\|^2} = \mathbf{u} - \mathbf{u} = \mathbf{0}$.

3. The gradient descent update is defined as $\mathbf{u}_{t+1} = \mathbf{u}_t - \alpha \nabla_{\mathbf{u}_t} L$, where $\alpha > 0$ is the learning rate. From the last question, we know that $\nabla_{\mathbf{u}_t} L$ is perpendicular to \mathbf{u}_t . Letting $c = \|\alpha \nabla_{\mathbf{u}_t} L\| / \|\mathbf{u}_t\|$, by Pythagorean theorem, we have

$$\|\mathbf{u}_{t+1}\| = \sqrt{\|\mathbf{u}_t\|^2 + \|\alpha \nabla_{\mathbf{u}_t} L\|^2} = \|\mathbf{u}_t\| \sqrt{1 + c^2}$$

which is strictly larger than $\|\mathbf{u}_t\|$ as long as $\alpha > 0$ and $\nabla_{\mathbf{u}_t} L$ has non-zero entries, hence the monotonic growth of $\|\mathbf{u}\|$. Furthermore, since $c \propto \frac{\alpha}{\|\mathbf{u}_t\|^2}$, the larger α is relative to $\|\mathbf{u}_t\|$, the faster the norm grows.

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.

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). More formally, you need to prove it using mathematical induction. You only need to prove the induction step in this question, assuming your expression holds for time step $t - 1$.
- *2. Let $\|\mathbf{A}\|$ denote the L_2 operator norm² of matrix \mathbf{A} ($\|\mathbf{A}\| := \max_{\mathbf{x}: \|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|$). Assume $\sigma(x)$ has bounded derivative, i.e. $|\sigma'| \leq \gamma$ for some $\gamma > 0$ and for all x . We denote as $\lambda_1(\cdot)$ the largest eigenvalue of a symmetric matrix. Show that if the largest eigenvalue of the weights is bounded by $\frac{\delta^2}{\gamma^2}$ for some $0 \leq \delta < 1$, gradients of the hidden state will vanish over time, i.e.

$$\lambda_1(\mathbf{W}^\top \mathbf{W}) \leq \frac{\delta^2}{\gamma^2} \implies \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| \rightarrow 0 \text{ as } T \rightarrow \infty$$

Use the following properties of the L_2 operator norm

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

3. What do you think will happen to the gradients of the hidden state if the condition in the previous question is reversed, i.e. if the largest eigenvalue of the weights is larger than $\frac{\delta^2}{\gamma^2}$? Is this condition *necessary* or *sufficient* for the gradient to explode? (Answer in 1-2 sentences).

Answer 5.

1. Let $\mathbf{g}_t := \sigma(\mathbf{h}_t)$. Assume it is also true that $\mathbf{g}_{t-1} = \sigma(\mathbf{h}_{t-1})$. Then

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

2. For consecutive units, the Jacobian is

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} = \mathbf{W} \frac{\partial \sigma(\mathbf{h}_{t-1})}{\partial \mathbf{h}_{t-1}}$$

Recall the following properties of 2-norm:

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

from which we have

$$\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right\| \leq \|\mathbf{W}\| \left\| \frac{\partial \sigma(\mathbf{h}_{t-1})}{\partial \mathbf{h}_{t-1}} \right\| \leq \frac{\delta}{\gamma} \gamma = \delta$$

which means the 2-norm is bounded by some $\delta \in [0, 1)$. Applying the sub-multiplicativity T times gives

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

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. By contraposition, gradients of the hidden state would not become arbitrarily large if the largest eigenvalue of weights is not larger than $\frac{\delta^2}{\gamma^2}$. Thus it is a necessary condition for gradient explosion. It is not sufficient: the product of the norms can be greater than the norm of the product. This can happen if the hidden state is orthogonal to the largest eigenvector of \mathbf{W} .

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. See Figure 2.

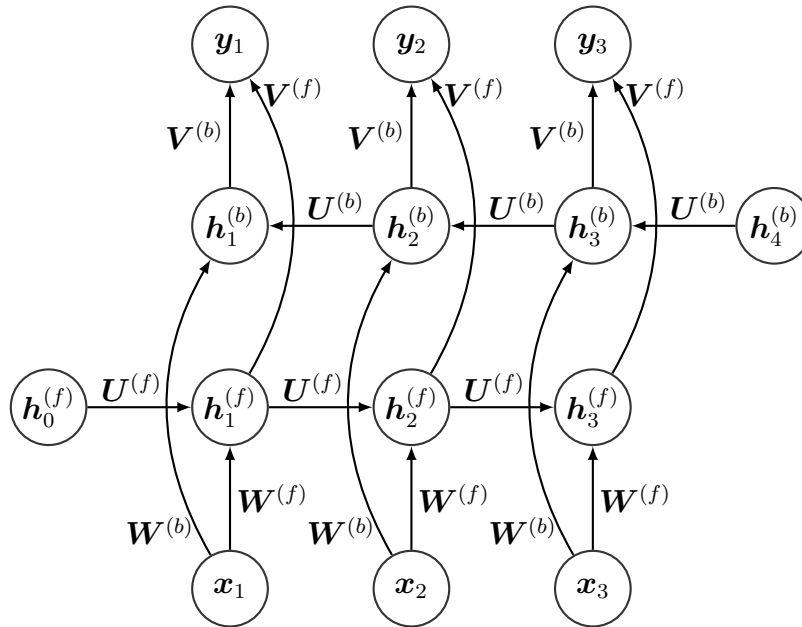


FIGURE 2 – Computational graph of the bidirectional RNN unrolled for three timesteps.

2. The gradients of L_t wrt $\mathbf{h}_t^{(f)}$ and $\mathbf{h}_t^{(b)}$ are

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

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

The Jacobian matrices of the recurrent units are

$$\frac{\partial \mathbf{h}_{t+1}^{(f)}}{\partial \mathbf{h}_t^{(f)}} = \text{diag}(\mathbf{h}_{t+1}^{(f)}(1 - \mathbf{h}_{t+1}^{(f)}))\mathbf{U}^{(f)}$$

$$\frac{\partial \mathbf{h}_{t-1}^{(b)}}{\partial \mathbf{h}_t^{(b)}} = \text{diag}(\mathbf{h}_{t-1}^{(b)}(1 - \mathbf{h}_{t-1}^{(b)}))\mathbf{U}^{(b)}$$

Combining the above, we can express the gradients of L wrt $\mathbf{h}_t^{(f)}$ and $\mathbf{h}_t^{(b)}$ using total derivative:

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

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

Let the subscript t denote the contribution of the weight matrices when computing the recurrent unit indexed by t , e.g. $\mathbf{W}_t^{(f)}$ and $\mathbf{U}_t^{(b)}$. Now, the gradients wrt the i 'th row of $\mathbf{W}_t^{(f)}$ and the i 'th row of $\mathbf{U}_t^{(b)}$ are

$$\nabla_{(\mathbf{W}_t^{(f)})_i} L = \left(\frac{\partial \mathbf{h}_t^{(f)}}{\partial (\mathbf{W}_t^{(f)})_i} \right)^\top \nabla_{\mathbf{h}_t^{(f)}} L = \left(\text{diag}(\mathbf{h}_t^{(f)}(1 - \mathbf{h}_t^{(f)}))(\mathbf{e}_i \mathbf{x}_t^\top) \right)^\top \nabla_{\mathbf{h}_t^{(f)}} L$$

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

where \mathbf{e}_i is a one-hot vector with its i 'th entry being 1. In total, in matrix form,

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