

Due Date: March 22nd 2021, 11pm EST

Question 1 (4-6-4). This question is about activation functions and vanishing/exploding gradients in recurrent neural networks (RNNs). Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be an activation function. When the argument is a vector, we apply σ element-wise. Consider the following recurrent unit:

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

1.1 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 $t - 1$.

If the relations hold for $t - 1$, then we have:

$$\begin{aligned}\mathbf{g}_{t-1} &= \sigma(\mathbf{W}\mathbf{g}_{t-2} + \mathbf{U}\mathbf{x}_{t-1} + \mathbf{b}) \\ \mathbf{h}_{t-1} &= \mathbf{W}\sigma(\mathbf{h}_{t-2}) + \mathbf{U}\mathbf{x}_{t-1} + \mathbf{b}\end{aligned}$$

Now replacing \mathbf{h}_{t-1} into the equation for \mathbf{h}_t we get:

$$\mathbf{h}_t = \mathbf{W}\sigma(\mathbf{h}_{t-1}) + \mathbf{U}\mathbf{x}_t + \mathbf{b} = \mathbf{W}\sigma(\mathbf{W}\sigma(\mathbf{h}_{t-2}) + \mathbf{U}\mathbf{x}_{t-1} + \mathbf{b}) + \mathbf{U}\mathbf{x}_t + \mathbf{b}$$

If we suppose that $\sigma(\mathbf{h}_t) = \mathbf{g}_t$, holds for $t - 1, t - 2$ we'll get:

$$\mathbf{h}_t = \mathbf{W}\sigma(\mathbf{W}\mathbf{g}_{t-2} + \mathbf{U}\mathbf{x}_{t-1} + \mathbf{b}) + \mathbf{U}\mathbf{x}_t + \mathbf{b}$$

Given the relation for $\mathbf{g}_t = \sigma(\mathbf{W}\mathbf{g}_{t-1} + \mathbf{U}\mathbf{x}_t + \mathbf{b})$, we can further simplify to obtain the following:

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

Now if we pass both sides through a sigmoid:

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

Now using the equation for \mathbf{g}_t for the right hand side, we will get:

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

Thus we have proved the inductive bias and shown that $\mathbf{g}_t = \sigma(\mathbf{h}_t)$.

*1.2 Let $\|\mathbf{A}\|$ denote the L_2 operator norm¹ of matrix \mathbf{A} ($\|\mathbf{A}\| := \max_{\mathbf{x}: \|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|$). Assume σ , seen as a function $\mathbb{R}^n \rightarrow \mathbb{R}^n$, has bounded differential, i.e. $\|D\sigma(\mathbf{x})\| \leq \gamma$ (here $\|\cdot\|$ is the L_2 operator norm) for some $\gamma > 0$ and for all \mathbf{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.

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

1. 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.

Use the following properties of the L_2 operator norm

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

From the previous part we have that $\mathbf{g}_t = \sigma(\mathbf{h}_t)$, also we have that $\mathbf{h}_t = \mathbf{W}\sigma(\mathbf{h}_{t-1}) + \mathbf{U}\mathbf{x}_t + \mathbf{b}$. Now taking the derivative w.r.t. \mathbf{h}_0 we arrive at the following (using the chain rule):

$$\frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} = \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \frac{\partial \mathbf{h}_{T-1}}{\partial \mathbf{h}_{T-2}} \frac{\partial \mathbf{h}_{T-2}}{\partial \mathbf{h}_{T-3}} \cdots \frac{\partial \mathbf{h}_1}{\partial \mathbf{h}_0}$$

But for any given $t \in [1, \dots, T]$ we have that:

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

Replacing it back to the equation for $\frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0}$:

$$\begin{aligned} \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} &= \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \frac{\partial \mathbf{h}_{T-1}}{\partial \mathbf{h}_{T-2}} \frac{\partial \mathbf{h}_{T-2}}{\partial \mathbf{h}_{T-3}} \cdots \frac{\partial \mathbf{h}_1}{\partial \mathbf{h}_0} \\ &= \mathbf{W} D\sigma(\mathbf{h}_{T-1}) \mathbf{W} D\sigma(\mathbf{h}_{T-2}) \cdots \mathbf{W} D\sigma(\mathbf{h}_0) \end{aligned}$$

Now applying the property $\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|$ multiple times:

$$\begin{aligned} \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| &= \left\| \mathbf{W} D\sigma(\mathbf{h}_{T-1}) \mathbf{W} D\sigma(\mathbf{h}_{T-2}) \cdots \mathbf{W} D\sigma(\mathbf{h}_0) \right\| \\ &\leq \|\mathbf{W}\| \|D\sigma(\mathbf{h}_{T-1})\| \cdots \|\mathbf{W}\| \|D\sigma(\mathbf{h}_0)\| \\ &\leq \|\mathbf{W}\|^T \gamma^T \end{aligned}$$

Where the last inequality comes from the assumption $\|D\sigma(\mathbf{x})\| \leq \gamma$. We can rewrite this in the following form:

$$\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| \leq (\|\mathbf{W}\|^2 \gamma^2)^{T/2}$$

But from the other matrix operator norm property $\|\mathbf{A}\| = \sqrt{\lambda_1(\mathbf{A}^\top \mathbf{A})}$, we know that $\|\mathbf{W}\|^2 = \lambda_1(\mathbf{W}^\top \mathbf{W})$. So if this largest eigenvalue is bounded by $\frac{\delta^2}{\gamma^2}$, then we have that $\|\mathbf{W}\|^2 \gamma^2 \leq \frac{\delta^2}{\gamma^2} \gamma^2 = \delta^2$. Plugging this back to the derivative inequality:

$$\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| \leq (\delta^2)^{T/2} = \delta^T$$

So the derivative is *upper bounded* by δ^T and this upper bound goes to 0 in the limit:

$$\lim_{T \rightarrow \infty} \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| \leq \lim_{T \rightarrow \infty} \delta^T = 0$$

Last inequality comes from the fact that $0 \leq \delta < 1$. So if the upper bound of the derivative tends to zero as $T \rightarrow \infty$, then consecutive derivatives cannot be larger than the upper bound so in the limit, those derivatives will vanish over time and we have the vanishing gradients problem.

1.3 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* and/or *sufficient* for the gradient to explode? (Answer in 1-2 sentences).

For the gradients to explode, the upper bound on the derivatives should tend to ∞ over time, if the upper bound asymptotically becomes a finite value, then the gradients cannot go beyond that value and we will not have any explosion. So the *necessary* condition is that the upper bound in the previous section becomes *unbounded* itself. Or it's necessary that the upper bound is not zero. Hence if the largest eigenvalue is not smaller than $\frac{\delta^2}{\gamma^2}$, then the upper bound won't tend to zero and gradient explosion is *possible*. It won't certainly explode but now it has the potential to explode so this condition is a *necessary* condition.

Question 2 (1-12-2-6). Suppose that we have a vocabulary containing N possible words, including a special token <BOS> to indicate the beginning of a sentence. Recall that in general, a language model with a full context can be written as

$$p(w_1, w_2, \dots, w_T \mid w_0) = \prod_{t=1}^T p(w_t \mid w_0, \dots, w_{t-1}).$$

We will use the notation $\mathbf{w}_{0:t-1}$ to denote the (partial) sequence (w_0, \dots, w_{t-1}) . Once we have a fully trained language model, we would like to generate realistic sequences of words from our language model, starting with our special token <BOS>. In particular, we might be interested in generating the most likely sequence $\mathbf{w}_{1:T}^*$ under this model, defined as

$$\mathbf{w}_{1:T}^* = \arg \max_{\mathbf{w}_{1:T}} p(\mathbf{w}_{1:T} \mid w_0 = \text{<BOS>}).$$

For clarity we will drop the explicit conditioning on w_0 , assuming from now on that the sequences always start with the <BOS> token.

2.1 How many possible sequences of length $T + 1$ starting with the token <BOS> can be generated in total? Give an exact expression, without the O notation. Note that the length $T + 1$ here includes the <BOS> token.

According to the discussions on piazza, <BOS> can appear any where in the middle of the sequence, so we got T possible spots and each spot has the possibility to be filled with any of the N words in the vocabulary (which includes <BOS> too). So this gives N^T possible sequences of length $T + 1$ starting with token <BOS>.

2.2 In this question only, we will assume that our language model satisfies the *Markov property*

$$\forall t > 0, p(w_t \mid w_0, \dots, w_{t-1}) = p(w_t \mid w_{t-1}).$$

Moreover, we will assume that the model is time-homogeneous, meaning that there exists a matrix $\mathbf{P} \in \mathbb{R}^{N \times N}$ such that $\forall t > 0, p(w_t = j \mid w_{t-1} = i) = [\mathbf{P}]_{i,j} = P_{ij}$.

2.2.a Let $\delta_t \in \mathbb{R}^N$ be the vector of size N defined for $t > 0$ as

$$\delta_t(j) = \max_{\mathbf{w}_{1:t-1}} p(\mathbf{w}_{1:t-1}, w_t = j),$$

where $\mathbf{w}_{1:0} = \emptyset$. Using the following convention for δ_0

$$\delta_0(j) = \begin{cases} 1 & \text{if } j = \text{<BOS>} \\ 0 & \text{otherwise,} \end{cases}$$

give the recurrence relation satisfied by δ_t (for $t > 0$).

Starting by the definition and using $P(A, B) = P(A | B)P(B)$:

$$\delta_t(j) = \max_{w_{1:t-1}} p(w_{1:t-1}, w_t = j) = \max_{w_{1:t-1}} p(w_t = j | w_{1:t-1})p(w_{1:t-1})$$

Now employing the Markov assumption on $p(w_t = j | w_{1:t-1})$ we get:

$$\delta_t(j) = \max_{w_{1:t-1}} p(w_t = j | w_{1:t-1})p(w_{1:t-1}) = \max_{w_{1:t-1}} p(w_t = j | w_{t-1})p(w_{1:t-1})$$

But we know that the product of a sequence of numbers is maximized when each of them is at its maximum, so $\max_{\text{args}} f(\text{args})g(\text{args}) = \max_{\text{args}} f(\text{args}) \max_{\text{args}} g(\text{args})$, so we can rewrite our maximization as follows:

$$\delta_t(j) = \max_{w_{1:t-1}} p(w_t = j | w_{t-1})p(w_{1:t-1}) = \max_{w_{1:t-1}} p(w_t = j | w_{t-1}) \max_{w_{1:t-1}} p(w_{1:t-1})$$

But notice that $\max_{w_{1:t-1}} p(w_t = j | w_{t-1})$ can only be affected by w_{t-1} among the max arguments $w_{1:t-1}$, so we can simplify that:

$$\begin{aligned} \delta_t(j) &= \max_{w_{1:t-1}} p(w_t = j | w_{t-1}) \max_{w_{1:t-1}} p(w_{1:t-1}) \\ &= \max_{w_{t-1}} p(w_t = j | w_{t-1}) \max_{w_{1:t-1}} p(w_{1:t-1}) \end{aligned}$$

Now notice how we can expand the second term above:

$$\begin{aligned} \max_{w_{1:t-1}} p(w_{1:t-1}) &= \max_{w_{1:t-2}, w_{t-1}} p(w_{1:t-2}, w_{t-1}) \\ &= \max_{w_{t-1}} \delta_{t-1}(w_{t-1}) \end{aligned}$$

Rewriting the δ_t again with this new observation and noting that $p(w_t = j | w_{t-1}) = P_{w_{t-1}w_t}$ and mixing the max operators (since they're the same):

$$\begin{aligned} \delta_t(w_t = j) &= \max_{w_{t-1}} p(w_t = j | w_{t-1}) \max_{w_{1:t-1}} p(w_{1:t-1}) \\ &= \max_{w_{t-1}} P_{w_{t-1}j} \delta_{t-1}(w_{t-1}) \end{aligned}$$

So after cleaning up a bit, we have arrived at our recurrence relation:

$$\delta_t(j) = \max_i P_{ij} \delta_{t-1}(i)$$

Also note this recurrence has to start with $t = 1$, so $\delta_1(j) = \max_i P_{ij} \delta_0(i)$, and we need an initialization for δ_0 . But $\delta_0(i)$ is given in the problem statement and is zero anywhere $i \neq \text{<BOS>}$, so the maximum value for $\max_i P_{ij} \delta_0(i)$ is $P_{\text{<BOS>}j} = \delta_1(j)$.

2.2.b Show that $w_T^* = \arg \max_j \delta_T(j)$.

w_T^* is the T^{th} element of the most likely sequence $\mathbf{w}_{1:T}^* = \arg \max_{\mathbf{w}_{1:T}} p(\mathbf{w}_{1:T})$ (which of course is conditioned on $w_0 = \langle \text{BOS} \rangle$ but the notation is dropped for convenience), but we know how to maximize the probability inside the $\arg \max$:

$$\max_{\mathbf{w}_{1:T}} p(\mathbf{w}_{1:T}) = \max_{w_T} \max_{\mathbf{w}_{1:T-1}} p(\mathbf{w}_{1:T}) = \max_{w_T} \max_{\mathbf{w}_{1:T-1}} p(\mathbf{w}_{1:T-1}, w_T)$$

But recall the definition of $\delta_t(j) = \max_{\mathbf{w}_{1:t-1}} p(\mathbf{w}_{1:t-1}, w_t = j)$, we can directly plug this back into the above equation to get:

$$\max_{\mathbf{w}_{1:T}} p(\mathbf{w}_{1:T}) = \max_{w_T} \max_{\mathbf{w}_{1:T-1}} p(\mathbf{w}_{1:T-1}, w_T) = \max_{w_T} \delta_T(w_T)$$

So we know that the probability for the most likely sequence is $\max_{w_T} \delta_T(w_T)$, and thus the sequence ending with the argument that maximizes this probability is the most likely sequence, and it ends with that argument, hence w_T^* is the argument which maximizes this term $\max_{w_T} \delta_T(w_T)$, so it's the $\arg \max$ of that or $w_T^* = \arg \max_j \delta_T(j)$.

2.2.c Let $\mathbf{a}_t \in \{1, \dots, N\}^N$ be the vector of size N defined for $t > 0$ as

$$a_t(j) = \arg \max_i P_{ij} \delta_{t-1}(i).$$

Show that $\forall 0 < t < T$, $w_t^* = a_{t+1}(w_{t+1}^*)$.

To solve this part, we need to be careful with the definitions. w_t^* is the t^{th} element of the most likely sequence $\mathbf{w}_{1:T}^*$. What is $\delta_t(j)$? It's the maximum possible probability that the t^{th} word is the j^{th} word of the vocabulary. $\arg \max_j \delta_T(j)$ gives the most likely word for ending the most likely sequence at position T . Now what is the most likely word at one position before or at $(T-1)^{th}$?

$$\delta_t(j) = \max_i P_{ij} \delta_{t-1}(i)$$

But note that any $\delta_t(j) = \max_i P_{ij} \delta_{t-1}(i)$ maximizes over all possible *previous words* to make the probability of having word j at the next position the highest. In particular, this holds for w_T^* as well so i that maximizes $\delta_T(w_T^*) = \max_i P_{iw_T^*} \delta_{T-1}(i)$, is the most likely word preceding w_T^* or equivalently we can say $w_{T-1}^* = \arg \max_i P_{iw_T^*} \delta_{T-1}(i)$. Considering this as the base case of an induction proof, we can prove the induction step as follows:

Suppose that for some $t < T-1$ we have $w_t^* = \arg \max_i P_{iw_{t+1}^*} \delta_t(i)$, then we need to show this holds for $t-1$ as well because in this case we're traversing backwards, so the induction step is also decreasing. Again using a similar argument, $\delta_t(j)$ is the maximum possible probability that the t^{th} word is the j^{th} word of the vocabulary, so if this probability is maximized when the t^{th} word is w_t^* , then this maximum value for $\delta_t(j)$ is achieved by $\delta_t(w_t^*) = \max_i P_{iw_t^*} \delta_{t-1}(i)$. Again this maximization is over all possible *previous words* to make the probability of having word w_t^* at the next position the highest so i that maximizes $\max_i P_{iw_t^*} \delta_{t-1}(i)$, is the most likely word preceding w_t^* or equivalently we can say $w_{t-1}^* = \arg \max_i P_{iw_t^*} \delta_{t-1}(i)$. Thus we have completed the induction step since we showed:

$$w_t^* = \arg \max_i P_{iw_{t+1}^*} \delta_t(i) \implies w_{t-1}^* = \arg \max_i P_{iw_t^*} \delta_{t-1}(i)$$

Now simply replacing the introduced operator $a_t(j) = \arg \max_i P_{ij} \delta_{t-1}(i)$, we get:

$$w_t^* = a_{t+1}(w_{t+1}^*) \implies w_{t-1}^* = a_t(w_t^*)$$

And our inductive prove holds for this too since it's just a rewritten version of our derivation, thus our proof concludes and we have $w_t^* = a_{t+1}(w_{t+1}^*)$.

- 2.2.d Using the previous questions, write the pseudo-code of an algorithm to compute the sequence $\mathbf{w}_{1:T}^*$ using dynamic programming. This is called *Viterbi decoding*.

Algorithm 1: Viterbi decoding

Input: A Markovian language model or matrix $\mathbf{P} \in \mathbb{R}^{N \times N}$ such that

$$\forall t > 0, p(w_t = j \mid w_{t-1} = i) = [\mathbf{P}]_{i,j} = P_{ij}$$

Output: The most likely sequences $\mathbf{w}_{1:T}$

Initialization: $\delta_1(j) \leftarrow P_{\text{BOS} \rightarrow j}$ for all $j \in \{1, \dots, N\}$

Do the forward pass, i.e.:

for $t = 1$ **to** T **do**

for $j = 1$ **to** N **do**

$$\quad \delta_t(j) \leftarrow \max_i P_{ij} \delta_{t-1}(i)$$

Compute $w_T^* \leftarrow \arg \max_i \delta_T^*(i)$

Create an empty list \mathbf{w} and add w_T^* at the end of that list

Recursively do the backward pass, i.e.:

for $t = T - 1$ **to** 1 **do**

$$\quad w_t^* \leftarrow a_{t+1}(w_{t+1}^*) \text{ (} w_{t+1}^* \text{ is the first element of } \mathbf{w} \text{ at any } t \text{)}$$

 Append w_t^* to the beginning of \mathbf{w}

Return the list \mathbf{w} as the desired sequence $\mathbf{w}_{1:T}^*$

- 2.2.e What is the time complexity of this algorithm ? Its space complexity ? Write the algorithmic complexities using the O notation. Comment on the efficiency of this algorithm compared to naively searching for $\mathbf{w}_{1:T}^*$ by enumerating all possible sequences (based on your answer to question 1).

The time complexity of Viterbi decoding:

- i. Initialization $\rightarrow O(N)$
- ii. Forward pass: Computing $P_{ij} \delta_{t-1}(i)$ for a given j and for all $i \in \{1, \dots, N\}$ requires $O(N)$ time. Finding the maximum over N elements has the worst case complexity of $O(N)$ because we can't find deterministically find the max without looking at all elements. Since the for loop is over all j, t , then the time complexity would be $O(N + N) \times O(TN) = O(TN^2)$
- iii. Computing w_T^* requires doing an $\arg \max$ over N possible words, which again like the max itself will take a worst case of $O(N)$ time.
- iv. Creating a list and appending w_T^* to that $\rightarrow O(1)$
- v. Backward pass: for each t , there's an $\arg \max$ taking $O(N)$ in the worst case, appending takes $O(1)$, and since there are $T - 1$ of such steps, the backward pass will contribute to a time complexity of $O(N) \times O(1) \times O(T) = O(NT)$

So the overall time complexity would be $O(N) + O(TN^2) + O(N) + O(1) + O(NT) = O(TN^2)$

The memory complexity of Viterbi decoding:

- i. keeping the matrix $\mathbf{P} \in \mathbb{R}^{N \times N} \rightarrow O(N^2)$
- ii. Initialization and keeping $\delta_1(j)$ for all $j \in \{1, \dots, N\} \rightarrow O(N)$
- iii. Forward pass: Computing $P_{ij}\delta_{t-1}(i)$ for a given j and for all $i \in \{1, \dots, N\}$ requires keeping $\delta_{t-1}(i)$ at any t for all words, but once computed, we do need to keep them for the backward pass and operator $a_t(j)$. Furthermore, we have to keep the products $P_{ij}\delta_{t-1}(i)$ which takes $O(N)$ memory, but unlike δ s, after finding the maximum, we only keep a single value for each j . So it takes $O(N)$ memory. This would be the overall memory complexity of the forward pass as the for loop doesn't affect it because we continuously replace those values and don't add anything to the memory. But for keeping all δ s we need $O(NT)$ memory. So the total memory complexity would be $O(NT)$ for δ s and a reusable $O(N)$ for finding the maximum for each j .
- iv. Computing w_T^* requires doing an arg max over N possible words, which again like the max itself will take a memory of $O(N)$ to keep all values and go over them.
- v. Creating a list and appending w_T^* to that $\rightarrow O(1)$ memory complexity
- vi. Backward pass: for each t , only having \mathbf{w} in addition to the δ s that we have kept is enough. The $a_t(j)$ operator requires keeping N products and finding the arg max over them takes another $O(N)$ in the worst case, but this won't add up for consecutive iterations and can be replaced. Appending takes $O(1)$. So in total, the backward pass will contribute to a memory complexity of $O(N)$

So the overall memory complexity would be $O(N^2) + O(N) + O(NT) + O(N) + O(N) + O(1) + O(N) = O(N^2)$. The major component contributing to the memory complexity is the matrix \mathbf{P} .

Using the naive approach, requires computing N^T probabilities for sequences of length T so if each probability computation takes $O(T)$, then the total computation for all the probabilities would give rise to an *exponential* time complexity of $O(TN^T)$. To compute the maximum, we can enumerate all of these N^T probabilities and at any time keep the highest value so far (taking another worst case $O(N^T)$ time complexity). But to find the most likely sequence, we need to keep all N^T probabilities (hence $O(N^T)$), and all possible sequences corresponding to those probabilities so it'll consume $O(N^T) + O(N^T \times T) = O(TN^T)$. (Each sequence requires a memory of $O(T)$.)

Comparing the approaches:

Method	Time Complexity	Memory Complexity
Naive	$O(TN^T)$	$O(TN^T)$
Viterbi	$O(TN^2)$	$O(N^2)$

So we observe a significant improvement from the exponential dependence on the very huge number N to a squared dependence to N which is a significant improvement. Also for the memory complexity, unless $T \not\geq N$ (which always is!), then the memory complexity doesn't have the T factor which could be another advantage. So Viterbi decoding is by far more efficient than the naive approach of enumerating all possibilities.

However one can define an ordering for constructing all N^T sequences given the index, and thus at each time computes a likelihood and just keeps the highest value so far and the index of that sequence, then the sequence of length T can be easily generated if we know

how to create it from the index. In this case, the space complexity of the naive approach would just be $O(T)$ so Viterbi will only be advantageous only in terms of time complexity.

- 2.2.f The component essential for the Viterbi algorithm is the time-homogeneous Markov assumption that constitutes those recursions we derived earlier. Without a Markov assumption, we need to enumerate all input possibilities at each step and use the multiplication of these conditionals to obtain the joint. So LSTM and GRU cannot exploit the Viterbi algorithm because they model the full conditional $p(w_t | w_{1:t-1})$ thus they won't satisfy any n^{th} order Markov property and cannot be used for Viterbi decoding as there will be no recursion. But a simple old RNN that satisfies some n^{th} order Markov property can yield a recursion (by defining a larger matrix \mathbf{P}) and then we can exploit the Viterbi algorithm with such language model.

For exact inference in LSTM/GRU, we need to enumerate all input possibilities at each step and use the multiplication of these conditionals to obtain the joint. Once we have that, we can take the maximum and find the most likely sequence. How to obtain these joints? This is possible but tremendously painful. Suppose we're interested in finding the most likely sequence of length T . At any time $t \leq T$, we should feed all possible sequences of length $t - 1$ to the LSTM since now it is responsible for returning $p(w_t | w_{1:t-1})$. But these are conditionals, what about the joints? At time 1 we have $p(w_1 | \langle \text{BOS} \rangle)$, multiplying $p(\langle \text{BOS} \rangle)p(w_1 | \langle \text{BOS} \rangle) = p(2_1, w_0)$. We should keep it for all words in the vocabulary. Then for $t = 2$ we should compute all possible conditionals (by feeding all possible combinations) $p(w_2 | w_{0:1})$, then multiply them by the corresponding probability $p(w_{0:1})$ to get all possible joints $p(w_{0:2})$. Continuing like that we can get all joints. Then we can take the maximum over the joint at $t = T$ to find the most likely sequence. It's more finding the exact solution of maximum likelihood by enumerating all possibilities, not the Viterbi algorithm.

Algorithm 2: Beam search decoding

Input: A language model $p(\mathbf{w}_{1:T} | w_0)$, the beam width B

Output: B sequences $\mathbf{w}_{1:T}^{(b)}$ for $b \in \{1, \dots, B\}$

Initialization: $w_0^{(b)} \leftarrow \langle \text{BOS} \rangle$ for all $b \in \{1, \dots, B\}$

Initial log-likelihoods: $l_0^{(b)} \leftarrow 0$ for all $b \in \{1, \dots, B\}$

for $t = 1$ **to** T **do**

for $b = 1$ **to** B **do**

for $j = 1$ **to** N **do**

$s_b(j) \leftarrow l_{t-1}^{(b)} + \log p(w_t = j | \mathbf{w}_{0:t-1}^{(b)})$

for $b = 1$ **to** B **do**

 Find (b', j) such that $s_{b'}(j)$ is the b -th largest score

 Save the partial sequence b' : $\tilde{\mathbf{w}}_{0:t-1}^{(b)} \leftarrow \mathbf{w}_{0:t-1}^{(b')}$

 Add the word j to the sequence b : $w_t^{(b)} \leftarrow j$

 Update the log-likelihood: $l_t^{(b)} \leftarrow s_{b'}(j)$

 Assign the partial sequences: $\mathbf{w}_{0:t-1}^{(b)} \leftarrow \tilde{\mathbf{w}}_{0:t-1}^{(b)}$ for all $b \in \{1, \dots, B\}$

- 2.3 What is the time complexity of Algorithm 2? Its space complexity? Write the algorithmic complexities using the O notation, as a function of T , B , and N . Is this a practical decoding algorithm when the size of the vocabulary is large?

The time complexity of Beam search decoding:

- Initialization of $w_0^{(b)}$ and $l_0^{(b)} \rightarrow O(B) + O(B) = O(B)$
- Computing $s_b(j)$ s for a given j, b takes $O(1)$ for computing the output of the language model $\log p(w_t = j \mid \mathbf{w}_{0:t-1}^{(b)})$, and $O(1)$ for addition, so for all $j \in \{1, \dots, N\}$ and all $b \in \{1, \dots, B\}$ it requires $O(BN)$ time.
- Finding the b -th largest score and the corresponding j requires a search over BN elements ($s_b(j)$) which takes a worst case of $O(BN \log BN)$ time complexity to be sorted using merge sort, then finding the b -th largest score takes $O(1)$.
- Saving and adding the word to a sequence take $O(1)$ time complexity.
- Updating the log-likelihood also takes $O(1)$ for a given t, b, j .
- So each iteration of the second loop over b takes $O(BN \log BN)$, thus the complete for loop take $O(B) \times O(BN \log BN) = O(B^2 N \log BN)$
- Assigning the partial sequence of length t takes $O(T)$ in the worst case, and doing that for all b will result in a factor of $O(B)$, hence $O(BT)$.

So for each outer loop or for each t , we have a time complexity of

$$O(B) + O(BN) + O(B^2 N \log BN) + O(1) + O(BT) = O(B(N + T + BN \log BN))$$

Now taking into account that there are T iterations, the overall time complexity would be larger by a factor of T : $O(BT(N + T + BN \log BN))$

The memory complexity of Beam search decoding:

- Initialization and keeping $w_0^{(b)}$ and $l_0^{(b)}$ for all $b \in \{1, \dots, B\} \rightarrow O(B) + O(B) = O(B)$
- At any t , we need to have kept all $l_{t-1}^{(b)}$ which takes $O(B)$ memory and we considered that previously when initializing $l_0^{(b)}$ s, but we don't need to keep them for all t , just the previous step suffices. But to keep the values of $s_b(j)$ we need $O(BN)$ memory which are reusable for each iteration of t , meaning that we don't need to allocate new memory for each iteration, we can discard previous values and store the new ones.
- Finding (b', j) requires the same $O(BN)$ allocated memory above, nothing extra.
- Saving the partial sequence in the worst case requires a memory of $O(BT)$ because in the worst case we need to store b beams of length T .
- Updating the log-likelihood doesn't need extra memory as we'll replace the previous values.

So the overall memory complexity would be $O(B) + O(BN) + O(BT) = O(B(N + T))$.

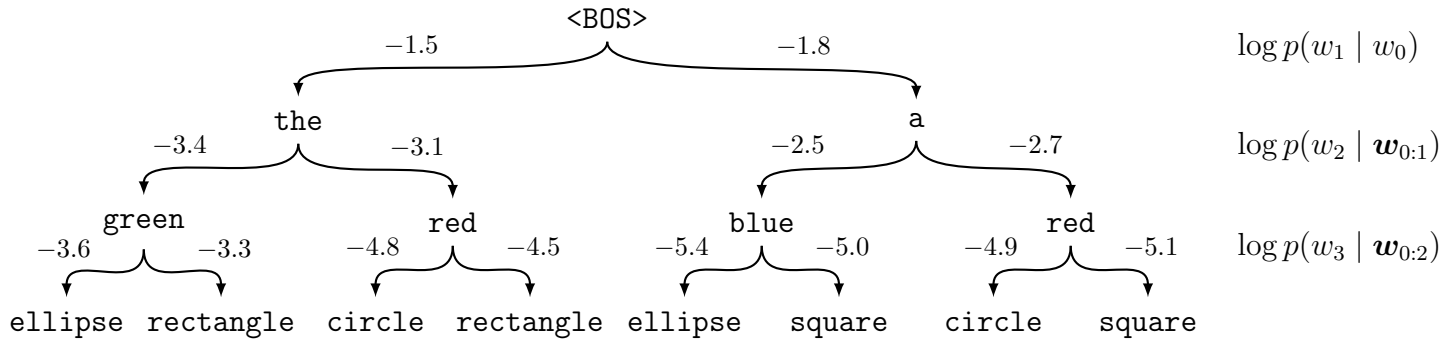
Now comparing the approaches:

Method	Time Complexity	Memory Complexity
Naive	$O(TN^T)$	$O(TN^T)$
Viterbi	$O(TN^2)$	$O(N^2)$
Beam Search	$O(BT(N + T + BN \log BN))$	$O(B(N + T))$

Thus indeed it's a practical approach compared to the naive and Viterbi decoding approaches specially when N is large. The reason is that its time complexity is not polynomial or exponential

in N but only in the worst case grows with $N \log N$ which is significantly better than N^2 and N^T that correspond to the Viterbi and naive approaches, respectively. The same goes for the memory complexity as well, instead of being exponential and polynomial N^T , N^2 , it only linearly grows with N which is greatly advantageous.

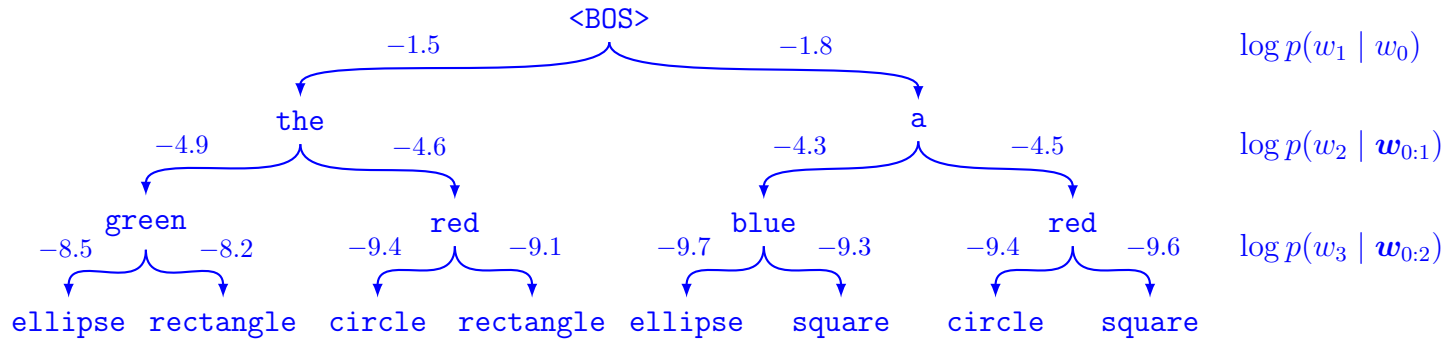
- 2.4 The different sequences that can be generated with a language model can be represented as a tree, where the nodes correspond to words and edges are labeled with the log-probability $\log p(w_t | \mathbf{w}_{0:t-1})$, depending on the path $\mathbf{w}_{0:t-1}$. In this question, consider the following language model (where the low probability paths have been removed for clarity)



- 2.4.a If you were given the whole tree, including the log-probabilities of all the missing branches (e.g. $\log p(w_2 = \text{a} | w_0 = \text{<BOS>, } w_1 = \text{red})$), could you apply Viterbi decoding from question 2 to this language model in order to find the most likely sequence $\mathbf{w}_{1:3}^*$? Why? Why not? Find $\mathbf{w}_{1:3}^*$, together with its corresponding log-likelihood $\log p(\mathbf{w}_{1:3}^*) = \max_{\mathbf{w}_{1:3}} \log p(\mathbf{w}_{1:3})$.

No, we can't, since the component essential for the Viterbi algorithm is the time-homogeneous Markov assumption. Here, the probability for every word depends on all previous words up to root so we cannot have any recursion and we need to do exact inference, keeping the track of all possible joints, and choosing the path from root to leaf having the highest joint log-likelihood. So here we need the joint probabilities not the conditionals. But we can easily obtain those joints, by traversing the tree and recursively adding the log probability of the parent conditional to that of the children and storing that as the new node attribute. This is because the log conditional for the root is $\log p(w_0 = \text{<BOS>}) = \log 1 = 0$, then for the first descendants of the root we have $\log p(w_1, w_0) = \log p(w_1 | w_0) + \log p(w_0) = \log p(w_1 | w_0) + 0 = \log p(w_1 | w_0)$. For next descendants, we need the joint $\log p(w_2, \mathbf{w}_{0:1}) = \log p(w_2 | \mathbf{w}_{0:1}) + \log p(w_1, w_0)$, but we got $\log p(w_1, w_0) = \log p(w_1 | w_0)$, so $\log p(w_2, \mathbf{w}_{0:1}) = \log p(w_2 | \mathbf{w}_{0:1}) + \log p(w_1 | w_0)$. The same goes for any descendant in the tree because the joint is obtainable by adding the log-conditional given previous words and the log-joint of the previous words.

Now that we have the joints, we can compute $p_t(w_t | w_{1:t-1})$ at each depth of the tree (corresponding to different t s) and then choose the path from root to some leaf having the highest summation of log-conditionals corresponding to the highest log-likelihood. Below is the same tree with node attributes being replaced by joints obtained by summing parent values to that of the children:



If we call the joints δ s, then, $w_3^* = \arg \max_j \delta_3(j)$, so let's see $\delta_3(j)$'s. At position 3, we have 4 possible words and the maximum value for the joint for each of them is (note that we have log of δ s instead of δ s which is not an issue since the log function is monotonic and a bijective):

$$\begin{aligned}\delta_3(\text{ellipse}) &= \max\{-8.5, -9.7\} = -8.5 \\ \delta_3(\text{rectangle}) &= \max\{-8.2, -9.1\} = -8.2 \\ \delta_3(\text{circle}) &= \max\{-9.4, -9.4\} = -9.4 \\ \delta_3(\text{square}) &= \max\{-9.3, -9.6\} = -9.3\end{aligned}$$

This leaves us with $w_3^* = \arg \max_j \delta_3(j) = \text{rectangle}$. Tracing back this leaf to the root, we get:

$$\mathbf{w}_{0:3}^* = \{\text{<BOS>, the, green, rectangle}\}$$

2.4.b *Greedy decoding* is a simple algorithm where the next word \bar{w}_t is selected by maximizing the conditional probability $p(w_t | \bar{\mathbf{w}}_{0:t-1})$ (with $\bar{\mathbf{w}}_0 = \text{<BOS>}$)

$$\bar{w}_t = \arg \max_{w_t} \log p(w_t | \bar{\mathbf{w}}_{0:t-1}).$$

Find $\bar{\mathbf{w}}_{1:3}$ using greedy decoding on this language model, and its log-likelihood $\log p(\bar{\mathbf{w}}_{1:3})$. Now it's even simpler and at each branch of the *log-conditional* tree (not the newly constructed *log-joint*) we just choose and traverse the one that is higher, thus we have:

$$\begin{aligned}\bar{w}_1 &= \arg \max_{w_1} \{-1.5, -1.8\} = \text{the} \\ \bar{w}_2 &= \arg \max_{w_2} \{-4.6, -4.5\} = \text{red} \\ \bar{w}_3 &= \arg \max_{w_3} \{-9.1, -9.3\} = \text{rectangle}\end{aligned}$$

So the sequence $\bar{\mathbf{w}}_{1:3}$ would be $\{\text{the, red, rectangle}\}$, and its log-likelihood is the summation of the numbers along the path which is -9.1 (from the previous log-joint tree). We see that at this computationally simpler greedy approach comes at the expense of returning a suboptimal solution because the log-likelihood for this sequence is lower than the one we derived using the Viterbi decoder so it's not really the most likely sequence overall. But this is the price we pay for going with a simpler algorithm.

2.4.c Apply beam search decoding (question 3) with a beam width $B = 2$ to this language model, and find $\mathbf{w}_{1:3}^{(1)}$ and $\mathbf{w}_{1:3}^{(2)}$, together with their respective log-likelihoods.

At each step t , we will keep two *beams* of words that are most likely based on log-likelihoods. Starting from the root, there are only two possibilities: $\{\mathbf{the}, \mathbf{a}\}$ with log-likelihoods $\{-1.5, -1.8\}$. Then going one step further, there are three possibilities and **red** appearing twice, so there are four candidates with their log-likelihoods calculated by adding the values for parents to childrens' (or using the log-joint tree constructed before): $\{\mathbf{red}^{(1)} : -4.5, \mathbf{blue} : -4.3, \mathbf{red}^{(2)} : -4.6, \mathbf{green} : -4.9\}$, discarding two beams and only keeping the top 2, we get the following beams:

$b_1: \mathbf{a} \quad \mathbf{blue} \rightarrow \text{log-likelihood} = -4.3$
 $b_2: \mathbf{a} \quad \mathbf{red} \rightarrow \text{log-likelihood} = -4.6$

Now we can go the last step and see where these beams take us. For the first beam there are two options: $\{\mathbf{ellipse}, \mathbf{square}\}$, and for the second beam there are two options: $\{\mathbf{circle}, \mathbf{rectangle}\}$, so below are all possible beams with their respective log-likelihoods:

$b_1: \mathbf{a}, \mathbf{blue} \quad \mathbf{ellipse} \rightarrow \text{log-likelihood} = -9.7$
 $\checkmark b_2: \mathbf{a}, \mathbf{blue} \quad \mathbf{square} \rightarrow \text{log-likelihood} = -9.3$
 $\checkmark b_3: \mathbf{a}, \mathbf{red} \quad \mathbf{circle} \rightarrow \text{log-likelihood} = -9.4$
 $b_4: \mathbf{a}, \mathbf{red} \quad \mathbf{square} \rightarrow \text{log-likelihood} = -9.6$

And the two top beams have been marked so to conclude:

$\mathbf{w}_{1:3}^{(1)}: \mathbf{a}, \mathbf{blue}, \mathbf{square} \rightarrow \text{log-likelihood} = -9.3$
 $\mathbf{w}_{1:3}^{(2)}: \mathbf{a}, \mathbf{red}, \mathbf{circle} \rightarrow \text{log-likelihood} = -9.4$

2.4.d Compare the behavior of these 3 decoding algorithms on this language model (in particular greedy decoding vs. maximum likelihood, and beam search decoding vs. the other two). How can you mitigate the limitations of beam search?

We have computed the complexity of the Viterbi algorithm, and the Beam search earlier. Finding the exact maximum likelihood solution requires a lot of computation even if the matrix \mathbf{P} is available, the greedy approach tries to find an approximation but is *very unstable* since it only looks one step forward and doesn't allow for any room to correct wrongly traversed paths, on the other hand, it's very cheap computationally, so we can run it several times and obtain multiple most likely sequences and decide based on that. However, beam search decoding can add the ingredient of robustness to our approach by not a significant computational overhead. It becomes more robust as we keep more and more beams by increasing B , and it reduces to the maximum likelihood solution if we keep all possible beams at each step. But even with a relatively small B we can keep several paths that potentially result in the maximum solution and cheaply traverse them as they're not too many. This way we've allowed some room for mistakes and have raised the probability of hitting the exact maximum likely solutions. In a sense beam search decoding introduces a tradeoff that avoids monitoring all possibilities but also avoids being unstable like the

greedy approach. So it seems to be the better approach. In this example the results are relatively different, but increasing the beam width could alleviate such discrepancies. In this problem it can be easily seen that a beam width of 4 could recover the exact solution.

Question 3 (2-3-4-4). Weight decay as L2 regularization In this question, you will reconcile the relationship between L2 regularization and weight decay for the Stochastic Gradient Descent (SGD) and Adam optimizers. Imagine you are training a neural network (with learnable weights θ) with a loss function $L(f(\mathbf{x}^{(j)}, \theta), \mathbf{y}^{(j)})$, under two different schemes. The *weight decay* scheme uses a modified SGD update rule: the weights θ decay exponentially by a factor of λ . That is, the weights at iteration $i + 1$ are computed as

$$\theta_{i+1} = \theta_i - \eta \frac{1}{m} \frac{\partial \sum_{j=1}^m L(f(\mathbf{x}^{(j)}, \theta), \mathbf{y}^{(j)})}{\partial \theta_i} + \lambda \theta_i$$

where η is the learning rate of the SGD optimizer. The *L2 regularization* scheme instead modifies the loss function (while maintaining the typical SGD or Adam update rules). The modified loss function is

$$L_{\text{reg}}(f(\mathbf{x}^{(j)}, \theta), \mathbf{y}^{(j)}) = L(f(\mathbf{x}^{(j)}, \theta), \mathbf{y}^{(j)}) + \gamma \|\theta\|_2^2$$

3.1 Prove that the *weight decay* scheme that employs the modified SGD update is identical to an *L2 regularization* scheme that employs a standard SGD update rule.

Single sample's and average over all samples' loss are:

$$\begin{aligned} L_{\text{reg}}^j &= L(f(\mathbf{x}^{(j)}, \theta), \mathbf{y}^{(j)}) + \gamma \|\theta\|_2^2 \\ L_{\text{reg}} &= \frac{1}{m} \sum_j L(f(\mathbf{x}^{(j)}, \theta), \mathbf{y}^{(j)}) + \frac{1}{m} \sum_j \gamma \|\theta\|_2^2 \\ &= \frac{1}{m} \sum_j L(f(\mathbf{x}^{(j)}, \theta), \mathbf{y}^{(j)}) + \gamma \|\theta\|_2^2 \end{aligned}$$

The standard SGD update rule would be $\theta_{t+1} = \theta_t - \eta \nabla_{\theta} L_{\text{reg}}$. Now taking the gradient:

$$\nabla_{\theta} L_{\text{reg}} = \frac{1}{m} \sum_j \frac{\partial L(f(\mathbf{x}^{(j)}, \theta), \mathbf{y}^{(j)})}{\partial \theta} + \gamma \frac{\partial}{\partial \theta} \|\theta\|_2^2$$

But $\|\theta\|_2^2 = \theta^\top \theta$, and thus $\frac{\partial}{\partial \theta} \|\theta\|_2^2 = \frac{\partial}{\partial \theta} \theta^\top \theta = 2\theta$, hence the overall update for the L2 regularized loss would be:

$$\theta_{t+1} = \theta_t - \eta \frac{1}{m} \sum_j \frac{\partial L(f(\mathbf{x}^{(j)}, \theta), \mathbf{y}^{(j)})}{\partial \theta} - 2\eta \gamma \theta_t$$

Comparing this with the modified SGD for weight decay, we find that they're identical if we set $\lambda = -2\eta\gamma$. This concludes the proof.

3.2 This question refers to the Adam algorithm as described in the lecture slide (also identical to Algorithm 8.7 of the deep learning book). It turns out that a one-line change to this algorithms gives us Adam with an L2 regularization scheme. Identify the line of the algorithm that needs to change, and provide this one-line modification.

One (naive) way would be to add the L2 penalty to the loss and this will modify the line that computes the gradients:

$$\mathbf{h} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)}) + \underbrace{\gamma \theta}_{\text{new term}}$$

We will soon see why it's naive, because it negates the motivation of L2 penalty to penalize all weights equally. Doing the modification here will result in imbalanced rate of changes for different weights.

3.3 Consider a “decoupled” weight decay scheme for the original Adam algorithm (see lecture slides, or equivalently, Algorithm 8.7 of the deep learning book) with the following two update rules.

- The **Adam-L2-reg** scheme computes the update by employing an L2 regularization scheme (same as the question above).
- The **Adam-weight-decay** scheme computes the update as $\Delta\theta = -\left(\epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \delta}} + \lambda\theta\right)$.

Now, assume that the neural network weights can be partitioned into two disjoint sets based on their magnitude: $\theta = \{\theta_{\text{small}}, \theta_{\text{large}}\}$, where each weight $\theta_s \in \theta_{\text{small}}$ has a much smaller gradient magnitude than each weight $\theta_l \in \theta_{\text{large}}$. Using this information provided, answer the following questions. In each case, provide a brief explanation as to why your answer holds.

- (a) Under the **Adam-L2-reg** scheme, which set of weights among θ_{small} and θ_{large} would you expect to be regularized (i.e., driven closer to zero) more strongly than the other? Why? Let's see how the updates will look like if we add the L2 penalty right to the loss and then take the gradient:

$$\theta_{t+1} = \theta_t - \eta \frac{\hat{\mathbf{m}}_t}{\sqrt{\hat{\mathbf{v}}_t} + \epsilon}$$

But what are $\hat{\mathbf{m}}_t, \hat{\mathbf{v}}_t$? They are unbiased estimates of the moving average of gradients and gradient squares.

$$\begin{aligned} \hat{\mathbf{m}}_t &= \frac{\mathbf{m}_t}{1 - \beta_1^t} = \frac{\beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t}{1 - \beta_1^t} \\ \hat{\mathbf{v}}_t &= \frac{\mathbf{v}_t}{1 - \beta_2^t} = \frac{\beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2}{1 - \beta_2^t} \end{aligned}$$

In above equations, \mathbf{g}_t is the gradient of the *regularized* loss, the one we modified in the previous section.

$$\mathbf{g}_t = \frac{1}{m} \nabla_{\theta_t} \sum_i L(f(\mathbf{x}^{(i)}; \theta_t), \mathbf{y}^{(i)}) + \gamma \theta_t$$

Now plugging all these equations to the update rule:

$$\begin{aligned}
 \theta_{t+1} &= \theta_t - \eta \frac{\beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) (\nabla_{\theta_t} L + \gamma \theta_t)}{(1 - \beta_1^t) (\sqrt{\hat{\mathbf{v}}_t} + \epsilon)} \\
 &= \theta_t - \eta \frac{\beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) (\nabla_{\theta_t} L + \gamma \theta_t)}{(1 - \beta_1^t) \left(\sqrt{\frac{\beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \nabla_{\theta_t}^2 L}{1 - \beta_2^t}} + \epsilon \right)} \\
 &= \theta_t - \eta \frac{\beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) (\nabla_{\theta_t} L)}{(1 - \beta_1^t) \left(\sqrt{\frac{\beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \nabla_{\theta_t}^2 L}{1 - \beta_2^t}} + \epsilon \right)} - \underbrace{\eta \frac{(1 - \beta_1) (\gamma \theta_t)}{(1 - \beta_1^t) \left(\sqrt{\frac{\beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \nabla_{\theta_t}^2 L}{1 - \beta_2^t}} + \epsilon \right)}}_{\text{regularization term}}
 \end{aligned}$$

From the last equation we can draw conclusions about $\{\theta_{\text{small}}, \theta_{\text{large}}\}$. We observe that the weight decay term ($\gamma \theta_t$) resulting from the L2 penalty, is involved in the moving averages, and for θ_{large} , the denominator would be huge (because of large $\nabla_{\theta_t}^2 L$) and thus results in less decay for this set of parameters, on the other hand, for θ_{small} , the gradients are smaller and thus the denominator is smaller and the regularization or decay is stronger for this set of parameters. This imbalanced strength of regularization for different set of parameters is undesired and contradicts the motivation of the L2 penalty in the first place: We want all weights to be decayed equally.

- (b) Would your answer change for the **Adam-weight-decay** scheme? Why/why not?

If we add the weight decay term (or equivalently the gradient of the regularizer) only after all moving averages have been taken care of, then the decay will not depend on any moving average and thus the magnitude of gradients for different set of parameters won't affect it and the effect/strength of the regularizer would be balanced for all parameters regardless of their gradient magnitude.

(Note: for the two sub-parts above, we are interested in the rate at which the weights are regularized, *relative* to their initial magnitudes.)

- 3.4 In the context of all of the discussion above, argue that weight decay is a better scheme to employ as opposed to L2 regularization; particularly in the context of adaptive gradient based optimizers. (Hint: think about how each of these schemes regularize each parameter, and also about what the overarching objective of regularization is).

We eluded to this problem several times earlier. The main objective of L2 regularizer is to decay all weights equally and the motivation for that comes from the observation that *networks with smaller weights or weights closer to zero* have shown *better generalization*. That is what we strive for by introducing the L2 penalty. Now in the context of adaptive gradient based optimizer, naively adding the regularization term to the loss and taking adaptive gradients, will bring in the gradient of *regularization term* into the play for all moving averages, and thus the strength or effect of the regularizer will depend on the magnitude of the gradient for each parameter which is undesirable. To alleviate this problem, using the alternative or equivalent of L2 penalty in a more appropriate part of the algorithm seems reasonable. That is, we add the weight decay to the last update where moving averages and bias corrections have been completed. Then we will use the weight decay and recover the same effect as balanced L2 penalty with an equal

strength over all parameters. That concludes the argument for why employing weight decay is a better scheme. Indeed they're the same, what matters is where you apply the regularization.

Question 4 (4-6-6). This question is about normalization techniques.

4.1 Batch normalization, layer normalization and instance normalization all involve calculating the mean μ and variance σ^2 with respect to different subsets of the tensor dimensions. Given the following 3D tensor, calculate the corresponding mean and variance tensors for each normalization technique: μ_{batch} , μ_{layer} , $\mu_{instance}$, σ_{batch}^2 , σ_{layer}^2 , and $\sigma_{instance}^2$.

$$\begin{bmatrix} \begin{bmatrix} 1, 3, 2 \\ 1, 2, 3 \end{bmatrix}, \begin{bmatrix} 3, 3, 2 \\ 2, 4, 4 \end{bmatrix}, \begin{bmatrix} 4, 2, 2 \\ 1, 2, 4 \end{bmatrix}, \begin{bmatrix} 3, 3, 2 \\ 3, 3, 2 \end{bmatrix} \end{bmatrix}$$

The size of this tensor is 4 x 2 x 3 which corresponds to the batch size, number of channels, and number of features respectively.

— Batch normalization:

Here we should take the mean over features *for each channel* (rows here). The mean is computed over all samples so we have ($\mu_{batch}^{(2 \times 1)}$):

$$\begin{aligned} \mu_{batch}^{(0)} &: \frac{1}{4 \times 3} \left((1 + 3 + 2) + (3 + 3 + 2) + (4 + 2 + 2) + (3 + 3 + 2) \right) = \frac{1}{12} (6 + 8 + 8 + 8) = 2.5 \\ \mu_{batch}^{(1)} &: \frac{1}{4 \times 3} \left((1 + 2 + 3) + (2 + 4 + 4) + (1 + 2 + 4) + (3 + 3 + 2) \right) = \frac{1}{12} (6 + 10 + 7 + 8) = 2.58 \end{aligned}$$

Now computing the variance is simple, summing the squared difference of elements with the mean for each channel and over all samples:

$$\begin{aligned} \sigma_{batch}^{2(0)} &= \frac{1}{4 \times 3} \left((1.5^2 + 0.5^2 + 0.5^2) + (0.5^2 + 0.5^2 + 0.5^2) + (1.5^2 + 0.5^2 + 0.5^2) \right. \\ &\quad \left. + (0.5^2 + 0.5^2 + 0.5^2) \right) = \frac{7}{12} \\ \sigma_{batch}^{2(1)} &= \frac{1}{4 \times 3} \left(((-1.58)^2 + (-0.58)^2 + 0.42^2) + ((-0.58)^2 + 1.42^2 + 1.42^2) \right. \\ &\quad \left. + ((-1.58)^2 + (-0.58)^2 + 1.42^2) + (0.42^2 + 0.42^2 + (-0.58)^2) \right) = \frac{12.92}{12} = 1.08 \end{aligned}$$

— Layer normalization:

Here we should take the mean over all features and channels *for each sample* (each matrix here). So the shape of μ_{layer} is 4 x 1:

$$\begin{aligned} \mu_{layer}^{(0)} &: \frac{1}{2 \times 3} \left((1 + 3 + 2) + (1 + 2 + 3) \right) = \frac{1}{6} (6 + 6) = 2 \\ \mu_{layer}^{(1)} &: \frac{1}{2 \times 3} \left((3 + 3 + 2) + (2 + 4 + 4) \right) = \frac{1}{6} (8 + 10) = 3 \\ \mu_{layer}^{(2)} &: \frac{1}{2 \times 3} \left((4 + 2 + 2) + (1 + 2 + 4) \right) = \frac{1}{6} (8 + 7) = 2.5 \\ \mu_{layer}^{(3)} &: \frac{1}{2 \times 3} \left((3 + 3 + 2) + (3 + 3 + 2) \right) = \frac{1}{6} (8 + 8) = 2.67 \end{aligned}$$

Now computing the variance is simple, summing the squared difference of elements with the mean for *each sample* (over all features and channels):

$$\sigma_{\text{layer}}^2{}^{(0)} = \frac{1}{2 \times 3} \left(((-1)^2 + 1^2 + 0^2) + ((-1)^2 + 0^2 + 1^2) \right) = \frac{4}{6} = 0.67$$

$$\sigma_{\text{layer}}^2{}^{(1)} = \frac{1}{2 \times 3} \left((0^2 + 0^2 + (-1)^2) + ((-1)^2 + 1^2 + 1^2) \right) = \frac{4}{6} = 0.67$$

$$\sigma_{\text{layer}}^2{}^{(2)} = \frac{1}{2 \times 3} \left((1.5^2 + (-0.5)^2 + (-0.5)^2) + ((-1.5)^2 + (-0.5)^2 + (1.5)^2) \right) = \frac{7.5}{6} = 1.25$$

$$\sigma_{\text{layer}}^2{}^{(3)} = \frac{1}{2 \times 3} \left((0.33^2 + 0.33^2 + (-0.67)^2) + (0.33^2 + 0.33^2 + (-0.67)^2) \right) = \frac{1.33}{6} = 0.22$$

— Instance normalization:

Here we should take the mean over all features *for each channel and sample*. So we will have $4 \times 2 = 8$ means and variances overall. The shape of each $\mu_{\text{instance}}, \sigma_{\text{instance}}^2$ is 1×1 (a scalar):

— Sample 1 $\rightarrow \begin{bmatrix} 1, 3, 2 \\ 1, 2, 3 \end{bmatrix}$

channel 0	$\mu_{\text{instance}}:$	$\frac{1}{3}(1 + 3 + 2) = \frac{1}{3}(6) = 2$
channel 1	$\mu_{\text{instance}}:$	$\frac{1}{3}(1 + 2 + 3) = \frac{1}{3}(6) = 2$
channel 0	$\sigma_{\text{instance}}^2:$	$\frac{1}{3}((-1)^2 + 1^2 + 0^2) = \frac{2}{3} = 0.67$
channel 1	$\sigma_{\text{instance}}^2:$	$\frac{1}{3}((-1)^2 + 0^2 + 1^2) = \frac{2}{3} = 0.67$

— Sample 2 $\rightarrow \begin{bmatrix} 3, 3, 2 \\ 2, 4, 4 \end{bmatrix}$

channel 0	$\mu_{\text{instance}}:$	$\frac{1}{3}(3 + 3 + 2) = \frac{8}{3} = 2.67$
channel 1	$\mu_{\text{instance}}:$	$\frac{1}{3}(2 + 4 + 4) = \frac{10}{3} = 3.33$
channel 0	$\sigma_{\text{instance}}^2:$	$\frac{1}{3}(0.33^2 + 0.33^2 + (-0.67)^2) = \frac{0.67}{3} = 0.22$
channel 1	$\sigma_{\text{instance}}^2:$	$\frac{1}{3}((-1.33)^2 + 0.67^2 + 0.67^2) = \frac{2.67}{3} = 0.89$

— Sample 3 $\rightarrow \begin{bmatrix} 4, 2, 2 \\ 1, 2, 4 \end{bmatrix}$

channel 0	$\mu_{\text{instance}}:$	$\frac{1}{3}(4 + 2 + 2) = \frac{8}{3} = 2.67$
channel 1	$\mu_{\text{instance}}:$	$\frac{1}{3}(1 + 2 + 4) = \frac{7}{3} = 2.33$
channel 0	$\sigma_{\text{instance}}^2:$	$\frac{1}{3}(1.33^2 + (-0.67)^2 + (-0.67)^2) = \frac{2.67}{3} = 0.89$
channel 1	$\sigma_{\text{instance}}^2:$	$\frac{1}{3}((-1.33)^2 + (-0.33)^2 + 1.67^2) = \frac{4.67}{3} = 1.56$

— Sample 4 $\rightarrow \begin{bmatrix} 3, 3, 2 \\ 3, 3, 2 \end{bmatrix}$

channel 0	$\mu_{\text{instance}}:$	$\frac{1}{3}(3 + 3 + 2) = \frac{8}{3} = 2.67$
channel 1	$\mu_{\text{instance}}:$	$\frac{1}{3}(3 + 3 + 2) = \frac{8}{3} = 2.67$
channel 0	$\sigma_{\text{instance}}^2:$	$\frac{1}{3}(0.33^2 + 0.33^2 + (-0.67)^2) = \frac{0.67}{3} = 0.22$
channel 1	$\sigma_{\text{instance}}^2:$	$\frac{1}{3}(0.33^2 + 0.33^2 + (-0.67)^2) = \frac{0.67}{3} = 0.22$

4.2 For the next two subquestions, 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} .

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

Starting by definitions:

$$\begin{aligned}\hat{y} &= \gamma \cdot \frac{y - \mu_y}{\sigma_y} + \beta \\ y &= \mathbf{u}^\top \mathbf{x}\end{aligned}$$

We thus need to compute μ_y, σ_y :

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

Last equality comes from the assumption of the problem ($\mathbb{E}[\mathbf{x}] = \mathbf{0}$) along with the fact that \mathbf{u} is independent of \mathbf{x} hence comes out of the expectation. Now for the variance:

$$\begin{aligned}\sigma_y^2 &= \mathbb{E}[(y - \mu_y)^2] \stackrel{\mu_y=0}{=} \mathbb{E}[(y)^2] = \mathbb{E}[(\mathbf{u}^\top \mathbf{x})^2] \\ &\stackrel{A^2=A^\top A}{=} \mathbb{E}[(\mathbf{u}^\top \mathbf{x})^\top \mathbf{u}^\top \mathbf{x}] = \mathbb{E}[(\mathbf{u}^\top \mathbf{x})^\top \mathbf{u}^\top \mathbf{x}]\end{aligned}$$

But $\mathbf{u}^\top \mathbf{x}$ is a scalar, so it's equal to its transpose $(\mathbf{u}^\top \mathbf{x})^\top$ and hence we can rewrite $(\mathbf{u}^\top \mathbf{x})^2 = \mathbf{u}^\top \mathbf{x}(\mathbf{u}^\top \mathbf{x})^\top$. Then we get:

$$\sigma_y^2 = \mathbb{E}[(\mathbf{u}^\top \mathbf{x})^2] = \mathbb{E}[\mathbf{u}^\top \mathbf{x}(\mathbf{u}^\top \mathbf{x})^\top] = \mathbb{E}[\mathbf{u}^\top \mathbf{x} \mathbf{x}^\top \mathbf{u}]$$

But note that $\text{Var}(\mathbf{x}) = \mathbf{I} = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])^2]$, but since $\mathbb{E}[\mathbf{x}] = \mathbf{0}$ by assumption, then $\text{Var}(\mathbf{x}) = \mathbf{I} = \mathbb{E}[(\mathbf{x})^2] = \mathbb{E}[\mathbf{x} \mathbf{x}^\top]$. In the equation above, \mathbf{u} is independent of expectation over \mathbf{x} so it can be pulled out and we can plug $\text{Var}(\mathbf{x}) = \mathbf{I} = \mathbb{E}[\mathbf{x} \mathbf{x}^\top]$ into that to get:

$$\sigma_y^2 = \mathbb{E}[\mathbf{u}^\top \mathbf{x} \mathbf{x}^\top \mathbf{u}] = \mathbf{u}^\top \mathbb{E}[\mathbf{x} \mathbf{x}^\top] \mathbf{u} = \mathbf{u}^\top \mathbf{I} \mathbf{u} = \|\mathbf{u}\|_2^2$$

So we got $\mu_y = 0$ and $\sigma_y = \|\mathbf{u}\|_2$, plugging these back to the equation for \hat{y} :

$$\hat{y} = \gamma \cdot \frac{y - \mu_y}{\sigma_y} + \beta = \gamma \cdot \frac{y - 0}{\|\mathbf{u}\|_2} + \beta = \gamma \cdot \frac{\mathbf{u}^\top}{\|\mathbf{u}\|_2} \mathbf{x} + \beta$$

And since $\mathbf{w} = \gamma \frac{\mathbf{u}}{\|\mathbf{u}\|} \rightarrow \mathbf{w}^\top = \gamma \frac{\mathbf{u}^\top}{\|\mathbf{u}\|}$, hence we arrive at the following and the proof concludes:

$$\hat{y} = \gamma \cdot \frac{\mathbf{u}^\top}{\|\mathbf{u}\|_2} \mathbf{x} + \beta = \mathbf{w}^\top \mathbf{x} + \beta$$

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

The loss is probably a distance between the post-activation and some ground truth output y_g . But we don't know the exact details, but we can use the chain rule and factor those details in s :

$$\nabla_{\mathbf{u}} L(\mathbf{u}, \gamma, \beta) = \frac{\partial L(\mathbf{u}, \gamma, \beta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{u}}$$

Now let's focus on the part we can compute $\frac{\partial \hat{y}}{\partial \mathbf{u}}$:

$$\begin{aligned} \frac{\partial \hat{y}}{\partial \mathbf{u}} &= \frac{\partial}{\partial \mathbf{u}} \left(\gamma \frac{\mathbf{u}^\top}{\|\mathbf{u}\|} \mathbf{x} + \beta \right) = \gamma \left(\frac{\partial}{\partial \mathbf{u}} \frac{\mathbf{u}^\top}{\|\mathbf{u}\|} \right) \mathbf{x} + \frac{\partial}{\partial \mathbf{u}} (\beta) \\ &= \gamma \left(\frac{\frac{\partial \mathbf{u}}{\partial \mathbf{u}} \|\mathbf{u}\| - \mathbf{u} \frac{\partial \|\mathbf{u}\|}{\partial \mathbf{u}}}{\|\mathbf{u}\|^2} \right) \mathbf{x} = \gamma \left(\frac{\mathbf{I} \|\mathbf{u}\| - \mathbf{u} \frac{\partial \|\mathbf{u}\|}{\partial \mathbf{u}}}{\|\mathbf{u}\|^2} \right) \mathbf{x} \end{aligned}$$

To compute $\frac{\partial \|\mathbf{u}\|}{\partial \mathbf{u}}$, note that $\|\mathbf{u}\| = \sqrt{\mathbf{u}^\top \mathbf{u}}$, then:

$$\frac{\partial \|\mathbf{u}\|}{\partial \mathbf{u}} = \frac{\partial (\mathbf{u}^\top \mathbf{u})^{1/2}}{\partial \mathbf{u}} = \frac{1}{2} (\mathbf{u}^\top \mathbf{u})^{-1/2} \frac{\partial (\mathbf{u}^\top \mathbf{u})}{\partial \mathbf{u}} = \frac{1}{2} (\mathbf{u}^\top \mathbf{u})^{-1/2} \times 2\mathbf{u}^\top$$

Inserting this result into the equation for $\frac{\partial \hat{y}}{\partial \mathbf{u}}$:

$$\begin{aligned} \frac{\partial \hat{y}}{\partial \mathbf{u}} &= \gamma \left(\frac{\mathbf{I} \|\mathbf{u}\| - \mathbf{u} \frac{\partial \|\mathbf{u}\|}{\partial \mathbf{u}}}{\|\mathbf{u}\|^2} \right) \mathbf{x} = \gamma \left(\frac{\mathbf{I} \|\mathbf{u}\| - \mathbf{u} (\mathbf{u}^\top \mathbf{u})^{-1/2} \mathbf{u}^\top}{\|\mathbf{u}\|^2} \right) \mathbf{x} = \gamma \left(\frac{\mathbf{I} \|\mathbf{u}\| - \mathbf{u} (\|\mathbf{u}\|^2)^{-1/2} \mathbf{u}^\top}{\|\mathbf{u}\|^2} \right) \mathbf{x} \\ &= \gamma \left(\frac{\mathbf{I} \|\mathbf{u}\|^2 - \mathbf{u} \mathbf{u}^\top}{\|\mathbf{u}\|^3} \right) \mathbf{x} = \left(\mathbf{I} - \frac{\mathbf{u} \mathbf{u}^\top}{\|\mathbf{u}\|^2} \right) \frac{\gamma \mathbf{x}}{\|\mathbf{u}\|} \end{aligned}$$

So finally we have:

$$\nabla_{\mathbf{u}} L(\mathbf{u}, \gamma, \beta) = \frac{\partial L(\mathbf{u}, \gamma, \beta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{u}} = \frac{\partial L(\mathbf{u}, \gamma, \beta)}{\partial \hat{y}} \left(\mathbf{I} - \frac{\mathbf{u} \mathbf{u}^\top}{\|\mathbf{u}\|^2} \right) \frac{\gamma \mathbf{x}}{\|\mathbf{u}\|}$$

but $\nabla_{\mathbf{w}} L = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{w}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial (\mathbf{w}^\top \mathbf{x} + \beta)}{\partial \mathbf{w}} = \frac{\partial L}{\partial \hat{y}} \mathbf{x}$, substituting this term in the equation for $\nabla_{\mathbf{u}} L(\mathbf{u}, \gamma, \beta)$ along with $\mathbf{W}^\perp = \left(\mathbf{I} - \frac{\mathbf{u} \mathbf{u}^\top}{\|\mathbf{u}\|^2} \right)$:

$$\nabla_{\mathbf{u}} L(\mathbf{u}, \gamma, \beta) = \frac{\partial L(\mathbf{u}, \gamma, \beta)}{\partial \hat{y}} \left(\mathbf{I} - \frac{\mathbf{u} \mathbf{u}^\top}{\|\mathbf{u}\|^2} \right) \frac{\gamma \mathbf{x}}{\|\mathbf{u}\|} = \frac{\gamma}{\|\mathbf{u}\|} \mathbf{W}^\perp \nabla_{\mathbf{w}} L$$

So if we set $s = \frac{\gamma}{\|\mathbf{u}\|}$, then we arrive at the desired result:

$$\nabla_{\mathbf{u}} L(\mathbf{u}, \gamma, \beta) = s \mathbf{W}^\perp \nabla_{\mathbf{w}} L$$

2. 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.