

**Due Date: March 23rd (11pm), 2022**

Instructions

- For all questions, show your work!
- Starred questions are **hard** questions, not **bonus** questions.
- Use LaTeX and the template we provide when writing your answers. You may reuse most of the notation shorthands, equations and/or tables. See the assignment policy on the course website for more details.
- 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 **Ankit Vani** and **Sai Aravind Sreeramadas**.

**Question 1** (6-9-6). This question is about activation functions and vanishing/exploding gradients in recurrent neural networks (RNNs). Let  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  be an activation function. When the argument is a vector, we apply  $\sigma$  element-wise. Consider the following recurrent unit:

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

- 1.1 Show that applying the activation function in this way results in an equivalent recurrence as 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  $t - 1$ .

**Answer.**

Note that here we are going to assume that  $\sigma(\mathbf{h}_{t-1}) = \mathbf{g}_{t-1}$  as suggested in the question. From that:

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

Therefore,

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

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

■.

- \*1.2 Let  $\|\mathbf{A}\|$  denote the  $L_2$  operator norm<sup>1</sup> of matrix  $\mathbf{A}$  ( $\|\mathbf{A}\| := \max_{\mathbf{x}: \|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|$ ). Assume  $\sigma(x)$  has bounded derivative, i.e.  $|\sigma'(x)| \leq \gamma$  for some  $\gamma > 0$  and for all  $x$ . We denote as  $\lambda_1(\cdot)$  the largest eigenvalue of a symmetric matrix. Show that if the largest eigenvalue of the weights is 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

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

**Answer.**

Using the chain rule:

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

Then:

$$\frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} = \frac{\partial}{\partial \mathbf{h}_{T-1}} \mathbf{W} \sigma(\mathbf{h}_{T-1}) + \mathbf{U} \mathbf{x}_T + \mathbf{b} = \mathbf{W} \frac{\partial \sigma(\mathbf{h}_{T-1})}{\partial \mathbf{h}_{T-1}}$$

Using given properties of  $L_2$  norm:

$$\begin{aligned} \left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \right\| &\leq \|\mathbf{W}\| \cdot \left\| \frac{\partial \sigma(\mathbf{h}_{T-1})}{\partial \mathbf{h}_{T-1}} \right\| \\ &\leq \sqrt{\lambda_1(\mathbf{W}^\top \mathbf{W})} \cdot \left\| \frac{\partial \sigma(\mathbf{h}_{T-1})}{\partial \mathbf{h}_{T-1}} \right\| \\ &\leq \frac{\delta}{\gamma} \gamma \leq \delta \end{aligned}$$

So:

$$\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| \leq \prod_{i=1}^T \left\| \frac{\partial \mathbf{h}_i}{\partial \mathbf{h}_{i-1}} \right\| \leq \delta^T$$

Therefore, as  $\delta \in [0, 1)$  as  $T \rightarrow \infty$ ,  $\delta^T \rightarrow 0$ . So as  $\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| \leq \delta^T$ ,  $\left\| \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_0} \right\| \rightarrow 0$ . ■

- 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).

**Answer.**

This is a necessary condition but not a sufficient one for the gradient to explode. This can be shown as we have  $\mathbf{W}$  such that  $\delta \leq \sqrt{\lambda_1(\mathbf{W}^\top \mathbf{W})} \gamma \leq 1$  than the gradient will not explode (it will vanish if  $\delta \leq \sqrt{\lambda_1(\mathbf{W}^\top \mathbf{W})} \gamma < 1$  and will be stabilized if  $\delta = \sqrt{\lambda_1(\mathbf{W}^\top \mathbf{W})} \gamma = 1$ ).

---

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.

**Question 2** (8-8-8). 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.

2.1 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  $(\alpha, \epsilon)$  as a function of  $(\beta, \delta)$ .

**Answer.**

**Recursive update expressions:**

- SGD with momentum:

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

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

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

Therefore, the two update rules are equivalent :

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

when  $\alpha = \beta$  and  $\epsilon = \delta(1 - \beta)$

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

**Answer.**

$$\begin{aligned}
 \Delta\theta_t &= -\delta\mathbf{v}_t = -\delta(\beta\mathbf{v}_{t-1} + (1-\beta)\mathbf{g}_t) \\
 &= -\delta\beta\mathbf{v}_{t-1} - \delta(1-\beta)\mathbf{g}_t \\
 &= -\delta\beta(\beta\mathbf{v}_{t-2} + (1-\beta)\mathbf{g}_{t-1}) - \delta(1-\beta)\mathbf{g}_t \\
 &= -\delta\beta^2\mathbf{v}_{t-2} - \delta\beta(1-\beta)\mathbf{g}_{t-1} - \delta(1-\beta)\mathbf{g}_t \\
 &= -\delta\beta^2(\beta\mathbf{v}_{t-3} + (1-\beta)\mathbf{g}_{t-2}) - \delta\beta(1-\beta)\mathbf{g}_{t-1} - \delta(1-\beta)\mathbf{g}_t \\
 &= -\delta\beta^3\mathbf{v}_{t-3} - \delta\beta^2(1-\beta)\mathbf{g}_{t-2} - \delta\beta(1-\beta)\mathbf{g}_{t-1} - \delta(1-\beta)\mathbf{g}_t \\
 &= -\delta\sum_{i=1}^t\beta^{i-1}(1-\beta)\mathbf{g}_{t-i+1} = \boxed{-\delta(1-\beta)\sum_{i=1}^t\beta^{t-i}\mathbf{g}_i}
 \end{aligned}$$

2.3 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.**

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

Therefore I just demonstrate that  $\mathbb{E}[\mathbf{v}_t]$  is biased as  $\mathbb{E}[\mathbf{v}_t] \neq \mathbb{E}[\mathbf{g}_t]$ .

We can fix that by using the same trick as used in the **ADAM algorithm**:

$$\mathbb{E}[\hat{\mathbf{v}}_t] = \mathbb{E}[\mathbf{g}_t] = \frac{\mathbb{E}[\mathbf{v}_t] - \delta}{(1-\beta^t)}$$

Where  $\mathbb{E}[\hat{\mathbf{v}}_t]$  is an unbiased exponential running average as  $\delta$  is very close to 0.

**Question 3** (8-8-6-9-3). In this question, you will analyze the performance of dot-product attention and derive an efficient approximation of it. Consider that *multi-head* dot-product attention for a sequence of length  $n$  is defined as follows:

$$\begin{aligned} \text{MultiHead}(\bar{\mathbf{Q}}, \bar{\mathbf{K}}, \bar{\mathbf{V}}) &= \text{Concat}(\text{head}_1, \dots, \text{head}_h) \mathbf{W}^O \\ \text{where } \text{head}_i &= \text{Attention}_{\text{std}}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \quad (\text{here, } \mathbf{Q} := \bar{\mathbf{Q}} \mathbf{W}_i^Q, \mathbf{K} := \bar{\mathbf{K}} \mathbf{W}_i^K, \mathbf{V} := \bar{\mathbf{V}} \mathbf{W}_i^V) \\ &= \text{softmax}_{\text{row}} \left( \frac{\mathbf{Q} \mathbf{K}^\top}{\sqrt{d_k}} \right) \mathbf{V} \end{aligned}$$

where  $\bar{\mathbf{Q}}, \bar{\mathbf{K}}, \bar{\mathbf{V}} \in \mathbb{R}^{n \times d_{\text{model}}}$  are the queries, keys, and values, and  $\mathbf{W}_i^Q, \mathbf{W}_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$ ,  $\mathbf{W}_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v} \forall i$ , and  $\mathbf{W}_O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$  are the weights. The softmax subscript “row” indicates that the softmax is computed along the rows, and the Attention subscript “std” indicates that this is the standard variant (we will see other variants later in the question). For this question, you can assume that  $d_k = d_v = d_{\text{model}}$  and call the value  $d$ .

For calculating the time and space complexities, you can also assume that matrix multiplications are performed naively. As an example, for  $\mathbf{C} = \mathbf{AB}$  where  $\mathbf{A} \in \mathbb{R}^{p \times q}$ ,  $\mathbf{B} \in \mathbb{R}^{q \times r}$ , and  $\mathbf{C} \in \mathbb{R}^{p \times r}$ , the time complexity is  $\Theta(pqr)$  due to the three nested loops, and the space complexity is  $\Theta(pq + qr + pr)$  from storing the inputs and the result.

3.1 What is the time and space complexity of the attention operation carried out by a single head in  $\Theta$ -notation in terms of  $n$  and  $d$ ? Use your answer to calculate the time and space complexity of multi-head dot-product attention in terms of  $n$ ,  $d$ , and  $h$ , assuming that the heads are computed sequentially. For very long sequences, where does the bottleneck lie?

For the remaining parts, let us focus on the attention operation carried out by a single head. Furthermore, you can omit the scaling factor  $\sqrt{d}$  without loss of generality by considering that  $\mathbf{Q}$  and  $\mathbf{K}$  can be scaled as desired.

**Answer.**

using  $n$  the sequence length,  $d$  the embedding dimension, and  $h$  the number of heads:

### Single-head time complexity

The time complexity for a single head is  $\mathcal{O}(n^2d)$  as  $\mathbf{Q} \cdot \mathbf{K}^\top$  multiplies a  $n \times d$  with a  $d \times n$  matrix leading to a complexity of  $\mathcal{O}(n^2d)$ . Then we need to compute the softmax of that resulting matrix ( $\mathbf{Q} \cdot \mathbf{K}^\top \in \mathbb{R}^{n \times n}$ ) which result in a time complexity of  $\mathcal{O}(n^2)$  as we need to compute the softmax for each element of the matrix. Finally, we multiply the matrix obtain after applying the softmax ( $\in \mathbb{R}^{n \times n}$ ) by  $\mathbf{V}$  ( $\in \mathbb{R}^{n \times d}$ ) leading to a final time complexity of  $\mathcal{O}(n^2d)$ .

### Single-head space complexity

The space complexity of a single head is  $\mathcal{O}(n^2 + nd)$  as we first need to store the result of  $\text{softmax}(\mathbf{Q} \cdot \mathbf{K}^\top)$  resulting in the storage of  $n^2$  values (as it is a  $n \times n$  matrix). After that, we need to store the result of  $\text{softmax}(\mathbf{Q} \cdot \mathbf{K}^\top) \cdot \mathbf{V}$  ( $\in \mathbb{R}^{n \times d}$ ) which need to store  $nd$  values. All together those storage needs leads to a space complexity of  $\mathcal{O}(n^2 + nd)$ .

### Multi-head time complexity

The time complexity of the multi-head dot-product is  $\mathcal{O}(n^2dh + nd^2h)$ . Firstly we can take the result from the single-head complexity which is  $\mathcal{O}(n^2d)$  but now we have  $h$  heads to compute **sequentially** leading to a time complexity of  $\mathcal{O}(n^2dh)$ . Then we need to concatenate those heads ( $\in \mathbb{R}^{n \times d}$ ) into one matrix ( $\in \mathbb{R}^{n \times dh}$ ) and multiply that matrix with  $\mathbf{W}^o \in \mathbb{R}^{dh \times d}$  leading to a complexity of  $\mathcal{O}(nd^2h)$ . Therefore we can plug to two computational part into the same complexity and get the final multi-head attention dot-product time complexity which is  $\mathcal{O}(n^2dh + nd^2h)$ .

### Multi-head space complexity

The space complexity of the multi-head dot-product is  $\mathcal{O}(n^2h + 2ndh + nd)$ . Firstly we can take the result from the single-head space complexity which is  $\mathcal{O}(n^2 + nd)$  but now we have  $h$  heads to compute **sequentially** leading to a multi-head space complexity of  $\mathcal{O}(n^2h + ndh)$ . Then we need to concatenate those heads ( $\in \mathbb{R}^{n \times d}$ ) into one matrix ( $\in \mathbb{R}^{n \times dh}$ ) which requires to store  $ndh$  values leading to a space complexity of  $\mathcal{O}(ndh)$  for the concatenation operation. Finally, we have to multiply the result of that operation ( $\in \mathbb{R}^{n \times dh}$ ) with  $\mathbf{W}^o \in \mathbb{R}^{dh \times d}$  leading to a space complexity of  $\mathcal{O}(nd)$ . Therefore we can put all those storage complexity together to finally get the multi-head dot-product space complexity which is  $\mathcal{O}(n^2h + 2ndh + nd)$ .

### For very long sequences, where does the bottleneck lie ?

The bottleneck lie on the fact that for both multi-head dot-product time and space complexity are dependant of  $n^2$ , the sequence length squared in their computation time or storage size. That means that we can't afford having arbitrarily long sequence length as it will quadratically increase both computation time and needed storage size.

If we want to localize what in the equation leads to such sequence-length dependant complexities we can guess that it is related to  $\mathbf{Q} \cdot \mathbf{K}^\top$  as it is always the step that create the  $n^2$  of both time and space complexities.

3.2 Let us consider an alternative form of attention, one that performs row-wise softmax on  $\mathbf{Q}$  and column-wise softmax on  $\mathbf{K}$  separately as follows:

$$\text{Attention}_{\text{separable}}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}_{\text{row}}(\mathbf{Q}) \text{softmax}_{\text{col}}(\mathbf{K})^\top \mathbf{V}.$$

Prove that  $\text{softmax}_{\text{row}}(\mathbf{Q}) \text{softmax}_{\text{col}}(\mathbf{K})^\top$  produces valid categorical distributions in every row, like  $\text{softmax}_{\text{row}}(\mathbf{Q} \mathbf{K}^\top)$ . If  $n \gg d$ , show that  $\text{Attention}_{\text{separable}}$  can be faster and requires less space than  $\text{Attention}_{\text{std}}$ . Is  $\text{Attention}_{\text{separable}}$  as expressive as  $\text{Attention}_{\text{std}}$ ?

(Hint: For a valid categorical distribution  $\mathbf{p} \in \mathbb{R}^d$  over  $d$  categories,  $p_i \geq 0 \forall i \in \{1, \dots, d\}$  and  $\sum_{i=1}^d p_i = 1$ .)

**Answer.**

**Proof.** We need to prove that  $\text{softmax}_{\text{row}}(\mathbf{Q}) \text{softmax}_{\text{col}}(\mathbf{K})^\top$  produces a valid categorical distribution over every row like  $\text{softmax}_{\text{row}}(\mathbf{Q} \mathbf{K}^\top)$ .

To do so, let fix  $\mathbf{A} = \text{softmax}_{\text{row}}(\mathbf{Q})\text{softmax}_{\text{col}}(\mathbf{K})^\top = \mathbf{BC}$

$$\mathbf{A} = \begin{bmatrix} B_{1,1} & \dots & B_{1,d} \\ \vdots & \ddots & \vdots \\ B_{n,1} & \dots & B_{n,d} \end{bmatrix} \begin{bmatrix} C_{1,1} & \dots & C_{1,n} \\ \vdots & \ddots & \vdots \\ C_{d,1} & \dots & C_{d,n} \end{bmatrix}$$

Notice that,

$$\sum_{j=1}^d B_{i,j} = 1 \quad \forall i \in \{1, \dots, n\}$$

$$\sum_{j=1}^n C_{i,j} = 1 \quad \forall i \in \{1, \dots, d\}$$

Therefore,  $\mathbf{A} = \begin{bmatrix} B_{1,\cdot} \cdot C_{\cdot,1} & \dots & B_{1,\cdot} \cdot C_{\cdot,n} \\ \vdots & \ddots & \vdots \\ B_{n,\cdot} \cdot C_{\cdot,1} & \dots & B_{n,\cdot} \cdot C_{\cdot,n} \end{bmatrix}$

Thus,

$$\begin{aligned} \sum_{j=1}^n A_{i,j} &= \sum_{j=1}^n B_{i,\cdot} \cdot C_{\cdot,j} \\ &= \sum_{j=1}^n \sum_{k=1}^d B_{i,k} C_{k,j} \\ &= \sum_{k=1}^d B_{i,k} \left( \sum_{j=1}^n C_{k,j} \right) \\ &= \sum_{k=1}^d B_{i,k} 1 \quad \text{using the softmax definition} \\ &= 1 \quad \forall i \in \{1, \dots, n\} \end{aligned}$$

So I have demonstrated that the sum of any row of matrix  $\mathbf{A}$  sums to 1. All element of  $\mathbf{A}$  are greather or equal to 0 as we applied the softmax to  $\mathbf{Q}$  and  $\mathbf{K}$  before the matrix multiplication leading to have all element of each of those matrices positives (or equal to zero). Therefore multiplying or summing positives elements with positives elements will never gives negatives elements. I have so demonstrated that the result of  $\text{Attention}_{\text{separable}}(\mathbf{Q}, \mathbf{K}, \mathbf{V})$  produces valid categorical distrubition over every row of  $\mathbf{A}$ .

**$\text{Attention}_{\text{separable}}$  can be faster and requires less space than  $\text{Attention}_{\text{std}}$  ?**

Let fix  $n \gg d$  with that constraint the  $\text{Attention}_{\text{separable}}$  will be faster to compute the attention as the softmax will consist of 2  $\mathcal{O}(nd)$  instead of  $\mathcal{O}(n^2)$  as in the  $\text{Attention}_{\text{std}}$  leading to compute the softmax much faster as  $n \gg d$ . It will also consume less storage as the result of two softmax consist of  $2nd$  values for  $\text{Attention}_{\text{separable}}$  instead of  $n^2$  values for  $\text{Attention}_{\text{std}}$  which will take less place in memory if  $n \gg d$ .

### 3.3 Verify that the standard attention can be written as

$$\text{Attention}_{\text{std}}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \mathbf{D}^{-1} \mathbf{A} \mathbf{V}$$

where  $\mathbf{A} = \exp(\mathbf{Q}\mathbf{K}^\top)$  and  $\mathbf{D} = \text{diag}(\mathbf{A}\mathbf{1})$ , where  $\exp$  is an element-wise operation,  $\text{diag}$  creates a diagonal matrix from a vector, and  $\mathbf{1}$  is a vector of ones. Note that you can store diagonal matrices in linear space and compute matrix multiplications with them in linear time. Let us now consider a variant  $\text{Attention}_{\text{approx}}$  where the elements  $a_{ij}$  of  $\mathbf{A}$  can be represented as  $a_{ij} = f(\mathbf{q}_i)^\top f(\mathbf{k}_j)$  for some  $f: \mathbb{R}^d \rightarrow \mathbb{R}_+^m$ , where  $\mathbf{q}_i$  and  $\mathbf{k}_j$  are the  $i$ th row of  $\mathbf{Q}$  and the  $j$ th row of  $\mathbf{K}$  respectively.

If  $n \gg m$  and  $n \gg d$ , how can you use this formulation to make attention efficient? What is the time and space complexity of  $\text{Attention}_{\text{approx}}$ ? You can assume that  $f$  takes  $\Theta(md)$  time and space.

(Hint: Decompose the matrix  $\mathbf{A}$ .)

**Answer.**

$$\text{Attention}_{\text{std}} = \text{softmax}(\mathbf{Q}\mathbf{K}^\top) \mathbf{V}$$

We can define  $\mathbf{C} = \mathbf{Q}\mathbf{K}^\top$

Therefore,

$$\begin{aligned} \text{Attention}_{\text{std}} &= \text{softmax}(\mathbf{Q}\mathbf{K}^\top) \mathbf{V} \\ &= \text{softmax}(\mathbf{C}) \mathbf{V} \\ &= \text{diag}\left(\left[ \frac{1}{\sum_{j=1}^n e^{C_{1,j}}} \quad \dots \quad \frac{1}{\sum_{j=1}^n e^{C_{n,j}}} \right]\right) e^{\mathbf{C}} \mathbf{V} \\ &= \text{diag}\left(\left[ \sum_{j=1}^n e^{C_{1,j}} \quad \dots \quad \sum_{j=1}^n e^{C_{n,j}} \right]\right)^{-1} e^{\mathbf{C}} \mathbf{V} \end{aligned}$$

Then, as  $\mathbf{C} = \mathbf{Q}\mathbf{K}^\top$  we have  $\mathbf{A} = e^{\mathbf{Q}\mathbf{K}^\top} = e^{\mathbf{C}}$ .

We also have  $\mathbf{D} = \text{diag}(\mathbf{A}\mathbf{1}) = \text{diag}\left(\left[ \sum_{j=1}^n e^{C_{1,j}} \quad \dots \quad \sum_{j=1}^n e^{C_{n,j}} \right]\right)$ .

Thus, we can show that:

$$\begin{aligned} \text{Attention}_{\text{approx}} &= \mathbf{D}^{-1} \mathbf{A} \mathbf{V} \\ &= \text{diag}\left(\left[ \sum_{j=1}^n e^{C_{1,j}} \quad \dots \quad \sum_{j=1}^n e^{C_{n,j}} \right]\right)^{-1} e^{\mathbf{C}} \mathbf{V} \\ &= \text{diag}\left(\left[ \sum_{j=1}^n e^{\mathbf{Q}\mathbf{K}_{1,j}^\top} \quad \dots \quad \sum_{j=1}^n e^{\mathbf{Q}\mathbf{K}_{n,j}^\top} \right]\right)^{-1} e^{\mathbf{Q}\mathbf{K}^\top} \mathbf{V} \\ &= \text{diag}\left(\left[ \frac{1}{\sum_{j=1}^n e^{\mathbf{Q}\mathbf{K}_{1,j}^\top}} \quad \dots \quad \frac{1}{\sum_{j=1}^n e^{\mathbf{Q}\mathbf{K}_{n,j}^\top}} \right]\right) e^{\mathbf{Q}\mathbf{K}^\top} \mathbf{V} \\ &= \text{softmax}_{\text{row}}(\mathbf{Q}\mathbf{K}^\top) \mathbf{V} \end{aligned}$$

Then we can compute the time and space of the  $\text{Attention}_{\text{approx}}$ . As we assume that  $n \gg m$  and  $n \gg d$ , we start to compute  $\mathbf{A}\mathbf{V}$  which have a time complexity  $\mathcal{O}(nd^2)$  as  $\mathbf{A} \in \mathbb{R}^{n \times n}$  and  $\mathbf{V} \in \mathbb{R}^{n \times d}$



\*3.4 Prove that in  $\text{Attention}_{\text{std}}$ ,

$$a_{ij} = \exp\left(\frac{-\|\mathbf{q}_i\|^2}{2}\right) \cdot \mathbb{E}_{\mathbf{x} \in \mathcal{N}(\mathbf{0}, \mathbf{I})} [\exp(\mathbf{x}^\top \mathbf{q}_i) \exp(\mathbf{x}^\top \mathbf{k}_j)] \cdot \exp\left(\frac{-\|\mathbf{k}_j\|^2}{2}\right).$$

Use this result to devise the function  $f : \mathbb{R}^d \rightarrow \mathbb{R}_+^m$  introduced in the previous part, such that  $\text{Attention}_{\text{approx}}$  approximates the expectation in  $\text{Attention}_{\text{std}}$  by sampling.

(Hint 1: If  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{I})$ ,  $p(\mathbf{x}) = (2\pi)^{-d/2} \exp(-\frac{1}{2}\|\mathbf{x} - \boldsymbol{\mu}\|^2)$  and  $\int_{\mathbf{x}} p(\mathbf{x}) d\mathbf{x} = 1$ .)

(Hint 2:  $\mathbf{x}^\top \mathbf{y} = -\frac{1}{2}(\mathbf{x}^\top \mathbf{x} - (\mathbf{x} + \mathbf{y})^\top (\mathbf{x} + \mathbf{y}) + \mathbf{y}^\top \mathbf{y})$ .)

**Answer.**

Based on the Performers's paper (<https://arxiv.org/abs/2009.14794>) and using the fact that  $\text{Attention}_{\text{std}}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \mathbf{D}^{-1} \mathbf{A} \mathbf{V}$  as proven in question 3.3 :

First we can show that:

$$\mathbf{A}_{i,j} = \exp(\mathbf{Q}_{i,:} \mathbf{K}_{:,j}^\top) = \exp(-\|\mathbf{Q}_{i,:}\|^2/2) \cdot \exp(\|\mathbf{q}_i + \mathbf{K}_{:,j}^\top\|^2/2) \cdot \exp(-\|\mathbf{K}_{:,j}^\top\|^2/2) \quad (1)$$

where  $\mathbf{Q}_{i,:}, \mathbf{K}_{:,j}^\top \in \mathbb{R}^d$ .

Therefore, with supposing that  $\mathbf{x} \in \mathbb{R}^d$ . Remind that for any  $\boldsymbol{\mu} \in \mathbb{R}^d$

$$(2\pi)^{-d/2} \int \exp(-\|\mathbf{x} - \boldsymbol{\mu}\|^2/2) d\mathbf{x} = 1$$

Thus,

$$\begin{aligned} \exp(\|\mathbf{q} + \mathbf{k}\|^2/2) &= (2\pi)^{-d/2} \exp(\|\mathbf{q} + \mathbf{k}\|^2/2) \int \exp(-\|\mathbf{x} - (\mathbf{q} + \mathbf{k})\|^2/2) d\mathbf{x} \\ &= (2\pi)^{-d/2} \int \exp(-\|\mathbf{x}\|^2/2 + \mathbf{x}^\top (\mathbf{q} + \mathbf{k}) - \|\mathbf{q} + \mathbf{k}\|^2/2 + \|\mathbf{q} + \mathbf{k}\|^2/2) d\mathbf{x} \\ &= (2\pi)^{-d/2} \int \exp(-\|\mathbf{x}\|^2/2 + \mathbf{x}^\top (\mathbf{q} + \mathbf{k})) d\mathbf{x} \\ &= (2\pi)^{-d/2} \int \exp(-\|\mathbf{x}\|^2/2) \cdot \exp(\mathbf{x}^\top \mathbf{q}) \cdot \exp(\mathbf{x}^\top \mathbf{k}) d\mathbf{x} \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)} [\exp(\mathbf{x}^\top \mathbf{q}) \cdot \exp(\mathbf{x}^\top \mathbf{k})] \end{aligned}$$

Finally, using (1) we can proof that in  $\text{Attention}_{\text{approx}}$ :

$$\mathbf{A}_{ij} = \exp\left(\frac{-\|\mathbf{q}_i\|^2}{2}\right) \cdot \mathbb{E}_{\mathbf{x} \in \mathcal{N}(\mathbf{0}, \mathbf{I})} [\exp(\mathbf{x}^\top \mathbf{q}_i) \exp(\mathbf{x}^\top \mathbf{k}_j)] \cdot \exp\left(\frac{-\|\mathbf{k}_j\|^2}{2}\right)$$

■.

3.5 Discuss the implications of the choice of  $m$  for  $\text{Attention}_{\text{approx}}$ . What are the trade-offs to think about ?

**Answer.**

**Here we are talking about choosing between accuracy of estimate of  $\mathbf{A}$  and computation time:** This is a trade-off in the sense that big  $m$  will leads to less variance on the estimation of  $\mathbf{A}$  but a higher cost for computation, where a low  $m$  leads to more variance on estimation of  $\mathbf{A}$  but much less computation time.

**Question 4 (4-5-6-6).** 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  $\theta$ ) with a loss function  $L(f(\mathbf{x}^{(i)}, \theta), \mathbf{y}^{(i)})$ , under two different schemes. The *weight decay* scheme uses a modified SGD update rule: the weights  $\theta$  decay exponentially by a factor of  $\lambda$ . That is, the weights at iteration  $i + 1$  are computed as

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

where  $\eta$  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

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

4.1 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.

**Answer.**

Firstly the derivative of the loss and the regularizer term with respect to parameters  $\theta_i$  is,

$$\begin{aligned} \frac{\partial}{\partial \theta_i} L(f(\mathbf{x}^{(i)}; \theta_i), \mathbf{y}^{(i)}) &= \frac{\partial}{\partial \theta_i} (L(f(\mathbf{x}^{(i)}; \theta_i), \mathbf{y}^{(i)}) + \gamma \|\theta_i\|_2^2) \\ &= \frac{\partial}{\partial \theta_i} L(f(\mathbf{x}^{(i)}; \theta_i), \mathbf{y}^{(i)}) + \frac{\partial}{\partial \theta_i} \gamma \|\theta_i\|_2^2 \\ &= \frac{\partial}{\partial \theta_i} L(f(\mathbf{x}^{(i)}; \theta_i), \mathbf{y}^{(i)}) + 2\gamma \theta_i \end{aligned}$$

Therefore,

$$\begin{aligned} \theta_{i+1} &= \theta_i - \eta \left( \frac{\partial L(f(\mathbf{x}^{(i)}; \theta_i), \mathbf{y}^{(i)})}{\partial \theta_i} + 2\gamma \theta_i \right) \\ &= \theta_i - \eta \frac{\partial L(f(\mathbf{x}^{(i)}; \theta_i), \mathbf{y}^{(i)})}{\partial \theta_i} - 2\eta \gamma \theta_i \end{aligned}$$

So, Weight decay scheme is equal to L2 regularization scheme when  $\lambda = 2\eta\gamma$ .

- 4.2 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 algorithm gives us Adam with an L2 regularization scheme. Identify the line of the algorithm that needs to change, and provide this one-line modification.

**Answer.**

The line that needs to be changed is the following:

$$\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)})$$

The modification consists of adding the L2 penalty term to the loss as following ( $g$  using the modification will become  $g'$  here):

$$\begin{aligned} \mathbf{g}' &\leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \left( \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)}) + \gamma \|\boldsymbol{\theta}\|_2^2 \right) \\ &\leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \left( \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)}) + \nabla_{\boldsymbol{\theta}} \gamma \|\boldsymbol{\theta}\|_2^2 \right) \\ &\leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \left( \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)}) + 2\gamma \boldsymbol{\theta} \right) \end{aligned}$$

- 4.3 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 gradient 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.

- Under the **Adam-L2-reg** scheme, which set of weights among  $\theta_{\text{small}}$  and  $\theta_{\text{large}}$  would you expect to be regularized (i.e., driven closer to zero) more strongly than the other? Why?
- Would your answer change for the **Adam-weight-decay** scheme? Why/why not?

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

**Answer.**

Adam  $L_2$  regularization scheme,  $\theta_{\text{small}}, \theta_{\text{large}}$ ?

$\theta_{\text{small}}$  will be pushed further 0 as  $r$  is computed using the gradient squared. Then as,  $\hat{r}$  is at the denominator of  $\Delta\theta$ , a larger  $\hat{r}$  will lead to small changes where small  $\hat{r}$  will change much more  $\theta$ .

Adam weight decay scheme,  $\theta_{small}, \theta_{large}$  ?

The answer to that question depends on the value of  $\theta$  as it is directly proportional to the weights updates  $\Delta\theta$  which is given by  $\Delta\theta = -\epsilon \frac{\hat{s}}{\sqrt{\hat{r} + \delta}} - \lambda\theta$ . Thus,  $\theta_{large}, \theta_{small}$  will change depending on the value of  $\theta$ .

- 4.4 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).

**Answer.**

Weight decay is better than L2 regularization as it avoid sparsity in the sense of it is not pushing  $\theta_{small}$  as L2 regularization does on Adam  $L_2$  regularization scheme. This is a good news because even if in some case sparsity can be interesting pushing  $\theta_{small}$  towards 0 leads to only keep high  $\theta_{large}$  which will create a lot of variance between weights that are set to 0 and the other ones. Furthermore, the goal of regularization is to avoid overfitting. For instance, by penalizing large weights which Adam using L2 regularization scheme doesn't seems to do in contrast to AdamW (Adam Weight Decay).