

Due Date: March 17th 23:00, 2020

Instructions

- For all questions, show your work!
- Starred questions are **hard** questions, not **bonus** questions.
- Please use a document preparation system such as LaTeX, unless noted otherwise.
- Unless noted that questions are related, assume that notation and definitions for each question are self-contained and independent
- All norms denote Euclidean norms unless otherwise specified.
- Submit your answers electronically via Gradescope.
- TAs for this assignment are **Jessica Thompson, Jonathan Cornford and Lluís Castrejon**.

Question 1 (4-4-4). In this question you will demonstrate that an estimate of the first moment of the gradient using an (exponential) running average is equivalent to using momentum, and is biased by a scaling factor. The goal of this question is for you to consider the relationship between different optimization schemes, and to practice noting and quantifying the effect (particularly in terms of bias/variance) of *estimating* a quantity.

Let \mathbf{g}_t be an unbiased sample of gradient at time step t and $\Delta\boldsymbol{\theta}_t$ be the update to be made. Initialize \mathbf{v}_0 to be a vector of zeros.

1. For $t \geq 1$, consider the following update rules:

- SGD with momentum:

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

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

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

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

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

Express the two update rules recursively ($\Delta\boldsymbol{\theta}_t$ as a function of $\Delta\boldsymbol{\theta}_{t-1}$). Show that these two update rules are equivalent; i.e. express (α, ϵ) as a function of (β, δ) .

Solution

Replacing $\Delta\boldsymbol{\theta}_t$ in the definition of the SGD momentum we obtain

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

Doing the same process on the SGD with running average, i.e replacing $\Delta\boldsymbol{\theta}_t = -\delta\mathbf{v}_t$ we arrive to

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

They are equivalent if $\alpha = \beta$ and $\epsilon = \delta(1 - \beta)$

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

Solution

$$\begin{aligned}\mathbf{v}_t &= \beta \mathbf{v}_{t-1} + (1 - \beta) \mathbf{g}_t \\ &= \beta^2 \mathbf{v}_{t-2} + \beta(1 - \beta) \mathbf{g}_{t-1} + (1 - \beta) \mathbf{g}_t \\ &= \beta^3 \mathbf{v}_{t-3} + \beta^2(1 - \beta) \mathbf{g}_{t-2} + \beta(1 - \beta) \mathbf{g}_{t-1} + (1 - \beta) \mathbf{g}_t\end{aligned}$$

So we can define the rule

$$\mathbf{v}_t = \beta^t \mathbf{v}_0 + (1 - \beta) \sum_{j=t}^0 \beta^j \mathbf{g}_{t-j}$$

Now we do the change of variable $t - j = i$ and assuming $\mathbf{v}_0 = 0$ we arrive to

$$\mathbf{v}_t = (1 - \beta) \sum_{i=1}^t \beta^{t-i} \mathbf{g}_i$$

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 .

Solution

The expected value of the last expression is given by

$$\mathbb{E}[\mathbf{v}_t] = (1 - \beta) \sum_{i=1}^t \beta^{t-i} \mathbb{E}[\mathbf{g}_i]$$

Considering $\mathbb{E}[\mathbf{g}_i]$ is independent of t and making the change of variable $t - i = j$ we obtain

$$\sum_{j=0}^{t-1} \beta^j = \frac{1 - \beta^t}{1 - \beta}$$

therefore

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

so such bias can be eliminated rescaling \mathbf{v}_t as $\mathbf{v}_t = \frac{\mathbf{v}_t}{1 - \beta^t}$

Question 2 (7-5-5-3). The point of this question is to understand and compare the effects of different regularizers (specifically dropout and weight decay) on the weights of a network. Consider a linear regression problem with input data $\mathbf{X} \in \mathbb{R}^{n \times d}$, weights $\mathbf{w} \in \mathbb{R}^{d \times 1}$ and targets $\mathbf{y} \in \mathbb{R}^{n \times 1}$. Suppose that dropout is applied to the input (with probability $1 - p$ of dropping the unit i.e. setting it to 0). Let $\mathbf{R} \in \mathbb{R}^{n \times d}$ be the dropout mask such that $\mathbf{R}_{ij} \sim \text{Bern}(p)$ is sampled i.i.d. from the Bernoulli distribution.

For a squared error loss function with dropout, we then have:

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

1. Let Γ be a diagonal matrix with $\Gamma_{ii} = (\mathbf{X}^\top \mathbf{X})_{ii}^{1/2}$. Show that the *expectation (over \mathbf{R})* of the loss function can be rewritten as $\mathbb{E}[L(\mathbf{w})] = \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1-p)\|\Gamma\mathbf{w}\|^2$. *Hint: Note we are trying to find the expectation over a squared term and use $\text{Var}(Z) = \mathbb{E}[Z^2] - \mathbb{E}[Z]^2$.*

Solution

Explicitly, the loss is done by

$$L(\mathbf{w}) = \sum_i \sum_j (y_i - x_{ij} R_{ij} w_j)^2$$

so using the hint

$$\mathbb{E}[L(\mathbf{w})] = \left(\mathbb{E} \left[\sum_i \sum_j y_i - x_{ij} R_{ij} w_j \right] \right)^2 + \text{Var} \left(\sum_i \sum_j y_i - x_{ij} R_{ij} w_j \right) \quad (1)$$

Firstly, we calculate the expectation

$$\left(\mathbb{E} \left[\sum_i \sum_j y_i - x_{ij} R_{ij} w_j \right] \right)^2 = \left(\sum_i \sum_j y_i - p x_{ij} w_j \right)^2 = \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 \quad (2)$$

where we used $\mathbb{E}[\mathbf{R}] = p$ because it is a Bernoulli distribution.

Now we calculate the variance of Eq.1

$$\text{Var} \left(\sum_i \sum_j y_i - x_{ij} R_{ij} w_j \right) = \text{Var} \left(\sum_i y_i - (\mathbf{X}_i \odot \mathbf{R}_i) \mathbf{w} \right)$$

By definition $\text{Var}(Z) = \mathbb{E}[(Z - \mathbb{E}[Z])^2]$ and using the result from Eq.2

$$\text{Var} \left(\sum_i y_i - (\mathbf{X}_i \odot \mathbf{R}_i) \mathbf{w} \right) = \mathbb{E} \left[\left(\sum_i y_i - (\mathbf{X}_i \odot \mathbf{R}_i) \mathbf{w} - \sum_i y_i + p \mathbf{X}_i \mathbf{w} \right)^2 \right] \quad (3)$$

$$= \mathbb{E} \left[\left(\sum_i p \mathbf{X}_i \mathbf{w} - (\mathbf{X}_i \odot \mathbf{R}_i) \mathbf{w} \right)^2 \right] \quad (4)$$

$$= \mathbb{E} \left[\left(\sum_i p \mathbf{X}_i \mathbf{w} - (\mathbf{X}_i \odot \mathbf{R}_i) \mathbf{w} \right)^T \left(\sum_i p \mathbf{X}_i \mathbf{w} - (\mathbf{X}_i \odot \mathbf{R}_i) \mathbf{w} \right) \right] \quad (5)$$

$$= \sum_i \mathbf{w}^T (-p^2 \mathbf{X}_i^T \mathbf{X}_i + \mathbf{X}_i^T \mathbf{X}_i \mathbb{E}[\mathbf{R}_i^T \mathbf{R}_i]) \mathbf{w} \quad (6)$$

Finally we define the value $\mathbb{E}[\mathbf{R}^T \mathbf{R}] = \text{Var}(\mathbf{R}) + \mathbb{E}[\mathbf{R}]^2 = p(1-p) + p^2 = p$ and replacing on Eq.6, we arrive to

$$\begin{aligned} \text{Var} \left(\sum_i y_i - (\mathbf{X}_i \odot \mathbf{R}_i) \mathbf{w} \right) &= p(1-p) \sum_i \mathbf{w}^T \mathbf{X}_i^T \mathbf{X}_i \mathbf{w} \\ &= p(1-p) \sum_i \mathbf{w}^T \Gamma_{ii} \mathbf{w} \\ &= p(1-p) \|\Gamma \mathbf{w}\|^2 \end{aligned}$$

Joining this result with Eq.2 we conclude that

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

2. Show that the solution $\mathbf{w}^{\text{dropout}}$ that minimizes the expected loss from question 2.1 satisfies

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

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

Solution From the solution above

$$E[L(\mathbf{w})] = (\mathbf{y} - p\mathbf{X}\mathbf{w})^T (\mathbf{y} - p\mathbf{X}\mathbf{w}) + p(1-p)(\Gamma\mathbf{w})^T (\Gamma\mathbf{w})$$

therefore

$$\frac{\partial \mathbb{E}}{\partial \mathbf{w}} = -(\mathbf{y}^T - p\mathbf{w}^T \mathbf{X}^T)(p\mathbf{X}) + p(1-p)(\mathbf{w}^T \Gamma^T) \Gamma$$

making the derivative to zero and doing the algebraic operations on the last equation we find that

$$\begin{aligned} p^2 \mathbf{w}^T \left[\mathbf{X}^T \mathbf{X} + \left(\frac{1-p}{p} \right) \Gamma^2 \right] &= p\mathbf{y}^T \mathbf{X} \Rightarrow \\ p\mathbf{w} &= \left[\mathbf{X}^T \mathbf{X} + \left(\frac{1-p}{p} \right) \Gamma^2 \right]^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

Finally

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

with $\lambda^{\text{drop}} \equiv \frac{1-p}{p}$. If $p = 1$ we do not have regularization, if p is very small the regularization very high.

3. Express the loss function for a linear regression problem without dropout and with L^2 regularization, with regularization coefficient λ^{L^2} . Derive its closed form solution \mathbf{w}^{L^2} .

Solution In this case, the lost is defined as

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

i.e

$$\begin{aligned} L(\mathbf{w}) &= (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda^{L^2} \mathbf{w}^T \mathbf{w} \Rightarrow \\ \frac{\partial L}{\partial \mathbf{w}} &= -(\mathbf{y}^T - \mathbf{w}^T \mathbf{X}^T) \mathbf{X} + \lambda^{L^2} \mathbf{w}^T \end{aligned}$$

Making the derivative equal to zero we obtain the condition

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

4. Compare the results of 2.2 and 2.3: identify specific differences in the equations you arrive at, and discuss qualitatively what the equations tell you about the similarities and differences in the effects of weight decay and dropout (1-3 sentences).

Solution In the standard L_2 the weights are penalized uniformly. In dropout we must use a mask so the weights are nor penalized uniformly.

Question 3 (6-10-2). The goal of this question is for you to understand the reasoning behind different parameter initializations for deep networks, particularly to think about the ways that the initialization affects the activations (and therefore the gradients) of the network. Consider the following equation for the t -th layer of a deep network:

$$\mathbf{h}^{(t)} = g(\mathbf{a}^{(t)}) \quad \mathbf{a}^{(t)} = \mathbf{W}^{(t)}\mathbf{h}^{(t-1)} + \mathbf{b}^{(t)}$$

where $\mathbf{a}^{(t)}$ are the pre-activations and $\mathbf{h}^{(t)}$ are the activations for layer t , g is an activation function, $\mathbf{W}^{(t)}$ is a $d^{(t)} \times d^{(t-1)}$ matrix, and $\mathbf{b}^{(t)}$ is a $d^{(t)} \times 1$ bias vector. The bias is initialized as a constant vector $\mathbf{b}^{(t)} = [c, \dots, c]^\top$ for some $c \in \mathbb{R}$, and the entries of the weight matrix are initialized by sampling i.i.d. from a Gaussian distribution $W_{ij}^{(t)} \sim \mathcal{N}(\mu, \sigma^2)$.

Your task is to design an initialization scheme that would achieve a vector of **pre-activations** at layer t whose elements are zero-mean and unit variance (i.e.: $\mathbb{E}[a_i^{(t)}] = 0$ and $\text{Var}(a_i^{(t)}) = 1$, $1 \leq i \leq d^{(t)}$) for the assumptions about either the activations or pre-activations of layer $t - 1$ listed below. Note we are not asking for a general formula; you just need to provide one setting that meets these criteria (there are many possibilities).

1. First assume that the activations of the previous layer satisfy $\mathbb{E}[h_i^{(t-1)}] = 0$ and $\text{Var}(h_i^{(t-1)}) = 1$ for $1 \leq i \leq d^{(t-1)}$. Also, assume entries of $\mathbf{h}^{(t-1)}$ are uncorrelated (the answer should not depend on g).

(a) Show $\text{Var}(XY) = \text{Var}(X)\text{Var}(Y) + \text{Var}(X)\mathbb{E}[Y]^2 + \text{Var}(Y)\mathbb{E}[X]^2$ when $X \perp Y$

Solution

By definition

$$\text{Var}(XY) = \mathbb{E}[X^2Y^2] - \mathbb{E}[XY]^2 = \mathbb{E}[X^2]\mathbb{E}[Y^2] - (\mathbb{E}[X]\mathbb{E}[Y])^2 \quad (7)$$

On the other hand, considering the definition of variance, we can say that

$$\begin{aligned} \text{Var}(X)\mathbb{E}[Y]^2 &= (\mathbb{E}[X^2] - \mathbb{E}[X]^2)\mathbb{E}[Y]^2 = \mathbb{E}[X^2]\mathbb{E}[Y]^2 - \mathbb{E}[X]^2\mathbb{E}[Y]^2 \\ \text{Var}(Y)\mathbb{E}[X]^2 &= (\mathbb{E}[Y^2] - \mathbb{E}[Y]^2)\mathbb{E}[X]^2 = \mathbb{E}[Y^2]\mathbb{E}[X]^2 - \mathbb{E}[Y]^2\mathbb{E}[X]^2 \end{aligned} \quad (8)$$

And using again the definition of variance, we have

$$\text{Var}(X)\text{Var}(Y) = \mathbb{E}[X^2]\mathbb{E}[Y^2] - \mathbb{E}[X^2]\mathbb{E}[Y]^2 - \mathbb{E}[X]^2\mathbb{E}[Y^2] + \mathbb{E}[X]^2\mathbb{E}[Y]^2 \quad (9)$$

Therefore adding Eq.9 with Eq.8 we obtain

$$\text{Var}(X)\text{Var}(Y) + \text{Var}(X)\mathbb{E}[Y]^2 + \text{Var}(Y)\mathbb{E}[X]^2 = \mathbb{E}[X^2]\mathbb{E}[Y^2] - \mathbb{E}[X]^2\mathbb{E}[Y]^2 \quad (10)$$

which is just Eq.7, therefore

$$\text{Var}(XY) = \text{Var}(X)\text{Var}(Y) + \text{Var}(X)\mathbb{E}[Y]^2 + \text{Var}(Y)\mathbb{E}[X]^2 \quad (11)$$

- (b) Write $\mathbb{E}[a_i^{(t)}]$ and $\text{Var}(a_i^{(t)})$ in terms of $c, \mu, \sigma^2, \text{Var}(h_i^{(t-1)}), \mathbb{E}[h_i^{(t-1)}]$.

Solution

We know that $a_i^{(t)} = \mathbf{W}_i^{(t)}\mathbf{h}^{(t-1)} + b_i$ so

$$\mathbb{E}[a_i^{(t)}] = \mathbb{E}[\mathbf{W}_i^{(t)}]\mathbb{E}[\mathbf{h}^{(t-1)}] + \mathbb{E}[b_i] \quad (12)$$

$$\mathbb{E}[a_i^{(t)}] = \mu \mathbf{1}_{d^{(t-1)}}^T \mathbb{E}[\mathbf{h}^{(t-1)}] + c \quad (13)$$

$$\mathbb{E}[a_i^{(t)}] = \mu d^{(t-1)} \mathbb{E}[h_i^{(t-1)}] + c \quad (14)$$

because $W_{ij}^{(t)} \sim \mathcal{N}(\mu, \sigma^2)$ and $\mathbb{E}[b_i]$ is a constant c . Remark that each component of the matrix \mathbf{W} has a normal distribution, so the expected value of a row \mathbf{W}_i is $\mu \mathbf{1}_{d^{(t-1)}}^T$, where $\mathbf{1}_{d^{(t-1)}}^T$ is a row vector of ones with the dimensions of the layer $t-1$. In the same way, we assume each component of the vector \mathbf{h} has the same distribution so $\mathbb{E}[\mathbf{h}] = \mathbf{1}_{d^{(t-1)}} \mathbb{E}[\mathbf{h}_i]$, with $\mathbf{1}_{d^{(t-1)}}$ a column vector of ones with the dimensions of the layer $t-1$. Using the Eq. 11 the variance of $a_i^{(t)}$ is

$$\text{Var}(a_i^{(t)}) = \text{Var}(\mathbf{W}_i^{(t)} \mathbf{h}^{(t-1)}) + \text{Var}(b_i^t) \quad (15)$$

$$\text{Var}(a_i^{(t)}) = \text{Var}(\mathbf{W}_i^{(t)}) \text{Var}(\mathbf{h}^{(t-1)}) + \text{Var}(\mathbf{W}_i^{(t)}) \mathbb{E}[\mathbf{h}^{(t-1)}]^2 + \text{Var}(\mathbf{h}^{(t-1)}) \mathbb{E}[\mathbf{W}_i^{(t)}]^2 \quad (16)$$

$$\text{Var}(a_i^{(t)}) = \sigma^2 d^{(t-1)} \text{Var}(\mathbf{h}_i^{(t-1)}) + \sigma^2 d^{(t-1)} \mathbb{E}[\mathbf{h}_i^{(t-1)}]^2 + \mu^2 d^{(t-1)} \text{Var}(\mathbf{h}_i^{(t-1)}) \quad (17)$$

In particular if $\mathbb{E}[\mathbf{h}_i^{(t-1)}] = 0$ and $\text{Var}(\mathbf{h}_i^{(t-1)}) = 1$ we have from Eq.14 and Eq.17

$$\mathbb{E}[a_i^{(t)}] = c \quad (18)$$

and

$$\text{Var}(a_i^{(t)}) = (\sigma^2 + \mu^2)(d^{(t-1)}) \quad (19)$$

- (c) Give values for c , μ , and σ^2 as a function of $d^{(t-1)}$ such that $\mathbb{E}[a_i^{(t)}] = 0$ and $\text{Var}(a_i^{(t)}) = 1$ for $1 \leq i \leq d^{(t)}$.

Solution

From Eq.18 $\mathbb{E}[a_i^{(t)}] = 0$ if $c = 0$ and a condition to have $\text{Var}(a_i^{(t)}) = 1$ assuming $\mu = 0$ from Eq.19 is

$$\sigma = \sqrt{\frac{1}{d^{(t-1)}}}$$

2. Now assume that the pre-activations of the previous layer satisfy $\mathbb{E}[a_i^{(t-1)}] = 0$, $\text{Var}(a_i^{(t-1)}) = 1$ and $a_i^{(t-1)}$ has a symmetric distribution for $1 \leq i \leq d^{(t-1)}$. Assume entries of $\mathbf{a}^{(t-1)}$ are uncorrelated. Consider the case of ReLU activation: $g(x) = \max\{0, x\}$.

- (a) Derive $\mathbb{E}[(h_i^{(t-1)})^2]$

Solution

$$\mathbb{E}[(h_i^{(t-1)})^2] = \mathbb{E}[g(a_i^{(t-1)})^2] = \mathbb{E}[\max\{0, a_i^{(t-1)}\}^2]$$

Using the definition of expectation and considering $f_p(a_i^{(t-1)})$ as the probability density of function, we have

$$\begin{aligned} \mathbb{E}[(h_i^{(t-1)})^2] &= \int_{-\infty}^{\infty} \max\{0, a_i^{(t-1)}\}^2 f_p(a_i^{(t-1)}) da_i^{(t-1)} \\ &= \int_0^{\infty} (a_i^{(t-1)})^2 f_p(a_i^{(t-1)}) da_i^{(t-1)} \\ &= \frac{1}{2} \int_{-\infty}^{\infty} (a_i^{(t-1)})^2 f_p(a_i^{(t-1)}) da_i^{(t-1)} * \\ &= \frac{1}{2} \mathbb{E}[(a_i^{(t-1)})^2] \end{aligned}$$

In the step * we consider that the argument of the integral is an even function. By definition

$$\mathbb{E}[(a_i^{(t-1)})^2] = \left[\text{Var}(a_i^{(t-1)}) + \mathbb{E}[(a_i^{(t-1)})]^2 \right] = 1$$

Therefore, we conclude that

$$\mathbb{E}[(h_i^{(t-1)})^2] = \frac{1}{2}$$

- (b) Using the result from (a), give values for c , μ , and σ^2 as a function of $d^{(t-1)}$ such that $\mathbb{E}[a_i^{(t)}] = 0$ and $\text{Var}(a_i^{(t)}) = 1$ for $1 \leq i \leq d^{(t)}$.

Solution

From Eq.12 we know that

$$\mathbb{E}[a_i^{(t)}] = \mathbb{E}[\mathbf{W}_i^{(t)}] \mathbb{E}[\mathbf{h}^{(t-1)}] + \mathbb{E}[b_i^{(t)}] \quad (20)$$

$$= \mathbb{E}[W_{ij}^{(t)}] \mathbf{1}_{d^{(t-1)}}^T \mathbb{E}[\mathbf{h}_i^{(t-1)}] \mathbf{1}_{d^{(t-1)}} + c \quad (21)$$

and considering $\mathbb{E}[a_i^{(t)}] = 0$ we arrive to

$$\mathbb{E}[\mathbf{h}_i^{(t-1)}] = -\frac{c}{\mathbb{E}[W_{ij}^{(t)}]} d^{(t-1)} \quad (22)$$

Now we will calculate

$$\text{Var}(a_i^{(t)}) = \text{Var}(\mathbf{W}_i^{(t)} \mathbf{h}^{(t-1)})$$

using the property shown in Eq.11

$$\text{Var}(a_i^{(t)}) = \text{Var}(\mathbf{W}_i^{(t)}) \text{Var}(\mathbf{h}^{(t-1)}) + \text{Var}(\mathbf{W}_i^{(t)}) \mathbb{E}[\mathbf{h}^{(t-1)}]^2 + \text{Var}(\mathbf{h}^{(t-1)}) \mathbb{E}[\mathbf{W}_i^{(t)}]^2 \quad (23)$$

Considering that $\mathbb{E}[(\mathbf{h}_i^{(t-1)})^2] = \frac{1}{2}$ (item a) and considering that

$$\text{Var}[\mathbf{h}_i^{(t-1)}] = \mathbb{E}[(h_i^{(t-1)})^2] - \mathbb{E}[(h_i^{(t-1)})]^2$$

we obtain

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

Replacing on Eq.23:

$$\begin{aligned} \text{Var}(a_i^{(t)}) = & \text{Var}(W_{ij}^{(t)}) \left(\frac{1}{2} - \mathbb{E}[h_i^{(t-1)}]^2 \right) d^{(t-1)} + \text{Var}(W_{ij}^{(t)}) \mathbb{E}[h_i^{(t-1)}]^2 d^{(t-1)} \\ & + \mathbb{E}[W_{ij}^{(t)}]^2 \left(\frac{1}{2} - \mathbb{E}[h_i^{(t-1)}]^2 \right) d^{(t-1)} \end{aligned}$$

and taking into account that $\text{Var}(a_i^{(t)}) = 1$ and Eq.22 we arrive to

$$1 = \frac{1}{2} d^{(t-1)} \left[\text{Var}(W_{ij}^{(t)}) + \mathbb{E}[W_{ij}^{(t)}]^2 \right] - \frac{c^2}{d^{(t-1)}} \quad (24)$$

If we have a Gaussian distribution $\text{Var}(W_{ij}^{(t)}) = \sigma^2$ and $\mathbb{E}[W_{ij}^{(t)}] = \mu$ Substituting on Eq.24 we achieve

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

Additionally, if $c = 0$, and $\mu = 0$ we obtain

$$\sigma = \sqrt{\frac{2}{d^{(t-1)}}}$$

(c) What popular initialization scheme has this form?

Solution The scheme above is known as He normal distribution.

(d) Why do you think this initialization would work well in practice? Answer in 1-2 sentences.

Solution

Using this scheme (which is the ReLu version of Glorot) we have more control of the number of neurons that are activated and the scale of activation.

3. For both assumptions (1,2) give values α, β for $W_{ij}^{(t)} \sim Uniform(\alpha, \beta)$ such that $\mathbb{E}[a_i^{(t)}] = 0$ and $\text{Var}(a_i^{(t)}) = 1$.

Solution

Assumption 1

As $a_i^{(t)} = \mathbf{W}_i^{(t)} \mathbf{h}^{(t-1)} + b_i^{(t)}$ we have

$$\begin{aligned}\mathbb{E}[a_i^{(t)}] &= \mathbb{E}[\mathbf{W}_i^{(t)}] \mathbb{E}[\mathbf{h}^{(t-1)}] + \mathbb{E}[b_i^{(t)}] \Rightarrow \\ \mathbb{E}[a_i^{(t)}] &= c\end{aligned}$$

Now

$$\begin{aligned}\text{Var}(a_i^{(t)}) &= \text{Var}(\mathbf{W}_i^{(t)}) \text{Var}(\mathbf{h}^{(t-1)}) + \text{Var}(\mathbf{W}_i^{(t)}) \mathbb{E}[\mathbf{h}^{(t-1)}]^2 + \text{Var}(\mathbf{h}^{(t-1)}) \mathbb{E}[\mathbf{W}_i^{(t)}]^2 \\ &= \text{Var}(W_{ij}^{(t)}) \mathbf{1}_{d^{(t-1)}}^T \text{Var}(h_i^{(t-1)}) \mathbf{1}_{d^{(t-1)}} + \mathbb{E}[W_{ij}^{(t)}]^2 \mathbf{1}_{d^{(t-1)}}^T \text{Var}(h_i^{(t-1)}) \mathbf{1}_{d^{(t-1)}}\end{aligned}$$

Assuming $\text{Var}(a_i^{(t)}) = 1$ we find the condition

$$1 = \text{Var}(W_{ij}^{(t)}) d^{(t-1)} + \mathbb{E}[W_{ij}^{(t)}]^2 d^{(t-1)} \quad (25)$$

For solving this equation we must calculate $\mathbb{E}[W_{ij}^{(t)}]$ to find $\text{Var}(W_{ij})$ using

$$\text{Var}(W_{ij}) = \mathbb{E}[(W_{ij}^{(t)})^2] - \mathbb{E}[W_{ij}^{(t)}]^2 \quad (26)$$

with $W_{ij} \sim U(\alpha, \beta)$ and

$$U(x) = \begin{cases} \frac{1}{\beta - \alpha}, & \text{if } \alpha \leq x \leq \beta \\ 0, & \text{otherwise} \end{cases}$$

In this order of ideas we can calculate the integrals

$$\mathbb{E}[W_{ij}^{(t)}] = \frac{1}{\beta - \alpha} \int_{\alpha}^{\beta} x \, dx = \frac{\beta + \alpha}{2} \quad (27)$$

$$\mathbb{E}[(W_{ij}^{(t)})^2] = \frac{1}{\beta - \alpha} \int_{\alpha}^{\beta} x^2 \, dx = \frac{\beta^2 + \alpha\beta + \alpha^2}{3} \quad (28)$$

Replacing Eq.27 and Eq.28 on Eq.26 we obtain

$$\text{Var}(W_{ij}) = \frac{(\beta - \alpha)^2}{12}$$

and finally substituting on Eq.25

$$\frac{1}{d^{(t-1)}} = \frac{(\beta - \alpha)^2}{12} + \frac{(\beta + \alpha)^2}{4}$$

$$\frac{3}{d^{(t-1)}} = \beta^2 + \alpha\beta + \alpha^2$$

A solution could be $\alpha = -\sqrt{\frac{3}{d^{(t-1)}}}$, $\beta = \sqrt{\frac{3}{d^{(t-1)}}}$ which is a glorot distribution.

Assumption 2

From the Eq.22 we know that $\mathbb{E}[a_i^{(t)}] = 0$ if $c = 0$ and from Eq.24

$$1 = \frac{1}{2}d^{(t-1)} \left[\text{Var}(W_{ij}^{(t)}) + \mathbb{E}[W_{ij}^{(t)}]^2 \right] - \frac{c^2}{d^{(t-1)}} \Rightarrow$$

$$1 = \frac{1}{2}d^{(t-1)} \left[\text{Var}(W_{ij}^{(t)}) + \mathbb{E}[W_{ij}^{(t)}]^2 \right] \Rightarrow$$

$$1 = \frac{1}{2}d^{(t-1)} \left[\frac{(\beta - \alpha)^2}{12} + \frac{(\beta + \alpha)^2}{4} \right] \Rightarrow$$

$$\frac{6}{d^{(t-1)}} = \beta^2 + \alpha\beta + \alpha^2$$

A solution could be $c = 0$, $\alpha = -\sqrt{\frac{6}{d^{(t-1)}}}$, $\beta = \sqrt{\frac{6}{d^{(t-1)}}}$ it is known as He uniform distribution.

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

1. Batch normalization, layer normalization and instance normalization all involve calculating the mean $\boldsymbol{\mu}$ and variance $\boldsymbol{\sigma}^2$ 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: $\boldsymbol{\mu}_{batch}$, $\boldsymbol{\mu}_{layer}$, $\boldsymbol{\mu}_{instance}$, $\boldsymbol{\sigma}_{batch}^2$, $\boldsymbol{\sigma}_{layer}^2$, and $\boldsymbol{\sigma}_{instance}^2$.

$$\left[\begin{bmatrix} 1, 3, 2 \\ 1, 2, 3 \end{bmatrix}, \begin{bmatrix} 3, 3, 2 \\ 2, 4, 4 \end{bmatrix}, \begin{bmatrix} 4, 2, 2 \\ 1, 2, 4 \end{bmatrix}, \begin{bmatrix} 3, 3, 2 \\ 3, 3, 2 \end{bmatrix} \right]$$

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.

Solution

- **Batch normalization:** Considering the tensor of the exercise as $\mathbf{X} \in \mathbb{R}^{T \times C \times W}$ where T is the number of batches (4), C the number of channels (2), and W the number of features (3) the mean batch normalization $\boldsymbol{\mu}$ is defined as

$$\mu_i = \frac{1}{TW} \sum_{t=1}^T \sum_{l=1}^W X_{til} = \frac{1}{12} \sum_{t=1}^4 \sum_{l=1}^3 X_{til}$$

doing the calculations we obtain $\boldsymbol{\mu}_{batch} = [2.5, 2.583]$. On the other hand, the batch variance $\boldsymbol{\sigma}_{batch}^2$ is defined as

$$\sigma_i^2 = \frac{1}{TW} \sum_{t=1}^T \sum_{l=1}^W (X_{til} - \mu_i)^2 = \frac{1}{12} \sum_{t=1}^4 \sum_{l=1}^3 (X_{til} - \mu_i)^2$$

obtaining $\boldsymbol{\sigma}_{batch}^2 = [0.583, 1.076]$

- **Layer normalization:** Using the same notation, mean layer normalization $\boldsymbol{\mu}_{layer}$ is defined as

$$\mu_t = \frac{1}{CW} \sum_{i=1}^C \sum_{l=1}^W X_{til} = \frac{1}{6} \sum_{i=1}^2 \sum_{l=1}^3 X_{til}$$

arriving to $\boldsymbol{\mu}_{layer} = [2.0, 3.0, 2.5, 2.67]$. The layer variance is defined as

$$\sigma_t^2 = \frac{1}{CW} \sum_{i=1}^C \sum_{l=1}^W (X_{til} - \mu_t)^2 = \frac{1}{6} \sum_{i=1}^2 \sum_{l=1}^3 (X_{til} - \mu_t)^2$$

obtaining $\boldsymbol{\sigma}_{layer}^2 = [0.67, 0.67, 1.25, 0.22]$

- **Instance normalization:** Finally the mean instance normalization is defined as

$$\mu_{ti} = \frac{1}{W} \sum_{l=1}^W X_{til} = \frac{1}{3} \sum_{l=1}^3 X_{til}$$

so the matrix $\boldsymbol{\mu}_{instance}$ is

$$\boldsymbol{\mu}_{instance} = \begin{bmatrix} 2.0, 2.0 \\ 2.67, 3.33 \\ 2.67, 3.33 \\ 2.67, 2.67 \end{bmatrix}$$

and the variance is defined as

$$\sigma_{ti} = \frac{1}{W} \sum_{l=1}^W (X_{til} - \mu_{ti})^2 = \frac{1}{3} \sum_{l=1}^3 (X_{til} - \mu_{ti})^2$$

achieving the result

$$\boldsymbol{\sigma}_{instance}^2 = \begin{bmatrix} 0.67, 0.67 \\ 0.22, 0.88 \\ 0.88, 1.55 \\ 0.22, 0.22 \end{bmatrix}$$

- For the next two subquestions, we consider the following parametrization of a weight vector \boldsymbol{w} :

$$\boldsymbol{w} := \gamma \frac{\boldsymbol{u}}{\|\boldsymbol{u}\|}$$

where γ is scalar parameter controlling the magnitude and \boldsymbol{u} is a vector controlling the direction of \boldsymbol{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 $y = \boldsymbol{u}^\top \boldsymbol{x}$. Assume the data \boldsymbol{x} (a random vector) is whitened ($\text{Var}(\boldsymbol{x}) = \boldsymbol{I}$) and centered at 0 ($\mathbb{E}[\boldsymbol{x}] = \mathbf{0}$). Show that $\hat{y} = \boldsymbol{w}^\top \boldsymbol{x} + \beta$.

Solution

We know that $y = \boldsymbol{u}^\top \boldsymbol{x}$ so $\text{Var}(y) = \text{Var}(\boldsymbol{u}^\top \boldsymbol{x})$, using the definition of variance

$$\text{Var}(\boldsymbol{u}^\top \boldsymbol{x}) = \mathbb{E}[(\boldsymbol{u}^\top \boldsymbol{x})^2] - \mathbb{E}[\boldsymbol{u}^\top \boldsymbol{x}]^2$$

but $\mathbb{E}[\mathbf{x}] = 0$ therefore

$$\text{Var}(\mathbf{u}^T \mathbf{x}) = \mathbb{E}[(\mathbf{u}^T \mathbf{x})^2] = \mathbb{E}[(\mathbf{u}^T \mathbf{x})(\mathbf{x}^T \mathbf{u})] = \|\mathbf{u}\|$$

On the other hand, we know that $\text{Var}(\mathbf{u}^T \mathbf{x}) = \sigma_y$ so $\sigma_y = \|\mathbf{u}\|$. By definition in the exercise $\hat{y} = \gamma \cdot \frac{y - \mu_y}{\sigma_y} + \beta$ so

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

but $\mu_y = \mathbb{E}[y] = 0$ and we conclude that

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

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 $\mathbf{W}^\perp = \left(\mathbf{I} - \frac{\mathbf{u} \mathbf{u}^T}{\|\mathbf{u}\|^2} \right)$. Note that $\mathbf{W}^\perp \mathbf{u} = \mathbf{0}$.

Solution

Applying chain rule $\nabla_{\mathbf{u}} L = \nabla_{\mathbf{w}} L \nabla_{\mathbf{u}} \mathbf{W}$. Furthermore, $\mathbf{W} = \gamma \frac{\mathbf{u}}{\|\mathbf{u}\|} = \gamma \frac{\mathbf{u}}{(\mathbf{u}^T \mathbf{u})^{1/2}}$ so

$$\begin{aligned} \nabla_{\mathbf{u}} \mathbf{W} &= \frac{\gamma}{(\mathbf{u}^T \mathbf{u})} \left[(\mathbf{u}^T \mathbf{u})^{1/2} - \frac{1}{2} \frac{\mathbf{u}}{(\mathbf{u}^T \mathbf{u})^{1/2}} d(\mathbf{u}^T \mathbf{u}) \right] \\ &= \frac{\gamma}{\|\mathbf{u}\|^2} \left[\|\mathbf{u}\| - \frac{1}{2} \frac{\mathbf{u}}{\|\mathbf{u}\|} 2\mathbf{u}^T \right] \\ &= \frac{\gamma}{\|\mathbf{u}\|} \left[\mathbf{I} - \frac{\mathbf{u} \mathbf{u}^T}{\|\mathbf{u}\|^2} \right] \end{aligned}$$

and we arrive to $\nabla_{\mathbf{u}} L = s \mathbf{W}^\perp \nabla_{\mathbf{w}} L$ with $s \equiv \frac{\gamma}{\|\mathbf{u}\|}$ and $\mathbf{W}^\perp = \left(\mathbf{I} - \frac{\mathbf{u} \mathbf{u}^T}{\|\mathbf{u}\|^2} \right)$

Question 5 (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. 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 $t-1$.

Solution We have that

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

if we assume $\mathbf{g}_{t-1} = \sigma(\mathbf{h}_{t-1})$ we obtain $\mathbf{g}_t = \sigma(\mathbf{h}_t)$

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

- *2. Let $\|\mathbf{A}\|$ denote the L_2 operator norm² of matrix \mathbf{A} ($\|\mathbf{A}\| := \max_{\mathbf{x}: \|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|$). Assume $\sigma(x)$ has bounded derivative, i.e. $|\sigma'| \leq \gamma$ for some $\gamma > 0$ and for all x . We denote as $\lambda_1(\cdot)$ the largest eigenvalue of a symmetric matrix. Show that if the largest eigenvalue of the weights is bounded by $\frac{\delta^2}{\gamma^2}$ for some $0 \leq \delta < 1$, gradients of the hidden state will vanish over time, i.e.

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

Use the following properties of the L_2 operator norm

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

Solution We begin with the condition $\mathbf{h}_t = \mathbf{W}\sigma(\mathbf{h}_{t-1}) + \mathbf{U}\mathbf{x}_t + \mathbf{b}$, so we can say that

$$\begin{aligned} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_0} &= \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \frac{\partial \mathbf{h}_{t-1}}{\partial \mathbf{h}_0} \\ &= \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \frac{\partial \mathbf{h}_{t-1}}{\partial \mathbf{h}_{t-2}} \frac{\partial \mathbf{h}_{t-2}}{\partial \mathbf{h}_0} \\ &= \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \frac{\partial \mathbf{h}_{t-1}}{\partial \mathbf{h}_{t-2}} \cdots \frac{\partial \mathbf{h}_1}{\partial \mathbf{h}_0} \end{aligned} \tag{29}$$

but every derivative in the equation above satisfies

$$\left\| \frac{\partial \mathbf{h}_{t-i}}{\partial \mathbf{h}_{t-i-1}} \right\| = \left\| \mathbf{W} \frac{\partial \sigma}{\partial \mathbf{h}_{t-i-1}} \right\| \leq \|\mathbf{W}\| \left\| \frac{\partial \sigma}{\partial \mathbf{h}_{t-i-1}} \right\|$$

but $\|\mathbf{W}\| = \sqrt{\lambda_1 \mathbf{W}^\top \mathbf{W}}$, and $\left\| \frac{\partial \sigma}{\partial \mathbf{h}_{t-i-1}} \right\| \leq \delta$, therefore using Eq.29 we obtain the result

$$\begin{aligned} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_0} &= \delta \times \delta \times \cdots \delta \text{ (} T \text{ times)} \\ &= \delta^T \end{aligned}$$

As $0 \leq \delta < 1$ we conclude that

$$\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_0} \right\| \rightarrow 0$$

when $T \rightarrow \infty$

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

Solution

From the last exercise we have the condition

$$\left\| \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_0} \right\| \leq \sqrt{\lambda_1 \mathbf{W}^\top \mathbf{W}} \gamma$$

But if $\lambda_1 > \frac{\delta^2}{\gamma^2}$ we do not have a conclusive bound condition for the norm of $\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_0}$, so this condition must be a necessary condition for gradient exploding but it is not sufficient.

2. The L_2 operator norm of a matrix \mathbf{A} is an *induced norm* corresponding to the L_2 norm of vectors. You can try to prove the given properties as an exercise.

Question 6 (4-8-8). Consider the following Bidirectional RNN:

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

where the superscripts f and b correspond to the forward and backward RNNs respectively and σ denotes the logistic sigmoid function. Let \mathbf{z}_t be the true target of the prediction \mathbf{y}_t and consider the sum of squared loss $L = \sum_t L_t$ where $L_t = \|\mathbf{z}_t - \mathbf{y}_t\|_2^2$.

In this question our goal is to obtain an expression for the gradients $\nabla_{\mathbf{W}^{(f)}} L$ and $\nabla_{\mathbf{U}^{(b)}} L$.

1. First, complete the following computational graph for this RNN, unrolled for 3 time steps (from $t = 1$ to $t = 3$). Label each node with the corresponding hidden unit and each edge with the corresponding weight. Note that it includes the initial hidden states for both the forward and backward RNNs.

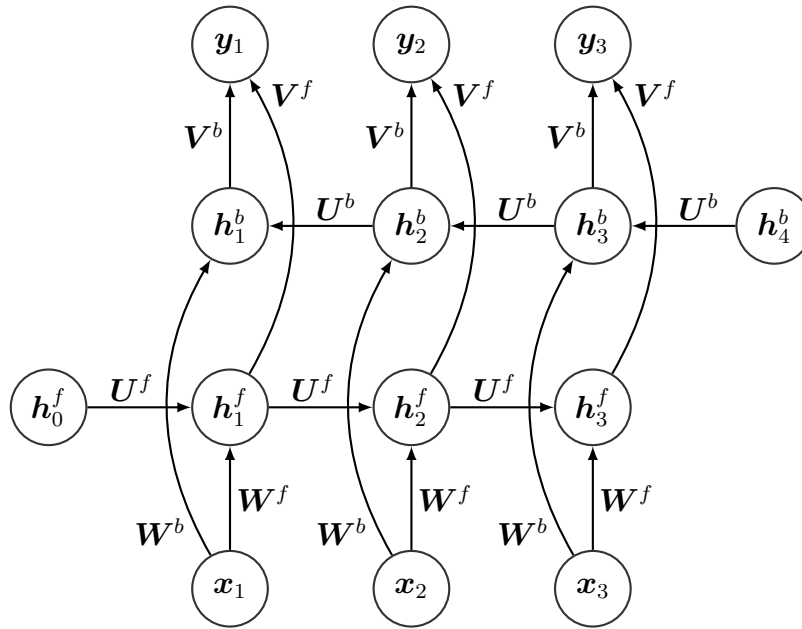


FIGURE 1 – Computational graph of the bidirectional RNN unrolled for three timesteps.

2. Using total derivatives we can express the gradients $\nabla_{\mathbf{h}_t^{(f)}} L$ and $\nabla_{\mathbf{h}_t^{(b)}} L$ recursively in terms of $\nabla_{\mathbf{h}_{t+1}^{(f)}} L$ and $\nabla_{\mathbf{h}_{t-1}^{(b)}} L$ as follows:

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

Derive an expression for $\nabla_{\mathbf{h}_t^{(f)}} L_t$, $\nabla_{\mathbf{h}_t^{(b)}} L_t$, $\frac{\partial \mathbf{h}_{t+1}^{(f)}}{\partial \mathbf{h}_t^{(f)}}$ and $\frac{\partial \mathbf{h}_{t-1}^{(b)}}{\partial \mathbf{h}_t^{(b)}}$.

Solution

For the derivatives we will use the set of equations:

$$\mathbf{h}_t^{(f)} = \sigma(\mathbf{W}^{(f)} \mathbf{x}_t + \mathbf{U}^{(f)} \mathbf{h}_{t-1}^{(f)}) \quad (30)$$

$$\mathbf{h}_t^{(b)} = \sigma(\mathbf{W}^{(b)} \mathbf{x}_t + \mathbf{U}^{(b)} \mathbf{h}_{t+1}^{(b)}) \quad (31)$$

$$\mathbf{y}_t = \mathbf{V}^{(f)} \mathbf{h}_t^{(f)} + \mathbf{V}^{(b)} \mathbf{h}_t^{(b)} \quad (32)$$

$$L_t = ||z_t - y_t||_2 \quad (33)$$

therefore

$$\begin{aligned} \nabla_{\mathbf{h}_t^{(f)}} L &= \frac{\partial L_t}{\partial \mathbf{h}_t^{(f)}} = \frac{\partial L_t}{\partial y_t} \left(\frac{\partial y_t}{\partial \mathbf{h}_t^{(f)}} \right)^T \\ &= -2(\mathbf{V}^{(f)})^T (z_t - y_t) \end{aligned}$$

$$\begin{aligned} \nabla_{\mathbf{h}_t^{(b)}} L &= \frac{\partial L_t}{\partial \mathbf{h}_t^{(b)}} = \frac{\partial L_t}{\partial y_t} \left(\frac{\partial y_t}{\partial \mathbf{h}_t^{(b)}} \right)^T \\ &= -2(\mathbf{V}^{(b)})^T (z_t - y_t) \end{aligned}$$

For making the derivative $\frac{\partial \mathbf{h}_{t+1}^{(f)}}{\partial \mathbf{h}_t^{(f)}}$ we must consider that $\frac{\partial \mathbf{h}_{t+1}^{(f)}}{\partial \mathbf{h}_t^{(f)}} = \frac{\partial \mathbf{h}_{t+1}^{(f)}}{\partial \sigma} \frac{\partial \sigma}{\partial \mathbf{h}_t^{(f)}}$ but

$$\frac{\partial \mathbf{h}_{t+1}^{(f)}}{\partial \sigma} = \sigma(\arg)(1 - \sigma(\arg)) = \mathbf{h}_{t+1}^{(f)}(1 - \mathbf{h}_{t+1}^{(f)})$$

with $\arg \equiv \mathbf{W}^{(f)} \mathbf{x}_t + \mathbf{U}^{(f)} \mathbf{h}_{t-1}^{(f)}$ and $\frac{\partial \sigma}{\partial \mathbf{h}_t^{(f)}} = \mathbf{U}^{(f)}$ summarizing the results we have

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

where diag indicates the diagonal matrix containing the elements $\mathbf{h}_{t+1}^{(f)}(1 - \mathbf{h}_{t+1}^{(f)})$. This is the Jacobian of the sigmoid function associated with the hidden unit i at time $t + 1$. In a similar way we obtain

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

3. Now derive $\nabla_{\mathbf{W}^{(f)}} L$ and $\nabla_{\mathbf{U}^{(b)}} L$ as functions of $\nabla_{\mathbf{h}_t^{(f)}} L$ and $\nabla_{\mathbf{h}_t^{(b)}} L$, respectively.

Hint: It might be useful to consider the contribution of the weight matrices when computing the recurrent hidden unit at a particular time t and how those contributions might be aggregated.

Solution

$$\begin{aligned} \nabla_{\mathbf{W}^{(f)}} L &= \frac{\partial L}{\partial \mathbf{W}^{(f)}} \\ &= \sum_t \frac{\partial L}{\partial \mathbf{h}_t^{(f)}} \left(\frac{\partial \mathbf{h}_t^{(f)}}{\partial \mathbf{W}^{(f)}} \right)^T \\ &= \text{diag}(\mathbf{h}_t^{(f)}(1 - \mathbf{h}_t^{(f)})) \nabla_{\mathbf{h}_t^{(f)}} L \mathbf{x}_t^T \end{aligned}$$

for achieving the last result we must remember that

$$\frac{\partial \mathbf{h}_t^{(f)}}{\partial \mathbf{W}^{(f)}} = \text{diag}(\mathbf{h}_t^{(f)}(1 - \mathbf{h}_t^{(f)}))\mathbf{x}^T$$

Finally

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

for achieving the last result we must remember that

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