

Due Date: March 22nd 23:59, 2019

Instructions

- For all questions, show your work!
- Starred questions are **hard** questions, not **bonus** questions.
- Please use a document preparation system such as LaTeX, unless noted otherwise.
- Unless noted that questions are related, assume that notation and definitions for each question are self-contained and independent
- Submit your answers electronically via Gradescope.
- **TAs for this assignment are David Krueger, Tegan Maharaj, and Chin-Wei Huang.**

Question 1 (6-10). The goal of this question is for you to understand the reasoning behind different parameter initializations for deep networks, particularly to think about the ways that the initialization affects the activations (and therefore the gradients) of the network. Consider the following equation for the t -th layer of a deep network:

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

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

For both of the assumptions (1 and 2) about the distribution of the inputs to layer t listed below, and for both (a) Gaussian, and (b) Uniform sampling, design an initialization scheme that would achieve preactivations with zero-mean and unit variance at layer t , i.e.: $\mathbb{E}[\mathbf{a}_i^{(t)}] = 0$ and $\text{Var}(\mathbf{a}_i^{(t)}) = 1$, for $1 \leq i \leq d^{(t)}$.

(Hint: if $X \perp Y$, $\text{Var}(XY) = \text{Var}(X)\text{Var}(Y) + \text{Var}(X)\mathbb{E}[Y]^2 + \text{Var}(Y)\mathbb{E}[X]^2$)

1. Assume $\mathbb{E}[\mathbf{h}_i^{(t-1)}] = 0$ and $\text{Var}(\mathbf{h}_i^{(t-1)}) = 1$ for $1 \leq i \leq d^{(t-1)}$. Assume entries of $\mathbf{h}^{(t-1)}$ are uncorrelated (the answer should not depend on g).
 - (a) Gaussian: give values for c , μ , and σ^2 as a function of $d^{(t-1)}$.
 - (b) Uniform: give values for c , α , and β as a function of $d^{(t-1)}$.
2. Assume that the preactivations of the previous layer satisfy $\mathbb{E}[\mathbf{a}_i^{(t-1)}] = 0$, $\text{Var}(\mathbf{a}_i^{(t-1)}) = 1$ and $\mathbf{a}_i^{(t-1)}$ has a symmetric distribution for $1 \leq i \leq d^{(t-1)}$. Assume entries of $\mathbf{a}^{(t-1)}$ are uncorrelated. Consider the case of ReLU activation: $g(x) = \max\{0, x\}$.
 - (a) Gaussian: give values for c , μ , and σ^2 as a function of $d^{(t-1)}$.
 - (b) Uniform: give values for c , α , and β as a function of $d^{(t-1)}$.
 - (c) What popular initialization scheme has this form?
 - (d) Why do you think this initialization would work well in practice? Answer in 1-2 sentences.

Answer 1.

Question 2 (4-6-4-4-3). The point of this question is to understand and compare the effects of different regularizers (specifically dropout and weight decay) on the weights of a network. Consider a linear regression problem with input data $\mathbf{X} \in \mathbb{R}^{n \times d}$, weights $\mathbf{w} \in \mathbb{R}^{d \times 1}$ and targets $\mathbf{y} \in \mathbb{R}^{n \times 1}$. Suppose that dropout is applied to the input (with probability $1 - p$ of dropping the unit i.e. setting it to 0). Let $\mathbf{R} \in \mathbb{R}^{n \times d}$ be the dropout mask such that $\mathbf{R}_{ij} \sim \text{Bern}(p)$ is sampled i.i.d. from the Bernoulli distribution.

1. For squared error loss, express the loss function $L(\mathbf{w})$ in matrix form (in terms of \mathbf{X} , \mathbf{y} , \mathbf{w} , and \mathbf{R}).
2. Let Γ be a diagonal matrix with $\Gamma_{ii} = (\mathbf{X}^\top \mathbf{X})_{ii}^{1/2}$. Show that the *expectation (over \mathbf{R})* of the loss function can be rewritten as $\mathbb{E}[L(\mathbf{w})] = \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1 - p)\|\Gamma\mathbf{w}\|^2$.
3. Show that the solution $\mathbf{w}^{\text{dropout}}$ that minimizes the expected loss from question 2.2 satisfies

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

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

4. Express the solution \mathbf{w}^{L^2} for a linear regression problem without dropout and with L^2 regularization, with regularization coefficient λ^{L^2} in closed form.
5. Compare the results of 2.3 and 2.4: identify specific differences in the equations you arrive at, and discuss qualitatively what the equations tell you about the similarities and differences in the effects of weight decay and dropout (1-3 sentences).

Answer 2.

1. The loss function is :

$$L(\mathbf{w}) = \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2$$

where $\hat{\mathbf{y}}$ is the output of the neural network, which is expressed as following using dropout mask \mathbf{R} on the input \mathbf{X} :

$$\hat{\mathbf{y}} = (\mathbf{X} \odot \mathbf{R}) \cdot \mathbf{w}$$

So, the loss is:

$$L(\mathbf{w}) = \|\mathbf{y} - (\mathbf{X} \odot \mathbf{R}) \cdot \mathbf{w}\|_2^2 \quad (1)$$

2. From equation 1 we get:

$$\begin{aligned} L(\mathbf{w}) &= (\mathbf{y} - (\mathbf{X} \odot \mathbf{R})\mathbf{w})^\top (\mathbf{y} - (\mathbf{X} \odot \mathbf{R})\mathbf{w}) \\ &= \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top (\mathbf{X} \odot \mathbf{R})\mathbf{w} - \mathbf{w}^\top (\mathbf{X} \odot \mathbf{R})^\top \mathbf{y} + \mathbf{w}^\top (\mathbf{X} \odot \mathbf{R})^\top (\mathbf{X} \odot \mathbf{R})\mathbf{w} \end{aligned} \quad (2)$$

The expected value of loss using equation 2 is:

$$\mathbb{E}_{\mathbf{R}}[L(\mathbf{w})] = \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbb{E}_{\mathbf{R}}[\mathbf{X} \odot \mathbf{R}]\mathbf{w} - \mathbf{w}^\top \mathbb{E}_{\mathbf{R}}[(\mathbf{X} \odot \mathbf{R})^\top]\mathbf{y} + \mathbf{w}^\top \mathbb{E}_{\mathbf{R}}[(\mathbf{X} \odot \mathbf{R})^\top (\mathbf{X} \odot \mathbf{R})]\mathbf{w} \quad (3)$$

Given:

$$\mathbb{E}_{\mathbf{R}}[\mathbf{X} \odot \mathbf{R}]_{ij} = \mathbb{E}_{\mathbf{R}}[X_{ij}R_{ij}] = pX_{ij} \quad (4)$$

And

$$\mathbb{E}_{\mathbf{R}}[(\mathbf{X} \odot \mathbf{R})^\top (\mathbf{X} \odot \mathbf{R})]_{ij} = \mathbb{E}_{\mathbf{R}}\left[\sum_{k=1}^n X_{ik}X_{kj}R_{ik}R_{kj}\right]$$

- If $i = j$:

$$\mathbb{E}_{\mathbf{R}}[(\mathbf{X} \odot \mathbf{R})^\top (\mathbf{X} \odot \mathbf{R})]_{ij} = p^2 \sum_{k=1}^n X_{ik} X_{kj} \mathbb{E}_{\mathbf{R}}(R_{ij}^2) \quad (5)$$

We can prove easily that $\mathbb{E}[X^2] = p$, for $X \sim \text{Bern}(p)$:

We know that $E(X) = p$, $\text{Var}(X) = p(1-p)$ and $\text{Var}(X) = E(X^2) - E(X)^2$, so:

$$E(X^2) = \text{Var}(X) + E(X)^2 = p - p^2 + p^2 = p$$

back to equation 5, we get:

$$\begin{aligned} \mathbb{E}_{\mathbf{R}}[(\mathbf{X} \odot \mathbf{R})^\top (\mathbf{X} \odot \mathbf{R})]_{ii} &= p \sum_{k=1}^n X_{ik}^2 \\ &= p(\mathbf{X}^\top \mathbf{X})_{ii} \\ &= p\Gamma_{ii}^2 \end{aligned} \quad (6)$$

- if $i \neq j$:

$$\begin{aligned} \mathbb{E}_{\mathbf{R}}[(\mathbf{X} \odot \mathbf{R})^\top (\mathbf{X} \odot \mathbf{R})]_{ij} &= p \sum_{k=1}^n X_{ik} X_{kj} \mathbb{E}_{\mathbf{R}}(R_{ik}) \mathbb{E}_{\mathbf{R}}(R_{kj}) \\ &= p^2 \sum_{k=1}^n X_{ik} X_{kj} \\ &= p^2 (\mathbf{X}^\top \mathbf{X})_{ij} \end{aligned} \quad (7)$$

Using equation 4, equation 3 become:

$$\begin{aligned} \mathbb{E}_{\mathbf{R}}[L(\mathbf{w})] &= \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top p \mathbf{X} \mathbf{w} - \mathbf{w}^\top p \mathbf{X}^\top \mathbf{y} + \mathbf{w}^\top \mathbb{E}_{\mathbf{R}}[(\mathbf{X} \odot \mathbf{R})^\top (\mathbf{X} \odot \mathbf{R})] \mathbf{w} \\ &= (\mathbf{y} - p \mathbf{X} \mathbf{w})^\top (\mathbf{y} - p \mathbf{X} \mathbf{w}) - p^2 \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} + \mathbf{w}^\top \mathbb{E}_{\mathbf{R}}[(\mathbf{X} \odot \mathbf{R})^\top (\mathbf{X} \odot \mathbf{R})] \mathbf{w} \\ &= \|\mathbf{y} - p \mathbf{X} \mathbf{w}\|^2 - \mathbf{w}^\top [p^2 \mathbf{X}^\top \mathbf{X} - \mathbb{E}_{\mathbf{R}}[(\mathbf{X} \odot \mathbf{R})^\top (\mathbf{X} \odot \mathbf{R})]] \mathbf{w} \end{aligned} \quad (8)$$

From equation 7, we have shown that $p^2 \mathbf{X}^\top \mathbf{X}$ and $\mathbb{E}_{\mathbf{R}}[(\mathbf{X} \odot \mathbf{R})^\top (\mathbf{X} \odot \mathbf{R})]$ have the same non-diagonal element, so :

$$(p^2 \mathbf{X}^\top \mathbf{X} - \mathbb{E}_{\mathbf{R}}[(\mathbf{X} \odot \mathbf{R})^\top (\mathbf{X} \odot \mathbf{R})])_{i \neq j} = 0$$

And from equation 6 we have get the following formula for the diagonal elements:

$$(p^2 \mathbf{X}^\top \mathbf{X} - \mathbb{E}_{\mathbf{R}}[(\mathbf{X} \odot \mathbf{R})^\top (\mathbf{X} \odot \mathbf{R})])_{i=j} = p^2 \Gamma_{ii}^2 - p \Gamma_{ii}^2 = p(p-1) \Gamma_{ii}^2$$

Now let's put all together, and replace the last equations back to equation 8:

$$\begin{aligned} \mathbb{E}_{\mathbf{R}}[L(\mathbf{w})] &= \|\mathbf{y} - p \mathbf{X} \mathbf{w}\|^2 - p(p-1) \mathbf{w}^\top \Gamma^2 \mathbf{w} \\ &= \|\mathbf{y} - p \mathbf{X} \mathbf{w}\|^2 + p(1-p) (\mathbf{w} \Gamma)^\top (\mathbf{w} \Gamma) \\ &= \|\mathbf{y} - p \mathbf{X} \mathbf{w}\|^2 + p(1-p) \|\Gamma \mathbf{w}\|^2 \end{aligned} \quad (9)$$

3. Let's calculate $\nabla_{\mathbf{w}} \mathbb{E}_{\mathbf{R}}(L(\mathbf{w}))$ from equation 9 and solve the equation $\nabla_{\mathbf{w}} \mathbb{E}_{\mathbf{R}}(L(\mathbf{w})) = 0$

$$\nabla_{\mathbf{w}} \mathbb{E}_{\mathbf{R}}(L(\mathbf{w})) = -2p\mathbf{X}^{\top}(\mathbf{y} - p\mathbf{X}\mathbf{w}) + 2p(1-p)\Gamma\mathbf{w} \quad (10)$$

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

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

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

Where $\lambda^{\text{dropout}} = \frac{1-p}{p}$

Thus, when $p \rightarrow 0$, meaning dropping all the input units, we get no solution, which means an infinit regularization that makes the model underfit the data.

If $p \rightarrow 1$, meaning there is no dropout, we get the usual analytical solution of the squared error loss with no regularization an intermediate value of p give the model a regularization equivalent to L_2 regularization.

4. The loss function without dropout and with L_2 regularization is:

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

which has the analytical solution (derived similarly as 2.3):

$$\mathbf{w}^{L_2} = (\mathbf{X}^{\top}\mathbf{X} + \lambda^{L_2}\mathbf{I})^{-1}\mathbf{X}^{\top}\mathbf{y}$$

5. From 2.3 and 2.4 we notice that dropout and L_2 regularization have similar analytical solution and thus dropout behave as a regularization method similarly as L_2 regularization. The main difference is on the amount of regularization that each method puts on the model. In fact, L_2 has a simple real coefficient λ^{L_2} that influence the effect of the regularization, instead of a more complexe effect of the dropout which depends on the data that is the value of Γ and the binary probability p of the dropout mask.

Question 3 (5-5-5). In this question you will demonstrate that an estimate of the first moment of the gradient using an (exponential) running average is equivalent to using momentum, and is biased by a scaling factor. The goal of this question is for you to consider the relationship between different optimization schemes, and to practice noting and quantifying the effect (particularly in terms of bias/variance) of *estimating* a quantity.

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

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

- SGD with momentum:

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

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

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

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

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

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

- Unroll the running average update rule, i.e. express \mathbf{v}_t as a linear combination of \mathbf{g}_i 's ($1 \leq i \leq t$).
- Assume \mathbf{g}_t has a stationary distribution independent of t . Show that the running average is biased, i.e. $\mathbb{E}[\mathbf{v}_t] \neq \mathbb{E}[\mathbf{g}_t]$. Propose a way to eliminate such a bias by rescaling \mathbf{v}_t .

Answer 3.

Question 4 (5-5-5). This question is about weight normalization. We consider the following parameterization of a weight vector \mathbf{w} :

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

where γ is scalar parameter controlling the magnitude and \mathbf{u} is a vector controlling the direction of \mathbf{w} .

- 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 $\hat{y} = \mathbf{w}^\top \mathbf{x} + \beta$.
- Show that the gradient of a loss function $L(\mathbf{u}, \gamma, \beta)$ with respect to \mathbf{u} can be written in the form $\nabla_{\mathbf{u}} L = s \mathbf{W}^\perp \nabla_{\mathbf{w}} L$ for some s , where $\mathbf{W}^\perp = \left(\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|^2} \right)$. Note that $\mathbf{W}^\perp \mathbf{u} = \mathbf{0}$.
- Figure 1 shows the norm of \mathbf{u} as a function of number of updates made to a two-layer MLP using gradient descent. Different curves correspond to models trained with different log-learning rate. Explain why (1) the norm is increasing, and (2) why larger learning rate corresponds to faster growth. (Hint: Use the Pythagorean theorem and the fact that $\mathbf{W}^\perp \mathbf{u} = \mathbf{0}$ from question 4.2).

Answer 4.

Question 5 (5-5-5). This question is about activation functions and vanishing/exploding gradients in recurrent neural networks (RNNs). Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be an activation function. When the argument is a vector, we apply σ element-wise. Consider the following recurrent unit:

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

- 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).
- 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. $|\sigma'(x)| \leq \gamma$ 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 upper-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

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

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. 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.

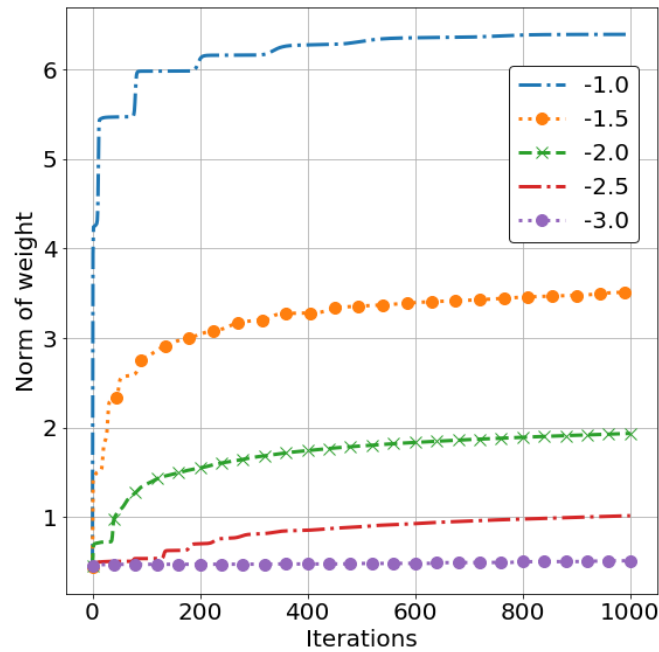


FIGURE 1 – Norm of parameters with different learning rate.

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

Answer 5.

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

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

where the superscripts f and b correspond to the forward and backward RNNs respectively.

1. Draw the computational graph for this RNN, unrolled for 3 time steps (from $t = 1$ to $t = 3$). Include and label the initial hidden states for both the forward and backward RNNs, $\mathbf{h}_0^{(f)}$ and $\mathbf{h}_4^{(b)}$ respectively. You may draw this by hand; you may also use a computer rendering package such as TikZ, but you are not required to do so. Label each node and edge with the corresponding hidden unit or weight.
- *2. Let \mathbf{z}_t be the true target of the prediction \mathbf{y}_t and consider the sum of squared loss $L = \sum_t L_t$ where $L_t = \|\mathbf{z}_t - \mathbf{y}_t\|_2^2$. Express the gradients $\nabla_{\mathbf{h}_t^{(f)}} L$ and $\nabla_{\mathbf{h}_t^{(b)}} L$ recursively (in terms of $\nabla_{\mathbf{h}_{t+1}^{(f)}} L$ and $\nabla_{\mathbf{h}_{t-1}^{(b)}} L$ respectively). Then derive $\nabla_{\mathbf{W}^{(f)}} L$ and $\nabla_{\mathbf{U}^{(b)}} L$.

Answer 6.