

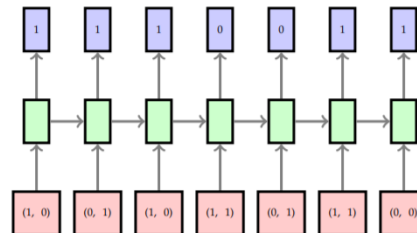
1 Binary Addition (2.5pts): In this problem, you will design a recurrent neural network to implement binary addition. The inputs are provided as binary sequences, starting with the least significant bit. The sequences are padded with at least one zero at the end. Here is an example:

$$0101101 + 0111010 = 1100111$$

where the inputs and outputs would be represented as:

- **Input 1:** 1, 0, 1, 1, 0, 1, 0
- **Input 2:** 0, 1, 0, 1, 1, 1, 0
- **Correct Output:** 1, 1, 1, 0, 0, 1, 1

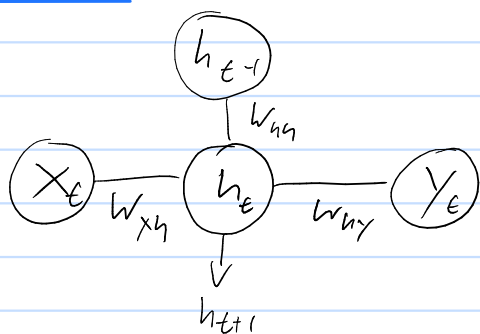
At each time instance, the RNN would have two inputs and one output. Therefore, the pattern of inputs and outputs for the above example would be:



Design, by hand, the weights and biases for a RNN that can perform binary addition. The RNN should consist of two inputs, three hidden nodes and one output. You can assume that all the nodes use the hard threshold activation function. In particular, find the weights W_{xh} , W_{hh} and W_{hy} , the bias vector b_h and scalar b_y .

Solution on next page!

Problem 1



$$x_t \in \left\{ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right\}$$

$$y_t = \Theta(w_{hy} h_t + b_y)$$

$$h_t = \Theta(w_{hh} h_{t-1} + w_{xh} x_t + b_h)$$

$$\text{where } \Theta(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases}$$

Example:

t	6	5	4	3	2	1	0
x_1	0	1	0	1	1	0	1
x_2	+	0	1	1	0	1	0
h_1	0	0	1	0	1	1	1
h_2	0	1	0	1	0	0	0
h_3	1	1	1	0	0	0	0
y	1	1	0	0	1	1	1

Idea: h_1, h_2 will be the bits of $x_1 + x_2$.
 h_3 will be the carry-over bit:

$$h_{t,3} = \begin{cases} 1, & h_{t-1,2} = 1 \\ 1, & h_{t-1,1} = 1 \text{ and } h_{t-1,3} = 1 \\ 0, & \text{otherwise} \end{cases}$$

We want $y = h_1 \text{ XOR } h_3 \Rightarrow w_{hy} = (-1 \ 0 \ -1)$
 $1 < b_y < 2$, let's say $b_y = 1.5$

We want $h_1 = x_1 \text{ XOR } x_2$
 and $h_2 = x_1 \text{ AND } x_2 \Rightarrow w_{xh} = \begin{pmatrix} -1 & -1 \\ 1 & 1 \\ 0 & 0 \end{pmatrix}, b_{y,1} = 1.5$
 $b_{y,2} = -1.5$
 $-2 < b_{y,2} < -1$

$$h_{t,3} = \begin{cases} 1, & h_{t-1,2} = 1 \\ 1, & h_{t-1,1} = 1 \text{ and } h_{t-1,3} = 1 \\ 0, & \text{otherwise} \end{cases}$$

$\Rightarrow w_{hh} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{pmatrix}, b_{y,3} = -1.5$
 $-2 < b_{y,2} < -1$
 or any other number $> |b_{y,3}|$

2 LSTM Gradient (3.5pts): In this problem you will derive the backpropagation-through-time equations for a univariate version of the Long-Term Short-Term Memory (LSTM) architecture we saw in class. Consider the case when the bias terms are assumed to be zero. So the LSTM computations are:

$$\begin{aligned} i^{(t)} &= \sigma(w_{xi}x^{(t)} + w_{hi}h^{(t-1)}) && \text{input} \\ f^{(t)} &= \sigma(w_{xf}x^{(t)} + w_{hf}h^{(t-1)}) && \text{forget} \\ o^{(t)} &= \sigma(w_{xo}x^{(t)} + w_{ho}h^{(t-1)}) && \text{output} \\ g^{(t)} &= \tanh(w_{xg}x^{(t)} + w_{hg}h^{(t-1)}) && \text{input activation} \\ c^{(t)} &= f^{(t)}c^{(t-1)} + i^{(t)}g^{(t)} && \text{cell state} \\ h^{(t)} &= o^{(t)}\tanh(c^{(t)}) && \text{output} \end{aligned}$$

1. (3pts) Derive the backpropagation-through-time equations for the activations and the gates:

$$\begin{aligned} \overline{h^{(t)}} &= \dots \\ \overline{c^{(t)}} &= \dots \\ \overline{g^{(t)}} &= \dots \\ \overline{o^{(t)}} &= \dots \\ \overline{f^{(t)}} &= \dots \\ \overline{i^{(t)}} &= \dots \end{aligned}$$

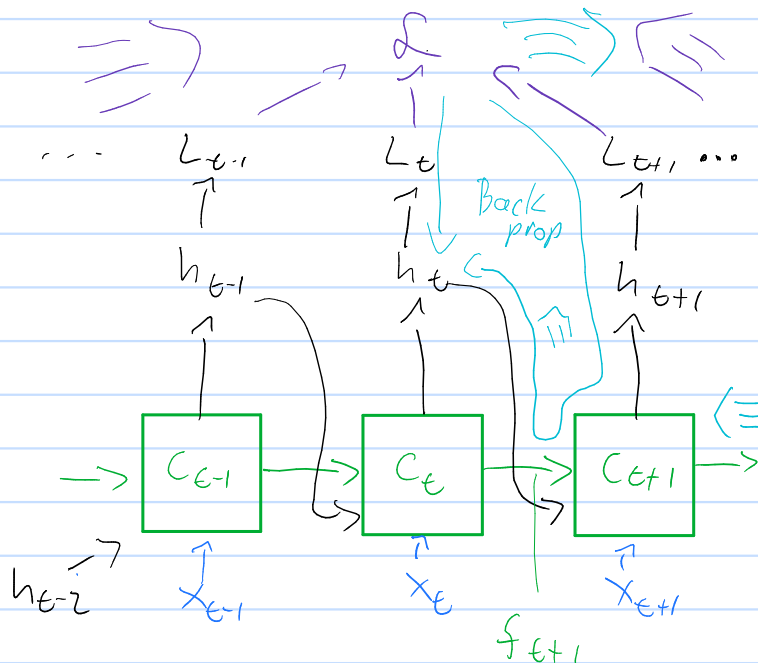
We simply use the regular backpropagation algorithm on the unfolded graph.

I prefer to write $i_e = i^{(e)}$, $o_e = o^{(e)}$, etc.

Before we start, recall the derivatives:

$$\frac{d}{dx} \tanh(x) = 1 - \tanh^2(x), \quad \frac{d}{dx} \sigma(x) = \sigma(x) \cdot (1 - \sigma(x))$$

We have $L = \sum_{t=1}^T L_t$ with some differentiable loss $L_t(h_t, y_t)$ where the y_t are labels.



$$\bar{h}_t = \frac{\partial \mathcal{L}}{\partial h_t} = \frac{\partial}{\partial h_t} \left(\sum_{k=1}^T L_k \right) = \sum_{k=t}^T \frac{\partial L_k}{\partial h_t} = \frac{\partial L_t}{\partial h_t} + \bar{c}_{t+1} \frac{\partial c_{t+1}}{\partial h_t}$$

only $\geq t$ contribute

$$= \frac{\partial L_t}{\partial h_t} + \bar{c}_{t+1} c_t \frac{\partial \sigma(w_{xf} x_{t+1} + w_{hf} h_t)}{\partial h_t}$$

$$= \frac{\partial L_t}{\partial h_t} + \bar{c}_{t+1} c_t w_{hf} f_{t+1} (1 - f_{t+1})$$

$$\bar{c}_t = \sum_{k=t}^T \frac{\partial L_k}{\partial c_t} = \bar{h}_t \frac{\partial h_t}{\partial c_t} + \sum_{k=t+1}^T \frac{\partial L_k}{\partial c_t} = \bar{h}_t o_t (1 - \tanh^2 c_t) + \bar{c}_{t+1} f_{t+1}$$

$$\bar{g}_t = \frac{\partial \mathcal{L}}{\partial c_t} \frac{\partial c_t}{\partial g_t} = \bar{c}_t \cdot \dot{c}_t$$

$$\bar{o}_t = \bar{h}_t \frac{\partial h_t}{\partial o_t} = \bar{h}_t \tanh(c_t)$$

$$\bar{f}_t = \bar{c}_t \frac{\partial c_t}{\partial f_t} = \bar{c}_t c_{t-1}$$

$$\dot{\bar{c}}_t = \bar{c}_t \frac{\partial c_t}{\partial \dot{c}_t} = \bar{c}_t g_t$$

2. (0.5pt) Derive the update equation for the weight w_{xi} :

$$\overline{w_{xi}} = \dots$$

$$\begin{aligned} \overline{w_{xi}} &= \sum_{t=1}^T \frac{\partial \mathcal{L}_t}{\partial w_{xi}} = \sum_{t=1}^T \sum_{k=1}^T \frac{\partial \mathcal{L}_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial w_{xi}} \\ &= \sum_{t=1}^T \overline{c_t} \frac{\partial \hat{c}_t}{\partial w_{xi}} = \sum_{t=1}^T \overline{c_t} g_t X_t \cdot \sigma'(\hat{c}_t) (1 - \sigma'(\hat{c}_t)) \\ &= \sum_{t=1}^T \underbrace{\overline{c_t} g_t X_t}_{\hat{c}_t} \hat{c}_t (1 - \hat{c}_t) \end{aligned}$$

gate activation
= $w_{xi} X_t + w_{hi} h_{t-1}$
↓

3. (Extra Credit: 0.5pt) Based on your answers above, explain why the gradient does not explode if the values of the forget gates are very close to 1 and the values of the input and output gates are very close to 0. **Hint:** Your answer should involve both $\overline{h^{(t)}}$ and $\overline{c^{(t)}}$.

$$f_t \approx 1, \quad i_t \ll 1, \quad o_t \ll 1$$

In this case $g_t X_t \hat{c}_t (1 - \hat{c}_t) \approx V_t$

where V_t is some small constant. It's roughly the same for each t : $V_t \approx V$

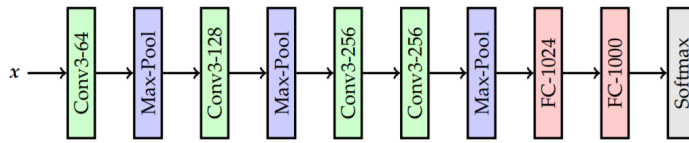
$$\begin{aligned} \text{So: } \overline{w_{xi}} &\approx V \sum_{t=1}^T \overline{c_t} \\ &= V \sum_{t=1}^T \left[\overline{h_t} \underbrace{o_t (1 - \tanh^2 c_t)}_{\text{small } \tilde{V}_t} + \overline{c_{t+1}} \underbrace{f_{t+1}}_{\approx 1} \right] \\ &\approx V \sum_{t=1}^T \overline{c_{t+1}} \approx VT \overline{c_T} \end{aligned}$$

Does not explode!

Because $\overline{c_t} \approx \overline{c_{t+1}}$, only the small contributions $\frac{\partial \mathcal{L}_t}{\partial h_t}$ in $\overline{h_t}$ accumulate.

3 Convolutional Neural Networks (2pts):

1. (1pt) Consider a CNN with 4 conv layers like in the diagram below. All 4 conv layers have kernel size of 3×3 . The number after the hyphen specifies the number of output channels or units of a layer (e.g. Conv3-64 layer has 64 output channels and FC-1024 has 1024 output units). All the Max Pool in the diagram has size of 2×2 . Assume zero padding of 1 for conv layers and stride 2 for Max Pool.



Size of the input image is 224×224 with 3 channels. Calculate the total number of parameters in the network including the bias units.

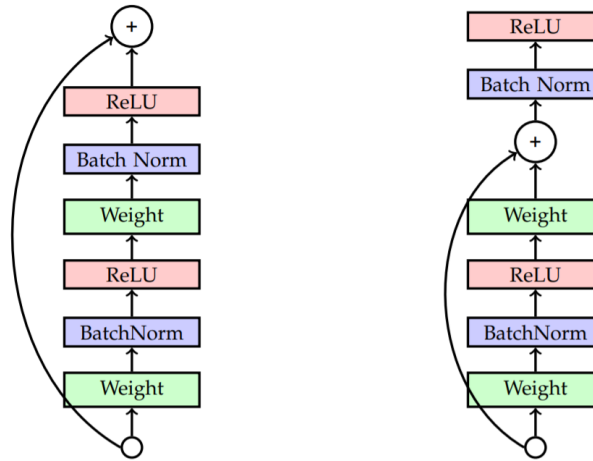
As usual: We want to figure out the number of parameters for each layer and then simply sum them up. I'll try to create a table, to keep track:

Layer	Output shape	#parameters
X	$(224, 224, 3)$	$\# \text{ filters} \rightarrow 0$
Conv3-64	$(224, 224, 64)$	$64 \cdot (3 \cdot 3 \cdot 3 + 1) = 1792$ \circ channels bias
Max Pool	$(112, 112, 64)$	\circ kernel
Conv3-128	$(112, 112, 128)$	$128 \cdot (64 \cdot 3 \cdot 3 + 1) = 73856$ \circ
Max Pool	$(56, 56, 128)$	\circ
Conv3-256	$(56, 56, 256)$	$256 \cdot (128 \cdot 3 \cdot 3 + 1) = 295168$
Conv3-256	$(56, 56, 256)$	295168
Max Pool	$(28, 28, 256) \Rightarrow 20704$	\circ
FC-1024	1024	$1024 \cdot (20704 + 1) = 21201920$
FC-1000	1000	$1000 \cdot (1024 + 1) = 1025000$
Softmax	1000	\circ

$$\text{So in total: } 1792 + 73856 + 2 \cdot 295168 + 21201920 + 1025000$$

$$= 22892904$$

2. (1pt) Consider the following two ResNet architectures for the placement of the batch norm and the residual connection. "Weight" layer can be either a matrix multiplication or convolution. Which architecture is easier to learn in terms of exploding / vanishing gradient? Provide a brief justification for your answer.



From the BN paper ([arXiv:1502.03167](https://arxiv.org/abs/1502.03167) [cs.LG]) :

Batch Normalization can be applied to any set of activations in the network. Here, we focus on transforms that consist of an affine transformation followed by an element-wise nonlinearity:

$$z = g(Wu + b)$$

where W and b are learned parameters of the model, and $g(\cdot)$ is the nonlinearity such as sigmoid or ReLU. This formulation covers both fully-connected and convolutional layers. We add the BN transform immediately before the nonlinearity, by normalizing $x = Wu + b$. We could have also normalized the layer inputs u , but since u is likely the output of another nonlinearity, the shape of its distribution is likely to change during training, and constraining its first and second moments would not eliminate the covariate shift. In contrast, $Wu + b$ is more likely to have a symmetric, non-sparse distribution, that is "more Gaussian" (Hyvärinen & Oja, 2000); normalizing it is likely to produce activations with a stable distribution.

Here we have: residual

$$z = w_+ g(BN(wu)) + x_+ \quad \text{left}$$

$$Vs \quad z = g(BN(w_+ wu + x_+)) \quad \text{right}$$

=> Therefore, the left architecture should be better suited to eliminate the internal covariate shift.

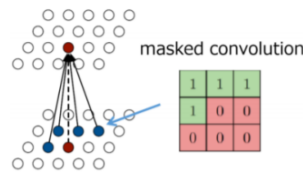
Indeed, this is what is widely used.

E.g. : [arXiv:1512.03385](https://arxiv.org/abs/1512.03385) [cs.CV] and <https://openreview.net/pdf?id=S1xU74med4>

However, this to me seems like something that requires specific case studies or complicated proof. It is not obvious which architecture will work better after all in practice.

4 Autoregressive Generative Models (2pt): In this question we will consider autoregressive models that model the distribution of images in an autoregressive manner. For simplicity, assume that each conditional distribution is a Gaussian with a fixed variance, so the model predicts the mean of the next pixel given all the previous pixels.

1. (1pt) Here we will consider PixelCNN, which models the distribution of images using a convolutional architecture. PixelCNN masks each convolutional filter to only see the pixels that appear before the current pixel in a raster scan order. See figure below for a visualization. Several of such layers are stacked sequentially.



- (a) Consider a d layer PixelCNN model, what is the total number of ~~connections~~ ^{weights}? Give your answers in terms of d, k, H, W . You only need to give $\mathcal{O}(\cdot)$, not an exact count.

Assuming that k is the size of the kernels, $H \times W$ the size of the image and d the number of sequentially stacked layers.

For a normal (unmasked) CNN we would have

$$\mathcal{O}(d k^2) \text{ weights.}$$

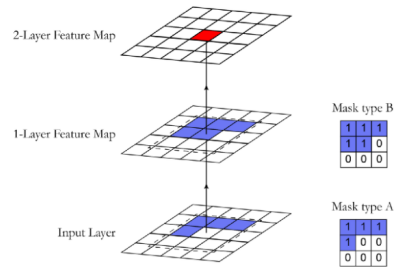
Here, the only difference is the mask, so instead of k^2 it's $\frac{k^2}{2}(-1)$, but $\mathcal{O}(\frac{k^2}{2} - 1) = \mathcal{O}(k^2)$.
only first layer

So it's still

$$\mathcal{O}(d k^2) \text{ weights.}$$

(b) Suppose that in each step we compute as many matrix-vector multiplications in parallel as we can, what is the minimum number of the sequential operations to compute the output of a Pixel CNN in terms of d, k, H, W . You only need to give $O(\cdot)$, not an exact count.

After training, when we actually want to generate, each pixel has to be predicted sequentially. Considering that Pixel CNN uses type B filters for deeper layers, which depend on the output of the previous layer:



Each filter position in each layer has to be evaluated
So $O(dHW)$ total operations.

During training, all pixels can be evaluated in parallel. However, all layers except the first (type A) have to be computed sequentially for each pixel.

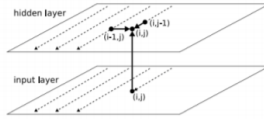
So during training: $O(d-1) = O(d)$ sequential operations

During evaluation, we can use the same parallel approach as I describe in 26 for the RNN.

So $O(\sqrt{dHW})$ sequential operations per layer.

For multiple layers, we can do the same thing in 3D, then: $O(\sqrt{dHW})$ sequential operations.

2. (1pt) Here we will consider Multidimensional RNN (MDRNN). This is like the RNNs we discussed in the lecture, except that instead of a 1-D sequence, we have a 2-D grid structure. Analogous to how ordinary RNNs have an input vector and a hidden vector for every time step, MDRNNs have an input vector and hidden vector for every grid square. Each hidden unit receives bottom-up connections from the corresponding input square, as well as recurrent connections from its north and west neighbors as shown in the figure below.



The activations are computed as:

$$h^{(i,j)} = \phi(W_{in}^T x^{(i,j)} + W_W^T x^{(i-1,j)} + W_N^T x^{(i,j-1)}) \quad (1)$$

Denote the number of input channels and the number of recurrent neurons at each layer to be k . The input image size is $H \times W$. For simplicity, we assume there are no bias parameters.

- (a) Consider a d layer MDRNN model, what is the total number of ^{connections} connections? Give your answers in terms of d, k, H, W . You only need to give $O(\cdot)$, not an exact count.

During training: x_{ij} are the input image pixels

Generation: $x_{ij} = y_{t-1}$, the predicted pixel from previous step.

$$W_k \in \mathbb{R}^{k \times k}, \quad x \in \mathbb{R}^k, \quad h \in \mathbb{R}^k$$

All 3 weight matrices are shared

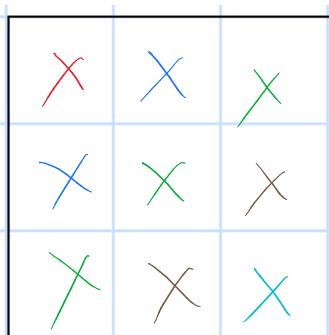
So we have $O(3k^2) = O(k^2)$ weights!

- (b) Suppose that in each step we compute as many matrix-vector multiplications in parallel as we can, what is the minimum number of the sequential operations to compute the output of a ~~Pixel~~ ~~CNN~~ in terms of d, k, H, W . You only need to give $O(\cdot)$, not an exact count.

mdrnn

There are $H \cdot W$ steps in total. However, we can compute the next step once the north and west neighbors have been computed. So we can traverse it diagonally and only $\approx 2\sqrt{HW} - 1$ sequential steps are needed. Exact for $H=W$

E.g. for a 3×3 grid:



1.

2.

3.

4.

5.

$$5 = 2\sqrt{2 \cdot 2} - 1$$

So: $O(\sqrt{HW})$ sequential steps!

3. (Extra Credit: 1pt) What is the benefit of using PixelCNN over MDRNN. Discuss the pros and cons of the two models in terms of their computational, memory complexity, parallelization potential and the size of their context windows.

PixelCNN

- + Fast to train in parallel
- Smaller context window due to blind spots (see Gated PixelCNN for improvement)
- Memory intensive

MDRNN

- Slow to train because of sequential dependency
- + Uses less memory
- + Maximum possible receptive field
- + \Rightarrow best performance