

## TTIC 31230 Fundamentals of Deep Learning, 2020

### Problems For Trainability.

In these problems, as in the lecture notes, capital letter indeces are used to indicate subtensors (slices) so that, for example,  $M[I, J]$  denotes a matrix while  $M[i, j]$  denotes one element of the matrix,  $M[i, J]$  denotes the  $i$ th row, and  $M[I, j]$  denotes the  $j$ th column.

Throughout these problems we assume a word embedding matrix  $e[W, I]$  where  $e[w, I]$  is the word vector for word  $w$ . We then have that  $e[w, I]^\top h[t, I]$  is the inner product of the word vector  $w[w, I]$  and the hidden state vector  $h[t, I]$ .

We will adopt the convention, similar to true Einstein notation, that repeated capital indeces in a product of tensors are implicitly summed. We can then write the inner product  $e[w, I]^\top h[t, I]$  simply as  $e[w, I]h[t, I]$  without the need for the (meaningless) transpose operation.

**Problem 1.** Consider a regression problem where we want to predict a scalar value  $y$  from a vector  $x$ . Consider the  $L$ -layer perceptron for this problem defined by the following equations which compute hidden layer vectors  $h_1[I], \dots, h_L[I]$  and predictions  $\hat{y}_1, \dots, \hat{y}_L$  where the prediction  $\hat{y}_\ell$  is done with a linear regression on the hidden vector  $h_\ell[I]$ .

$$\begin{aligned} h_0[i] &= x[i] \\ &\vdots \\ h_{\ell+1}[i] &= \sigma(W_{\ell+1}^{h,h}[i, I]h_\ell[I] - B_{\ell+1}^{h,h}[i]) \\ \hat{y}_{\ell+1} &= W_{\ell+1}^{h,y}[I]h_{\ell+1}[I] - B_{\ell+1}^{h,y} \\ &\vdots \\ \text{Loss} &= \sum_{\ell=1}^L (y - \hat{y}_\ell)^2 \end{aligned}$$

Each term  $(y - \hat{y}_\ell)^2$  is called a “loss head” and defines a loss on each prediction  $\hat{y}_\ell$ . Note, however, that there is only one scalar loss minimized by SGD which is the sum of the losses of each loss head.

**(a) 7 pts** Explain why these multiple loss terms might improve the ability of SGD to find a useful  $L$ -layer MLP regression  $\hat{y}_L$  when  $L$  is large.

**Solution:** SGD on deep networks with the loss term only occurring at the final layer is not generally effective because the lower layers never get meaningful gradients. Placing loss functions near the lower layers will cause the lower hidden layers to have meaningful gradients and produce informative features.

(b) **7 pts** As a function of  $L$  (ignoring the dimension size  $I$ ) what is the order of run time for the backpropagation procedure. Explain your answer.

**Solution:** It is  $O(L)$  — linear in  $L$ . Backpropagation loops over the assignments of the program and takes time proportional to the size of the program. Backpropagation over the final sum of losses produces a gradient for each prediction  $\hat{y}_\ell$  which can be used as we back-propagate over the earlier assignments.

(c) **11 pts** Rewrite the above MLP equations to use residual connections rather than multiple heads. There are multiple correct solutions differing in minor details. Pick one that seems good to you.

**Solution:**

$$\begin{aligned}
 h_0[i] &= x[i] \\
 &\vdots \\
 \tilde{h}_{\ell+1}[i] &= \sigma(W_{\ell+1}^{h,h}[i, I]h_\ell[I] - B_{\ell+1}^{h,h}[i]) \\
 h_{\ell+1}[i] &= \tilde{h}_{\ell+1}[i] + h_\ell[i] \\
 &\vdots \\
 \hat{y} &= W^{h,y}[I]h_L[I] - B^{h,y} \\
 \text{Loss} &= (y - \hat{y})^2
 \end{aligned}$$

## Problem 2. Gated CNNs

A UGRNN is defined by the following equations.

$$\tilde{R}_t[b, j] = \left( \sum_i W^{h,R}[j, i]h_{t-1}[b, i] \right) + \left( \sum_k W^{x,R}[j, k]x_t[b, k] \right) - B^R[j]$$

$$R_t[b, j] = \tanh(\tilde{R}_t[b, j])$$

$$\tilde{G}_t[b, j] = \left( \sum_i W^{h,G}[j, i]h_{t-1}[b, i] \right) + \left( \sum_k W^{x,G}[j, k]x_t[b, k] \right) - B^G[j]$$

$$G_t[b, j] = \sigma(\tilde{G}_t[b, j])$$

$$h_t[b, j] = G_t[b, j]h_{t-1}[b, j] + (1 - G_t[b, j])R_t[b, j]$$

Modify these to form a data-dependent data-flow CNN for vision — an Update-Gate CNN (UGCNN). More specifically, give equations analogous to those for

UGRNN for computing a CNN “box”  $L_{\ell+1}[b, x, y, j]$  from  $L_\ell[b, x, y, i]$  (stride 1) using a computed “gate box”  $G_{\ell+1}[b, x, y, j]$  and an “update box”  $R_{\ell+1}[b, x, y, j]$ .

$$\Phi = (W_{\ell+1}^{L,R}[\Delta x, \Delta y, j, j'], B_{\ell+1}^R[j], W_{\ell+1}^{L,G}[\Delta x, \Delta y, j, j'], B_{\ell+1}^G[j])$$

**Solution:**

$$\begin{aligned} R_{\ell+1}[b, x, y, j] &= \tanh \left( \left( \sum_{\Delta x, \Delta y, j'} W_{\ell+1}^{L,R}[\Delta x, \Delta y, j', j] L_\ell[b, x + \Delta x, y + \Delta y, j'] \right) - B_{\ell+1}^R[j] \right) \\ G_{\ell+1}[b, x, y, j] &= \sigma \left( \left( \sum_{\Delta x, \Delta y, i} W_{\ell+1}^{L,G}[\Delta x, \Delta y, i, j] L_\ell[b, x + \Delta x, y + \Delta y, i] \right) - B_{\ell+1}^G[j] \right) \\ L_{\ell+1}[b, x, y, j] &= G_{\ell+1}[b, x, y, j] L_\ell[b, x, y, j] + (1 - G_{\ell+1}[b, x, y, j]) R_{\ell+1}[b, x, y, j] \end{aligned}$$

### Problem 3. RNNs

Below are the equations defining the update cell of the UGRNN which takes as data inputs  $h[B, t-1, I]$  and  $x[B, t, K]$  and produces  $h[B, t, J]$ .

$$R[b, t, j] = \tanh (W^{h,R}[j, I]h[b, t-1, I] + W^{x,R}[j, K]x[b, t, K] - B^R[j])$$

$$G^h[b, t, j] = \sigma (W^{h,G}[j, I]h[b, t-1, I] + W^{x,G}[j, K]x[b, t, K] - B^G[j])$$

$$h[b, t, j] = G^h[b, t, j]h[b, t-1, j] + (1 - G^h[b, t, j])R[b, t, j]$$

Here I have written the gate as  $G^h$  to emphasize that it is a gate used to define a convex combination of the  $h_{t-1}$  and  $x_t$  inputs to  $h_t$ . Modify these equations to use a second gate  $G^R$  which gates inputs to  $R$  so that the input to the activation function (threshold function) producing  $R[b, t, j]$  is a convex combination of

$$W^{h,R}[j, I]h[b, t-1, I] - B^{h,R}[j]$$

and

$$W^{x,R}[j, K]x[b, t, K] - B^{x,R}[j]$$

where the weighting in the convex combination is given by your new gate  $G^R[b, t, j]$ . This gated RNN is similar to, but different from, a GRU which also has two gates.

**Solution:**

$$G^R[b, t, j] = \sigma (W^{h, GR}[j, I]h[b, t-1, I] + W^{x, GR}[j, K]x[b, t, K] - B^{GR}R[j])$$

$$R[b, t, j] = \tanh \left( \begin{aligned} &G^R[b, t, j] (W^{h, R}[j, I]h[b, t-1, I] - B^{h, R}[j]) \\ &+ (1 - G^R[b, t, j]) (W^{x, R}[j, K]x[b, t, K] - B^{x, R}[j]) \end{aligned} \right)$$

$$G^h[b, t, j] = \sigma (W^{h, Gh}[j, I]h[b, t-1, I] + W^{x, Gh}[j, K]x[b, t, K] - B^{Gh}[j])$$

$$h[b, t, j] = G^h[b, t, j]h[b, t-1, j] + (1 - G^h[b, t, j])R[b, t, j]$$

**Problem 4. RNN parallel run time.** Consider an autoregressive RNN neural language model with  $P_\Phi(w_{t+1}|w_1, \dots, w_t)$  defined by

$$P_\Phi(w_t|w_1, \dots, w_{t-1}) = \text{softmax}_{w_{t+1}} e[w_t, I]h[t-1, I]$$

Here  $e[w, I]$  is the word vector for word  $w$ ,  $h[t, I]$  is the hidden state vector at time  $t$  of a left-to-right RNN, and as described above  $e[w, I]h[t, I]$  is the inner product of these two vectors where we have assumed that they have the same dimension. For the first word  $w_1$  we have an externally provided initial hidden state  $h[0, I]$  and  $w_1, \dots, w_0$  denotes the empty string. We train the model on the full loss

$$\begin{aligned} \Phi^* &= \underset{\Phi}{\operatorname{argmin}} E_{w_1, \dots, w_T \sim \text{Train}} - \ln P_\Phi(w_1, \dots, w_T) \\ &= \underset{\Phi}{\operatorname{argmin}} E_{w_1, \dots, w_T \sim \text{Train}} \sum_{t=1}^T - \ln P_\Phi(w_t|w_1, \dots, w_{t-1}) \end{aligned}$$

What is the order of run time as a function of sentence length  $T$  for the back-propagation for this model run on a sentence  $w_1, \dots, w_T$ ? Explain your answer.

**Solution:** The backpropagation takes  $O(T)$  time (not  $O(T^2)$ ). The model consists of  $O(T)$  objects each of which performs a single forward operation and a single backward operation. As the backpropagation proceeds more of the loss terms in the sum over  $t$  get incorporated.

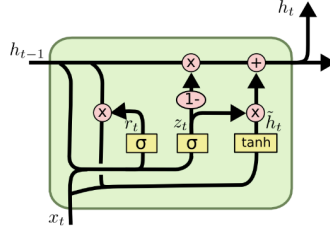
**Problem 5. Translating diagrams into equations.** A UGRNN cell for computing  $h[b, t, J]$  from  $h[b, t-1, J]$  and  $x[b, t, J]$  can be written as

$$G[b, t, j] = \sigma (W^{h,G}[j, I]h[b, t-1, I] + W^{x,G}[j, K]x[b, t, K] - B^G[j])$$

$$R[b, t, j] = \tanh (W^{h,R}[j, I]h[b, t-1, I] + W^{x,R}[j, K]x[b, t, K] - B^R[j])$$

$$h[b, t, j] = G[b, t, j]h[b, t-1, j] + (1 - G[b, t, j])R[b, t, j]$$

Modify the above equations so that they correspond to the following diagram for a Gated Recurrent Unit (GRU).



[Christopher Olah]

**Solution:**

$$G_1[b, t, j] = \sigma (W^{h,G_1}[j, I]h[b, t-1, I] + W^{x,G_1}[j, K]x[b, t, K] - B^{G_1}[j])$$

$$h'[b, t, j] = G_1[b, t, j]h[b, t-1, j]$$

$$G_2[b, t, j] = \sigma (W^{h,G_2}[j, I]h[b, t-1, I] + W^{x,G_2}[j, K]x[b, t, K] - B^{G_2}[j])$$

$$R[b, t, j] = \tanh (W^{h,R}[j, I]h'[b, t, j] + W^{x,R}[j, K]x[b, t, K] - B^R[j])$$

$$h[b, t, j] = (1 - G_2[b, t, j])h[b, t-1, j] + G_2[b, t, j]R[b, t, j]$$