

# TTIC 31230 Fundamentals of Deep Learning, Autumn 2021

## Exam 1

**Problem 1: (25 pts).** This problem is on softmax, entropy and cross entropy. Let the variable  $i$  range over the non-negative integers (0 to  $\infty$ ). Consider the following softmax distribution where  $\beta > 0$  is an inverse temperature parameter.

$$P(i) = \text{softmax}_i [-\beta i]$$

(a) Give an expression for  $P(i)$  in terms of  $i$  without using a softmax or infinite sum. You can use the equality that for  $0 < \lambda < 1$  we have

$$\sum_{n=0}^{\infty} \lambda^n = \frac{1}{1-\lambda}$$

(b) Solve for the entropy of this distribution. You can use the equality that for  $0 < \lambda < 1$  we have

$$\sum_{n=0}^{\infty} n\lambda^n = \frac{\lambda}{(1-\lambda)^2}$$

(c) Consider a one-hot population distribution defined by  $\text{Pop}(k) = 1$  and  $\text{Pop}(i) = 0$  for  $i \neq k$ . What is the cross-entropy  $H(\text{Pop}, P)$  and KL-divergence  $KL(\text{Pop}, P)$ ? What is the cross-entropy  $H(P, \text{Pop})$  and the KL divergence  $KL(P, \text{Pop})$ ?

**Problem 2. (25 pts)** 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  $h_\ell[I]$ .

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

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

**Problem 3. (25 pts) 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

$$W^{x,R}[j, K]x[b, t-1, K] - B^{x,R}[j]$$

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.

**Problem 4. (25 pts)** The self-attention in the transformer is computed by the following equations.

$$\text{Query}_{\ell+1}[k, t, i] = W_{\ell+1}^Q[k, i, J]L_\ell[t, J]$$

$$\text{Key}_{\ell+1}[k, t, i] = W_{\ell+1}^K[k, i, J]L_\ell[t, J]$$

$$\alpha_{\ell+1}[k, t_1, t_2] = \text{softmax}_{t_2} \left[ \frac{1}{\sqrt{I}} \text{Query}_{\ell+1}[k, t_1, I] \text{Key}_{\ell+1}[k, t_2, I] \right]$$

Notice that here the shape of  $W^Q$  and  $W^K$  are both  $[K, I, J]$ . We typically have  $I < J$  which makes the inner product in the last line an inner product of lower dimensional vectors.

(a) Give an equation computing a tensor  $\tilde{W}^Q[K, J, J]$  computed from  $W^Q$  and  $W^K$  such that the attention  $\alpha(k, t_1, t_2)$  can be written as

$$\alpha_{\ell+1}(k, t_1, t_2) = \text{softmax}_{t_2} \left[ L_\ell[t_1, J_1] \tilde{W}^Q[k, J_1, J_2] L_\ell[t_1, J_2] \right]$$

For a fixed  $k$  we have that  $W^Q[k, I, J]$  and  $W^K[k, I, J]$  are matrices. We want a matrix  $\tilde{W}^Q[k, J, J]$  such that the attention can be written in matrix notation as  $h_1^\top \tilde{W}^Q h_2$  where  $h_1$  and  $h_2$  are vectors and  $\tilde{W}^Q$  is a matrix. You need write this matrix  $\tilde{W}^Q$  in terms of the matrices for  $W^Q$  and  $W^K$ . But write your final answer in Einstein notation with  $k$  as the first index.

(b) Part (a) shows that we can replace the key and query matrix with a single query matrix without any loss of expressive power. If we eliminate the key matrix in this way what is the resulting number of query matrix parameters for a given layer and how does this compare to the number of key-query matrix parameters for a given layer in the original transformer version.