

ML HW 7 - Deep Learning

Thursday, January 14, 2021 12:41 PM

Problem 3: In machine learning you often come across problems which contain the following quantity

$$y = \log \sum_{i=1}^N e^{x_i} \leftarrow$$

For example if we want to calculate the log-likelihood of neural network with a softmax output we get this quantity due to the normalization constant. If you try to calculate it naively, you will quickly encounter underflows or overflows, depending on the scale of x_i . Despite working in log-space, the limited precision of computers is not enough and the result will be ∞ or $-\infty$.

To combat this issue we typically use the following identity:

$$y = \log \sum_{i=1}^N e^{x_i} = a + \log \sum_{i=1}^N e^{x_i - a}$$

for an arbitrary a . This means, you can shift the center of the exponential sum. A typical value is setting a to the maximum ($a = \max_i x_i$), which forces the greatest value to be zero and even if the other values would underflow, you get a reasonable result.

Your task is to show that the identity holds.

This is called the *log-sum-exp trick* and is often used in practice.

$$\begin{aligned} y &= \log \sum_{i=1}^N e^{x_i} \\ e^y &= \sum_{i=1}^N e^{x_i} & | \exp(\cdot) \\ e^{-a} e^y &= e^{-a} \sum_{i=1}^N e^{x_i} & | e^{-a} \\ e^{y-a} &= \sum_{i=1}^N e^{-a} e^{x_i} \\ y - a &= \log \sum_{i=1}^N e^{x_i - a} \\ y &= a + \log \sum_{i=1}^N e^{x_i - a} \end{aligned}$$

$$e^a \cdot e^b = e^{a+b}$$

Problem 4: Similar to the previous exercise we can compute the output of the softmax function $\pi_i = \frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}}$ in a numerically stable way by shifting by an arbitrary constant a :

$$\frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}} = \frac{e^{x_i - a}}{\sum_{j=1}^N e^{x_j - a}}$$

often chosen $a = \max_i x_i$. Show that the above identity holds.

For some arbitrary constant C we have

$$\frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}} = \frac{C e^{x_i}}{C \sum_{j=1}^N e^{x_j}} = \frac{e^{x_i + \log(C)}}{\sum_{j=1}^N e^{x_j + \log(C)}}$$

Since C is just an arbitrary constant, we can replace $\log(C) = -a$ and get $\frac{e^{x_i - a}}{\sum_{j=1}^N e^{x_j - a}}$.

$$C = e^{-a}$$

$$C \cdot e^a = e^{\log C} \cdot e^a = e^{\log C + a}$$

forward:

X $[N, D]$ W $[D, H]$ b $[H]$ \mapsto A $[N, H]$

x_i $[D]$ $a_i = W^T x_i + b$ $[H]$

$A_{ij} = x_{i1} \cdot w_{1j} + x_{i2} \cdot w_{2j} + \dots + x_{iD} \cdot w_{Dj} + b_j$
 $= X[i, :] @ W[:, j] + b[j]$

$A = X @ W + b$ \leftarrow forward
 $[N, H]$ $[H]$

$\frac{\partial E}{\partial b}$ $[1, H]$

backward

d_out $[N, H]$ \mapsto d_inputs $[N, D]$ d_weight $[D, H]$ d_bias $[H]$

$d_out[i, j] = \frac{\partial E}{\partial A_{ij}}$ $d_inputs[k, l] = \frac{\partial E}{\partial x_{kl}}$ $d_weight[k, l] = \frac{\partial E}{\partial w_{kl}}$ $d_bias[k] = \frac{\partial E}{\partial b_k}$

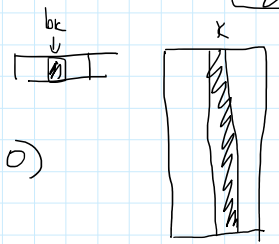
$\frac{\partial E}{\partial x_{kl}} = \sum_{i=1}^N \frac{\partial E}{\partial A_{ij}} \cdot \frac{\partial A_{ij}}{\partial x_{kl}} = \sum_{i=1}^N d_out[i, j] \cdot w_{il}$

$$d_bias[k] = \frac{\partial E}{\partial b_k} = \sum_{i=1}^N \sum_{j=1}^H \frac{\partial E}{\partial A_{ij}} \cdot \frac{\partial A_{ij}}{\partial b_k} = \sum_{i=1}^N \sum_{j=1}^D \frac{\partial E}{\partial A_{ij}} \cdot \mathbb{I}(j=k) = \sum_{i=1}^N \frac{\partial E}{\partial A_{ik}}$$

$$\frac{\partial A_{ij}}{\partial b_k} = \begin{cases} 0 & \text{if } j \neq k \\ 1 & \text{if } j = k \end{cases}$$

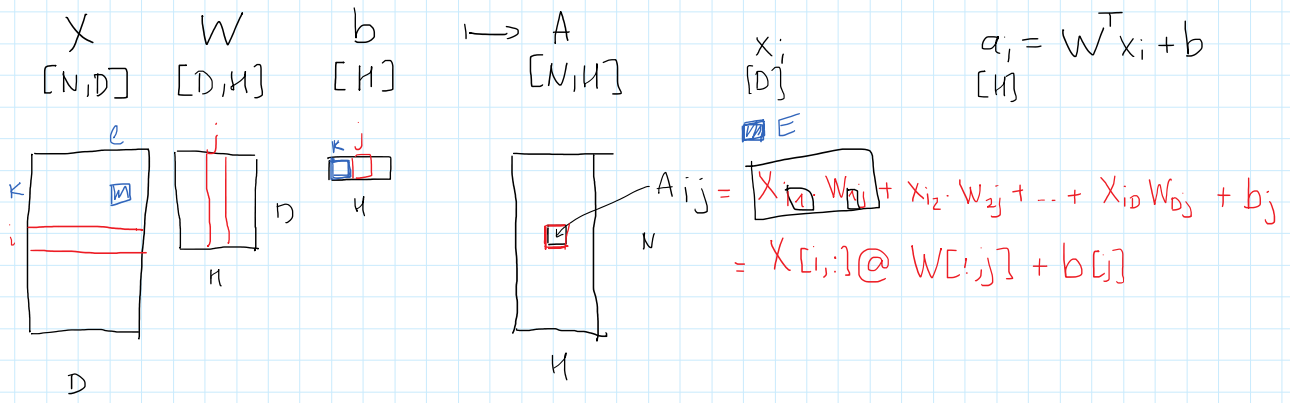
$$\frac{\partial E}{\partial A_{ik}} \cdot \mathbb{I}(k=1) + \frac{\partial E}{\partial A_{ik}} \cdot \mathbb{I}(k=2) + \dots$$

$$d_bias[k] = \sum_{i=1}^N \frac{\partial E}{\partial A_{ik}} \quad d_out = \left(\frac{\partial E}{\partial A} \right)^T$$



$$d_bias = np.sum(d_out, axis=0)$$

forward:



$$d_inputs[k,l] = \frac{\partial E}{\partial x_{kl}} = \sum_{i=1}^N \sum_{j=1}^D \frac{\partial E}{\partial A_{ij}} \cdot \frac{\partial A_{ij}}{\partial x_{kl}} = \sum_{i=1}^N \sum_{j=1}^H \frac{\partial E}{\partial A_{ij}} \cdot w_{lj} \cdot \mathbb{I}(i=k) = \sum_{j=1}^H \frac{\partial E}{\partial A_{kj}} \cdot w_{lj}$$

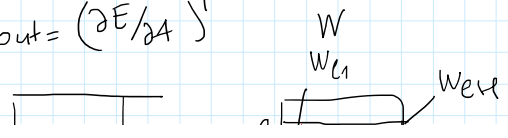
$$\frac{\partial A_{ij}}{\partial x_{kl}} = \begin{cases} 0 & \text{if } i \neq k \\ w_{lj} & \text{if } i = k \end{cases}$$

$$= w_{lj} \cdot \mathbb{I}(i=k)$$

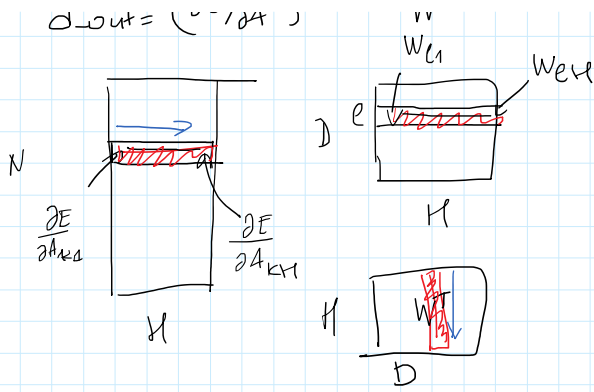
$$\frac{\partial A_{ij}}{\partial x_{il}} = w_{lj}$$

$$\frac{\partial E}{\partial x_{kl}} = \sum_{j=1}^H \frac{\partial E}{\partial A_{kj}} \cdot w_{lj}$$

$$d_out = \left(\frac{\partial E}{\partial A} \right)^T$$



$$d_inputs[k,l] = \sum_{j=1}^H \frac{\partial E}{\partial A_{kj}} \cdot w_{lj}$$



$$d_{out}[k, :] @ W[l, :] = \frac{\partial E}{\partial x_{kl}}$$

$$d_{out} @ W.T \rightarrow d_{input}$$

$[N, H] \quad [H, D] \quad [N, D]$

$$y = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases} \Rightarrow y = x \cdot \mathbb{I}(x \geq 0) =$$

$$\frac{\partial y}{\partial x} = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases} \Rightarrow \frac{\partial y}{\partial x} = \mathbb{I}(x \geq 0)$$

$$\frac{\partial E}{\partial x} = \frac{\partial E}{\partial y} \frac{\partial y}{\partial x} = \frac{\partial E}{\partial y} \mathbb{I}(x \geq 0)$$

Cross Ent $(y, A) \mapsto E$

forward:

y (labels) A (logits)
 $[N, K] \quad [N, K]$

0	0	1	0
1	0	0	0
0	1	0	0
1	0	0	0

$$E = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{ik} \cdot \log p(y_i = k | x_i)$$

$$= \frac{1}{N} \sum_{i,k} y_{ik} \cdot \log \frac{e^{A_{ik}}}{\sum_{c=1}^K e^{A_{ic}}} \leftarrow \text{softmax } A[i, :]$$

$\log \frac{a}{b} = \log a - \log b$

$$= \frac{1}{N} \sum_{i,k} y_{ik} (A_{ik} - \log(\sum_c e^{A_{ic}}))$$

log-sum-exp

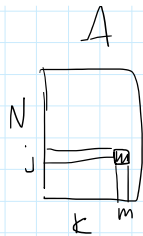
$$\frac{1}{N} \text{np.sum}(y \odot \log P)$$

$[N, K]$

$$\log P_{ik} \rightarrow (A_{ik} - m_i - \log(\sum_c e^{A_{ic} - m_i}))$$

$A[i, :]$

$$m_i = \max A_{ic}$$



$$m_i = \max_c A_{ic}$$

backward:

$$d_{\text{out}} \mapsto d_{\text{logits}} [j, m] = \frac{\partial E}{\partial A_{jm}}$$

[N, K]

$$E = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{ik} (A_{ik} - \log(\sum_{c=1}^K e^{A_{ic}}))$$

$$= \frac{1}{N} \sum_{i=1}^N \left[\left(\sum_{k=1}^K y_{ik} \cdot A_{ik} \right) - \left(\sum_{k=1}^K y_{ik} \cdot \log(\sum_c e^{A_{ic}}) \right) \right]$$

$$E = \frac{1}{N} \sum_i \left[\left(\sum_k y_{ik} A_{ik} \right) - \log(\sum_c e^{A_{ic}}) \cdot \sum_k y_{ik} \right]$$

010001

$$\frac{\partial E}{\partial A_{jm}} = \frac{1}{N} \frac{\partial}{\partial A_{jm}} \left[\left(\sum_k y_{jk} A_{jk} \right) - \log(\sum_c e^{A_{jc}}) \right]$$

$$= \frac{1}{N} \left[y_{jm} - \frac{\partial}{\partial A_{jm}} \log(\sum_c e^{A_{jc}}) \right]$$

$$\frac{\partial}{\partial A_{jm}} \log(e^{A_{j1}} + e^{A_{j2}} + \dots + e^{A_{jK}})$$

$$\frac{1}{\sum_c e^{A_{jc}}} \cdot e^{A_{jm}} = \frac{e^{A_{jm}}}{\sum_c e^{A_{jc}}}$$

$$d_{\text{logits}} = \frac{1}{N} (y - p)$$

$$= \frac{1}{N} \cdot (\text{labels} - \text{probas}) = \frac{1}{N} \left[y_{jm} - \frac{e^{A_{jm}}}{\sum_c e^{A_{jc}}} \right]$$

$$= \frac{1}{N} [y_{jm} - p_{jm}]$$

