

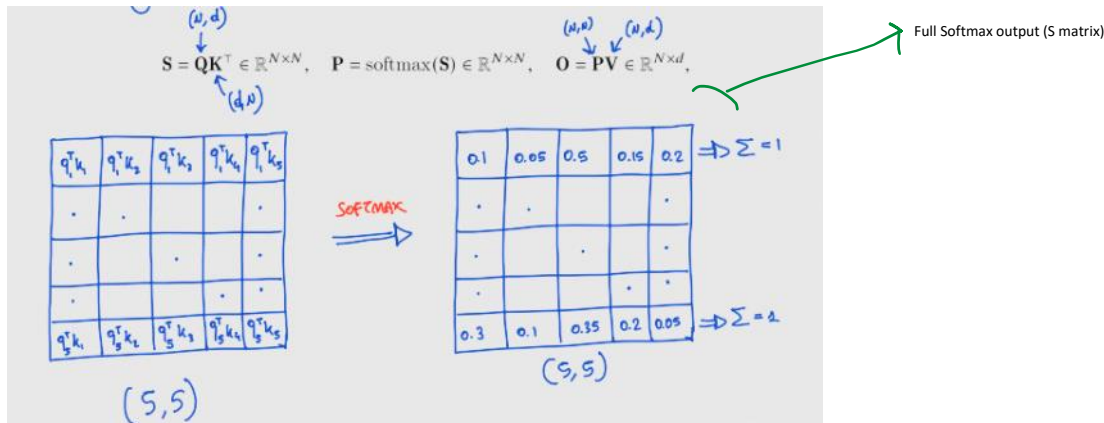
Online Softmax

Thursday, October 23, 2025 4:12 PM

Process of turning: Unsafe Naive Softmax --> Safe Naive Softmax --> Online Softmax

Making Softmax "safe":

Considerations:



If x_i is too large, e^{x_i} explodes! (eg: e^{100} , e^{200})

Making Softmax "safe" (numerically stable)

Given a vector $x \in \mathbb{R}^N$, the softmax is defined as:

$$\text{softmax}(x_i) = \frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}}$$

But there's a problem! If the values of the vector are large, the exponential will explode!

Numerically unstable = cannot be represented with a float 32 or float 16

If exponential softmax will not be stored. Reason: outside representable range of standard formats (fp16, fp32)

If e^{x_i} explodes, the softmax output cannot be stored. Reason: outside representable range of standard formats (fp16, fp32)

Normalization Factor: sum of e^{x_j} , where j is each row vector element

Luckily, we have a solution:

$$\begin{aligned} \frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}} &= \frac{c \cdot e^{x_i}}{c \cdot \sum_{j=1}^N e^{x_j}} = \frac{c e^{x_i}}{\sum_{j=1}^N c e^{x_j}} = \frac{e^{x_i + \log(c)}}{\sum_{j=1}^N e^{x_j + \log(c)}} \\ &= \frac{e^{x_i - k}}{\sum_{j=1}^N e^{x_j - k}} \quad \text{where } k = -\log(c) \end{aligned}$$

So we can "sneak in" a constant in the exponential to decrease its argument and make it numerically stable.

We will choose $k = \max_i(x_i)$

Largest x_i among all x_i elements making up THAT row vector on the S matrix

(S : the full softmax output matrix)

Largest x_i among all x_i elements making up THAT row vector on the S matrix (S : the full softmax output matrix)

So effectively, we have two cases:

- $x_i = k$: Causes $x_i - k = 0$, giving $e^{x_i - k} = 1$. "1" is easily representable in standard formats (fp16, fp32)
- $x_i < k$: Causes $x_i - k < 0$, causing $e^{x_i - k}$ (negative value) to be easily representable in standard formats (fp16, fp32)

So effectively, we can have two cases:

- $x_i = k$: Causes $x_i - k = 0$, causing $e^{x_i - k} = 1$

Causes $x_i - k = 0$, causing $e^*0 = 1$
 = '1': easily representable in standard formats (fp16, fp32)
 o $x_i < k$:
 = $e^{x_i - k} < 1$, causing e^{**} (negative value)
 = e^{**} (negative value): easily representable in standard formats (fp16, fp32)
 - Finally, the Softmax output can be represented safely! (does not explode to infinity or zero)

Safe Softmax:

$$\text{softmax}(x_i) = \frac{e^{x_i - x_{\max}}}{\sum_{j=1}^N e^{x_j - x_{\max}}}$$

- So effectively, we can have two cases:

- o $x_i == k$:
 - Causes $x_i - k = 0$, causing $e^*0 = 1$
 - '1': easily representable in standard formats (fp16, fp32)
 - o $x_i < k$:
 - Causes $x_i - k < 0$, causing e^{**} (negative value)
 - 'e**'(negative value)': easily representable in standard formats (fp16, fp32)
- Finally, the Softmax output can be represented safely! (does not explode to infinity)

2

But, this is still "naïve"
 Why? Look on ahead:

Let's review the algorithm:

$$\text{softmax}(x_i) = \frac{e^{x_i - x_{\max}}}{\sum_{j=1}^N e^{x_j - x_{\max}}}$$

given a $N \times N$ matrix, for each row.

STEP 1 Find the max value among all elements
 Time complexity: $O(N)$
 Memory reads: $O(N)$

STEP 2 Calculate the normalization factor
 Time complexity: $O(N)$
 Memory reads: $O(N)$

STEP 3 Apply the softmax to each element of the vector
 Time Complexity: $O(N)$
 Memory reads: $O(N)$

Pseudocode:

```

m0 = -∞
for i = 1 to N
  m_i = max(m_{i-1}, x_i)
l0 = 0
for j = 1 to N
  l_j = l_{j-1} + e^{x_j - m_N}
for k = 1 to N
  x_k ← e^{x_k - m_N} / l_N
  
```

STEP 1: Finding x_{\max} by looping over all x_i in the corresponding row vector, and accumulating max

STEP 2: Computing Normalization Factor (l_j)
 (can only be done AFTER Step 1, as we need it for STEP 2)

Again, takes $O(n)$ time

Takes $O(n)$ time

STEP 3: Compute Softmax over each row vector element x_i
 (can only be done AFTER STEPS 1 & 2, as we need calculations from both for STEP 3)

One more time!
 Takes $O(n)$ time

TOO SLOW!!!
 3 sequential $O(n)$ operations here

3

So, what do we do to make this FASTER
 (cutting time complexity without exploding exponential, i.e. keeping it "safe")?

New pseudocode (Online Softmax)

```

m0 = -∞
l0 = 0
for i = 1 to N
  m_i = max(m_{i-1}, x_i)
  l_i = l_{i-1} * e^{m_{i-1} - m_i} + e^{x_i - m_i}
for k = 1 to N
  x_k ← e^{x_k - m_N} / l_N
  
```

- Fusing STEP 1 and STEP 2 (by using a "Correction Factor")
- So, as we loop over each row vector element, we calculate $\max x_i$ (SO FAR) AND the Normalization Factor (SO FAR), in the same loop.
- This is opposed to Safe Naïve Softmax, where we went looping over row vector to first get $\max x_i$, then again looping over row vector to get Normalization Factor

Correction Factor
 $(m_{i-1} : \max x_i \text{ until the 'i-1'th row vector element (looking at all } x_i \text{ from the start, the } \max x_i \text{ so far),}$
 $m_i : \max x_i \text{ at the 'i'th row vector element})$

STEP 3 (remains as it was)

2 sequential $O(n)$ operations vs the 3 sequential $O(n)$ earlier

$$X = [3, 2, 5, 1]$$

We can take this row vector X and run the Fused Steps 1 & 2 from the above to see if the math works out, turns out it does!!
Finally, we can do a 'Proof by Induction' for this, turns out that works as well.