# The problem

The problem statement is: can we find a way to compute the softmax of a vector without going through the vector 3 times, but also preventing the exponential from exploding?

$$X = \begin{bmatrix} 3, & 2, 5, & 1 \end{bmatrix}$$

Pseudocode:

$m_0 = -\infty$
for $i = 1$ to $N$
  $m_i = \max(m_{i-1}, x_i)$

$l_0 = 0$
for $J = 1$ to $N$
  $l_J = l_{J-1} + e^{x_J - m_N}$

for $k = 1$ to $N$
  $x_k \leftarrow \dfrac{e^{x_k - m_N}}{l_N}$

Let's try fusing the computation of the max element with the normalization constant

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

## STEP 1:

$$x = [3, \quad 2, 5, \quad 1]$$

$$\text{max}_1 = 3$$

$$l_1 = e^{3-3}$$

## STEP 2

$$\text{max}_2 = \text{max}(3, 2) = 3$$

$$l_2 = l_1 + e^{2-3}$$

If our vector only contained the first two elements, then the max element and the normalization factor we've computed would be correct.

However, things change at position 3...

$$X = \begin{bmatrix} 3, & 2, 5, & 1 \end{bmatrix}$$

STEP 3

$$\max_3 = \max(3, 5) = 5$$

$$\ell_3 = \ell_2 + e^{5-5} = e^{3-3} + e^{2-3} + e^{5-5}$$

The $\ell_3$ we computed is wrong!

$$\ell_3 = \underbrace{e^{3-3} + e^{2-3}}_{} + e^{5-5}$$

Here it is wrong!

Can we fix it ON THE FLY? YES!

$$\ell_3 = \ell_2 \cdot \boxed{e^{3-5}} + e^{5-5} = \left( e^{3-3} + e^{2-3} \right) e^{3-5} + e^{5-5}$$

correction factor

$$= e^{3-3+3-5} + e^{2-3+3-5} + e^{5-5} = e^{3-5} + e^{2-5} + e^{5-5}$$

CORRECT!

So, every time we encounter a number bigger than the current maximum, we can "fix" the normalization constant computed so far!

$$x = \begin{bmatrix} 3, & 2, 5, & 1 \end{bmatrix}$$

STEP 4

$$\max_4 = \max(5, 1) = 5$$

$$l_4 = l_3 \cdot \boxed{e^{5-5}} + e^{1-5}$$

IN THIS CASE
WE DO NOT NEED
TO FIX ANYTHING

New pseudocode

$$m_0 = -\infty$$
$$l_0 = 0$$
for $i = 1$ to $N$
$$\quad m_i = \max(m_{i-1}, x_i)$$
$$\quad l_i = l_{i-1} \cdot e^{m_{i-1} - m_i} + e^{x_i - m_i}$$
for $u = 1$ to $N$
$$\quad x_u \leftarrow \frac{e^{x_u - m_N}}{l_N}$$

Can we prove this algorithm is correct?

Let's do it!

$m_0 = -\infty$

$l_0 = 0$

for $i = 1$ to $N$

$\quad m_i = \max(m_{i-1}, x_i)$

$\quad l_i = l_{i-1} \cdot e^{m_{i-1} - m_i} + e^{x_i - m_i}$

for $k = 1$ to $N$

$\quad x_k \leftarrow \dfrac{e^{x_k - m_N}}{l_N}$

We want to prove that at the end of this loop:

$$m_N = \max_i (x_i) = x_{MAX}$$

$$l_N = \sum_{J=1}^{N} e^{x_J - x_{MAX}}$$

We will prove it by induction:

1) Prove that it holds for a vector of size

$N=1$

$$m_1 = \max\left(-\infty, x_1\right) = x_1 = \max_i \left(x_i\right) = x_{MAX}$$

$$l_1 = 0 \times e^{-\infty} + e^{x_1 - x_1} = \sum_{J=1}^{N} e^{x_i - x_{MAX}}$$

2) If we assume it holds for a vector of size $N$, does it hold for a vector of size $N+1$ ?

$$m_{N+1} = \max\left(m_N, x_{N+1}\right) = \max_i \left(x_i\right)$$

$$l_{N+1} = l_N e^{m_N - m_{N+1}} + e^{x_{N+1} - m_{N+1}} =$$

$$= \left(\sum_{J=1}^{N} e^{x_J - m_N}\right) e^{m_N - m_{N+1}} + e^{x_{N+1} - m_{N+1}}$$

$$= \sum_{J=1}^{N} e^{x_J - m_{N+1}} + e^{x_{N+1} - m_{N+1}}$$

$$= \sum_{J=1}^{N+1} e^{x_J - m_{N+1}}$$