# Making softmax safe

$$\mathbf{S} = \mathbf{Q}\mathbf{K}^\top \in \mathbb{R}^{N \times N}, \quad \mathbf{P} = \text{softmax}(\mathbf{S}) \in \mathbb{R}^{N \times N}, \quad \mathbf{O} = \mathbf{P}\mathbf{V} \in \mathbb{R}^{N \times d},$$

(N, d) over $\mathbf{Q}$, (d, N) over $\mathbf{K}^\top$

(N, N) over $\mathbf{P}$, (N, d) over $\mathbf{V}$

| $q_1^T k_1$ | $q_1^T k_2$ | $q_1^T k_3$ | $q_1^T k_4$ | $q_1^T k_5$ |
|---|---|---|---|---|
| . | . | . | | . |
| . | . | . | | . |
| . | | . | . | . |
| $q_5^T k_1$ | $q_5^T k_2$ | $q_5^T k_3$ | $q_5^T k_4$ | $q_5^T k_5$ |

(5,5)

**SOFTMAX** ⟹

| 0.1 | 0.05 | 0.5 | 0.15 | 0.2 | ⟹ $\Sigma = 1$ |
|---|---|---|---|---|---|
| . | . | | | . | |
| . | | . | | . | |
| . | | | . | . | |
| 0.3 | 0.1 | 0.35 | 0.2 | 0.05 | ⟹ $\Sigma = 1$ |

(5,5)

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

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

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

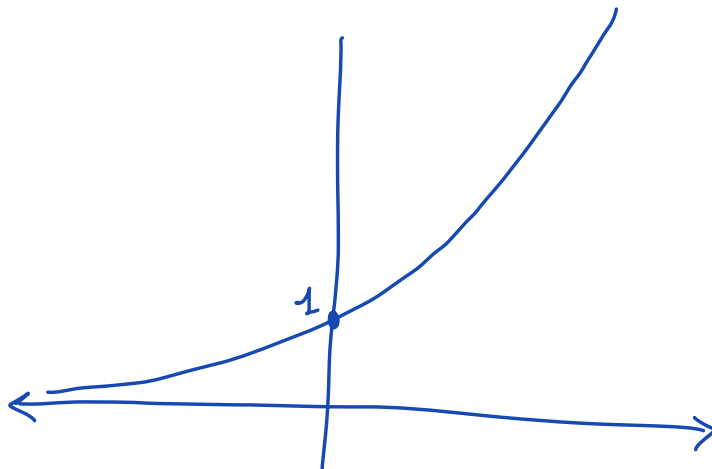Numerically unstable = cannot be represented with a float 32 or float 16

Luckily, we have a solution:

$$\frac{e^{x_i}}{\sum\limits_{j=1}^{N} e^{x_j}} = \frac{c \cdot e^{x_i}}{c \cdot \sum\limits_{j=1}^{N} e^{x_j}} = \frac{c e^{x_i}}{\sum\limits_{j=1}^{N} c e^{x_j}} = \frac{e^{\log(c)} e^{x_i}}{\sum\limits_{j=1}^{N} e^{\log(c)} e^{x_j}} =$$

$$= \frac{e^{x_i + \log(c)}}{\sum\limits_{j=1}^{N} e^{x_j + \log(c)}} = \frac{e^{x_i - k}}{\sum\limits_{j=1}^{N} e^{x_j - k}} \quad \text{where } k = -\log(c)$$

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

We will choose $k = \max\limits_{i} (x_i)$

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.

1) Find the max value among all elements

Time complexity: $O(N)$
Memory reads : $O(N)$

2) Calculate the normalization factor

Time complexity: $O(N)$
Memory reads: $O(N)$

3) Apply the softmax to each element of the vector

Time Complexity: $O(N)$
Memory reads: $O(N)$

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

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 see a practical example

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

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

1) $X_{max} = 5$

2) $e^{3-5} + e^{2-5} + e^{5-5} + e^{1-5} = e^{-2} + e^{-3} + e^{0} + e^{-4} = \ell$

3) $X_1 = \dfrac{e^{3-5}}{\ell}$   $X_2 = \dfrac{e^{2-5}}{\ell}$   $X_3 = \dfrac{e^{5-5}}{\ell}$   $X_4 = \dfrac{e^{1-5}}{\ell}$

To apply the softmax to a $N \times N$ matrix, we need to load each of its elements 3 times, and it must be done sequentially...

Is there a better way?