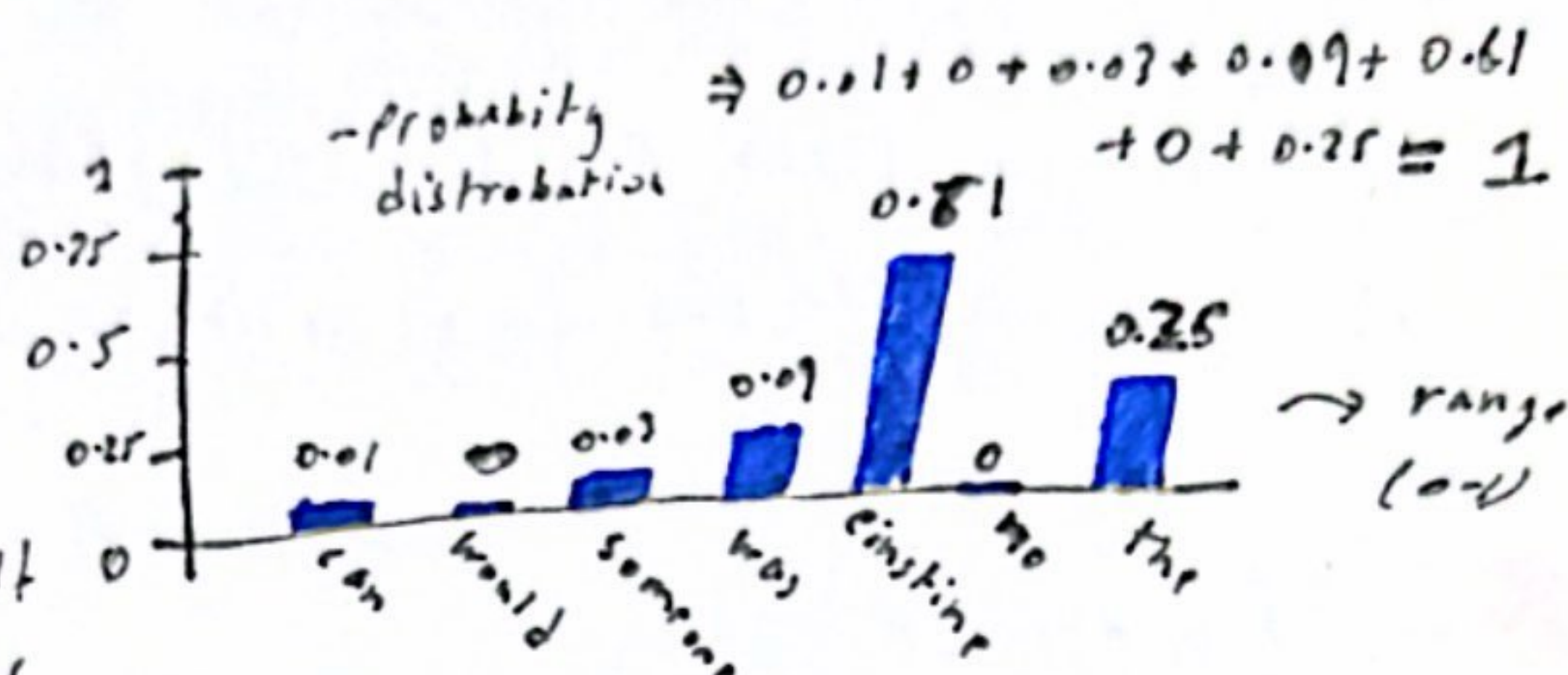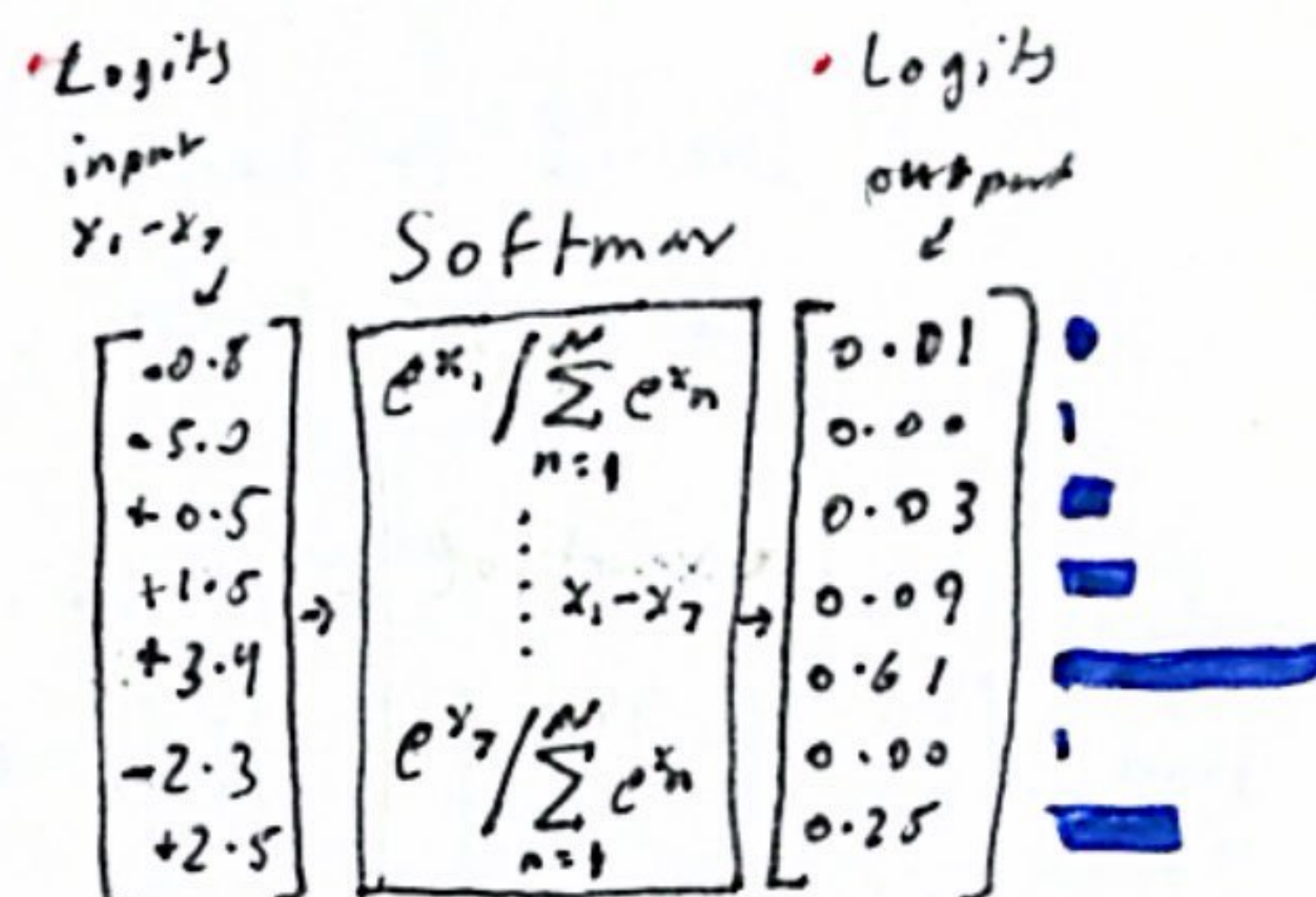# Softmax (LLM notes)

- if you want a sequence of numbers to act as a probability distribution ex a distribution over all next words, then each value has to be between 0 and 1 and for all of them to add up to 1

- but since we have alot of dot products and matrix multiplication the output we have by default dont follow these rules of probability distrabution



-Probability distribution

$\Rightarrow 0.01 + 0 + 0.03 + 0.09 + 0.61$
$+ 0 + 0.25 = 1$

$\rightarrow$ range (0-1)

- Soft max is a way to turn a list of arbitrary numbers into a valid distrabution is such a way that the largest values end up closer to 1 and the smaller value end up closer to 0

- this is known as a normalization step

- Logits input $x_1 - x_j$

$$Softmax \quad \begin{bmatrix} -0.8 \\ -5.0 \\ +0.5 \\ +1.5 \\ +3.4 \\ -2.3 \\ +2.5 \end{bmatrix} \rightarrow \begin{bmatrix} e^{x_1}/\sum_{n=1}^{N} e^{x_n} \\ \vdots \\ x_1 - x_7 \\ e^{x_7}/\sum_{n=1}^{N} e^{x_n} \end{bmatrix} \rightarrow$$

- Logits output

$$\begin{bmatrix} 0.01 \\ 0.00 \\ 0.03 \\ 0.09 \\ 0.61 \\ 0.00 \\ 0.25 \end{bmatrix}$$

$$Softmax (x)_i = \frac{e^{x_i}}{\sum_{n=1}^{N} e^{x_n}}$$

- in this way if one input is meaningfully bigger (3.4) then it dominates the distribution and vise versa so it we were sampling from this we would almost alway pick the largest and most meaning full one (best prediction)

## -Temprature (Creativity)

- we can add a constant T to the denominator of softmax

$$Softmax \text{ with } (x)_i = \frac{e^{x_i}/T}{\sum e^{x_n}/T}$$

- if T is larger more weight is given to the lower values boosting them up and vise versa at T=0 all weight = 1 goes to the max input

$$T=9 \begin{bmatrix} 6 \\ -5 \\ 4 \\ 1.5 \\ 1.9 \\ -2.3 \\ 4 \end{bmatrix} \rightarrow Softmax \rightarrow \begin{bmatrix} 0.18 \\ 0.05 \\ 0.15 \\ 0.11 \\ 0.27 \\ 0.08 \\ 0.15 \end{bmatrix}$$
With Temperature

$$T=2 \begin{bmatrix} 6 \\ -5 \\ 4 \\ 1.5 \\ 1.9 \\ -2.3 \\ 4 \end{bmatrix} \rightarrow Softmax \rightarrow \begin{bmatrix} 0.11 \\ 0.04 \\ 0.01 \\ 0.79 \\ 0 \\ 0.04 \end{bmatrix}$$
with temperature

53