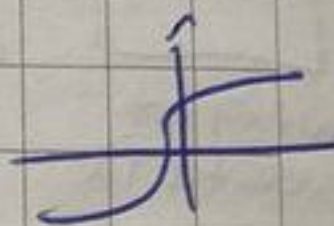


Логистическая
функция

$$f(x) = \frac{1}{1+e^{-x}} = \frac{e^x}{1+e^x}$$

$x \in \mathbb{R}$



Обобщение на высокие размерности:

Softmax

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

$$S(x) = \begin{pmatrix} \frac{e^{x_1}}{\sum_{j=1}^n e^{x_j}} \\ \vdots \\ \frac{e^{x_n}}{\sum_{j=1}^n e^{x_j}} \end{pmatrix} \in \mathbb{R}^n$$

Свойства

$$1) \|S(x)\|_1 = \sum_i \frac{e^{x_i}}{\sum_j e^{x_j}} = 1$$

\Rightarrow вероятностная интерпретация

2) вычислительная стабильность
(numerical stability) Softmax'a

— это
$$\frac{e^{x_j+D}}{\sum_k e^{x_k+D}} = \frac{e^D \cdot e^{x_j}}{e^D \cdot \sum_k e^{x_k}} = \frac{e^{x_j}}{\sum_k e^{x_k}}$$

Свойство исп. при "нормализации"
данных, потому что в Python

$$\max(\text{float64}) = 10^{38}$$

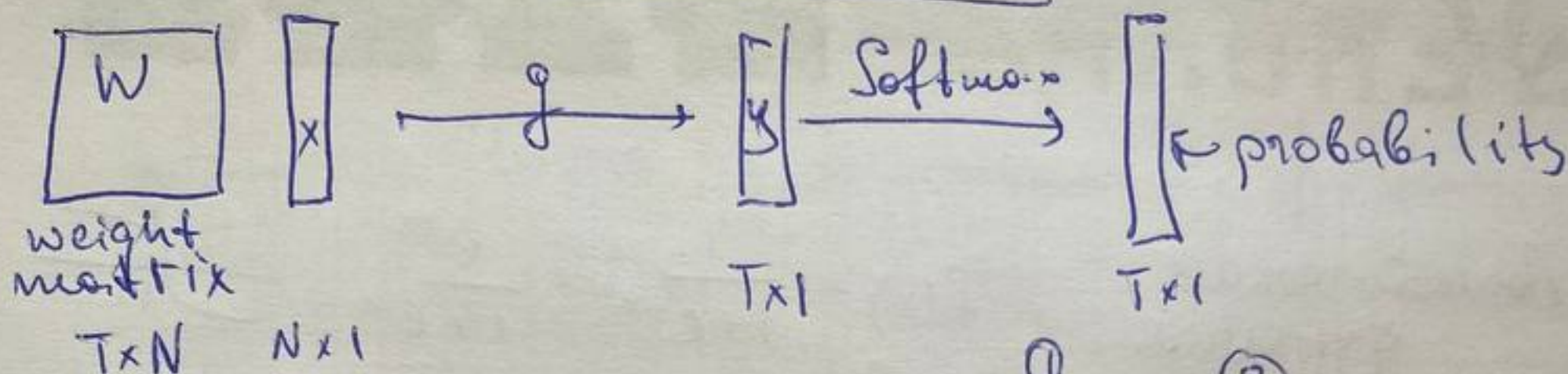
$$\Rightarrow \text{softmax}(1000) \sim e^{1000} = \text{NaN}$$

Нормализуем: $e^{x_j} \rightarrow e^{x_j - \max_i(x_i)}$

(в предположении, что данные
сгруппированы)

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

Вычисление градиента (сложной ф-ции)



$$\nabla S(g(W)) = \nabla S(g) \overset{(1)}{\nabla g} \overset{(2)}{W}$$

① given Softmax

$$\nabla S = \begin{bmatrix} \nabla_1 S_1 & \dots & \nabla_N S_1 \\ \vdots & & \vdots \\ \nabla_1 S_N & \dots & \nabla_N S_N \end{bmatrix} \quad \text{Jacobian}$$

$$\nabla_i S_j = \frac{\partial S_j}{\partial x_i} = \frac{\partial}{\partial x_i} \left(\frac{e^{x_j}}{\sum_{k=1}^N e^{x_k}} \right)$$

$$= \begin{cases} i \neq j: & \frac{-e^{x_i} e^{x_j}}{(\sum x_k)^2} = -S_i \cdot S_j \\ i = j & \frac{e^{x_i} (\sum x_k) - e^{2x_i}}{(\sum x_k)^2} = \frac{e^{x_i}}{\sum} \cdot \frac{\sum - e^{x_i}}{\sum} = S_i (1 - S_i) \end{cases}$$

$$= \begin{cases} -S_i S_j, & i \neq j \\ S_i (1 - S_i), & i = j \end{cases} = S_j (\delta_{ij} - S_i)$$

$$\delta_{ij} = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases}$$

$$\textcircled{2} \quad \nabla g = \begin{bmatrix} \nabla_1 g_1 & \dots & \nabla_N g_1 \\ \vdots & & \vdots \\ \nabla_1 g_T & \dots & \nabla_N g_T \end{bmatrix}$$

$g_i = w_{i1}x_1 + \dots + w_{iN}x_N$ (матричное умножение)

$$\nabla_{ij} g_t = \frac{\partial (w_{i1}x_1 + \dots)}{\partial w_{ij}} = \begin{cases} x_i, & i = t \\ 0, & i \neq t \end{cases}$$

Чтого же сложной ф-ции:

$$\begin{aligned} \nabla_{ij} (S_t(g(w))) &= \sum_k \nabla_k S_t \underbrace{\nabla_{ij} g_k}_{=0, i \neq k} = \\ &= \nabla S_t \underbrace{\nabla_{ij} g_i}_{=x_j} = \underbrace{\nabla_i S_t}_{\textcircled{1}} x_j = S_t(\delta_{ti} - s_i) x_j \end{aligned}$$

Чтого : • градиент от Softmax'a -
это сам Softmax (потти)

• Градиент линейной функции
от Softmax'a - это сам Softmax (потти)

⇒ НЕТ МАТРИЧНОГО УМНОЖЕНИЯ! ♥