

①  $\forall x \in \mathbb{R}^n, y \in \{1, 2, 3, 4, \dots, k\}$ . 数据集是一个由  $m \times n$  的  $x/y$  组成的

② 任务: 设计一种 softmax 机制来实现标签  $y$  的预测任务  
即寻找  $f: x \rightarrow y$ ,

③ 引入线性变化和 softmax:

$$\begin{array}{ccc} x_{11} & x_{12} & \dots & x_{1n} & & \theta_{11} & \theta_{12} & \dots & \theta_{1k} & & h_{11} & h_{12} & \dots & h_{1k} \\ x_{21} & x_{22} & \dots & x_{2n} & \times & \theta_{21} & \theta_{22} & \dots & \theta_{2k} & = & h_{21} & h_{22} & \dots & h_{2k} \\ \vdots & & & & & \vdots & & & & & \vdots & & & \\ x_{m1} & x_{m2} & \dots & x_{mn} & & \theta_{n1} & \theta_{n2} & \dots & \theta_{nk} & & h_{m1} & h_{m2} & \dots & h_{mk} \end{array}$$

$$X_{mn} \times \theta_{nk} \longrightarrow H_{mk}$$

引入 softmax ( $H$ )  $\longrightarrow$

$$\text{其中 } z_{ij} = \frac{e^{h_{ij}}}{\sum_{a=1}^k e^{h_{ia}}}$$

$$\begin{bmatrix} z_{11} & z_{12} & \dots & z_{1k} \\ z_{21} & z_{22} & \dots & z_{2k} \\ \vdots & \vdots & \vdots & \vdots \\ z_{m1} & z_{m2} & \dots & z_{mk} \end{bmatrix}$$

④  $z_{ij}$  可以看作  $P(\text{label} = j | x_i)$ , 一种预测类别为  $j$  的  $P$

定义损失函数  $\text{loss}(H; x_i, y_i) = -\log(p_{\text{label} = j | x_i})$

loss  $\uparrow$

图 2.1 损失函数



那么整体 
$$L_{\text{loss}} = \frac{1}{m} \sum_{i=1}^m -\log p(\text{label} = y^i | x_i, y_i)$$

⑤ 现在开始研究 Loss 和  $\theta$  的关系

$$L_{\text{loss}} = \frac{1}{m} \sum_{i=1}^m \left( -\log \frac{e^{h_i y_i}}{\sum_{a=1}^K e^{h_i a}} \right) = \frac{1}{m} \sum_{i=1}^m \left( -h_i y_i + \log \sum_{a=1}^K e^{h_i a} \right)$$

Loss 与  $h$  有关,  $h$  与  $\theta$  有关, 那么引入偏导

$$\frac{d L_{\text{loss}}}{d \theta_{r1}} = \frac{1}{m} \cdot \frac{d \sum_{i=1}^m (-h_i y_i + \log \sum_{a=1}^K e^{h_i a})}{d \theta_{r1}}$$

⑥  $\theta_{r1}$  会影响哪些  $h$ ?  $h_{a0} = x_{a1} \theta_{10} + x_{a2} \theta_{20} + \dots + x_{an} \theta_{n0}$   
观察发现  $\theta_{r1}$  和  $h_{11}, h_{21}, h_{31}, \dots, h_{m1}$  有关

$$\begin{cases} h_{11} = x_{11} \theta_{11} + x_{12} \theta_{21} + \dots + x_{1r} \theta_{r1} + \dots + x_{1n} \theta_{n1} \\ h_{21} = x_{21} \theta_{11} + x_{22} \theta_{21} + \dots + x_{2r} \theta_{r1} + \dots + x_{2n} \theta_{n1} \\ \vdots \\ h_{m1} = x_{m1} \theta_{11} + x_{m2} \theta_{21} + \dots + x_{mr} \theta_{r1} + \dots + x_{mn} \theta_{n1} \end{cases}$$

⑦ 记  $f_1 = -h_i y_i$ ,  $f_2 = \log \sum_{a=1}^K e^{h_i a}$

$$\frac{d \text{Loss}}{d \theta_{r1}} = \frac{1}{m} \sum_{i=1}^m \left( \frac{df_1}{d \theta_{r1}} + \frac{df_2}{d \theta_{r1}} \right)$$

$$\frac{df_1}{d \theta_{r1}} = \begin{cases} -x_{ir} & y^i = 1 \\ 0 & y^i \neq 1 \end{cases} \quad (h_{i1} = x_{i1}\theta_{11} + x_{i2}\theta_{21} + \dots + x_{ir}\theta_{r1} + \dots + x_{ik}\theta_{k1})$$

$$\begin{aligned} \frac{df_2}{d \theta_{r1}} &= \frac{df_2}{dh_{i1}} \cdot \frac{dh_{i1}}{d \theta_{r1}} = \frac{1}{e^{h_{i1}} + e^{h_{i2}} + \dots + e^{h_{ik}}} \cdot e^{h_{i1}} \cdot x_{ir} \\ &= \frac{e^{h_{i1}}}{\sum_{a=1}^k e^{h_{ia}}} \cdot x_{ir} \\ &= z_{i1} \cdot x_{ir} \end{aligned}$$

$$\text{所以 } \frac{d \text{Loss}}{d \theta_{r1}} = \begin{cases} \frac{1}{m} \sum_{i=1}^m (-x_{ir} + z_{i1} \cdot x_{ir}) & y^i = 1 \\ \frac{1}{m} \sum_{i=1}^m (0 + z_{i1} \cdot x_{ir}) & y^i \neq 1 \end{cases}$$

④ 把  $\theta_{r1}$  推广到  $\theta_{rq}$ ,  $q \in \{1, 2, \dots, k\}$

$$\frac{d \text{Loss}}{d \theta_{rq}} = \frac{1}{m} \sum_{i=1}^m (x_{ir} (z_{iq} - I_{iq}))$$

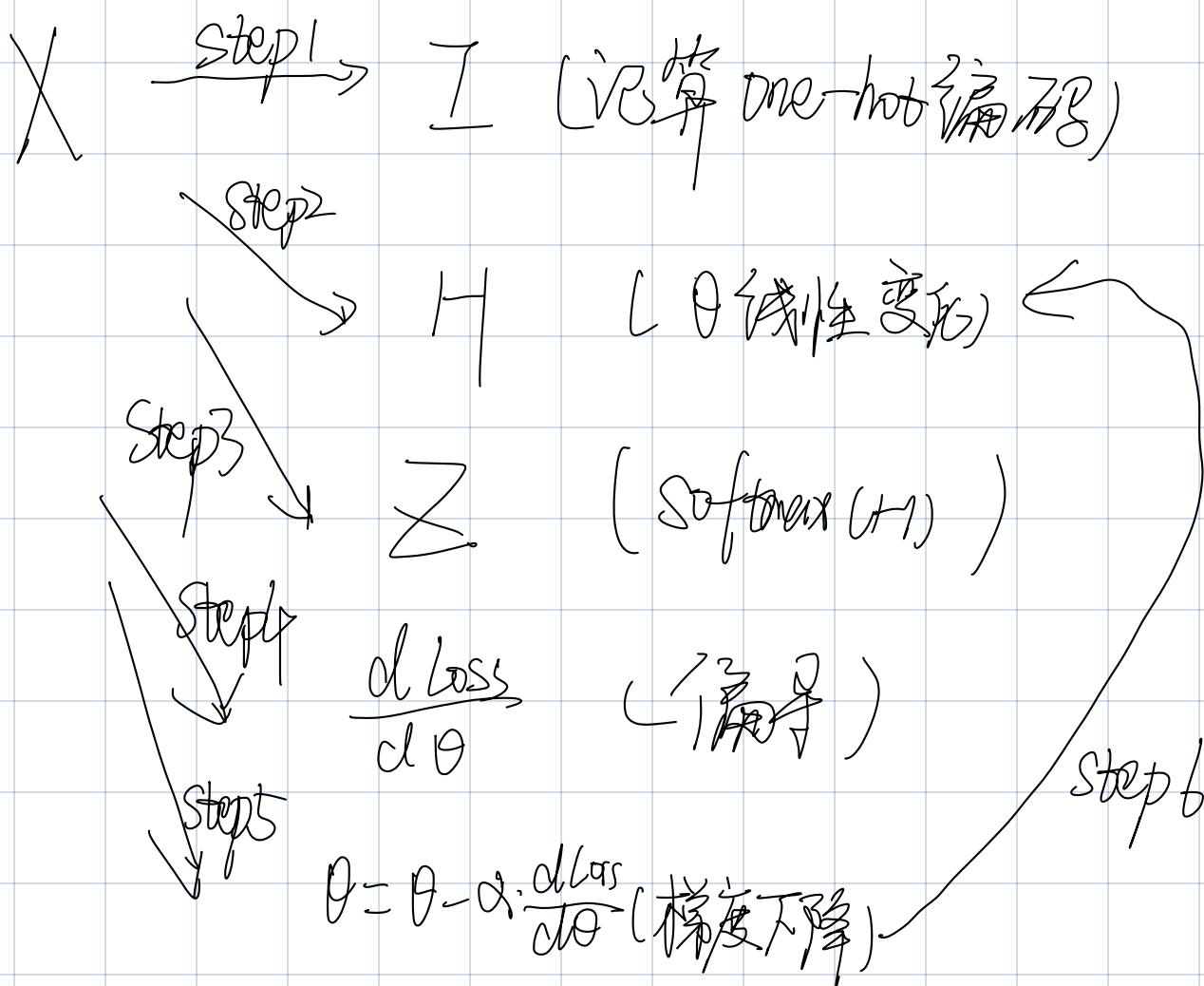
$$\text{其中 } I_{iq} = \begin{cases} 0 & y^i \neq q \\ 1 & y^i = q \end{cases}$$

⑨ 我们再把  $\frac{d \text{loss}}{d \theta}$  整理成矩阵形式

$$\frac{d \text{Loss}}{d \theta} = \frac{1}{n} X^T (Z - I)$$

$n \times k$                        $n \times m$      $m \times k$      $m \times k$

⑩ 现在有了偏导, everything is okay!



六步求解  $\theta$  参数!

