

约定：特征： $X(n, m)$ (n : 批量数 m : 特征数)

标签： $y_s(n, 1)$ 或 $\hat{y}(n, c)$ (c : 类别数)

参数： $\theta(m, 1)$ 或 $\theta(m, c)$

可学习

二分类

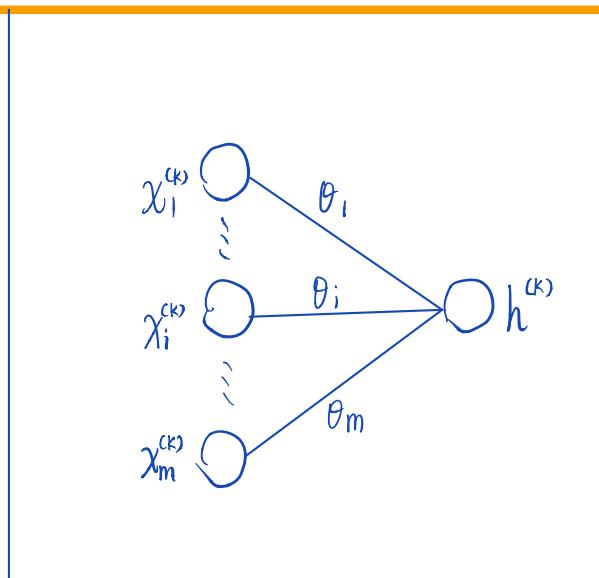
one-hot ("独热" 编码)

多分类

Linear Regression (线性回归)

假设： $\pi\theta$

损失：均方差 (MSE)



$$h^{(k)} = \sum_{i=1}^m x_i^{(k)} \cdot \theta_i$$

$$L^{(k)} = \frac{1}{2} \cdot (h^{(k)} - y^{(k)})^2$$

$$\frac{\partial L^{(k)}}{\partial \theta_i} = \frac{\partial L^{(k)}}{\partial h^{(k)}} \cdot \frac{\partial h^{(k)}}{\partial \theta_i} = (h^{(k)} - y^{(k)}) \cdot x_i^{(k)}$$

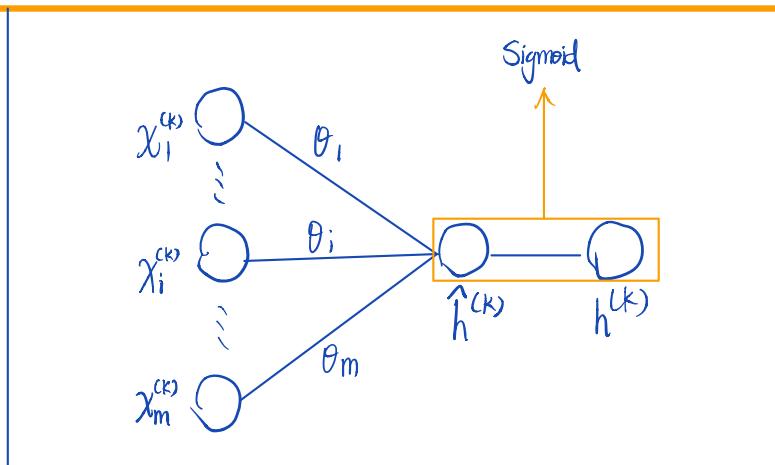
$$= x^T \cdot (h - y)$$

Logistic Regression (逻辑回归) (二分类)

①

假设：Sigmoid ($\sigma(\theta)$)

损失：交叉熵 (CE)



$$\hat{h}^{(k)} = \sum_{i=1}^m x_i^{(k)} \cdot \theta_i$$

$$h^{(k)} = \text{Sigmoid}(\hat{h}^{(k)}) = \frac{1}{1 + e^{-\hat{h}^{(k)}}}$$

$$L^{(k)} = -\log \left(\underbrace{h^{(k)} y^{(k)}}_{y^{(k)} \log h^{(k)}} \cdot \underbrace{(1-h^{(k)})}_{(1-y^{(k)}) \log (1-h^{(k)})} \right) = - \left[y^{(k)} \log h^{(k)} + (1-y^{(k)}) \log (1-h^{(k)}) \right]$$

似然(最大化)

$$\frac{\partial L^{(k)}}{\partial \theta_i} = \frac{\partial L^{(k)}}{\partial h^{(k)}} \cdot \frac{\partial h^{(k)}}{\partial \hat{h}^{(k)}} \cdot \frac{\partial \hat{h}^{(k)}}{\partial \theta_i}$$

$$= - \left(\frac{y^{(k)}}{h^{(k)}} - \frac{1-y^{(k)}}{1-h^{(k)}} \right) \cdot (h^{(k)} \cdot (1-h^{(k)})) \cdot x_i$$

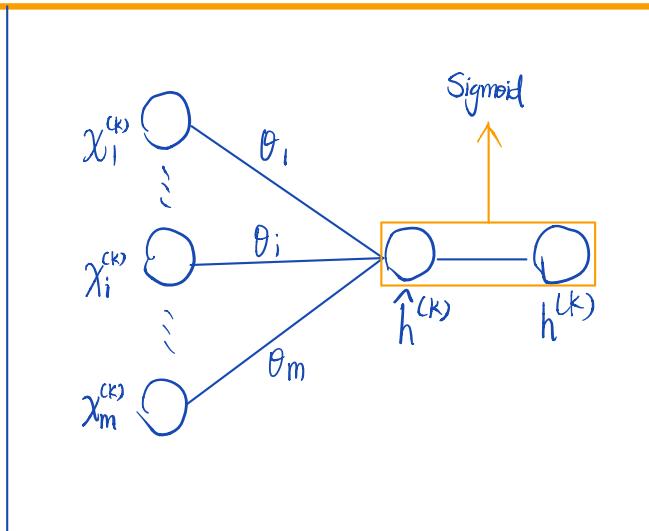
$$= (h^{(k)} - y^{(k)}) \cdot x_i^{(k)}$$

$$= x^T \cdot (h - y)$$

(2)

假设 $\hat{y} = \text{sigmoid}(x\theta)$

损失 \hat{L} 均方差 (mse)



$$\hat{h}^{(k)} = \sum_{i=1}^m x_i^{(k)} \cdot \theta_i$$

$$h^{(k)} = \text{Sigmoid}(\hat{h}^{(k)}) = \frac{1}{1 + e^{-\hat{h}^{(k)}}}$$

$$L^{(k)} = \frac{1}{2} (h^{(k)} - y^{(k)})^2$$

$$\frac{\partial L^{(k)}}{\partial \theta_i} = \frac{\partial L^{(k)}}{\partial h^{(k)}} \cdot \frac{\partial h^{(k)}}{\partial \hat{h}^{(k)}} \cdot \frac{\partial \hat{h}^{(k)}}{\partial \theta_i}$$

$$= (h^{(k)} - y^{(k)}) \cdot (h^{(k)} \cdot (1 - h^{(k)})) \cdot x_i$$

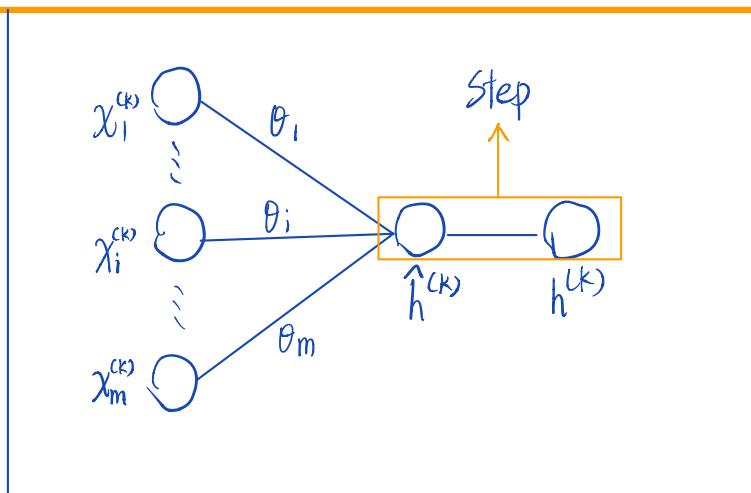
$$= x^T \cdot [(h - y) \times h \times (1 - h)]$$

感知机 (Perceptron) (二分类)

假设流:

$$\begin{cases} 1 & \chi \theta \geq 0 \\ 0 & \chi \theta < 0 \end{cases}$$

损失: $\sum_{v \in M_0} \chi \theta - \sum_{s \in M_1} \chi \theta \quad ("本被判1" - "本被判0")$



$$\hat{h}^{(k)} = \sum_{i=1}^m x_i^{(k)} \cdot \theta_i$$

$$h^{(k)} = \begin{cases} 1 & \sum_{i=1}^m x_i^{(k)} \theta_i \geq 0 \\ 0 & \sum_{i=1}^m x_i^{(k)} \theta_i < 0 \end{cases}$$

$$\begin{aligned} L^{(k)} &= (1 - y^{(k)}) \cdot h^{(k)} \cdot \hat{h}^{(k)} - y^{(k)} \cdot (1 - h^{(k)}) \cdot \hat{h}^{(k)} \\ &= (h^{(k)} - y^{(k)}) \cdot \hat{h}^{(k)} \end{aligned}$$

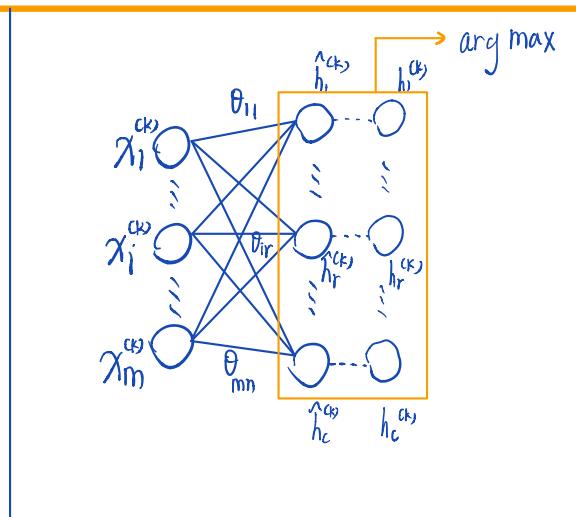
$$\frac{\partial L^{(k)}}{\partial \theta_i} = (h^{(k)} - y^{(k)}) \cdot x_i$$

$$= \chi^T \cdot (h - y)$$

多分类感知机 (Multi-class Perceptron)

假设说: $h(x) = \arg \max_{r=1, \dots, c} \sum_{i=1}^m x_i \theta_{ir}$

损失: $L = \max_{r=1, \dots, c} \sum_{i=1}^m x_i \theta_{ir} - \sum_{i=1}^m x_i \theta_{ic^*}$



$$\hat{h}_r^{(k)} = \sum_{i=1}^m x_i^{(k)} \theta_{ir}$$

$$h_r^{(k)} = \begin{cases} \hat{h}_r^{(k)} & r = \arg \max_{d=1, \dots, c} \hat{h}_d^{(k)} \\ 0 & \text{其它} \end{cases}$$

$$L^{(k)} = h_r^{(k)} - \hat{h}_{c^*}^{(k)} \quad (\text{注: } y_{c^*}^{(k)} = 1, \text{ 其它 } y_{\neq c^*}^{(k)} = 0)$$

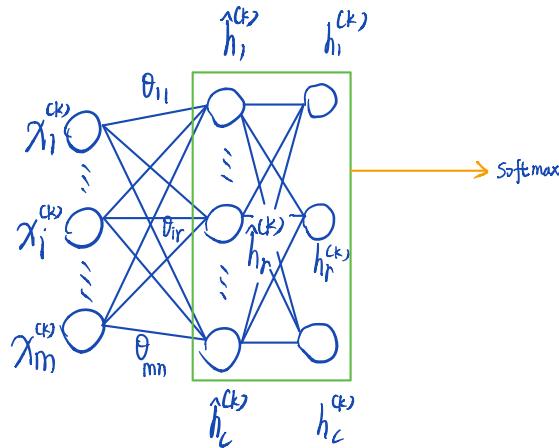
$$\frac{\partial L^{(k)}}{\partial \theta_{ir}} = (1\{r=r^*\} - 1\{r=c^*\}) \cdot x_i^{(k)}$$

$$= x^T \cdot (h - y)$$

Softmax (多分类)

假设: $h_\theta(x) = \text{softmax}(\mathbf{x}\theta)$

损失: 交叉熵 (CE)



方式一

$$\hat{h}_r^{(k)} = \sum_{i=1}^m x_i^{(k)} \theta_{ir}$$

$$h_r^{(k)} = \text{softmax}(\hat{h}_r^{(k)}) = \frac{e^{\hat{h}_r^{(k)}}}{\sum_{f=1}^c e^{\hat{h}_f^{(k)}}}$$

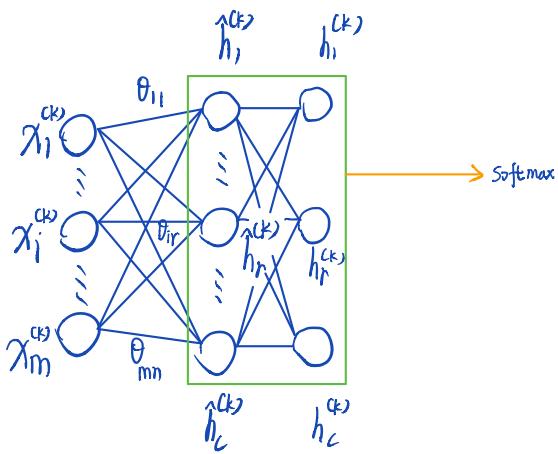
$$L^{(k)} = -\log \left(\prod_{r=1}^c h_r^{(k)} y_r^{(k)} \right) = -\sum_{r=1}^c y_r^{(k)} \log h_r^{(k)} = -y_{c^*}^{(k)} \log h_{c^*}^{(k)}$$

$$\frac{\partial L^{(k)}}{\partial \theta_{ir}} = \frac{\partial L^{(k)}}{\partial h_{c^*}^{(k)}} \cdot \frac{\partial h_{c^*}^{(k)}}{\partial \hat{h}_r^{(k)}} \cdot \frac{\partial \hat{h}_r^{(k)}}{\partial \theta_{ir}}$$

$$= \begin{cases} r=c^* \text{ 时} & \left(-\frac{\sum_{f=1}^c e^{\hat{h}_f^{(k)}}}{e^{\hat{h}_r^{(k)}}} \right) \cdot \left(\frac{e^{\hat{h}_r^{(k)}} (\sum_{f=1}^c e^{\hat{h}_f^{(k)}}) - e^{\hat{h}_r^{(k)}} \cdot e^{\hat{h}_r^{(k)}}}{(\sum_{f=1}^c e^{\hat{h}_f^{(k)}})^2} \right) \cdot x_i^{(k)} = (h_r^{(k)} - 1) \cdot x_i^{(k)} \\ r \neq c^* \text{ 时} & \left(-\frac{\sum_{f=1}^c e^{\hat{h}_f^{(k)}}}{e^{\hat{h}_r^{(k)}}} \right) \cdot \left(\frac{0 - e^{\hat{h}_r^{(k)}} \cdot e^{\hat{h}_r^{(k)}}}{(\sum_{f=1}^c e^{\hat{h}_f^{(k)}})^2} \right) \cdot x_i^{(k)} = (h_r^{(k)} - 0) \cdot x_i^{(k)} \end{cases}$$

$$= (h_r^{(k)} - y_r^{(k)}) \cdot x_i^{(k)}$$

$$= \mathbf{x}^T \cdot (\mathbf{h} - \mathbf{y})$$



方 式 二

$$\hat{h}_r^{(k)} = \sum_{i=1}^m x_i^{(k)} \theta_{ir}$$

$$h_r^{(k)} = \text{softmax}(\hat{h}_r^{(k)}) = \frac{e^{\hat{h}_r^{(k)}}}{\sum_{i=1}^c e^{\hat{h}_i^{(k)}}}$$

$$L^{(k)} = -\log \left(\prod_{r=1}^c p_r^{y_r^{(k)}} \right) = -\sum_{r=1}^c y_r^{(k)} \log h_r^{(k)}$$

$$\frac{\partial L^{(k)}}{\partial \theta_{ir}} = \sum_{d=1}^c \frac{\partial L^{(k)}}{\partial h_d^{(k)}} \cdot \frac{\partial h_d^{(k)}}{\partial \hat{h}_r^{(k)}} \cdot \frac{\partial \hat{h}_r^{(k)}}{\partial \theta_{ir}}$$

$$= \left(-\frac{y_r^{(k)} \left(\sum_{i=1}^c e^{\hat{h}_i^{(k)}} \right)}{e^{\hat{h}_r^{(k)}}} \right) \cdot \left(\frac{e^{\hat{h}_r^{(k)}} \left(\sum_{i=1}^c e^{\hat{h}_i^{(k)}} \right) - e^{\hat{h}_r^{(k)}} \cdot e^{\hat{h}_r^{(k)}}}{\left(\sum_{i=1}^c e^{\hat{h}_i^{(k)}} \right)^2} \right) \cdot x_i^{(k)} + \sum_{d \neq r}^c \left(-\frac{y_d^{(k)} \left(\sum_{i=1}^c e^{\hat{h}_i^{(k)}} \right)}{e^{\hat{h}_d^{(k)}}} \right) \cdot \left(\frac{0 - e^{\hat{h}_d^{(k)}} \cdot e^{\hat{h}_r^{(k)}}}{\left(\sum_{i=1}^c e^{\hat{h}_i^{(k)}} \right)^2} \right) \cdot x_i^{(k)}$$

$$= \left(-y_r^{(k)} (1 - h_r^{(k)}) \cdot z_s^{(k)} \right) + \sum_{d \neq r}^c \left(-y_d^{(k)} (0 - h_r^{(k)}) \cdot z_s^{(k)} \right)$$

$$= (h_r^{(k)} - y_r^{(k)}) \cdot z_s^{(k)} \longrightarrow \text{注} \begin{cases} ① \text{如果 } y_r^{(k)} = 1, \text{ 则 } y_{d \neq r}^{(k)} \text{ 全 } = 0, \text{ 第二项化掉} \\ ② \text{如果 } y_r^{(k)} = 0, \text{ 则 } y_{d \neq r}^{(k)} \text{ 中有 } 1, \text{ 第一项化掉, 第二项累加只保留一项} \end{cases}$$

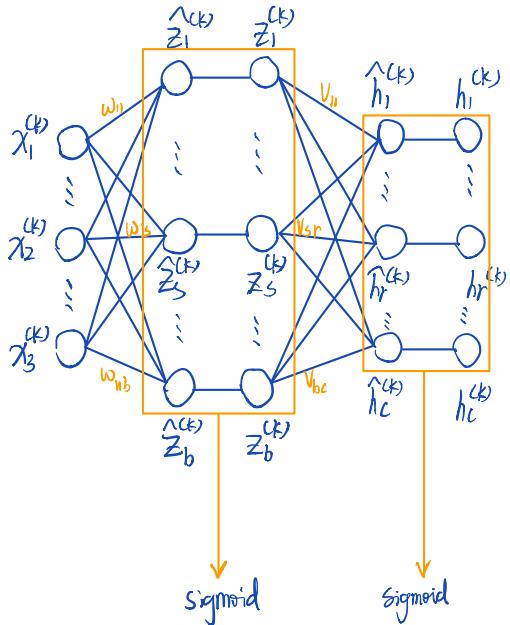
$$= x^T \cdot (h - y)$$

人口神经网络(ANN)、多层感知机(MLP)

①

假设说: $xw \rightarrow \text{Sigmoid} \rightarrow zv \rightarrow \text{Sigmoid}$

损失: mean square error (mse)



$$\begin{aligned}
 \hat{z}_s^{(k)} &= \sum_{i=1}^m x_i^{(k)} \cdot w_{is} \\
 \hat{z}_s^{(k)} &= \text{Sigmoid}(\hat{z}_s^{(k)}) = \frac{1}{1 + e^{-\hat{z}_s^{(k)}}} \\
 \hat{h}_r^{(k)} &= \sum_{s=1}^b \hat{z}_s^{(k)} \cdot v_{sr} \\
 \hat{h}_r^{(k)} &= \text{Sigmoid}(\hat{h}_r^{(k)}) = \frac{1}{1 + e^{-\hat{h}_r^{(k)}}} \\
 L^{(k)} &= \frac{1}{2} \sum_{r=1}^c (\hat{h}_r^{(k)} - y_r^{(k)})^2 \quad \text{仅有 r 动力, 其余全加口} \\
 \frac{\partial L^{(k)}}{\partial v_{sr}} &= \frac{\partial L^{(k)}}{\partial \hat{h}_r^{(k)}} \cdot \frac{\partial \hat{h}_r^{(k)}}{\partial \hat{h}_r^{(k)}} \cdot \frac{\partial \hat{h}_r^{(k)}}{\partial v_{sr}} \\
 &= (\hat{h}_r^{(k)} - y_r^{(k)}) \cdot (\hat{h}_r^{(k)} \cdot (1 - \hat{h}_r^{(k)})) \cdot \hat{z}_s^{(k)} \quad \rightarrow \text{OutputLayer } L_r^{(k)} \\
 &= \hat{z}^T \cdot [(h - y) \times (h \times (1 - h))] \\
 \end{aligned}$$

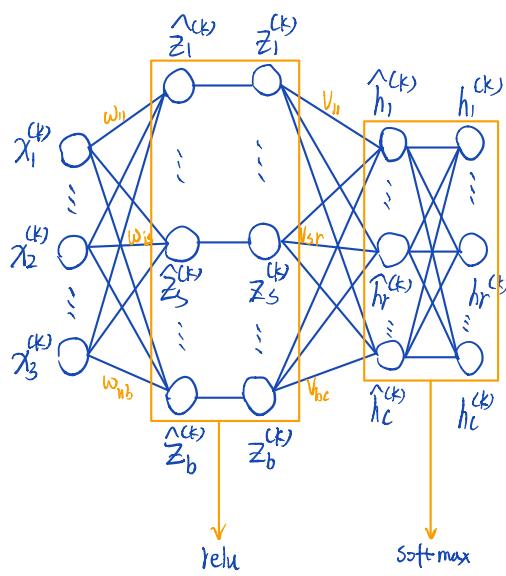
$$\begin{aligned}
 \frac{\partial L^{(k)}}{\partial w_{is}} &= \sum_{r=1}^c \frac{\partial L^{(k)}}{\partial \hat{h}_r^{(k)}} \cdot \frac{\partial \hat{h}_r^{(k)}}{\partial \hat{h}_r^{(k)}} \cdot \frac{\partial \hat{h}_r^{(k)}}{\partial \hat{z}_s^{(k)}} \cdot \frac{\partial \hat{z}_s^{(k)}}{\partial \hat{z}_s^{(k)}} \cdot \frac{\partial \hat{z}_s^{(k)}}{\partial w_{is}} \\
 &= \sum_{r=1}^c \text{OutputLayer } L_r^{(k)} \cdot \frac{\partial \hat{h}_r^{(k)}}{\partial \hat{z}_s^{(k)}} \cdot \frac{\partial \hat{z}_s^{(k)}}{\partial \hat{z}_s^{(k)}} \cdot \frac{\partial \hat{z}_s^{(k)}}{\partial w_{is}} \\
 &= \sum_{r=1}^c \underbrace{\text{OutputLayer } L_r^{(k)} \cdot v_{sr} \cdot (\hat{z}_s^{(k)} \cdot (1 - \hat{z}_s^{(k)})) \cdot x_i^{(k)}}_{(n, c) \quad (h, c) \quad (n, h) \quad (n, m)} \quad \rightarrow \text{HiddenLayer } L_s^{(k)}
 \end{aligned}$$

$$= X^T \cdot \left(\left((h - y) \times h \times (1 - h) \right) \cdot V^T \right) \times Z \times (1 - Z)$$

2.1

假设说: $xw \rightarrow \text{relu} \rightarrow zv \rightarrow \text{softmax}$

损失 = cross entropy



$$\begin{aligned}\hat{z}_s^{(k)} &= \sum_{i=1}^m x_i^{(k)} \cdot w_{is} \\ \hat{z}_s^{(k)} &= \text{relu}(\hat{z}_s^{(k)}) = \hat{z}_s^{(k)} \times \text{mask}_s \\ \hat{h}_r^{(k)} &= \sum_{s=1}^b \hat{z}_s^{(k)} \cdot v_{sr} \\ h_r^{(k)} &= \text{softmax}(\hat{h}_r^{(k)}) = \frac{e^{\hat{h}_r^{(k)}}}{\sum_{f=1}^c e^{\hat{h}_f^{(k)}}} \\ L^{(k)} &= -\log \left(\prod_{r=1}^c h_r^{(k)} \right) = -\sum_{r=1}^c y_r^{(k)} \log h_r^{(k)} = -y_{r^*}^{(k)} \log h_{r^*}^{(k)} \quad (y_{r^*}^{(k)} = 1, y_{r \neq r^*}^{(k)} \text{ 全=0})\end{aligned}$$

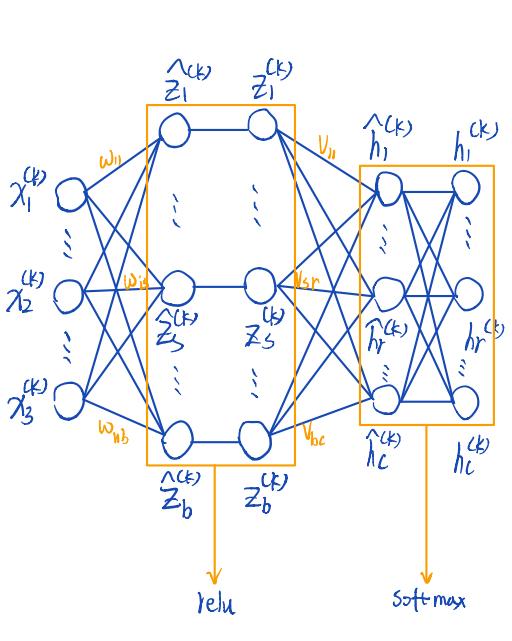
$$\begin{aligned}\frac{\partial L^{(k)}}{\partial v_{sr}} &= \frac{\partial L^{(k)}}{\partial h_{r^*}^{(k)}} \cdot \frac{\partial h_{r^*}^{(k)}}{\partial \hat{h}_r^{(k)}} \cdot \frac{\partial \hat{h}_r^{(k)}}{\partial v_{sr}} \\ &\xrightarrow{r=c^*} -\left[\left(\frac{\left(\sum_{f=1}^c e^{\hat{h}_f^{(k)}} \right)}{e^{\hat{h}_{r^*}^{(k)}}} \right) \cdot \left(\frac{e^{\hat{h}_{r^*}^{(k)}} \cdot \left(\sum_{f=1}^c e^{\hat{h}_f^{(k)}} \right) - e^{\hat{h}_{r^*}^{(k)}} \cdot e^{\hat{h}_{r^*}^{(k)}}}{\left(\sum_{f=1}^c e^{\hat{h}_f^{(k)}} \right)^2} \right) \cdot \hat{z}_s^{(k)} \right] \\ &\quad \xrightarrow{r \neq c^*} -\left[\left(\frac{\left(\sum_{f=1}^c e^{\hat{h}_f^{(k)}} \right)}{e^{\hat{h}_{r^*}^{(k)}}} \right) \cdot \left(\frac{0 - e^{\hat{h}_{r^*}^{(k)}} \cdot e^{\hat{h}_{r^*}^{(k)}}}{\left(\sum_{f=1}^c e^{\hat{h}_f^{(k)}} \right)^2} \right) \cdot \hat{z}_s^{(k)} \right] \\ &= -(1 - h_{r^*}^{(k)}) \cdot \hat{z}_s^{(k)} \\ &\quad \xrightarrow{\text{OutputLayer } L_r^{(k)}} \underbrace{(h_{r^*}^{(k)} - y_{r^*}^{(k)}) \cdot \hat{z}_s^{(k)}}_{=} \\ &= \mathbf{x}^T \cdot (\mathbf{h} - \mathbf{y})\end{aligned}$$

$$\begin{aligned}\frac{\partial L^{(k)}}{\partial w_{is}} &= \sum_{r=1}^c \frac{\partial L^{(k)}}{\partial h_r^{(k)}} \cdot \frac{\partial h_r^{(k)}}{\partial \hat{h}_r^{(k)}} \cdot \frac{\partial \hat{h}_r^{(k)}}{\partial \hat{z}_s^{(k)}} \cdot \frac{\partial \hat{z}_s^{(k)}}{\partial z_s^{(k)}} \cdot \frac{\partial z_s^{(k)}}{\partial \hat{z}_s^{(k)}} \cdot \frac{\partial \hat{z}_s^{(k)}}{\partial w_{is}} \\ &= \sum_{r=1}^c \text{OutputLayer } L_r^{(k)} \cdot \frac{\partial h_r^{(k)}}{\partial \hat{z}_s^{(k)}} \cdot \frac{\partial \hat{z}_s^{(k)}}{\partial \hat{z}_s^{(k)}} \cdot \frac{\partial \hat{z}_s^{(k)}}{\partial w_{is}} \\ &= \sum_{r=1}^c \text{OutputLayer } L_r^{(k)} \cdot v_{sr} \cdot \text{mask}_s \cdot x_i^{(k)} \\ &\quad \xrightarrow{\text{HiddenLayer } L_s^{(k)}} \underbrace{\mathbf{x}^T \cdot ((\mathbf{h} - \mathbf{y}) \cdot \mathbf{V}^T) \times \text{mask}}_{=}\end{aligned}$$

2.2

假设说: $xw \rightarrow \text{relu} \rightarrow zv \rightarrow \text{softmax}$

损失 = cross entropy



$$\begin{aligned}\hat{z}_s^{(k)} &= \sum_{i=1}^m x_i^{(k)} \cdot w_{is} \\ \hat{z}_s^{(k)} &= \text{relu}(\hat{z}_s^{(k)}) = \hat{z}_s^{(k)} \times \text{mask}_s \\ \hat{h}_r^{(k)} &= \sum_{s=1}^b \hat{z}_s^{(k)} \cdot v_{sr} \\ \hat{h}_r^{(k)} &= \text{softmax}(\hat{h}_r^{(k)}) = \frac{e^{\hat{h}_r^{(k)}}}{\sum_{f=1}^c e^{\hat{h}_f^{(k)}}} \\ L^{(k)} &= -\log \left(\prod_{r=1}^c y_r^{(k)} \hat{h}_r^{(k)} \right) = -\sum_{r=1}^c y_r^{(k)} \log \hat{h}_r^{(k)}\end{aligned}$$

做一个归一化，其余全为0 (one-hot)

$$\begin{aligned}\frac{\partial L^{(k)}}{\partial v_{sr}} &= \sum_{d=1}^c \frac{\partial L^{(k)}}{\partial h_d^{(k)}} \cdot \frac{\partial h_d^{(k)}}{\partial \hat{h}_r^{(k)}} \cdot \frac{\partial \hat{h}_r^{(k)}}{\partial v_{sr}} \\ &= - \left[\left(\frac{y_r^{(k)} \cdot \left(\sum_{f=1}^c e^{\hat{h}_f^{(k)}} \right)}{e^{\hat{h}_r^{(k)}}} \right) \cdot \left(\frac{e^{\hat{h}_r^{(k)}} \cdot \left(\sum_{f=1}^c e^{\hat{h}_f^{(k)}} \right) - e^{\hat{h}_r^{(k)}} \cdot e^{\hat{h}_r^{(k)}}}{\left(\sum_{f=1}^c e^{\hat{h}_f^{(k)}} \right)^2} \cdot \hat{z}_s^{(k)} \right) \right] - \sum_{d \neq r} \left[\left(\frac{y_d^{(k)} \cdot \left(\sum_{f=1}^c e^{\hat{h}_f^{(k)}} \right)}{e^{\hat{h}_d^{(k)}}} \right) \cdot \left(\frac{0 - e^{\hat{h}_d^{(k)}} \cdot e^{\hat{h}_r^{(k)}}}{\left(\sum_{f=1}^c e^{\hat{h}_f^{(k)}} \right)^2} \cdot \hat{z}_s^{(k)} \right) \right] \\ &= - \left[y_r^{(k)} (1 - \hat{h}_r^{(k)}) \cdot \hat{z}_s^{(k)} \right] - \sum_{d \neq r} \left[y_d^{(k)} (0 - \hat{h}_d^{(k)}) \cdot \hat{z}_s^{(k)} \right]\end{aligned}$$

注
 ①如果 $y_r^{(k)} = 1$, 则 $y_{d \neq r}^{(k)}$ 全 = 0, 第二项化掉
 ②如果 $y_r^{(k)} = 0$, 则 $y_{d \neq r}^{(k)}$ 中将有一个 = 1, 第一项化掉, 第二项里加中只会保留一项

$$\begin{aligned}&= (\hat{h}_r^{(k)} - y_r^{(k)}) \cdot \hat{z}_s^{(k)} \\ &\quad \xrightarrow{\text{OutputLayer } L_r^{(k)}}\end{aligned}$$

$$\begin{aligned}\frac{\partial L^{(k)}}{\partial w_{is}} &= \sum_{r=1}^c \frac{\partial L^{(k)}}{\partial h_r^{(k)}} \cdot \frac{\partial h_r^{(k)}}{\partial \hat{h}_r^{(k)}} \cdot \frac{\partial \hat{h}_r^{(k)}}{\partial \hat{z}_s^{(k)}} \cdot \frac{\partial \hat{z}_s^{(k)}}{\partial \hat{z}_s^{(k)}} \cdot \frac{\partial \hat{z}_s^{(k)}}{\partial w_{is}} \\ &= \sum_{r=1}^c \text{OutputLayer } L_r^{(k)} \cdot \frac{\partial \hat{h}_r^{(k)}}{\partial \hat{z}_s^{(k)}} \cdot \frac{\partial \hat{z}_s^{(k)}}{\partial \hat{z}_s^{(k)}} \cdot \frac{\partial \hat{z}_s^{(k)}}{\partial w_{is}} \\ &= \sum_{r=1}^c \text{OutputLayer } L_r^{(k)} \cdot v_{sr} \cdot \text{mask}_s \cdot x_i^{(k)} \\ &\quad \xrightarrow{\text{HiddenLayer } L_s^{(k)}} \\ &= x^T \cdot ((h - y) \cdot v^T) \times \text{mask}\end{aligned}$$

① 好坏的初值非常重要的，一开始直接 $\text{np.random.randn}(x, x)$ ，太大会使 L_2 正则化效果显著，但 loss 太高而无法训练

建议先初始化为 $0.01 * \text{np.random.randn}(x, x)$ ，网络此用效果会好很多。
 L_2 效果不明显，且此时损失 (Loss) 可能会上升，因此 之后损失让权值上升大于 L_2 损失

② drop out 完成细节

(1) 每个神经元的输出： $p \downarrow x + (1-p) \cdot 0$

$\begin{cases} \text{法一：} & \left\{ \begin{array}{l} \text{训练时补偿：乘以 } \frac{1}{p} \\ \text{测试时不用除} \end{array} \right. \\ \text{法二：} & \left\{ \begin{array}{l} \text{训练时不用乘} \\ \text{测试时补偿：乘以 } p \end{array} \right. \end{cases}$

(2) $\begin{cases} \text{随机剔除神经元：} & \text{drop_index} = \text{np.random.rand}(x.shape[1]) \quad \text{drop_index} = (\text{drop_index} > \text{ratio}) \quad \text{dropout.mask[:, drop_index]} = 0 \\ \text{随机样本剔除随机特征：} & \text{dropout.mask} = (\text{np.random.rand}(x.shape) < \text{ratio}) \\ \text{(经测试，更优)} & \end{cases}$