

(1) Regresión

(1)



* Modelo:

$$Y_i | \bar{X}_i = \bar{x}_i \sim \mathcal{N}(\mu_i, 1)$$

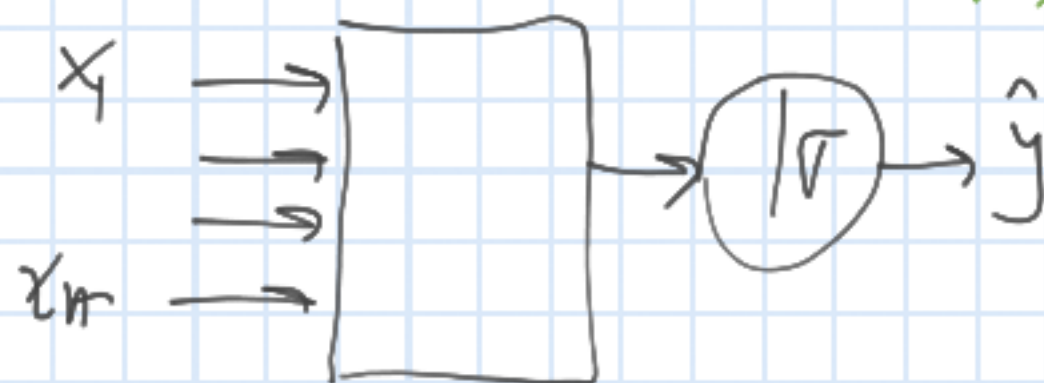
$$\mu_i = f(\bar{x}_i, \theta)$$

* Costo: $J_{\theta} = \frac{1}{n} \sum_{i=1}^n (y_i - f(\bar{x}_i, \theta))^2$

* Optimizador
↳ Mini-Batch

(2) Clasificación Bin

(2)



* Modelo: $Y_i | \bar{X}_i = \bar{x}_i \sim \text{Ber}(p_i)$
 $p_i = \sigma(f(\bar{x}_i, \theta))$

* Aplicaciones MV

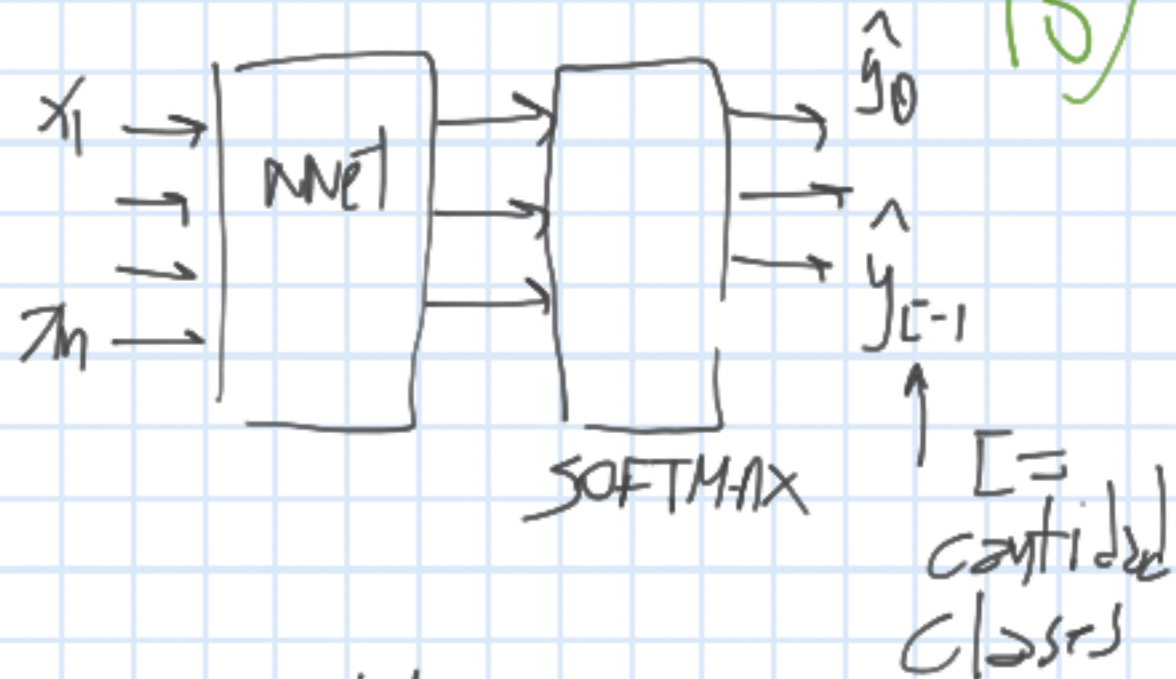
$$J_{\theta} = -\frac{1}{n} \sum_{i=1}^n y_i \log(p_i) + (1 - y_i) \log(1 - p_i)$$

Labels z'_0

* Optimizador → SGD / Mini-batch

(3) Clasificación MultiClase

(3)



* Modelo

$$Y_i | \bar{X}_i = \bar{x}_i \sim \text{Multinomial}$$

$$P(Y_i = 0 | \bar{x}_i) = p_0$$

$$P(Y_i = 1 | \bar{x}_i) = p_1$$

$$P(Y_{[c-1]} = 1 | \bar{x}_i) = p_{[c-1]}$$

¿Costo?

(2) Probar que optimizar $D_{KL}(y_i || \hat{y}_i) \equiv$ optimizar MV

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n D_{KL}(y_i || \hat{y}_i) \quad (\text{clasificación binaria})$$

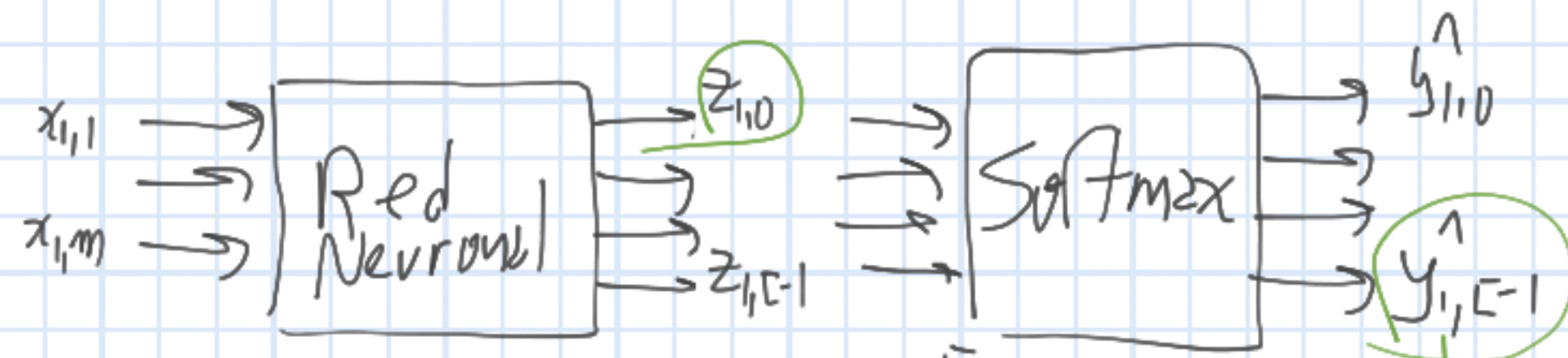
$$\frac{1}{n} \sum_{i=1}^n y_i \log\left(\frac{y_i}{p_i}\right) + (1-y_i) \log\left(\frac{1-y_i}{1-p_i}\right) \rightarrow \text{Depende de } \theta$$

$$\frac{1}{n} \sum_{i=1}^n \left(-y_i \ln(p_i) - (1-y_i) \ln(1-p_i) \right) + \left(y_i \ln(y_i) + (1-y_i) \ln(1-y_i) \right)$$

cte para
el problema
de optimización

$$D_{KL}(y_i || \hat{y}_i) = \underline{H(y_i, \hat{y}_i)} - \cancel{H(y_i)}$$

(3) MV aplicada al modelo multinomial



One hot encoding

VS
$$\begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}^y$$

$$\begin{array}{c|c|c} x_{1,1} & \dots & x_{1,m} \\ \vdots & & \vdots \\ x_{m,1} & \dots & x_{m,m} \end{array} \left| \begin{array}{l} y_1 = \{0, \dots, E-1\} \\ \vdots \\ y_m = \{0, \dots, E-1\} \end{array} \right| \begin{array}{c} \underbrace{[0 \ 0 \ 1 \ \dots \ 0]}_{E-1} \\ [1 \ 0 \ 0 \ \dots \ 0] \\ \vdots \\ [0 \ 0 \ 1 \ \dots \ 1] \end{array}$$

¿Que pasa cuando tengo muchas clases?
 $\sum_{k=0}^{E-1} e^{z_k}$ → Tarda mucho en calcularse! \gg
 ↳ Negative sampling

$$g(z_{1,0}) = \frac{e^{z_{1,0}}}{\sum_{k=0}^{E-1} e^{z_{1,k}}}$$

$$g(z_{1,E-1}) = \frac{e^{z_{1,E-1}}}{\sum_{k=0}^{E-1} e^{z_{1,k}}}$$

sum() = 1

¿Cómo queda la verosimilitud?

arg min_w
$$-\frac{1}{n} \sum_{i=1}^n \sum_{k=0}^{L-1} \delta_{k, y_i} \log \left(\frac{e^{z_k}}{\sum_{w=0}^{L-1} e^{z_w}} \right)$$

$$\delta_{k, y_i} = \begin{cases} 1 & y_i = k \\ 0 & y_i \neq k \end{cases}$$

J_{θ}

Cross Entropy