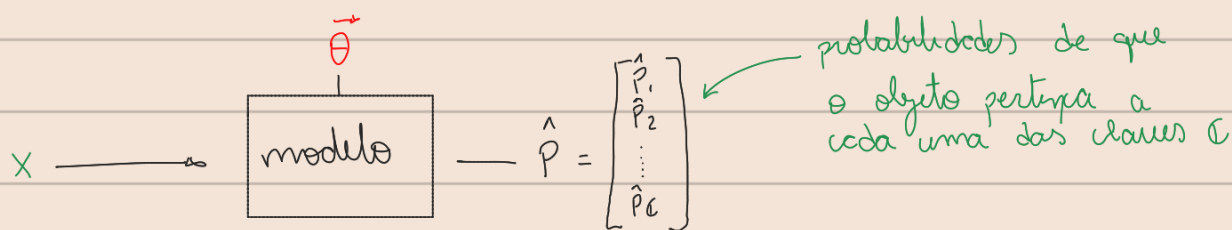
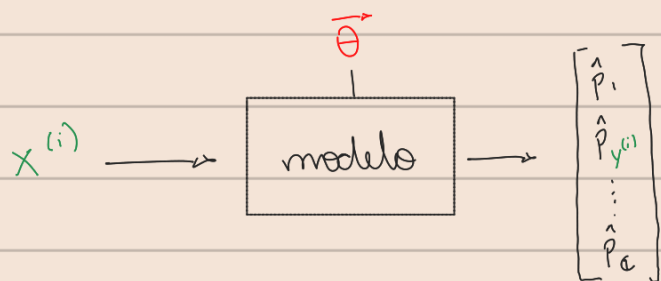


# Treinando um classificador



Para um exemplo  $(x^{(i)}, y^{(i)})$  qual é a chance, de acordo com o modelo, de que o objeto de features  $x^{(i)}$  pertença a classe  $y^{(i)}$ ?



$$\hat{p}_{y^{(i)}} = \text{Prob}(Y = y^{(i)} \mid X = x^{(i)}, \vec{\theta})$$

um modelo ruim **discorda** da realidade:

$$\rightarrow \text{Prob}(Y = y^{(i)} \mid X = x^{(i)}, \vec{\theta}_{\text{RUIM}})$$

vai ser **baixa**

um modelo bom **concorda** com a realidade:

$$\rightarrow \text{Prob}(Y = y^{(i)} \mid X = x^{(i)}, \vec{\theta}_{\text{BOM}})$$

vai ser **alta**

IDEIA: Achar parâmetros  $\vec{\theta}$  que maximizam a OPINIÃO DO MODELO acerca dos observados

$$\vec{\theta}_{ML} = \underset{\vec{\theta}}{\text{argmax}} \underbrace{\text{Prob} \left[ \begin{array}{c|c} \text{targets do exemplo de treino} & \text{features do exemplo de treino} \end{array} \mid \vec{\theta} \right]}_{\text{verossimilhança (likelihood)}} = \text{Prob}(\text{Dados} \mid \vec{\theta})$$

Construindo a formula de

$$\text{Prob} [\text{targets} \mid \text{features}, \theta]$$

• Premissas

1. os  $(x^{(i)}, y^{(i)})$  vem da mesma distribuição  $P(x, y)$  (iid)
2.  $(x^{(i)}, y^{(i)})$  não independentes

$$\begin{aligned} & \text{Prob} [Y_1 = y^{(1)} \text{ e } Y_2 = y^{(2)} \text{ e } Y_m = y^{(m)} \mid X_1 = \vec{x}^{(1)} \text{ e } X_2 = \vec{x}^{(2)} \text{ e } X_m = \vec{x}^{(m)} \text{ e } \theta] \\ &= \text{Prob} [Y = y^{(1)} \mid X = \vec{x}^{(1)}, \vec{\theta}] \cdot \text{Prob} [Y = y^{(2)} \mid X = \vec{x}^{(2)}, \vec{\theta}] \cdot \text{Prob} [Y = y^{(m)} \mid X = \vec{x}^{(m)}, \vec{\theta}] \\ &= \text{Prob} [\text{targets} \mid \text{features}, \theta] \\ &= \prod_{i=1}^m \hat{p}_{y^{(i)} \mid x^{(i)}, \vec{\theta}} \end{aligned}$$

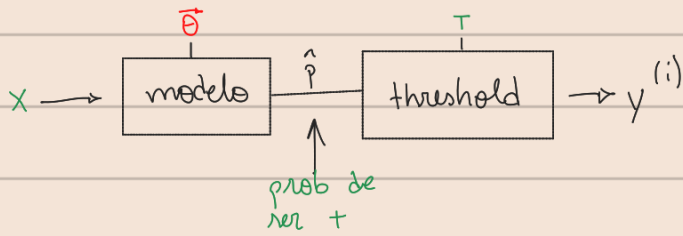
$$\begin{aligned} \vec{\theta}_{ML} &= \underset{\vec{\theta}}{\text{argmin}} \text{Prob} [\text{targets} \mid \text{features}, \theta] \\ &= \underset{\vec{\theta}}{\text{argmax}} \log \text{Prob} [\text{targets} \mid \text{features}, \theta] \\ &= \sum_{i=1}^m \log \hat{p}_{y^{(i)}}^{(i)} \\ &= \sum_{i=1}^m \sum_{k=1}^C \underbrace{[y^{(i)} == k]}_{\text{expressão booleana}} \log (\hat{p}_k^{(i)}(x^{(i)}, \theta)) \\ &\quad \underbrace{\hspace{10em}}_{\text{Notação de Iverson}} \end{aligned}$$

## Entropia cruzada

maximum likelihood:

$$\begin{aligned} \vec{\theta}^{\text{ótimo}} &= \underset{\vec{\theta}}{\text{argmax}} \left\{ \sum_{i=1}^m \sum_{k=1}^C [y^{(i)} == k] \log (\hat{p}_k^{(i)}(x^{(i)}, \theta)) \right\} \\ &\quad \Updownarrow \\ \vec{\theta}^{\text{ótimo}} &= \underset{\vec{\theta}}{\text{argmin}} \left\{ -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^C [y^{(i)} == k] \log (\hat{p}_k^{(i)}(x^{(i)}, \theta)) \right\} \quad * \text{cross entropy} \end{aligned}$$

## Exemplo: classificador linear



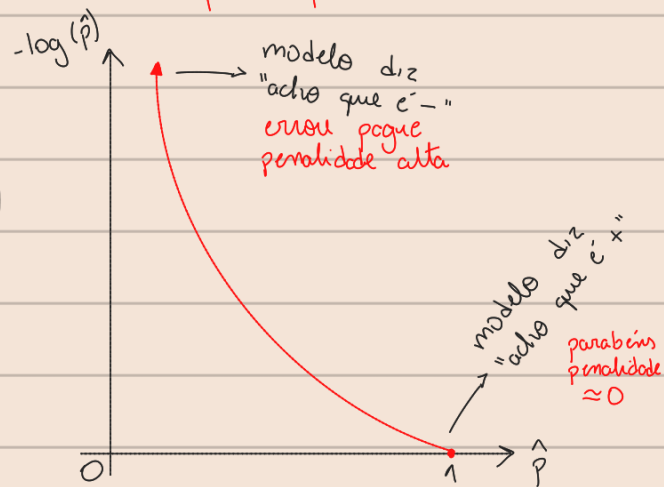
$$\begin{aligned} \text{Loss} &= -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K [\underline{y}^{(i)} = k] \log(\hat{p}_k^{(i)}) \\ &= -\frac{1}{m} \sum_{i=1}^m [\underline{y}^{(i)} = 0] \log(\hat{p}_0^{(i)}) \\ &\quad + -\frac{1}{m} \sum_{i=1}^m [\underline{y}^{(i)} = 1] \log(\hat{p}_1^{(i)}) \end{aligned}$$

$$= -\frac{1}{m} [(1 - y^{(i)}) \cdot \log(1 - \hat{p}^{(i)}) + y^{(i)} \cdot \log(\hat{p}^{(i)})]$$

$$= \frac{1}{m} \sum_{i=1}^m \underbrace{[(1 - y^{(i)}) (-\log(1 - \hat{p}^{(i)})) + y^{(i)} (-\log(\hat{p}^{(i)}))]}_{\text{perda por elemento}}$$

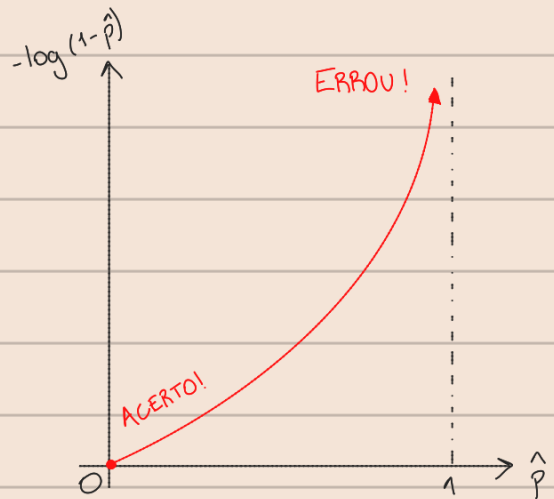
Caso 1:  $y^{(i)} = 1$

$$\begin{aligned} &(1 - y^{(i)}) (-\log(1 - \hat{p}^{(i)})) + y^{(i)} (-\log(\hat{p}^{(i)})) \\ &= (1 - 1) (-\log(1 - \hat{p}^{(i)})) + 1 \cdot (-\log(\hat{p}^{(i)})) \\ &= -\log(\hat{p}^{(i)}) \end{aligned}$$

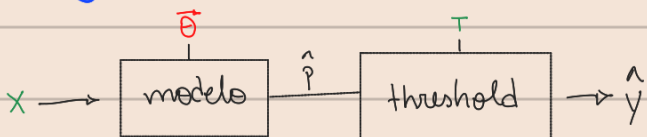


Caso 2:  $y^{(i)} = 0$

$$\begin{aligned} &(1 - y^{(i)}) (-\log(1 - \hat{p}^{(i)})) + y^{(i)} (-\log(\hat{p}^{(i)})) \\ &= (1 - 0) (-\log(1 - \hat{p}^{(i)})) + 0 \cdot (-\log(\hat{p}^{(i)})) \\ &= -\log(1 - \hat{p}^{(i)}) \end{aligned}$$



## Regressão Logística



$$\hat{p} = \sigma(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n)$$

função sigmoide

