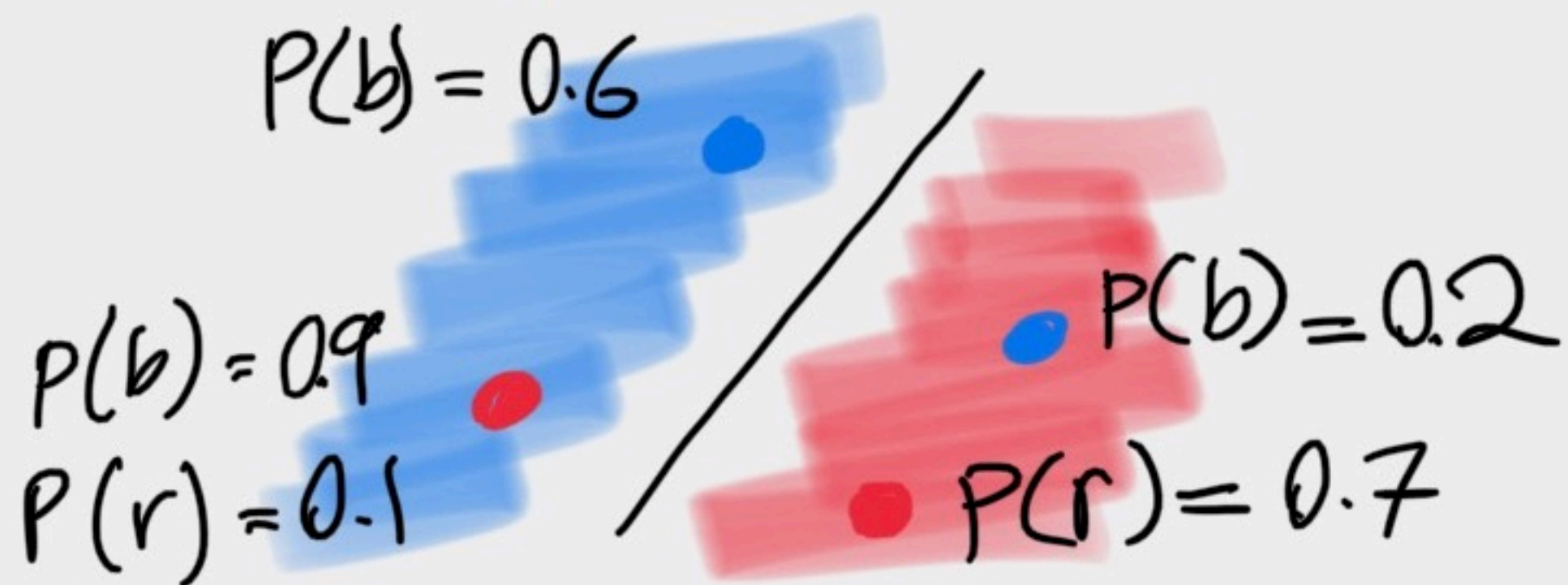


# Maximum Likelihood

Pick the model that gives the existing labels the highest probability.

So we try to maximize this probability.

example:

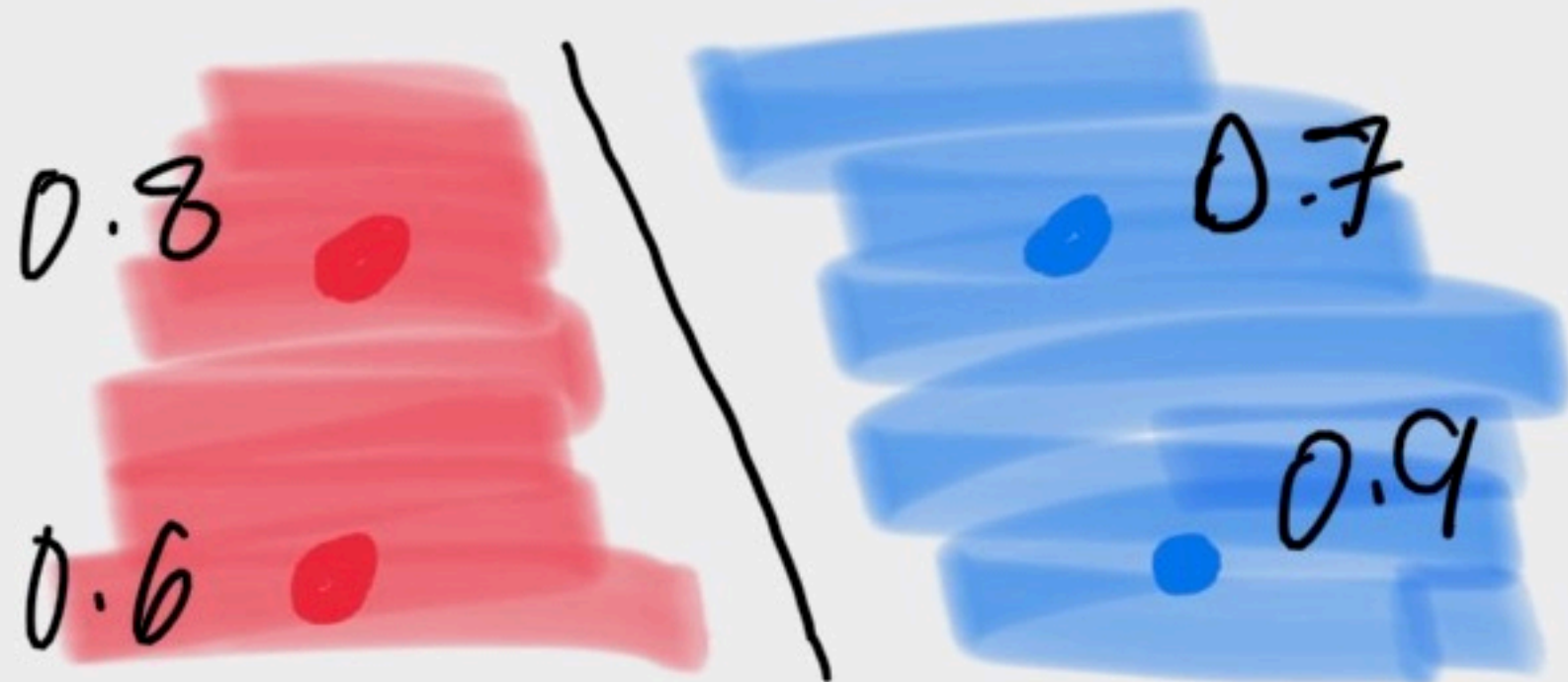


$$\hat{y} = \sigma(Wx + b) = P(\text{blue})$$

We assume the points are iid:

$$P(\text{all}) = \underbrace{0.1 \times 0.7}_{\text{red dots}} \times \underbrace{0.6 \times 0.2}_{\text{blue dots}} = 0.0084$$

$\downarrow$   
small



$$P(\text{all}) = 0.3024$$

better than the other model.

Maximum likelihood  
tries to maximize this.

When dealing with maximum likelihood, we prefer sums to products.  $\rightarrow \log_e() == \ln()$

$\ln(x) < 0$  if  $x \in [0, 1]$ . So we use  $[-\ln(x)]$ .



Cross Entropy of good models is smaller.

we can think of it as error

Error Function = Cross Entropy

So we now change the goal to minimizing cross entropy.

Probabilities —

|            | G   | R   | B   |
|------------|-----|-----|-----|
| P(gift)    | 0.8 | 0.7 | 0.1 |
| P(no gift) | 0.2 | 0.3 | 0.9 |

most likely event?  
choosing the largest  
prob. in each column.

$$0.8 \times 0.7 \times 0.9$$

$$\text{Cross Entropy} = - \sum_{i=1}^m y_i \ln(P_i) + (1 - y_i) \ln(1 - P_i)$$

↓  
it can tell us

how close two

vectors are. e.g.,

$$\text{CE}[(1, 1, 0), (0.8, 0.7, 0.1)] = 0.69$$

$$\text{CE}[(0, 0, 1), (0.8, 0.7, 0.1)] = 5.12$$

likely to  
happen



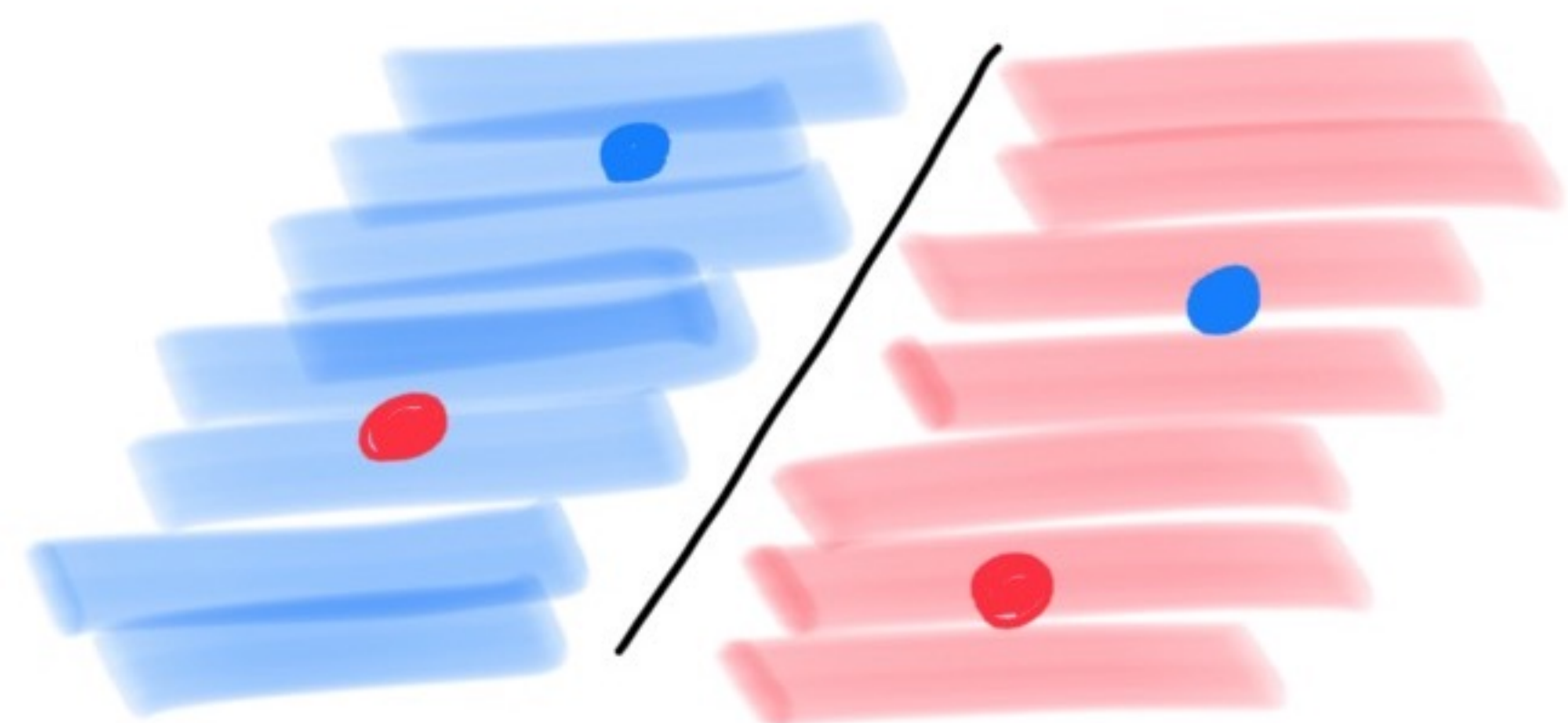
## Multi-class Entropy:

number  
of classes

$$CE = - \sum_{i=1}^n \sum_{j=1}^m y_{ij} \ln(P_{ij})$$

probability of  $i$  and  $j$

$\rightarrow 0/1$  to indicate which event we are talking about.



$$\left\{ \begin{array}{l} \text{if } y=1 \rightarrow P(\text{blue}) = \hat{y} \rightarrow \text{Error} = -\ln \hat{y} \end{array} \right.$$

$$\left\{ \begin{array}{l} \text{if } y=0 \rightarrow P(\text{red}) = 1 - P(\text{blue}) = 1 - \hat{y} \rightarrow \text{Error} = -\ln(1 - \hat{y}) \end{array} \right.$$

$$\text{Error function: } \frac{1}{m} \sum_{i=1}^m (1 - y_i) (\ln(1 - \hat{y}_i)) + y_i \ln \hat{y}_i$$

Same as:

binary  
classification:

$$E(w, b) = \frac{1}{m} \sum_{i=1}^m (1 - y_i) \ln(1 - \sigma(wx_i + b)) + y_i \ln(\sigma(wx_i + b))$$

Multi-class classification:

$$\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n y_{ij} \ln \hat{y}_{ij}$$