

$$J_m = \sum_{m=1}^M w_m^{(m)} I(y_m(x_n) \neq t_n) \quad I(e) = \begin{cases} 1 & e = \text{true} \\ 0 & \text{otherwise} \end{cases}$$

2. evaluate  $\epsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} I(y_m(x_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}}$   $\alpha_m = \ln\left(\frac{1-\epsilon_m}{\epsilon_m}\right)$

3. update

$$w_n^{(m+1)} \leftarrow w_n^{(m)} e^{\alpha_m I(y_m(x_n) \neq t_n)}$$

3. output  $y_m = \arg \left( \sum_{m=1}^M \alpha_m y_m(x) \right)$

ML EXAMS 04/11/2018

EX. 1

1. in a binary classification problem the target function is  $f: X \rightarrow Y$  where  $X = \{FR, NR, NK\}$  and  $Y = \{Y, NY\}$  the dataset is  $D = \{(x_i, y_i)\}_{i=1}^N$

2. if we use ID3 algorithm for choose which attribute to use we compute the gain

$$g_{\text{gain}} = \text{entropy}(D) - \sum_{V \in V} \frac{|S_V|}{|S|} \text{entropy}(S_V) \quad \text{where } V = \text{values of attribute}$$

$$\text{entropy} = -\sum_{i=1}^n p_i \log_2(p_i)$$

3.

~~good gain~~

$$\text{entropy}(F_Y) = -\frac{1}{2} \log\left(\frac{1}{2}\right) - \frac{1}{2} \log\left(\frac{1}{2}\right) = 1$$

$$\text{entropy}(FN) = -\frac{2}{3} \log\left(\frac{2}{3}\right) - \frac{1}{3} \log\left(\frac{1}{3}\right) = 0.916$$

$$\text{entropy}(NR_3) = -\frac{2}{3} \log\left(\frac{2}{3}\right) - \frac{1}{3} \log\left(\frac{1}{3}\right) = 0.916$$

$$\text{entropy}(NR_4) = -\frac{2}{2} \log\left(\frac{2}{2}\right) = 0$$

$$\text{entropy}(NK_4) = -1 \log(1) = 0$$

$$\text{entropy}(NKN) = -\frac{2}{4} \log\left(\frac{2}{4}\right) - \frac{1}{2} \log\left(\frac{1}{2}\right) = 1$$

$$\text{entropy}(S) = -\frac{3}{5} \log\left(\frac{3}{5}\right) - \frac{2}{5} \log\left(\frac{2}{5}\right) = 0.971$$

$$g_{\text{aim}}(S, F) = \text{entropy}(S) - \frac{2}{5} \text{entropy}(F_y) - \frac{3}{5} \text{entropy}(F_N) = 0.0184$$

$$g_{\text{aim}}(S, N_R) = \text{entropy}(S) - \frac{3}{5} \text{entropy}(N_3) - \frac{2}{5} \text{entropy}(N_4) = 0.4214$$

$$g_{\text{aim}}(S, N_C) = \text{entropy}(S) - \frac{1}{5} \text{entropy}(N_{C_y}) - \frac{4}{5} \text{entropy}(N_{C_N}) = 0.372$$

NR highest  $g_{\text{aim}}$

EX.2

$$1. h_{\text{MAP}} = \underset{\substack{\text{argmax} \\ P(D|h)}}{P(D|h)} \quad \text{if we assume } P(h_i) = P(h_j)$$

$$h_{\text{ML}} = \underset{\text{argmax}}{P(D|h)}$$

2. this is false if we introduce 3 hypothesis:

$$h_1(x) = + \quad h_2(x) = - \quad h_3(x) = - \quad P(D|h_1) = 0.4 \quad P(D|h_2) = 0.3$$

$$P(D|h_3) = 0.3$$

$$h_{\text{ML}} = \underset{\{0.4, 0.3, 0.3\}}{\text{argmax}} = 0.4 \rightarrow \text{class +}$$

now I introduce Bayesian Optimal classifier

$$V_{\text{BOC}} = \underset{\substack{\text{argmax} \\ h_i \in H}}{\sum} P(v|x, h_i) P(h_i|D)$$

is an optimal classifier that return the optimal solution

$$\cancel{P(+|x, h_1) = 1} \quad P(-|x, h_1) = 0$$

$$P(+|x, h_2) = 0 \quad P(-|x, h_2) = 1$$

$$P(+|x, h_3) = 0 \quad P(-|x, h_3) = 1$$

$$V_{\text{BOC}} = \underset{\{1 \cdot 0.4 + 0 + 0, 0 + 1 \cdot 0.3 + 1 \cdot 0.3\}}{\text{argmax}} = -$$

EX.3

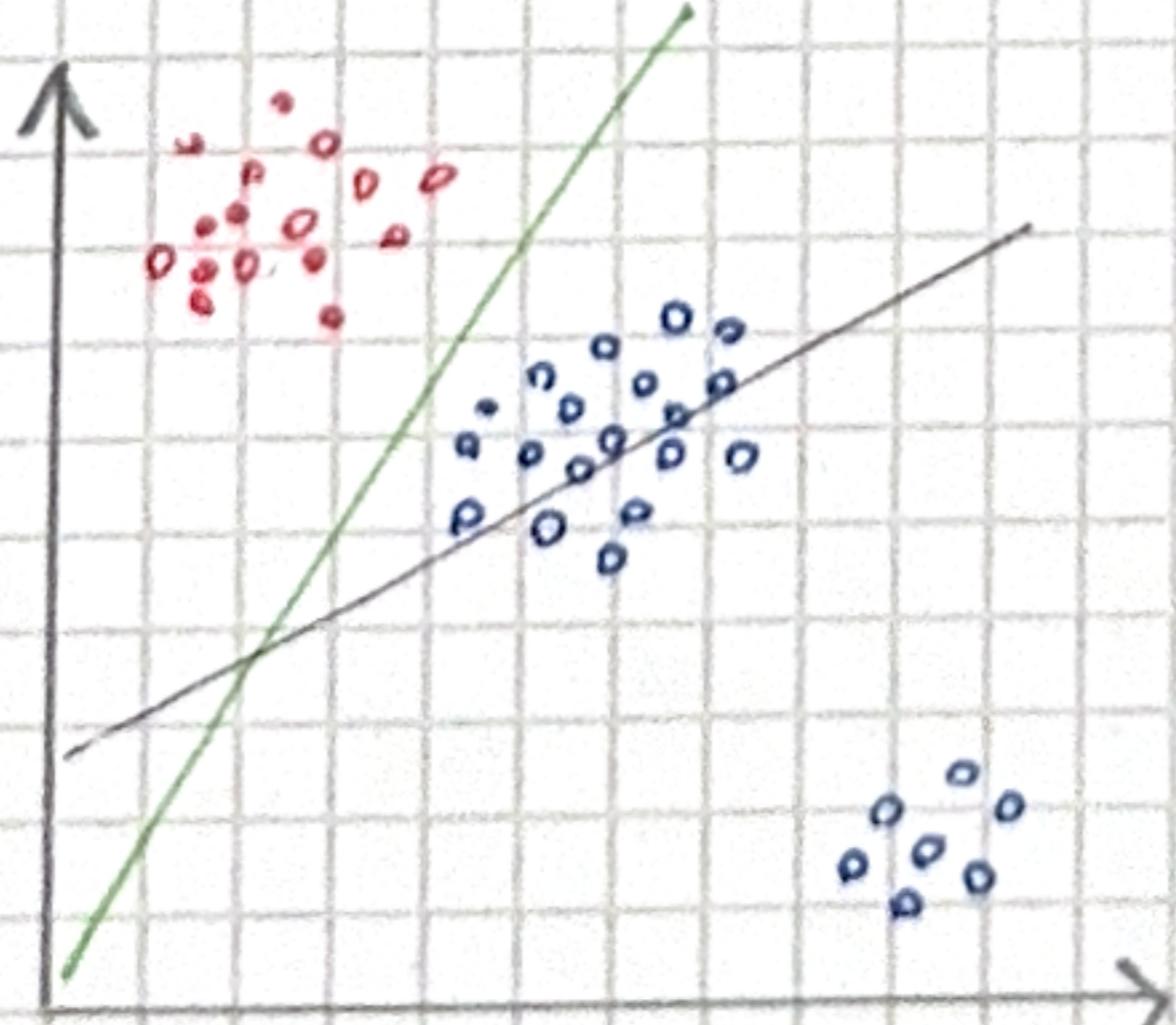
Learned  $\omega$  is a linear model for classification. In this method we want minimize an error function that is called sum of squares

$$J(\tilde{\omega}) = \frac{1}{2} \text{tr}((\tau - \tilde{x}\tilde{\omega})^T (\tau - \tilde{x}\tilde{\omega}))$$

We call this sum of squares because the trace is a sum of elements and the product of one matrix for his transpose is the square

$$\tilde{\omega}^* = \tilde{x}^T \tau \text{ and } y = \tilde{\omega}^T \tilde{x} = \tau^T (\tilde{x}^T)^T \tilde{x}$$

is not robust to outliers because is based on distance



#### EX.4

1. the gram matrix is defined as  $H = \mathbf{x}\mathbf{x}^T$ . if we have a kernel function  $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}'\mathbf{x}$  and a model  $y(\mathbf{x}, \mathbf{w}^*) = \sum_{m=1}^N w_m x_m^T \mathbf{x}$

the gram matrix is  $H = \begin{bmatrix} x_1^T x_1 & \dots & x_1^T x_M \\ \vdots & \ddots & \vdots \\ x_N^T x_1 & \dots & x_N^T x_M \end{bmatrix}$

if we have a generic kernel function  $k(\mathbf{x}, \mathbf{x}') \Rightarrow y(\mathbf{x}, \mathbf{w}^*) = \sum_{n=1}^N w_n k(x_n, \mathbf{x})$

and gram matrix will be  $H = \begin{bmatrix} k(x, x_1) & \dots & k(x, x_M) \\ \vdots & \ddots & \vdots \\ k(x_N, x_1) & \dots & k(x_N, x_M) \end{bmatrix}$ .

2. we use a modified error function  $J(\mathbf{w}) = \sum_{n=1}^N \alpha_n x_n^T \mathbf{x}$  now we can use this function to solve

$$y(\mathbf{x}, \mathbf{w}^*) = \sum_{n=1}^N \alpha_n k(x_n, \mathbf{x}) \Rightarrow \text{computationally expensive}$$

#### EX.5

1. I will use linear regression model for solve this task

$$y(\mathbf{x}, \mathbf{w}) = \sum_{n=1}^N w_n \phi(x_n)$$

that is linear in  $w_n$  the goal is minimize

the error function:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(x_n))^2$$

we use a numerical algorithm to find a solution

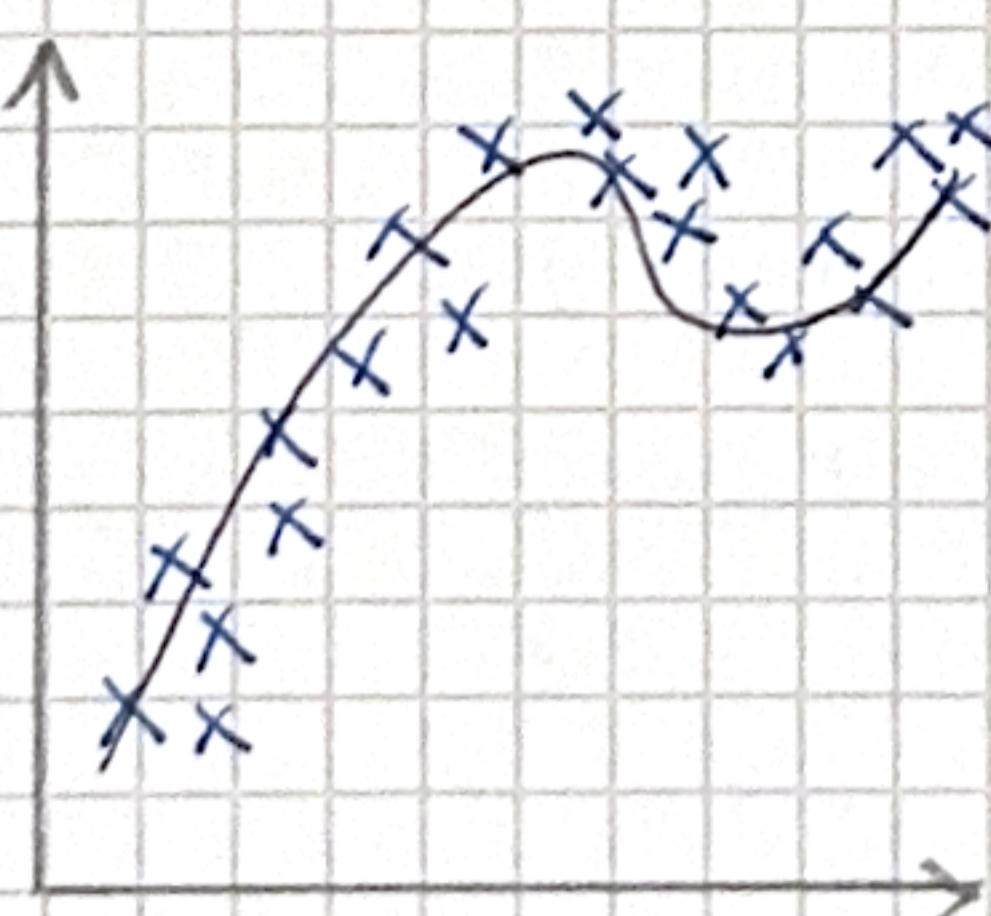
$$\hat{\mathbf{w}} \leftarrow \hat{\mathbf{w}} + \eta (t_n - \mathbf{w}^T \phi(x_n)) \phi(x_n)$$

2. the overfitting can be reduced by applying a regularization factor

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{arg\,min}} E(\mathbf{w}) + \lambda E_w(\mathbf{w}) \text{ with}$$

$$E_w(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

3.



EX. 6

1. given a dataset  $D = \{(x_n, t_n)\}_{n=1}^N$  and a target function  $f: X \rightarrow Y$ .  
kNN is an imitative model model.

1. find  $k$  nearest neighbors of new instance  $x$

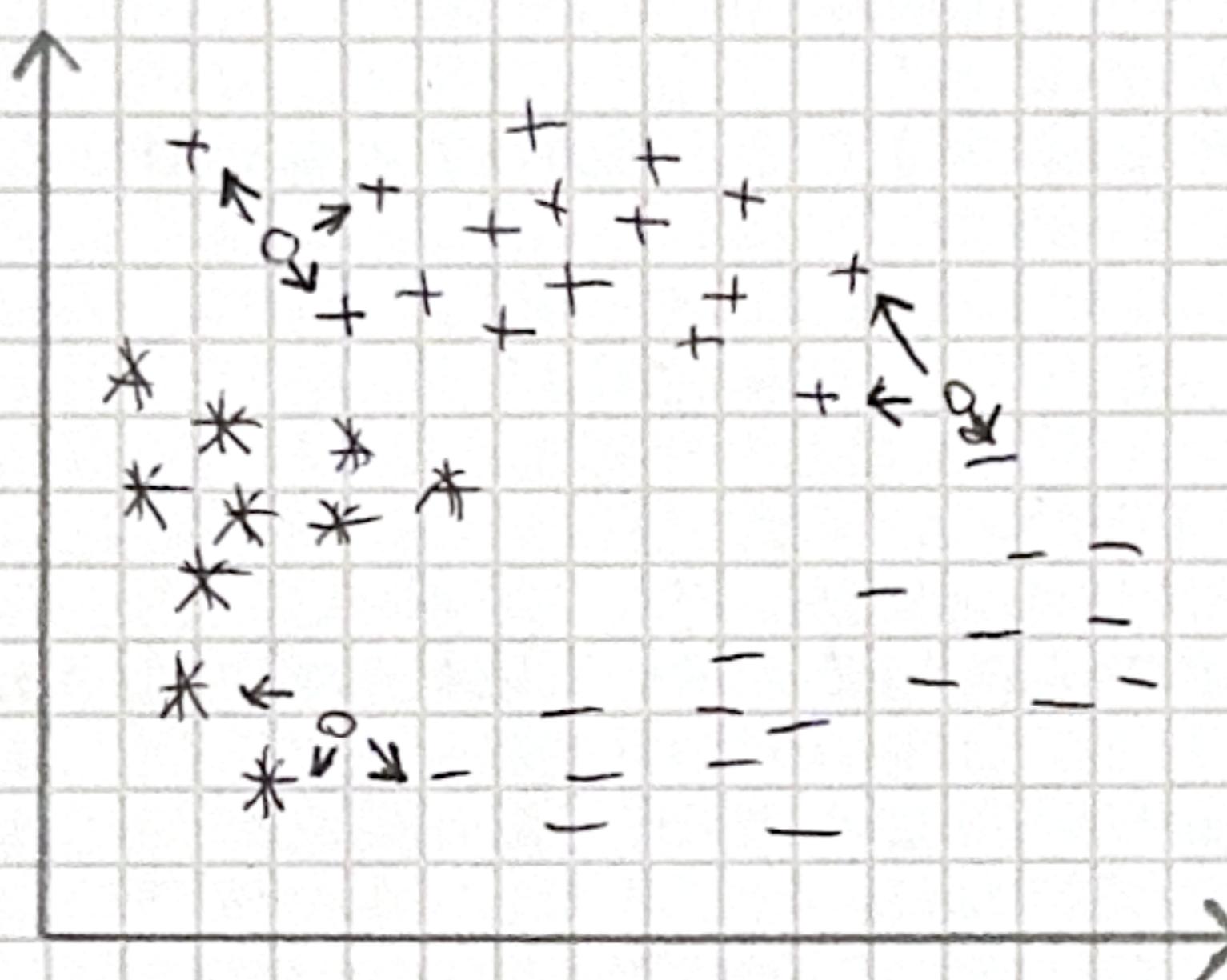
2. assign to  $x$  the most common label among the majority of neighbors

likelihood  $\pi$  of class of new instance  $x'$

$$P(c|x, D, k) = \frac{1}{k} \sum_{x_n \in N_k(x, D)} \mathbb{I}(t_n=c) \text{ when } k > 0$$

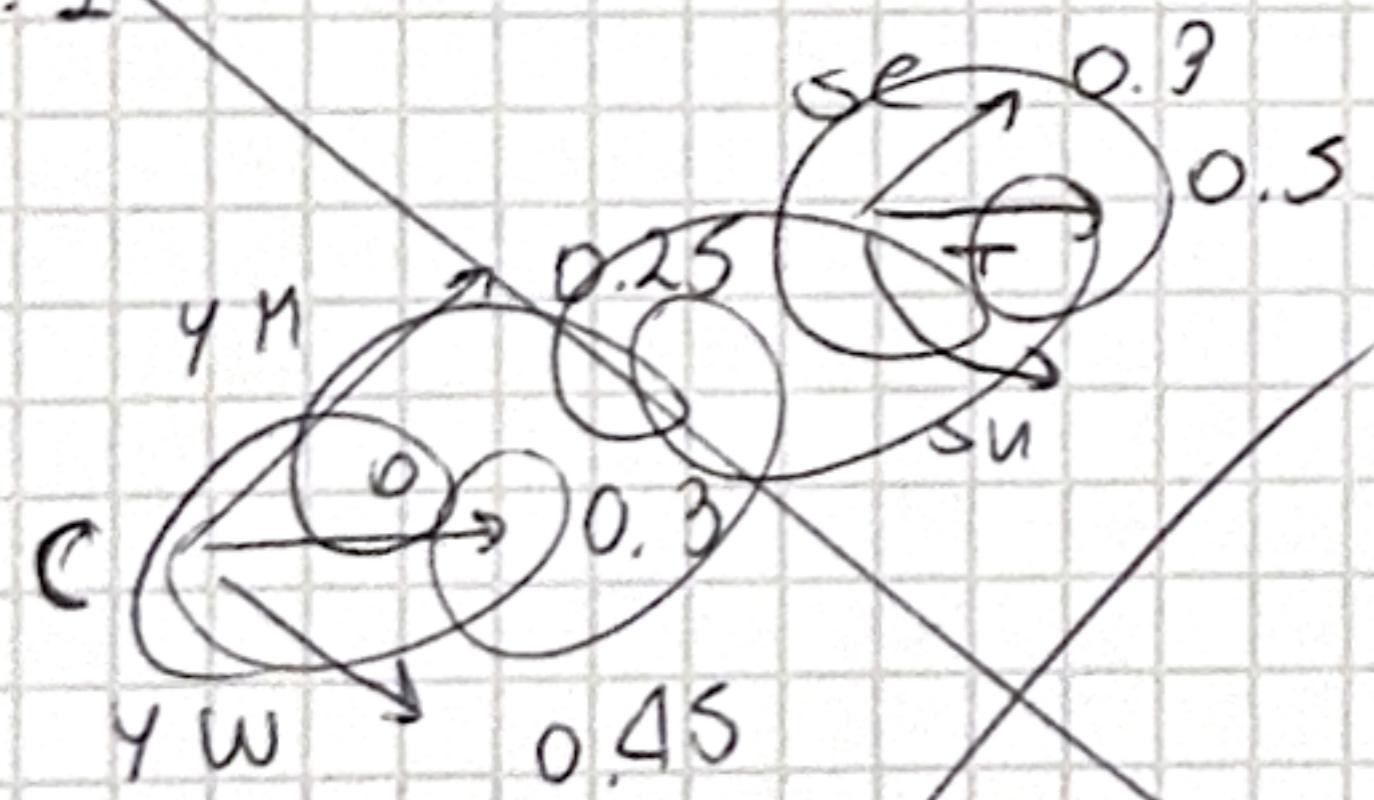
$$\mathbb{I}(c) = \begin{cases} 1 & \text{if } c = \text{true} \\ 0 & \text{otherwise} \end{cases}$$

2



ML EXAMS 2020 20/01/2020

EX. 1



$$P(Y_M) = 0.25 \quad P(Y_W) = 0.45 \quad P(O_I) = 0.3$$

$$P(Y_M|S) = 0.3 \quad P(Y_W|T) = 0.5 \quad P(Y_M|R) = 0.2$$

$$P(Y_W|S) = 0.5 \quad P(Y_W|T) = 0.3 \quad P(Y_W|R) = 0.2$$

$$P(O_I|S) = 0.3 \quad P(O_I|T) = 0.3 \quad P(O_I|R) = 0.4$$

$V^* = \text{argmax}$