

EXAM TEST 04.11.19

Question 1

1). This is a binary classification task. The target function is $f: X \rightarrow Y$ with $X = \{F \times MR \times NK\}$ and $Y = \{Y, N\}$. The dataset is $D = \{(x_i, y_i)\}_{i=1}^n$.

2). If we use the ID3 algorithm we choose the attribute with the highest information gain: $IG(S, a) = Entropy(S) - \sum_{v \in V} \frac{|S_v|}{|S|} \text{entropy}(S_v)$ with $V = \text{values } \{e\}$

With the entropy = $\sum_{i=1}^K -p_i \log_2 p_i$ that measures the impurity of the information

$$3). ent(S) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.943$$

$$IG(S, F) = ent(S) - \frac{2}{5} ent(F_F) - \frac{3}{5} ent(F_N) = 0.0196$$

$$ent(F_F) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.916 \quad ent(F_N) = 1$$

$$IG(S, MR) = ent(S) - \frac{2}{5} ent(MR_4) - \frac{3}{5} ent(MR_3) = 0.4216 \rightarrow \text{HIGHEST IG!}$$

$$ent(MR_3) = 0.916 \quad ent(MR_4) = 0$$

$$IG(S, NK) = ent(S) - \frac{1}{5} ent(NK_Y) - \frac{4}{5} ent(NK_N) = 0.141$$

$$ent(NK_Y) = 0 \quad ent(NK_N) = 1$$

Question 2

1). Given a classification task with a dataset D and an hypotheses space H . We are interested in $p(H|D)$. We know that $p(H|D) = \frac{p(D|h)p(h)}{p(D)}$

We want to find the Maximum A Posteriori (MAP) hypothesis: $h_{MAP} = \underset{h \in H}{\operatorname{argmax}} \frac{p(D|h)p(h)}{p(D)}$

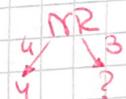
$= \underset{h \in H}{\operatorname{argmax}} p(D|h)p(h)$. If we assume that $p(h_i) = p(h_j)$ we can further simplify and we obtain the maximum likelihood (ML) hypothesis $h_{ML} = \underset{h \in H}{\operatorname{argmax}} p(D|h)p(h)$

2). This is false. We introduce 3 hypotheses:

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

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} \{0.4, 0.3, 0.3\} = +$$

Now we introduce Bayes Optimal Classifier:



We continue using this technique.

$$v_{\text{acc}} = \underset{v \in V}{\operatorname{argmax}} \sum_{h \in H} p(v|x, h) p(h|D).$$

Bayes optimal classifier is an optimal classifier, in fact it returns always the optimal solution.

$$p(+|x, h_1) = 1 \quad p(+|x, h_2) = 0 \quad p(+|x, h_3) = 0$$

$$p(-|x, h_1) = 0 \quad p(-|x, h_2) = 1 \quad p(-|x, h_3) = 1$$

$$r_{\text{NB}} = \operatorname{argmax} \{ 1 \cdot 0.4 + 0 + 0, 0 + 1 \cdot 0.3 + 1 \cdot 0.3 \} = -$$

This confirms that the hypothesis is false.

Question 3

Least squares is a linear model for classification. In this method we want to minimize an error function that is also called sum of squares.

$$J(\tilde{w}) = \frac{1}{2} \operatorname{Tr} \{ (\tilde{T} - \tilde{X}\tilde{w})^T (\tilde{T} - \tilde{X}\tilde{w}) \}$$

We call this function sum of squares because

the trace is simply the sum of the elements in the main diagonal and the product of one matrix for the transpose is simply the square.

$$\tilde{w}^* = \tilde{X}^T \tilde{T} \text{ and } y = \tilde{W}^T \tilde{X} = \tilde{T}^T (\tilde{X}^T)^T \tilde{w}$$

Question 4

1). The gram matrix is $K = X X^T$. If we have a kernel function $K(x, x') = x^T x'$ and a model $y(x, w^*) = \sum_{n=1}^N w_n x_n^T X$ the gram matrix is $K = \begin{pmatrix} x_1^T x_1 & x_1^T x_N \\ x_N^T x_1 & x_N^T x_N \end{pmatrix}$

If we have a generic kernel function $K(x, x') \rightarrow y(x, w^*) = \sum_{n=1}^N w_n K(x_n, x)$ and the gram matrix K is $K = \begin{pmatrix} K(x_1, x_1) & K(x_1, x_N) \\ K(x_N, x_1) & K(x_N, x_N) \end{pmatrix}$

2). We use the modified error function $J(w) = \sum_{n=1}^N E(y_n, t_n) + \lambda \|w\|^2$

We know that $E(y_n, t_n) = (y_n - t_n)^2$ and $y(x, w^*) = \sum_{n=1}^N \alpha_n x_n^T X$. Now we kernelize this solution: $y(x, w^*) = \sum_{n=1}^N \alpha_n K(x_n, x)$. This is a solution for our problem but is really computationally expensive, in fact requires $|D|^2$ matrix multiplications.

QUESTION 5

1). We can use a linear regression model. We know that $y(n, w) = \sum_{j=1}^n w_j \phi_j(x)$.

The target value is $t = y(n, w) + \varepsilon$ with ε additive gaussian noise.

$p(\varepsilon | \beta) = N(\varepsilon | 0, \beta^{-1})$. With this hypothesis we obtain: $p(t | n_1, \dots, n_r, w, \beta) = N(t | y(n, w), \beta^{-1})$

Using the iid hypothesis: $p(\{t_1, \dots, t_n\} | n_1, \dots, n_r, w, \beta) = \prod_{n=1}^N N(t_n | w^\top \phi(x_n), \beta^{-1}) =$

$$= \prod_{n=1}^N \ln N(t_n | w^\top \phi(x_n), \beta^{-1}) = -\beta \cdot \frac{1}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2 - \frac{n}{2} \ln(2\pi\beta^{-1})$$

The error function is $E(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2$. We use sequential learning

to find the solution: $\hat{w} \leftarrow \hat{w} + \eta [(t_n - w^\top \phi(x_n)) \phi(x_n)]$

2). We can reduce overfitting using regularization:

$$w^* = \underset{w}{\operatorname{argmin}} E_D(w) + \lambda E_R(w) \quad \text{with } \lambda > 0 \text{ regularization factor.}$$

$$E_R(w) = \frac{1}{2} w^\top w r.$$

QUESTION 6

1). We compute the k nearest neighbors of a new instance $x \notin D$ and then we assign to x the most common label among the majority of neighbors.

$$p(c | n, k, D) = \frac{1}{k} \sum_{x_n \in N_k(x, D)} \mathbb{I}(t_n = c) \quad \text{with } \mathbb{I}(e) = \begin{cases} 1 & \text{if } e \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$