

EXAM 23.03.18

Question 1

1). Given a training set \mathcal{D} and an hypotheses space H : we overfits training data if exists $h' \in H$ such that $\text{errors}(h) < \text{errors}(h')$ and $\text{errors}(h) > \text{errors}(h')$.

2). In decision tree is very probable to have solutions that are overfitted, for example if we have two different solutions for a problem and one tree is deeper than another for sure we have overfitting. We have two methods to reduce overfitting:

- Reduced error pruning: we split our training set into training set and validation set. We evaluate the impact on validation set of pruning each possible node and we greedily remove the node that most improves validation accuracy.

- Rule post pruning: we generate all the tree. We move our tree into a logical form and we compute the accuracy. We start pruning, if the accuracy increases we continue, otherwise we stop.

Question 2

1). The Naive Bayes classifier is based on the property of conditional independence.

X is conditionally independent from Y given Z if $p(X, Y | Z) = p(X|Z)p(Y|Z)$.

Given a target function $f: X \rightarrow V$ with $V = \{V_1, V_2, \dots, V_k\}$ knowing that each sample x can be represented as $\langle o_1, o_2, \dots, o_n \rangle$

$$p(V|x, \mathcal{D}) = p(V|o_1, \dots, o_n, \mathcal{D}) = \frac{p(o_1, \dots, o_n | V, \mathcal{D})p(V|\mathcal{D})}{p(o_1, \dots, o_n | \mathcal{D})}$$

$V_{NB} = \underset{V \in V}{\operatorname{argmax}} p(o_1, \dots, o_n | V, \mathcal{D})p(V|\mathcal{D})$. Now we introduce the conditional independence

$V_{NB} = \underset{V \in V}{\operatorname{argmax}} p(o_1, \dots, o_n | V, \mathcal{D}) = \prod_{n=1}^N p(o_n | V, \mathcal{D})$; each sample is mutually conditionally independent from another one. $V_{NB} = \underset{V \in V}{\operatorname{argmax}} \prod_{n=1}^N p(V|\mathcal{D})p(o_n | V, \mathcal{D})$.

2). The target function is $f: X \rightarrow Y$ with $X = \{\text{title}, \text{abstract}, \text{author}, \text{site}\}$ and

$Y = \{HL, KR, PL\}$. The dataset is $\mathcal{D} = \{(d_i, y_i)\}_{i=1}^N$.

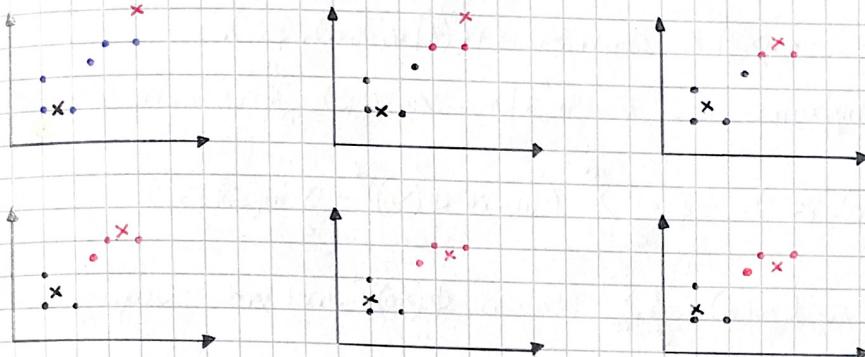
$V_{NB} = \underset{V \in V}{\operatorname{argmax}} p(V|\mathcal{D}) \prod_{i=1}^N p(d_i | V, \mathcal{D})$ knowing that $p(d_i | V, \mathcal{D}) = \prod_{j=1}^{\text{length}(d)} p(d_{ij} = v_{ij} | V, \mathcal{D})$

$$\hat{p}(d_{ij} | V, \mathcal{D}) = \frac{T_{ij} + 1}{T_{ij} + 2}$$

Question 3

- 1). In an unsupervised learning task we have a target function $f: X \rightarrow Y$ and a dataset $D = \{(x_n)\}_{n=1}^N$. It's possible to see that in this case we don't have labels, so the dataset is composed only by input values.

2):



- 3). It's possible to see that initially one instance is classified as a point of the "blue centroid" but after the first iteration the instance will be classified as red.

Question 4

- 1). We compute the k -nearest neighbors of a new instance $x \notin D$.

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

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

with $N_k(x_n, D)$ k -nearest neighbors of x_n .

Question 5

- 1). Backpropagation (BackProp) is an algorithm to compute the gradient, is not a learning algorithm and is used in ANNs. In this mode we propagate the gradient through all the net. In the forward step we compute the values of the parameter of the functions while in the backward step we compute the gradient.

SGD is a learning algorithm that is used to find the minimum or the maximum of a function, is an iterative method and is an optimization strategy very common.

FORWARD STEP

Requires $W^{(l)}$ weights matrices, n input value, $b^{(l)}$ bias parameters and t target value.

$$h^{(0)} = n$$

for $l=1 \dots L$ do:

$$d^{(l)} = b^{(l)} + W^{(l)} h^{(l-1)}$$

$$h^{(l)} = f(d^{(l)})$$

$$h^{(L)} = y$$

$$J = J(t, y)$$

BACKWARD STEP

$$g \leftarrow \nabla_y J = \nabla_y J(t, y)$$

for $l=L$ to 1 :

$$g \leftarrow \nabla_{d^{(l)}} J = g \circ f'(d^{(l)})$$

$$\nabla_{b^{(l)}} J = g$$

$$\nabla_{w^{(l)}} J = g (h^{(l-1)})^\top$$

$$g \leftarrow \nabla_{w^{(l)}} J = (W^{(l)})^\top g$$

Question 6

1). We have a target function $f: X \rightarrow Y$ with $X \subseteq \mathbb{R}^d$ and $Y = \mathbb{R}$. Our model is:

$$y(x, w) = \sum_{j=1}^n w_j \phi_j(x) \text{ that is linear in } w \text{ and } \mathcal{D} = \{(x_n, t_n)\}_{n=1}^N.$$

We know that $t = y(x, w) + \varepsilon$ with ε additive Gaussian noise. If we know that:

$$\varphi(\varepsilon | \beta) = N(\varepsilon | 0, \beta^{-1}) \rightarrow \varphi(t | x_1, \dots, x_N, w, \beta) = N(t | y(x, w), \beta^{-1}).$$

Now we use the iid hypothesis: $\varphi(\{t_1, \dots, t_N\} | x_1, \dots, x_N, w, \beta) = \prod_{n=1}^N N(t_n | w^\top \phi(x_n), \beta^{-1}) =$

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

The error function is $E_\beta(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2$ and $w^* = \operatorname{argmin}_w E_\beta(w)$.

2). We can use Least Squares: $E_\beta(w) = \frac{1}{2} (t - \phi w)^\top (t - \phi w)$ and $w^* = \phi^\top t$.

We can also use a sequential algorithm: $\hat{w} \leftarrow \hat{w} + \eta [t_n - w^\top \phi(x_n)] \phi(x_n)$