

EXAH 12.06.19

Question 1

- Given a dataset D and an hypotheses space H : H overfits training data if exists another hypothesis $h' \in H$ such that $\text{errors}(h) < \text{errors}(h')$ and $\text{errors}(h) > \text{errors}(h)$
- In decision trees is very common to have problems of overfitting. If we have two solutions for a problem and one tree is deeper than the other for sure we have overfitting. We have 2 possible techniques to avoid 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. We prune the node that most improves validation accuracy.
 - Rule post pruning: we generate all the tree and we move to a logical representation of the tree. We compute the accuracy and we start pruning. If the accuracy after pruning is greater than the accuracy before pruning we have overfitting and we can prune.

Question 2

- 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 $V = \{V_1, V_2, \dots, V_n\}$ knowing that each instance x can be described as $\langle o_1, o_2, \dots, o_n \rangle$

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

Now we introduce the conditional independence hypothesis:

$$p(o_1, \dots, o_n | V, D) = \prod_{i=1}^n p(o_i | V, D) \quad \text{Each sample is mutually conditionally independent from another.}$$

$$V_{NB} = \underset{V \in V}{\operatorname{argmax}} \prod_{i=1}^n p(o_i | V, D)$$

Is possible to see that this solution, differently from the BoC is not optimal but can be used also if the hypotheses space is large or if we don't have analytical solutions.

- This is a classification task. The target function is $f: X \rightarrow Y$ with $X = \{\text{title}, \text{author}, \text{x_abs}, \text{x_site}\}$ and $Y = \{\text{ML}, \text{UR}, \text{PL}\}$. The dataset is $D = \{(d_i, y_i)\}_{i=1}^N$.

$$V_{NB} = \underset{V \in V}{\operatorname{argmax}} \prod_{i=1}^N p(y_i | V, D)$$

We know that $\hat{p}(d_i | v, \Delta) = \prod_{j=1}^L \hat{p}(d_i = w_j | v, \Delta)$ with $L = \text{length}(d_i)$ and $\hat{p}(d_i = w_j | v, \Delta)$ probability that the word w_j occurs in the document i .

$$\hat{p}(d_i | v, \Delta) = \frac{\sum_j p(d_i = w_j | v, \Delta)}{L+1}$$

Question 3

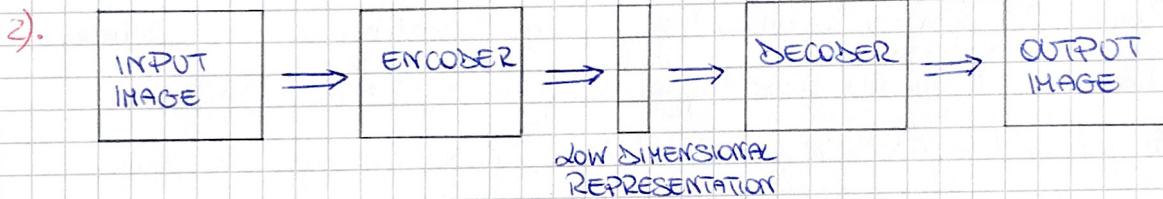
$$1). p(\hat{p}(\tilde{w}) = \prod_{n=1}^N y_n p_n (1-y_n)^{1-p_n}) \text{ with } y_n = \sigma(\tilde{w}^\top \tilde{w}_n)$$

We want to minimize the cross-entropy error function $E(w) = -\sum_{n=1}^N [p_n \log(y_n) + (1-p_n) \log(1-y_n)]$

The parameters of the model that have to be learned are the numbers of hours o_i that the student i has attended a course and the number of hours s_i that the student has studied for the exam.

Question 4

1). An autoencoder is a combination of two NNs, a encoder and a decoder. The training is based on reconstruction loss and the hidden layers of the nets have a reduced size (bottlenecks). An autoencoder takes as input and as output the same sample x_n from the dataset $\{x_n\}$ and uses it for training.



Question 5

1). In a dynamic system we have the property of fully observability if we can see the state resulting from the execution of one action. Also if we have non deterministic action we can not predict the state resulting from the execution but we can observe it.

2). We can define an MDP as: $MDP = \langle X, A, \delta, r \rangle$. X is the set of states, A is the set of actions, $\delta: X \times A \rightarrow X$ is the transition function and $r: X \times A \rightarrow \mathbb{R}$ is the reward function. We can define a POMDP as: $POMDP = \langle X, A, Z, \delta, r, o \rangle$. The differences with the previous model are: Z is the set of observations, δ is the probability of the transition function and o is the probability of the observations and $p(x_0)$ is the probability of the initial state.

2).

