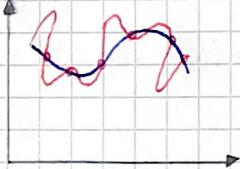


## EXAM II.02.20

### Question 1

- 3). In classification we have a target function  $f: X \rightarrow Y$  with  $X \subseteq \mathbb{R}^d$  and  $Y = \{y_1, y_2\}$ .
- In regression we have a target function  $f: X \rightarrow Y$  with  $X \subseteq \mathbb{R}^d$  and  $Y = \mathbb{R}$ .
- 2). In linear regression we have a model  $y(n, w) = \sum_{j=1}^N w_j \phi_j(x) = w^\top \phi(x)$ . The target value is  $t = y(n, w) + \varepsilon$  with  $\varepsilon$  additive gaussian noise. Since the noise is gaussian  $p(\varepsilon | \beta) = N(\varepsilon | 0, \beta^{-1}) \Rightarrow p(t | n_1, \dots, n_N, w, \beta) = N(t | y(n, w), \beta^{-1})$ .
- If we use the iid hypothesis:  $p(t_1, \dots, t_N | n_1, \dots, n_N, w, \beta) = \prod_{n=1}^N N(t_n | w^\top \phi(x_n), \beta^{-1}) =$
- $$= \prod_{n=1}^N \ln N(t_n | w^\top \phi(x_n), \beta^{-1}) = -\beta \cdot \frac{1}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2 - \frac{N}{2} \ln(2\pi\beta^{-1})$$
- We want to find  $w^* = \underset{w}{\operatorname{argmin}} E_D(w)$  with  $E_D(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2$
- We can find  $w^*$  with least squares  $\Rightarrow E_D(w) = \frac{1}{2} (t - \phi w)^\top (t - \phi w) \Rightarrow w^* = \phi^\top t$ .
- 3).
- 
- It is possible to see that this solution is overfitted.
- We can avoid overfitting applying regularization.
- $\underset{w}{\operatorname{argmin}} E_D(w) + \lambda E_w(w)$  with  $\lambda > 0$  regularization factor
- and  $E_w(w) = \frac{1}{2} w^\top w$ .

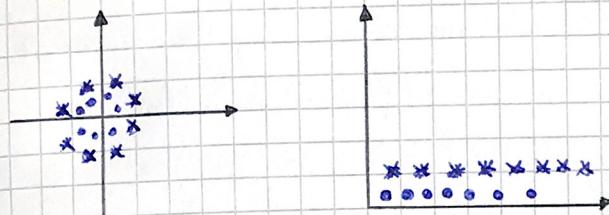
### Question 2

- 1). Logistic regression is a probabilistic discriminative model based on maximum likelihood.
- $p(t|\tilde{w}) = \prod_{n=1}^N y_n^{t_n} (1-y_n)^{1-t_n}$ . We want to minimize the error function (cross-entropy):
- $$E(\tilde{w}) = -\sum_{n=1}^N t_n \log(y_n) + (1-t_n) \log(1-y_n)$$
- We can apply Newton-Raphson to find  $\tilde{w}^* = \underset{\tilde{w}}{\operatorname{argmin}} E_D(\tilde{w}) \Rightarrow \tilde{w} \leftarrow \tilde{w} + H(\tilde{w})^{-1} \nabla E(\tilde{w})$  with  $\nabla E(\tilde{w}) = \sum_{n=1}^N (y_n - t_n) \tilde{x}_n$ .
- 2). The second solution fits the data better. We have to do the product  $w_i x_i$  and then we apply the sigmoid function. If  $\sigma(w_i x_i) = t_i$  the prediction is correct otherwise is wrong. It is possible to see that in the first case we have 0 correct predictions, in the second case 3.

### Question 3

- 1). If the input vector appears only in the form of an inner product  $x^T x'$ , we can replace the inner product with a kernel function  $K(x, x')$ . For example if we have  $\phi(x)\phi(x')$  we can replace this product with  $K(\phi(x), \phi(x'))$ .

2).



Use a polynomial kernel function:

$$K(x, x') = (Bx^T x' + y)^d, \text{ selecting for example } d=3.$$

### Question 5

- 1). In supervised learning we have a target function  $f: X \rightarrow Y$  and a dataset  $D = \{(x_i, y_i)\}_{i=1}^n$ , and we want to find an approximated function  $\hat{f}$  that returns values similar to  $f$ , specially for samples not in the dataset.

In reinforcement learning we have a dataset  $D = \{x_1, o_1, r_1, \dots, x_n, o_n, r_n\}$  and we have an agent that acts in an environment; the goal of RL is to find a behaviour function  $\Pi: X \rightarrow A$ .

- 2). In MDPs we have the property of full observability, states are fully observable. If we have non-deterministic actions we cannot predict the outcome of one action but we can observe the resulting state after the execution of this action.

### Question 6

- 1). K-Means takes as input the dataset  $\{x_i\}$  and the number of clusters  $K$  and returns as output  $y_1, \dots, y_K$ . The algorithm works in the following way:

① We select the value of  $K = \# \text{clusters}$ .

② We have to split our samples in clusters. We can do this randomly or systematically as follow:

- We take the first  $K$  samples and we assign each sample to a single element cluster;
- We take the remaining  $N - K$  samples and we assign each sample to the cluster with the nearest centroid.

③ We select each sample in sequence and we compute the distance from all the clusters. If the sample is not correctly classified we move it to the correct cluster and we recompute the centroids of the modified clusters. We repeat this step until convergence.