

# Machine Learning from Data – IDC – 2022

## HW5 – Theory + SVM

### 1. Kernels and mapping functions (25 pts)

- a. (20 pts) Let  $K(x, y) = (x \cdot y + 1)^3$  be a function over  $\mathbb{R}^2 \times \mathbb{R}^2$  (i.e.,  $x, y \in \mathbb{R}^2$ ).

Find  $\psi$  for which  $K$  is a kernel. (It may help to first expand the above term on the right-hand side).

- b. (2 pts) What did we call the function  $\psi$  in class if we remove all coefficients?
- c. (3 pts) How many multiplication operations do we save by using  $K(x, y)$  versus  $\psi(x) \cdot \psi(y)$ ?

### 2. Lagrange multipliers (25 pts)

Let  $f(x, y) = 2x - y$ . Find the minimum and the maximum points for  $f$  under the constraint  $g(x, y) = \frac{x^2}{4} + y^2 = 1$ .

### 3. PAC Learning (25 pts)

Let  $X = \mathbb{R}^2$ . Let vectors  $u = \left(\frac{\sqrt{3}}{2}, \frac{1}{2}\right)$ ,  $w = \left(\frac{\sqrt{3}}{2}, -\frac{1}{2}\right)$ ,  $v = (0, -1)$

$$\text{and } C = H = \left\{ h(r) = \left\{ (x_1, x_2) \left| \begin{array}{l} (x, y) \cdot u \leq r, \\ (x, y) \cdot v \leq r, \\ (x, y) \cdot w \leq r \end{array} \right. \right\}, \text{ for } r > 0, \right\}$$

the set of all origin-centered upright equilateral triangles.

Describe a polynomial sample complexity algorithm  $L$  that learns  $C$  using  $H$ . State the time complexity and the sample complexity of your suggested algorithm. Prove all your steps.

4. (15 pts) A business manager at your ecommerce company asked you to make a model to predict whether a user is going to proceed to checkout or abandon their cart. You created the model using, and reported 20% error on your test set of size 1000 samples. In the business manager's presentation to upper management, he presented your

model and stated that the company can expect 20% error when deploying the model live on the website.

Luckily, you realize that this is a mistaken assumption, and you correct the statement to say that with 95% confidence, the true error they can expect is up to what percentage? (Just state the error percentage).

5. SVM (10 pts)

See the notebook in the homework files and follow the instructions there.

Take a **screenshot** of your resulting graph near the bottom of the notebook (titled “My Graph”) and paste into your submission PDF along with your answers to the theoretical questions. Do **NOT** submit your code.

**Omer Wachman – 312332372**  
**Adam Vinestock - 209795624**  
**HW No.5**

**Question 1:**

Let  $K(x, y) = (x \cdot y + 1)^3$  be a function over  $\mathbb{R}^2 \times \mathbb{R}^2$  (i.e.,  $x, y \in \mathbb{R}^2$ ).

**a. Find  $\psi$  for which  $K$  is a kernel:**

$$\begin{aligned} (x \cdot y + 1)^3 &= \\ &= (x \cdot y)^3 + (x \cdot y)^2 + (x \cdot y) + 1 = (x_1 y_1 + x_2 y_2)^3 + (x_1 y_1 + x_2 y_2)^2 + (x_1 y_1 + x_2 y_2) + 1 \\ &= x_1^3 y_1^3 + 3x_1^2 y_1^2 x_2 y_2 + 3x_1 y_1 x_2^2 y_2^2 + x_2^3 y_2^3 + 3x_1^2 y_1^2 + 6x_1 y_1 x_2 y_2 + 3x_2^2 y_2^2 \\ &\quad + 3x_1 y_1 + 3x_2 y_2 + 1 \\ \Rightarrow \psi(x_1, x_2) &= \langle x_1^3, \sqrt{3}x_1^2 x_2, \sqrt{3}x_1 x_2^2, x_2^3, \sqrt{3}x_1^2, \sqrt{6}x_1 x_2, \sqrt{3}x_2^2, \sqrt{3}x_1, \sqrt{3}x_2, 1 \rangle \end{aligned}$$

**b. What did we call the function  $\psi$  in class if we remove all coefficients?**

We have called this function – cubic mapping.

**c. How many multiplication operations do we save by using  $K(x, y)$  versus  $\psi(x) \cdot \psi(y)$ ?**

When using the mapping function of  $\psi$  we map from  $\mathbb{R}^2$  to  $\mathbb{R}^{10}$ , therefore when applying the operation  $\psi(x_1, x_2) \cdot \psi(y_1, y_2)$  we use 10 multiplications.

When applying the kernel  $K(x, y) = (x_1 \cdot y_1 + x_2 \cdot y_2 + 1)^3$  we use 4 multiplications.

Therefore, in total we save 6 multiplication operations.

**Question 2:**

Let  $f(x, y) = 2x - y$ . Find the minimum and the maximum points for  $f$  under the constraint  $g(x, y) = x^2 + y^2 = 1$ .

Solution:

Lagrange equation:

$$\exists \lambda \text{ s.t., } L(x, y) = f(x, y) + \lambda g(x, y)$$

Find the partial derivatives and equal to zero

$$\frac{\partial}{\partial x} L(x, y) = 2 + \frac{1}{2} \lambda x = 0 \Rightarrow x = -\frac{4}{\lambda}$$

$$\frac{\partial}{\partial y} L(x, y) = -1 + 2\lambda y = 0 \Rightarrow y = \frac{1}{2\lambda}$$

$$\frac{\partial}{\partial \lambda} L(x, y) = \frac{x^2}{4} + y^2 - 1 = 0 \Rightarrow \frac{\left(-\frac{4}{\lambda}\right)^2}{4} + \left(\frac{1}{2\lambda}\right)^2 - 1 = \frac{16}{4\lambda^2} + \frac{1}{4\lambda^2} - 1 = \frac{17}{4\lambda^2} - 1 \Rightarrow$$

$$\Rightarrow \frac{17}{4\lambda^2} - 1 = 0 \Rightarrow (\lambda_1, \lambda_2) = \left(+\frac{\sqrt{17}}{4}, -\frac{\sqrt{17}}{4}\right)$$

Because we squared  $\lambda$  after the substitution, we must consider that the resulted value of  $\lambda$  can be also negative. Hence, we got two different  $\lambda$ .

Substitute  $(\lambda_1, \lambda_2)$  into X and Y values:

$$x = -\frac{4}{\lambda} \Rightarrow (x_1, x_2) = \left(-\frac{8}{\sqrt{17}}, +\frac{8}{\sqrt{17}}\right)$$

$$y = \frac{1}{2\lambda} \Rightarrow (y_1, y_2) = \left(+\frac{1}{\sqrt{17}}, -\frac{1}{\sqrt{17}}\right)$$

Now we have to different point to substitute in  $f(x, y)$ :

$$1. f(x_1, y_1) = f\left(-\frac{8}{\sqrt{17}}, +\frac{8}{\sqrt{17}}\right) = -\frac{16}{\sqrt{17}} - \frac{1}{\sqrt{17}} = -\frac{17}{\sqrt{17}} = -\sqrt{17} \Rightarrow \textbf{Minimum point}$$

$$2. f(x_2, y_2) = f\left(+\frac{1}{\sqrt{17}}, -\frac{1}{\sqrt{17}}\right) = \frac{16}{\sqrt{17}} - \left(-\frac{1}{\sqrt{17}}\right) = \frac{17}{\sqrt{17}} = \sqrt{17} \Rightarrow \textbf{Maximum point}$$

### Question 3:

PAC Learning

Let  $X = \mathbb{R}^2$ . Let vectors  $u = (\sqrt{3}, 1)$ ,  $w = (\sqrt{3}, -1)$ ,  $v = (0, -1)$

and  $C = H = \{(x, y) \mid (x, y) \cdot v \leq r, (x, y) \cdot w \leq r, (x, y) \cdot u \leq r\}$ , for  $r > 0$

the set of all origin-centered upright equilateral triangles.

Describe a polynomial sample complexity algorithm  $L$  that learns  $C$  using  $H$ . State the time complexity and the sample complexity of your suggested algorithm. Prove all your steps.

We will define an algorithm that learns the  $r$  parameter according to given data.

```

r = 0 # Initializing r
for every instance  $\vec{x} = (x_1, x_2) \in D$  with “True” label:
    if  $r < \max(\vec{x} \cdot \mathbf{u}, \vec{x} \cdot \mathbf{v}, \vec{x} \cdot \mathbf{w})$ :
         $r = \max(\vec{x} \cdot \mathbf{u}, \vec{x} \cdot \mathbf{v}, \vec{x} \cdot \mathbf{w})$ 
return r

```

such that:

$$h = L(D) = \{(x_1, x_2) \mid (x, y) \cdot v \leq r, (x, y) \cdot w \leq r, (x, y) \cdot u \leq r\}$$

$r$  represents the distance from the origin to each of the triangle vertices.

Given some  $\delta \in [0, 0.5]$ , we can bound our sample size, such that our learning algorithm is not epsilon-bad with confidence  $1 - \delta$ .

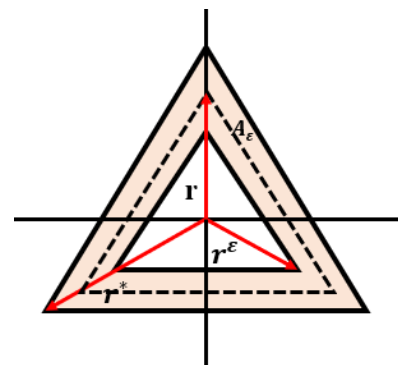
Note that our learning algorithm is a consistent learner and by taking a large number of examples we get closer to the concept.

Let  $c \in \mathcal{C}$  be a concept.

Let  $r$  be the value that defines the hypothesis from the learning algorithm  $L(D)$  shown above.

Where  $A_{\mathcal{C}}$  is the area between  $L(D)$  and the concept.

Then  $\forall x \in X$ , the probability that  $x \in A_\varepsilon$  is less than  $\varepsilon$ .



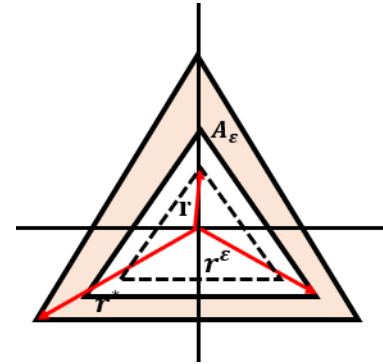
2<sup>nd</sup> case:  $r^\epsilon \geq r$

Given  $m$  samples to the learning algorithm, the probability of an instance  $x \in X$  being misclassified is  $(1 - \epsilon)^m \leq e^{-\epsilon m}$

Now, let's bound it with  $\delta$ :

$$e^{-\epsilon m} \leq \delta \Rightarrow \frac{1}{\delta} \leq e^{\epsilon m} \Rightarrow \ln\left(\frac{1}{\delta}\right) \leq \ln(e^{\epsilon m}) \Rightarrow m \geq \frac{\ln\left(\frac{1}{\delta}\right)}{\epsilon}$$

So with confidence  $\delta$  in order for hypothesis  $L(D)$  to be  $\epsilon$ -bad we need at least  $\frac{\ln\left(\frac{1}{\delta}\right)}{\epsilon}$  instances.



### Time complexity:

For each sample we perform  $O(1)$  calculations.

In total the time complexity is linear to the sample size -  $O(m)$

### Question 4:

We can say that with 95% confidence, the true error they can expect is in the interval:

$$(\hat{p} - 2se, \hat{p} + 2se)$$

$$\hat{p} = \text{true error} = 20\% = 0.2$$

$$n = \text{number of samples in test set} = 1000$$

$$se = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = \sqrt{\frac{0.2(1-0.2)}{1000}} = \sqrt{\frac{0.2 \cdot 0.8}{1000}} = \sqrt{0.00016} \approx 0.01265$$

$\Rightarrow$  we can say with confidence of 95% that the true error between:

$$(0.2 - 1.96 * 0.01265, 0.2 + 1.96 * 0.01265) = (0.1752, 0.2247) \approx (17.52\%, 22.47\%)$$

### Question 5:

