

# NYCU Introduction to Machine Learning, Homework 4

110550014, 吳權祐

The screenshot and the figures we provided below are just examples. **The results below are not guaranteed to be correct.** Please make sure your answers are clear and readable, or no points will be given. Please also remember to convert it to a pdf file before submission. **You should use English to answer the questions.** After reading this paragraph, you can delete this paragraph.

## Part. 1, Coding (50%):

- (10%) Show the accuracy score of the testing data using *linear\_kernel*. Your accuracy score should be higher than 0.8.

Accuracy of using linear kernel (C = 0.1): 0.83

- (20%) Tune the hyperparameters of the *polynomial\_kernel*. Show the accuracy score of the testing data using *polynomial\_kernel* and the hyperparameters you used.

Accuracy of using polynomial kernel (C = 0.1, degree = 4): 0.99

- (20%) Tune the hyperparameters of the *rbf\_kernel*. Show the accuracy score of the testing data using *rbf\_kernel* and the hyperparameters you used.

Accuracy of using rbf kernel (C = 1.5, gamma = 0.5): 0.98

## Part. 2, Questions (50%):

- (20%) Given a valid kernel  $k_1(x, x')$ , prove that the following proposed functions are or are not valid kernels. If one is not a valid kernel, give an example of  $k(x, x')$  that the corresponding  $K$  is not positive semidefinite and shows its eigenvalues.

a.  $k(x, x') = k_1(x, x') + \exp(x^T x')$

a. Let  $\phi(x) = x \Rightarrow K_1(x, x') = \langle \phi(x), \phi(x') \rangle = \phi(x)^T \phi(x') = x^T \cdot x' \Rightarrow$  is valid kernel  
(b.16)  $\exp(x^T \cdot x')$  is a valid kernel  
(b.17)  $k(x, x') = k_1(x, x') + \exp(x^T \cdot x')$  is a valid kernel

b.  $k(x, x') = k_1(x, x') - 1$

b. Let  $k_1(x, x') = x^T \cdot x'$  is a valid kernel and  $x_1 = 0, x_2 = 1$   
 $K(x_1, x_2) = \begin{bmatrix} -1 & 0 \\ 0 & 0 \end{bmatrix} \Rightarrow \begin{cases} \lambda_1 = -1 < 0 \\ \lambda_2 = 0 \end{cases} \Rightarrow k(x, x') = k_1(x, x') - 1$  is not valid

c.  $k(x, x') = \exp(\|x - x'\|^2)$

c. Let  $x_1 = 1, x_2 = 2$   
 $K(x_1, x_2) = \begin{bmatrix} 1 & e \\ e & 1 \end{bmatrix} \Rightarrow \begin{cases} \lambda_1 = 1+e \\ \lambda_2 = 1-e < 0 \end{cases} \Rightarrow k(x, x') = \exp(\|x - x'\|^2)$  is not valid

d.  $k(x, x') = \exp(k_1(x, x')) - k_1(x, x')$

d.  $\exp(k_1(x, x')) = \lim_{n \rightarrow \infty} \left( 1 + k_1(x, x') + \frac{1}{2!} k_1^2(x, x') + \frac{1}{3!} k_1^3(x, x') + \dots + \frac{1}{n!} k_1^n(x, x') \right)$   
 $k(x, x') = \exp(k_1(x, x')) - k_1(x, x') = \lim_{n \rightarrow \infty} \left( 1 + \frac{1}{2!} k_1^2(x, x') + \frac{1}{3!} k_1^3(x, x') + \dots + \frac{1}{n!} k_1^n(x, x') \right)$   
 (6.13) & (6.17) & (6.18)  $\Rightarrow k_2(x, x') = \lim_{n \rightarrow \infty} \left( \frac{1}{2!} k_1^2(x, x') + \frac{1}{3!} k_1^3(x, x') + \dots + \frac{1}{n!} k_1^n(x, x') \right)$  is valid  
 $\phi_1(x)$  is a feature map of  $K_2$ . Let  $\phi: x \mapsto \begin{bmatrix} \phi_1(x) \\ 1 \end{bmatrix}$   
 $\Rightarrow \langle \phi(x), \phi(x') \rangle = \langle \phi_1(x), \phi_1(x') \rangle + 1 = K_2(x, x') + 1 = \exp(k_1(x, x')) - k_1(x, x') + 1 = \exp(k_1(x, x'))$   
 $\Rightarrow k(x, x') = \exp(k_1(x, x')) - k_1(x, x')$  is a valid kernel  $\#$

2. (15%) One way to construct kernels is to build them from simpler ones. Given three possible “construction rules”: assuming  $K_1(x, x')$  and  $K_2(x, x')$  are kernels then so are

- (scaling)  $f(x)K_1(x, x')f(x')$ ,  $f(x) \in \mathbb{R}$
- (sum)  $K_1(x, x') + K_2(x, x')$
- (product)  $K_1(x, x')K_2(x, x')$

Use the construction rules to build a normalized cubic polynomial kernel:

$$K(x, x') = \left( 1 + \left( \frac{x}{\|x\|} \right)^T \left( \frac{x'}{\|x'\|} \right) \right)^3$$

You can assume that you already have a constant kernel  $K_0(x, x') = 1$  and a linear kernel  $K_1(x, x') = x^T x'$ . Identify which rules you are employing at each step.

2.  $K_0(x, x') = 1$   $K_1(x, x') = x^T x'$   
 (scaling)  $K_2(x, x') = f(x)K_1(x, x')f(x')$ ;  $f(x) = \frac{1}{\|x\|} \Rightarrow K_2(x, x') = \left( \frac{x}{\|x\|} \right)^T \left( \frac{x'}{\|x'\|} \right)$   
 (sum)  $K_3(x, x') = K_0(x, x') + K_2(x, x') = 1 + \left( \frac{x}{\|x\|} \right)^T \left( \frac{x'}{\|x'\|} \right)$   
 (product)  $K_4(x, x') = K_3(x, x') \cdot K_3(x, x') \cdot K_3(x, x')$   
 $= \left( 1 + \left( \frac{x}{\|x\|} \right)^T \left( \frac{x'}{\|x'\|} \right) \right)^3 \Rightarrow$  is a valid kernel  $\#$

3. (15%) A social media platform has posts with text and images spanning multiple topics like news, entertainment, tech, etc. They want to categorize posts into these topics using SVMs. Discuss two multi-class SVM formulations: 'One-versus-one' and 'One-versus-the-rest' for this task.

I use OvO to represent One-versus-one and OvR to represent One-versus-the rest in the following questions.

- a. The formulation of the method [how many classifiers are required]

If there are  $n$  classes need to be classified, in OvO formulation, we need to train  $n(n-1) / 2$  classifiers, and in OvR, there are only  $n$  classifiers need to be trained.

- b. Key trade offs involved (such as complexity and robustness).

Because of the number of classifiers need to be trained, OvO has higher complexity on training. But since the classifiers in OvO could distinguish two specific classes better than classifiers in OvR, OvO should be more robust.

- c. If the platform has limited computing resources for the application in the inference phase and requires a faster method for the service, which method is better.

If there are some computing resources limitation, and require a faster method for the service, OvR is a better method, since it has lower cost for training, and lower number of classifiers. This could have faster process time, when it needs to do classifying.