

# NYCU Introduction to Machine Learning, Homework 4

110550130 劉秉驊

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

For this coding assignment, you are required to implement some fundamental parts of the [Support Vector Machine Classifier](#) using only NumPy. After that, train your model and tune the hyperparameter on the provided dataset and evaluate the performance on the testing data.

## (50%) Support Vector Machine

### Requirements:

- Implement the *gram\_matrix* function to compute the [Gram matrix](#) of the given data with an argument **kernel\_function** to specify which kernel function to use.
- Implement the *linear\_kernel* function to compute the value of the linear kernel between two vectors.
- Implement the *polynomial\_kernel* function to compute the value of the [polynomial kernel](#) between two vectors with an argument **degree**.
- Implement the *rbf\_kernel* function to compute the value of the [rbf kernel](#) between two vectors with an argument **gamma**.

### Tips:

- Your functions will be used in the SVM classifier from [scikit-learn](#) like the code below.

```
svc = SVC(kernel='precomputed')
svc.fit(gram_matrix(X_train, X_train, your_kernel), y_train)
y_pred = svc.predict(gram_matrix(X_test, X_train, your_kernel))
```
- For hyperparameter tuning, you can use any third party library's algorithm to automatically find the best hyperparameter, such as [GridSearch](#). In your submission, just give the best hyperparameter you used and do not import any additional libraries/packages.

### Criteria:

1. (10%) Show the accuracy score of the testing data using *linear\_kernel*. Your accuracy score should be higher than 0.8.
2. (20%) Tune the hyperparameters of the *polynomial\_kernel*. Show the accuracy score of the testing data using *polynomial\_kernel* and the hyperparameters you used.

3. (20%) Tune the hyperparameters of the *rbf\_kernel*. Show the accuracy score of the testing data using *rbf\_kernel* and the hyperparameters you used.

(next page)

The following table is the grading criteria for question 2 and 3:

Points	Testing Accuracy
20 points	$0.98 \leq \text{acc}$
15 points	$0.90 \leq \text{acc} < 0.98$
10 points	$0.85 \leq \text{acc} < 0.90$
5 points	$0.8 \leq \text{acc} < 0.85$
0 points	$\text{acc} < 0.8$

Ans:

```
root@MSI:/mnt/c/Users/user/Desktop/Class/3up/ml/hw4# python3 110550130_Hw4.py
Accuracy of using linear kernel (C = 1.5): 0.83
Accuracy of using polynomial kernel (C = 1, degree = 3): 0.98
Accuracy of using rbf kernel (C = 1, gamma = 0.8): 0.99
```

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

1. (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.

- $k(x, x') = k_1(x, x') + \exp(x^T x')$
- $k(x, x') = k_1(x, x') - 1$
- $k(x, x') = \exp(\|x - x'\|^2)$
- $k(x, x') = \exp(k_1(x, x')) - k_1(x, x')$

Ans:

a.

$$1. \quad k(x, x') = k_1(x, x') + e^{(x^T x')}$$

let basic function  $\phi(x) = x$ .

$$\begin{aligned} \text{then } k_2(x, x') &= \phi^T(x) \phi(x') \\ &= x^T x' \end{aligned}$$

$\rightarrow k_2(x, x')$  is valid

$$k(x, x') = k_1(x, x') + e^{k_2(x, x')}$$

$\therefore k_2(x, x')$  is valid. and by 6.16

$\therefore e^{k_2(x, x')}$  is valid

We know  $k_1(x, x')$  is valid

By 6.17.  $k(x, x') = k_1(x, x') + e^{k_2(x, x')}$  is valid  $\neq$

b.

2.

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

$$\text{suppose } k_1(x, x') = x - x'$$

$$x_1 = 1, \quad x_2 = 2$$

$$\begin{aligned} K(x_i, x_j) &= \begin{bmatrix} k(x_1, x_1) & k(x_2, x_1) \\ k(x_1, x_2) & k(x_2, x_2) \end{bmatrix} \\ &= \begin{bmatrix} k_1(x_1, x_1) - 1 & k_1(x_2, x_1) - 1 \\ k_1(x_1, x_2) - 1 & k_1(x_2, x_2) - 1 \end{bmatrix} \\ &= \begin{bmatrix} 0 & 1 \\ 1 & 3 \end{bmatrix} \end{aligned}$$

$$(0 - \lambda)(3 - \lambda) - 1 = 0$$

$$\rightarrow \begin{cases} \lambda_1 = \frac{3 + \sqrt{13}}{2} \\ \lambda_2 = \frac{3 - \sqrt{13}}{2} \end{cases}$$

$$\therefore \lambda_2 < 0$$

$\rightarrow$  not valid  $\#$

c.

3.

$$k(x, x') = e^{-\|x - x'\|^2}$$

let  $K(x_1, x_2)$  be kernel matrix for  $k(x_1, x_2)$

then choose  $x_1 = 0, x_2 = 1$

$$K(x_1, x_2) = \begin{bmatrix} k(x_1, x_1) & k(x_2, x_1) \\ k(x_1, x_2) & k(x_2, x_2) \end{bmatrix}$$
$$= \begin{bmatrix} 1 & e \\ e & 1 \end{bmatrix}$$

$$(1-\lambda)(1-\lambda) - e^2 = 0$$

$$\rightarrow \begin{cases} \lambda_1 = 1 + e \\ \lambda_2 = 1 - e \end{cases}$$

$$\therefore \lambda_2 < 0$$

$\therefore$  not valid  $\nexists$

d.

$$4. \quad k(x, x') = e^{k_1(x, x')} - k_1(x, x')$$

$$e^{k_1(x, x')} = 1 + k_1(x, x') + \frac{1}{2!} k_1^2(x, x') + \frac{1}{3!} k_1^3(x, x') + \dots$$

$$\text{then } k(x, x') = 1 + \frac{1}{2!} k_1^2(x, x') + \frac{1}{3!} k_1^3(x, x') + \dots$$

$$\text{By 6.13, 6.17, 6.18, and let } k_2(x, x') = k_1(x, x') - 1.$$

then  $k_2(x, x')$  is valid

By 6.15, let  $g(x) = x + 1$ , for  $g$  is a polynomial with non-negative coefficients

$$\text{then } k(x, x') = g(k_2(x, x')) = k_2(x, x') + 1 \text{ is valid}$$

$$\Rightarrow k(x, x') = e^{k_1(x, x')} - k_1(x, x') \text{ is valid}$$

(next page)

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.

Ans:

2.

1. By scaling, let  $f(x) = \frac{1}{\|x\|} \in \mathbb{R}$

$$k_2(x, x') = \frac{1}{\|x\|} k_1(x, x') \frac{1}{\|x'\|}$$

$$= \frac{1}{\|x\|} (x^T x') \frac{1}{\|x'\|}$$

$$= \left( \frac{x}{\|x\|} \right)^T \cdot \left( \frac{x'}{\|x'\|} \right)$$

2. By sum, let

$$k_3(x, x') = k_0(x, x') + k_2(x, x')$$

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

3. By product, let

$$k_4(x, x') = k_3(x, x') \cdot k_3(x, x')$$

$$= \left( 1 + \left( \frac{x}{\|x\|} \right)^T \left( \frac{x'}{\|x'\|} \right) \right)^2$$

4. By product, let

$$k_5(x, x') = k_4(x, x') \cdot k_3(x, x')$$

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

~~✗~~

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.
- The formulation of the method [how many classifiers are required]
  - Key trade offs involved (such as complexity and robustness).
  - 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.

Ans:

a. Formulation:

In one versus one, we need to train  $n$  classes with  $n*(n-1)/2$  classifiers, and then choose the one classifier with the most votes.

In one versus the rest, we need to train  $n$  classes with  $n$  classifiers, and then each class predicts whether the testing data belongs to itself, that's the highest confidence.

b. Key trade off:

For complexity, we can see one versus one trains  $n*(n-1)/2$  classifiers much more than one versus the rest having  $n$  classifiers. One versus one costs more.

For robustness, since one versus one compares each other classifiers and one versus the rest is hard to handle overlapped classes, thus one versus one has better robustness.

c. Resource constraints:

As we can see, one versus one has  $n * (n - 1) / 2$  classifiers that are much more than one versus the rest having  $n$  classifiers to train, which one versus the rest computes less, and is better at facing resource constraints.