

ML Homework 05-2: SVM

- Student ID: 309553002
- Name: 林育愷

ML Homework 05-2: SVM

Code with Detailed Explanations

Prerequisites

Input Data

Part 1

Part 2

Part 3

Experiments Settings and Results

Part 1

Part 2

Part 3

Observations and Discussion

References

Code with Detailed Explanations

Prerequisites

I use Python 3.6 for this implementation based on platform Ubuntu 18.04, with the following packages

- NumPy,
- SciPy,
- Matplotlib,
- seaborn, and
- libsvm (<https://github.com/cjlin1/libsvm>, embedded in this project).

Input Data

We use MNIST as our dataset. As it is of csv format, we can simply import them with NumPy.

```
1 import numpy as np
2
3
4 # image_filename = 'X_train.csv' or 'X_test.csv'
5 # label_filename = 'Y_train.csv' or 'Y_test.csv'
6 def load_data(image_filename, label_filename):
7     images = np.genfromtxt(image_filename, delimiter=',')
8     labels = np.genfromtxt(label_filename, delimiter=',').flatten()
9     return images, labels
```

Part 1

According to the documentation ¹, LIBSVM supports the following types of kernel functions:

- linear: $u^T v$,
- polynomial: $(\gamma u^T v + c_0)^d$,
- radial basis function (RBF): $\exp(-\gamma \|u - v\|^2)$, and
- sigmoid: $\tanh(\gamma u^T v + c_0)$.

In the first part we are going to perform linear, polynomial, and RBF, so we define ones corresponding arguments.

```
1 import enum
2
3
4 class KernelType(enum.Enum):
5     LINEAR = '-t 0'
6     POLYNOMIAL = '-t 1'
7     RADIAL_BASIS_FUNCTION = '-t 2'
8     PRECOMPUTED_KERNEL = '-t 4'
```

Later we feed them into solver and find out the corresponding results.

```
1 import time
2 import libsvm.python.svmutil as svmutil
3
4
5 for ktype in KernelType:
6     # skip custom mode
7     if ktype == KernelType.PRECOMPUTED_KERNEL:
8         continue
9
10    param = ktype.value + ' -q'
11    start = time.perf_counter()
12    model = svmutil.svm_train(train_label, train_image, param)
13    end = time.perf_counter()
14    duration = end - start
15
16    _, accuracy, _ = svmutil.svm_predict(test_label, test_image, model)
17    print('[%s] accuracy: %.2f%%, mse: %.2f, time: %.2f sec' %
18          (ktype.name, accuracy[0], accuracy[1], duration))
```

Part 2

In the second part we are going to fine-tune the parameters of C-SVC with RBF kernel function. It means there are C and γ that we are going to figure out.

Note that each run takes time and only a processor is used, so I use multiprocessing to boost the running time.

```
1 import multiprocessing as mp
2 import numpy as np
3
4
5 NUMBER_OF_CPUS = mp.cpu_count()
```

```

6
7 ktype = KernelType.RADIAL_BASIS_FUNCTION
8 cs = np.exp(np.arange(-5, 3))
9 gammas = np.exp(np.arange(-9, 0))
10 accuracy_matrix = np.zeros(cs.shape + gammas.shape)
11
12 for i, c in enumerate(cs):
13     args_list = []
14     for j, gamma in enumerate(gammas):
15         param = '-q %s -v 3 -s 0 -c %f -g %f' % (ktype.value, c, gamma)
16         svm_args = (train_label, train_image, param)
17         args_list.append(svm_args)
18
19     with mp.Pool(NUMBER_OF_CPUS) as pool:
20         results = pool.starmap(svmutil.svm_train, args_list)
21
22     accuracy_matrix[i] = results

```

Part 3

In the final part we are going to build a custom kernel function. In this assignment the kernel function is set to be

$$k(x, x') = k_{\text{linear}}(x, x') + k_{\text{RBF}}(x, x').$$

For efficiency we separate distance function and kernel function in RBF so as to improve the running time. Also by documentation the format for each row should be `[i, k(xi, x1), k(xi, x2), ..., k(xi, xn)]`.

```

1 import numpy as np
2 from scipy.spatial.distance import cdist
3
4
5 def calc_feature_distance(x1, x2):
6     return cdist(x1, x2, 'sqeuclidean')
7
8
9 def linear_plus_rbf_kernel(x1, x2, feature_distance, gamma):
10     linear = x1 @ x2.T
11     # https://www.csie.ntu.edu.tw/~cjlin/libsvm/
12     rbf = np.exp(-gamma * feature_distance)
13     param = np.hstack(
14         (np.arange(x1.shape[0])[:, np.newaxis] + 1, linear + rbf))
15     return param

```

Again we use multiprocessing to run multiple cases of γ so as to find a good configuration.

```

1 import multiprocessing as mp
2 import numpy as np
3
4
5 ktype = KernelType.PRECOMPUTED_KERNEL
6 gammas = np.exp(np.arange(-9, 0))
7 args_list = []
8
9 train_feature_distance = calc_feature_distance(train_image,

```

```

10                                     train_image)
11 test_feature_distance = calc_feature_distance(test_image, test_image)
12
13 for j, gamma in enumerate(gammas):
14     param = '-q %s' % (ktype.value)
15     train_kernel = linear_plus_rbf_kernel(train_image, train_image,
16                                         train_feature_distance,
17                                         gamma)
18     test_kernel = linear_plus_rbf_kernel(test_image, test_image,
19                                         test_feature_distance, gamma)
20
21     svm_args = (train_label, train_kernel, test_label, test_kernel,
22                param)
23     args_list.append(svm_args)
24
25 with mp.Pool(NUMBER_OF_CPUS) as pool:
26     accuracy_matrix = pool.starmap(train_and_predict, args_list)
27
28 accuracy_matrix = np.array([accuracy_matrix])

```

Experiments Settings and Results

For visualization we use `seaborn` to colorize the accuracy matrix.

```

1 import seaborn as sns
2
3
4 sns.heatmap(accuracy_matrix, xticklabels=np.log(gammas), annot=True)

```

Part 1

We see that RBF and linear have similar accuracy of classification problem but polynomial kernel function get poor result and takes lots of time.

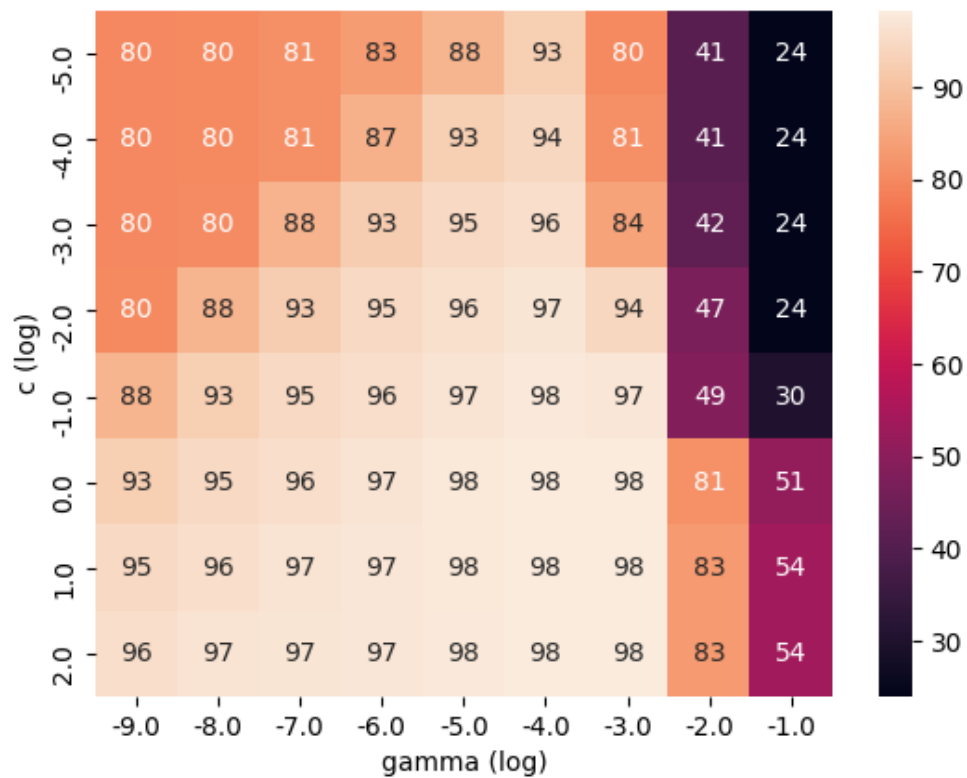
```

1 PART I
2 -----
3 [LINEAR] accuracy: 95.08%, mse: 0.14, time: 1.80 sec
4 [POLYNOMIAL] accuracy: 34.68%, mse: 2.62, time: 21.00 sec
5 [RADIAL_BASIS_FUNCTION] accuracy: 95.32%, mse: 0.15, time: 3.92 sec

```

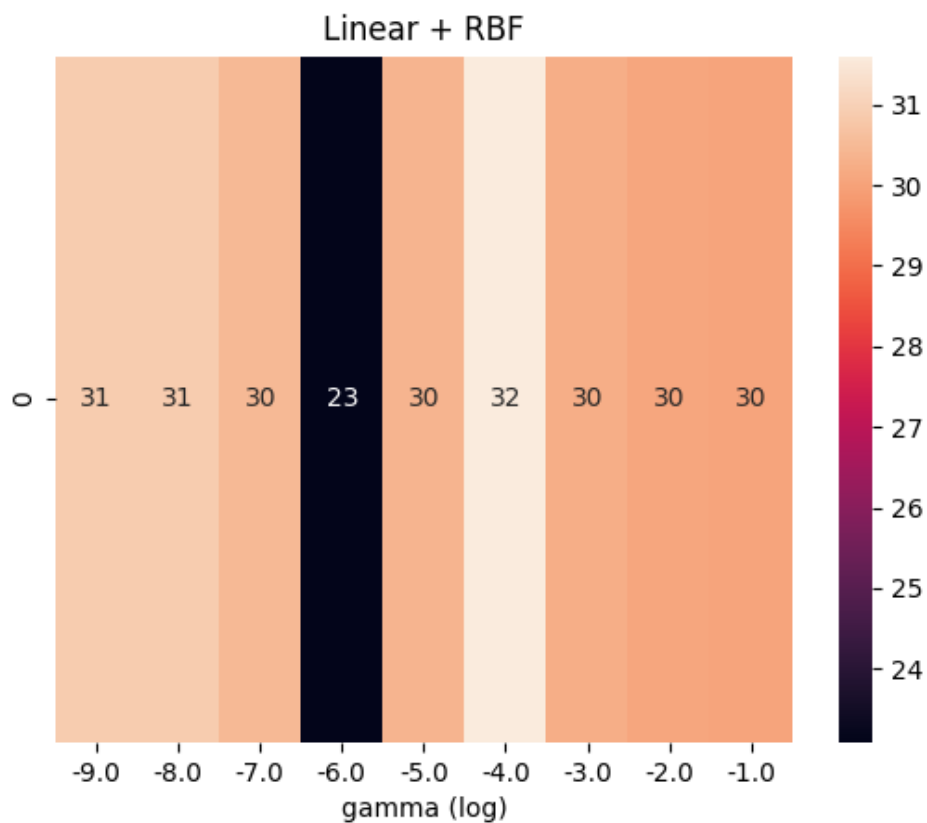
Part 2

In part 2 we find that with $(C, \gamma) = (\exp(2), \exp(-4))$ a better result is obtained (with cross-validation 98.34%).



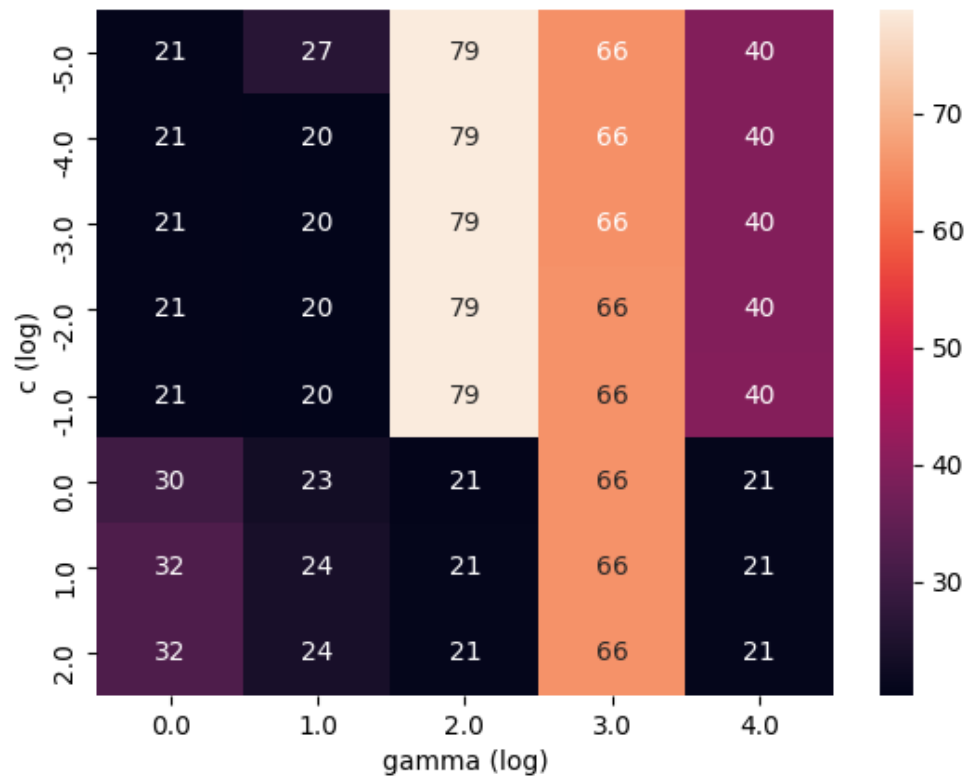
Part 3

We find that the given kernel function for this problem cannot get a better result.



Observations and Discussion

In part 2 we find that with larger C we have better performance. In general the larger γ (greater than 1), the poor result, but there are some exceptions.



References

1. <https://github.com/cjlin1/libsvm> ↗