

ML Homework 5

Bo-Cheng Huang, 309652022

1 Gaussian Process

1.1 Code with detailed explanations

For Part 1, I use a rational quadratic kernel to compute similarities between different points. It is parameterized by a length scale parameter $l > 0$ and a scale mixture parameter $\alpha > 0$. The formula of a rational quadratic kernel is

$$k(x_i, x_j) = \left(1 + \frac{d(x_i, x_j)^2}{2\alpha l^2}\right)^{-\alpha}$$

where d is the Euclidean distance.

```
1 def kernel(X1,X2,l,alpha):
2     '''
3     Impletation of rational quadratic kernel
4     X1: (n)-array
5     X2: (m)-array
6     return a (n,m)-array
7     '''
8     sqddist=np.power(X1.reshape(-1,1)-X2.reshape(1,-1),2)
9     return np.power((1+sqddist/(2*alpha*l**2)),-alpha)
```

Listing 1: Kernel

Hence I can predict the distribution of f by rational quadratic kernel. For $\mathbf{f} = [f(x_1), \dots, f(x_N)]^T$ and $\mathbf{y} = [y_1, \dots, y_N]^T$ let

$$p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \beta^{-1}\mathbf{I}_N), \quad p(\mathbf{f}) = \mathcal{N}(\mathbf{0}, \mathbf{K}), \quad p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{C}),$$

where $\mathbf{C}(\mathbf{x}_n, \mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m) + \beta^{-1}\delta_{nm}$. Then the conditional distribution $p(y^*|\mathbf{y})$ is a Gaussian distibution with:

$$\begin{aligned}\mu(\mathbf{x}^*) &= k(\mathbf{x}, \mathbf{x}^*)^T \mathbf{C}^{-1} \mathbf{y}, \\ \sigma^2(\mathbf{x}^*) &= k^* - k(\mathbf{x}, \mathbf{x}^*)^T \mathbf{C}^{-1} k(\mathbf{x}, \mathbf{x}^*), \\ k^* &= k(\mathbf{x}^*, \mathbf{x}^*) + \beta^{-1}.\end{aligned}$$

So I can do the Gaussian process of predicting the distribution of f in Listing 2 and Listing 3 is the code for data preprocessing and visualization:

```

1  def predict_GP(X_old,X_new,y,C,beta,l,alpha):
2      '''
3          Use train data (old) and test data (new) to compute new mean
4          and variance.
5          X_old: (n)-array, train data
6          X_new: (m)-array, test data
7          y: (n)-array, ground truth
8          K: (n,n)-array, covariance matrix C
9          beta, l, alpha: parameters
10         return (len(X_new),1)-array, (len(X_new),1)-array
11         '''
12         ker_on=kernel(X_old,X_new,l,alpha)
13         ker_nn=kernel(X_new,X_new,l,alpha)
14         mean=ker_on.transpose()@inv(C)@y.reshape(-1,1)
15         ker_star=ker_nn+1/beta*np.identity(len(ker_nn))
16         var=ker_star-ker_on.transpose()@inv(C)@ker_on
17         return mean,var

```

Listing 2: Gaussian Process

```

1  X=np.zeros((34,1))
2  Y=np.zeros((34,1))
3  m=0
4  input_file=open('input.data','rb')
5  line=input_file.readline()
6  while line:
7      s=line.decode("utf-8")
8      x,y=s.split(' ')
9      X[m][0]=float(x)
10     Y[m][0]=float(y)
11     line=input_file.readline()
12     m+=1
13
14  def print_graph():
15     plt.plot(xlin,mean,'black')
16     plt.xlim(-60,60)
17     plt.fill_between(xlin,mean+2*vars,mean-2*vars,color='aquamarine')
18     plt.scatter(X,Y)
19     plt.show()

```

Listing 3: Data preprocessing and visualization

To get the visualization, let `xlin` be a sample set which contains data points in the range $[-60, 60]$, then I can do prediction between origin data and sample to get the graph.

```

1  beta=5
2  C=kernel(X,X,1,1)+1/beta*np.identity(len(X))
3  xlin=np.linspace(-60,60,200).transpose()
4  mean,vars=predict_GP(X,xlin,Y,C,beta,1,1)
5  mean=mean.reshape(-1)
6  vars=np.sqrt(np.diag(vars))
7  print_graph()

```

Listing 4: Main function of Part 1

For Part 2, I need to optimize the kernel parameters by minimizing negative marginal log-likelihood (for rational quadratic kernel, $\theta = [l, \alpha]$),

$$p(\mathbf{y}|\theta) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{C}_\theta),$$

$$-\ln p(\mathbf{y}|\theta) = \frac{1}{2}\ln|\mathbf{C}_\theta| + \frac{1}{2}\mathbf{y}^T \mathbf{C}_\theta^{-1} \mathbf{y} + \frac{N}{2}\ln(2\pi).$$

Using cholesky from numpy.linalg to get the cholesky decomposition of C . Then use L to compute the determinant of C . $\text{inv}(C)$ is the inverse matrix of C .

```

1  def nll_fn(X_train,Y_train,beta):
2      '''
3      X_train: (n)-array
4      Y_train: (n)-array
5      beta: parameter
6      return the value of negative log-likelihood in GP
7      '''
8      Y_train = Y_train.ravel()
9      def obj(theta):
10         C=kernel(X_train,X_train,theta[0],theta[1])\
11         +1/beta*np.identity(len(X_train))
12         L=cholesky(C)
13         value=0.5*np.sum(np.log(np.diagonal(L)))\
14         +0.5*Y_train.transpose()@inv(C)@Y_train\
15         +0.5*len(X_train)*np.log(2*np.pi)
16         return value
17     return obj

```

Listing 5: Negative marginal log-likelihood

Next, I use minimize from scipy.optimize to get the parameters l, α . Using different initial value of l and α in $[0.01, 0.1, 1, 10, 100]$ to avoid some bigger local minimum. And set `Lopt` and `alpha_opt` to predict y .

```

1  value=1e9
2  for i in range(-2,3):
3      for j in range(-2,3):
4          res=minimize(nll_fn(X,Y,beta),[10**i, 10**j],\
5                      bounds=((1e-5, 1e5),(1e-5, 1e5)))
6          if res.fun<value:
7              value=res.fun
8              l_opt,alpha_opt=res.x
9
10 beta=5
11 C=kernel(X,X,l_opt,alpha_opt)+1/beta*np.identity(len(X))
12 xlin=np.linspace(-60,60,200).transpose()
13 mean,vars=predict_GP(X,xlin,Y,C,beta,l_opt,alpha_opt)
14 mean=mean.reshape(-1)
15 vars=np.sqrt(np.diag(vars))
16 print_graph()

```

Listing 6: Main function of Part 2

1.2 Experiments settings and results

In Part 1, I set xlin equal to $\text{np.linspace}(-60,60,100)$, $\text{np.linspace}(-60,60,200)$ and $\text{np.linspace}(-60,60,500)$. And the visualization of them are below:

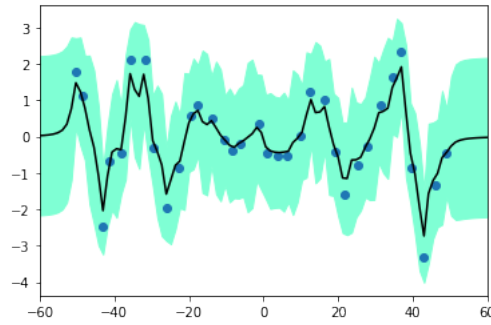


Figure 1: $\text{xlin}=\text{np.linspace}(-60,60,100)$

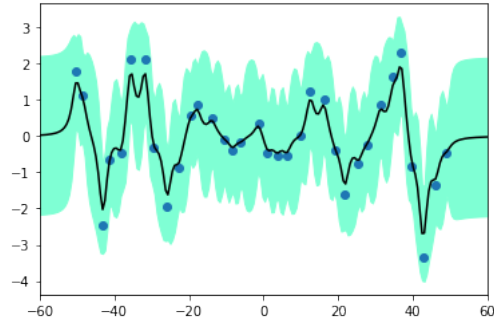


Figure 2: $xlin=np.linspace(-60,60,200)$

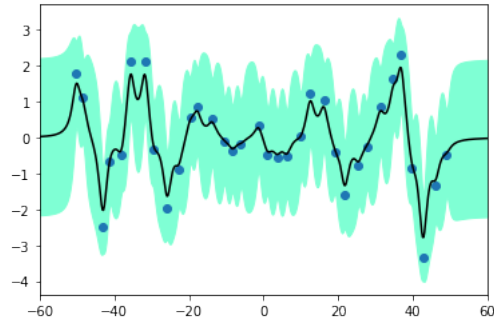


Figure 3: $xlin=np.linspace(-60,60,500)$

In Part 2, I found that $l_{opt}=2.4468811091753806$ and $\alpha_{opt}=901.503494930027$. And the visualization of them are below:

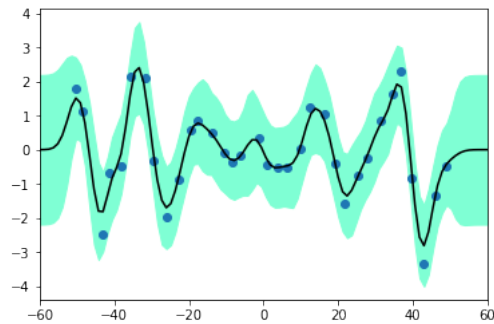


Figure 4: $xlin=np.linspace(-60,60,100)$

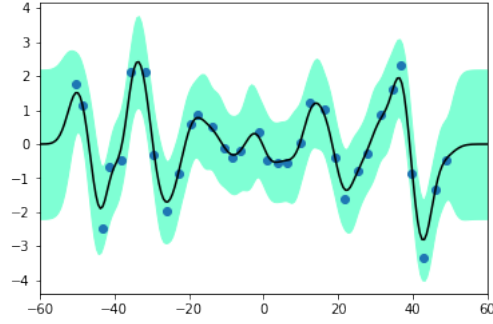


Figure 5: `xlin=np.linspace(-60,60,200)`

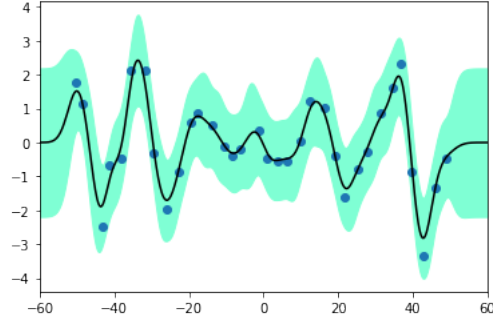


Figure 6: `xlin=np.linspace(-60,60,500)`

1.3 Observations and discussion

From the visualization of different `xlin`, when the parameters $l = \alpha = 1$, I think that more sample points would make the distribution more smoothing. After using `l_opt` and `alpha_opt`, the variance of some sample points is smaller than the case $l = \alpha = 1$.

2 SVM

2.1 Code with detailed explanations

First we need to transform data into array.

```
1 def handle_with_csv(x,y,str):
2     data=np.zeros((x,y))
3     with open(str,newline='') as csvfile:
4         rows = csv.reader(csvfile)
5         m=0
6         for row in rows:
7             for i in range(y):
8                 data[m][i]=float(row[i])
9             m+=1
10    return data
11
12 train_x=handle_with_csv(5000,784,'X_train.csv')
13 train_y=handle_with_csv(5000,1,'Y_train.csv')
14 test_x=handle_with_csv(2500,784,'X_test.csv')
15 test_y=handle_with_csv(2500,1,'Y_test.csv')
```

Listing 7: Data preprocessing

In LIBSVM library, I use `svm_train` and `svm_predict` to do Support Vector Machine. The default type of SVM in LIBSVM is C-SVC and there are 4 types of kernel for users: linear, polynomial, radial basis function and sigmoid.

- linear: $K(x_i, x_j) = x_i^T x_j$.
- polynomial: $K(x_i, x_j) = (\gamma x_i^T x_j + r)^d$, $\gamma > 0$.
- radial basis function: $K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$, $\gamma > 0$.
- sigmoid: $K(x_i, x_j) = \tanh(\gamma x_i^T x_j + r)$.

For Part 1, I use `svm_train` with options `-t` to test different kernels and use `svm_predict` to get the accuracy of each model.

```
1 K_type={'linear':'-t 0','polynomial':'-t 1','RBF':'-t 2'}
2 for parameters in K_type:
3     model=svm_train(train_y.reshape(-1),train_x,K_type[parameters])
4     p_label,p_acc,p_val=svm_predict(test_y.reshape(-1),test_x,model)
5     print('%0.2f'%(p_acc[0]),end=' ')
```

Listing 8: Main function of Part 1

For Part 2, to do grid search for finding parameters of the best performing model, I decide some ranges of the parameters of each kernel function. C-SVC is soft-margin SVM, it contains a parameter C . In polynomial and RBF kernel, there is a parameter γ . (I think that the parameters $r = 0$, $d = 3$ in polynomial kernel have less effect to SVM.)

```

1  for parameters in K_type:
2      print(parameters)
3      if parameters=='linear':
4          for c in range(-9,1):
5              model=svm_train(train_y.reshape(-1),train_x,\
6                  '-c '+str(2**c)+' '+K_type[parameters])
7              p_label,p_acc,p_val=svm_predict(test_y.reshape(-1),test_x,model)
8              print('%.2f'%(p_acc[0]),end=' ')
9      elif parameters=='polynomial':
10         for c in range(-5,5):
11             for gamma in range(-9,-3):
12                 model=svm_train(train_y.reshape(-1),train_x,\
13                     '-c '+str(2**c)+' -g '+str(2**gamma)+' '+K_type[parameters])
14                 p_label,p_acc,p_val=svm_predict(test_y.reshape(-1),test_x,model)
15                 print('%.2f'%(p_acc[0]),end=' ')
16             print()
17         elif parameters=='RBF':
18             for c in range(-3,6):
19                 for gamma in range(-11,-3):
20                     model=svm_train(train_y.reshape(-1),train_x,\
21                         '-c '+str(2**c)+' -g '+str(2**gamma)+' '+K_type[parameters])
22                     p_label,p_acc,p_val=svm_predict(test_y.reshape(-1),test_x,model)
23                     print('%.2f'%(p_acc[0]),end=' ')
24                 print()
25         print()

```

Listing 9: Main function of Part 2

In Part 3, I construct a user-defined kernel by linear kernel + RBF kernel. After adding two kernels, it remains to add a column of indexes from 1 to len(X) in front of kernel (the first column).


```

1 def computed_ker(X,gamma):
2     '''
3     X: (n,m)-array
4     gamma: parameter
5     return (n,n)-array defined by linear+RBF
6     '''
7     linear=X@X.transpose()
8     RBF=dis.squareform(np.exp(-gamma*dis.pdist(X,'sqeuclidean')))
9     ker=linear+RBF
10    ker=np.hstack((np.arange(1,len(X)+1).reshape(-1,1),ker))
11    return ker

```

Listing 10: User-defined kernel

I select $\gamma = 2^{-4}$ with good performance in RBF kernel. Then use `svm_problem` to transform `train_ker` into a legal type to run `svm_train`.

```

1 train_ker=computed_ker(train_x,2**-4)
2 add=svm_problem(train_y.reshape(-1),train_ker,isKernel=True)
3 p=svm_parameter('-t 4')
4 model=svm_train(add,p)
5 test_ker=computed_ker(test_x,2**-4)
6 p_label,p_acc,p_val=svm_predict(test_y.reshape(-1),test_ker,model)

```

Listing 11: Main function of Part 3

2.2 Experiments settings and results

In Part 1, with default parameters, the score of linear, polynomial and RBF are 95.08, 34.68, 95.32, respectively. It has bad performance on default polynomial kernel.

In Part 2, for linear kernel, I select C from 2^{-9} to 2^0 , then I observe that the best score is 96.04 when $C = 2^{-4}$.

C	2^{-9}	2^{-8}	2^{-7}	2^{-6}	2^{-5}	2^{-4}	2^{-3}	2^{-2}	2^{-1}	2^0
	95.16	95.52	95.84	95.92	96.00	96.04	95.92	95.80	95.52	95.08

Table 1: Hyperparameters in linear kernel

For polynomial kernel, I select C from 2^{-5} to 2^5 and γ from 2^{-10} to 2^{-4} , then I observe that the best score is 97.84 where $C = 2^{-5}$ and $\gamma = 2^{-4}$; $C = 2^{-2}$ and $\gamma = 2^{-5}$; ... Besides, I observe a better score 97.92 where degree $d = 2$ with $C = 2^{-2}$ and $\gamma = 2^{-5}$.

C, γ	2^{-10}	2^{-9}	2^{-8}	2^{-7}	2^{-6}	2^{-5}	2^{-4}
2^{-5}	28.88	28.88	33.56	74.88	92.08	97.04	97.84
2^{-4}	28.88	28.88	43.88	83.40	93.64	97.52	97.48
2^{-3}	28.88	28.88	61.92	88.84	95.16	97.80	97.48
2^{-2}	28.88	33.56	74.88	92.08	97.04	97.84	97.48
2^{-1}	28.88	43.88	83.40	93.64	97.52	97.48	97.48
2^0	28.88	61.92	88.84	95.16	97.80	97.48	97.48
2^1	33.56	74.88	92.08	97.04	97.84	97.48	97.48
2^2	43.88	83.40	93.64	97.52	97.48	97.48	97.48
2^3	61.92	88.84	95.16	97.80	97.48	97.48	97.48
2^4	74.88	92.08	97.04	97.84	97.48	97.48	97.48

Table 2: Hyperparameters in polynomial kernel

For RBF kernel, I select C from 2^{-3} to 2^5 and γ from 2^{-11} to 2^{-4} , then I observe that the best score is 98.52 where $C \geq 1$ and $\gamma = 2^{-4}$.

C, γ	2^{-11}	2^{-10}	2^{-9}	2^{-8}	2^{-7}	2^{-6}	2^{-5}	2^{-4}
2^{-3}	90.80	92.64	94.04	94.96	95.48	96.44	96.60	88.20
2^{-2}	92.80	94.08	94.92	95.44	96.12	97.12	97.32	93.44
2^{-1}	94.04	94.68	95.32	95.72	96.72	97.52	98.04	96.40
2^0	94.68	95.24	95.60	96.48	97.32	98.08	98.52	97.36
2^1	95.16	95.48	96.24	96.88	97.64	98.28	98.52	97.44
2^2	95.52	95.96	96.32	97.24	97.92	98.40	98.52	97.44
2^3	95.92	96.12	96.80	97.44	98.04	98.40	98.52	97.44
2^4	96.08	96.32	97.04	97.64	98.04	98.44	98.52	97.44
2^5	96.20	96.52	97.24	97.56	98.04	98.44	98.52	97.44

Table 3: Hyperparameters in RBF kernel

In Part 3, the score of linear+RBF kernel with $\gamma = 2^{-4}$ is 28.12. It is not good for this data set with respect to other kernels.

2.3 Observations and discussion

In Part 1, this model has bad performance on default polynomial kernel. In Part 2, when doing the grid search for finding parameters, I have set an initial range (e.g. $C = 2^{-5}, 2^{-3}, \dots, 2^{15}$ and $\gamma = 2^{-15}, 2^{-13}, \dots, 2^3$ in "A Practical Guide to Support Vector Classification") and using the results to tight a new range. In linear kernel, there is less difference from the selection of C . In polynomial kernel, by the results of grid search, there may exist some ratio with γ and C . So there are a few of combinations of γ and C with same score. In RBF kernel, the performance of smaller γ is better than $\gamma \geq 1$ and the score converges for bigger C with some γ (e.g. $2^{-4}, 2^{-5}, \dots$). In Part 3,

although $\gamma = 2^{-4}$ has good performance in RBF kernel, it may not be good in linear+RBF kernel. After trying some parameters of linear+RBF kernel, it still has bad performance for this data set ($\leq 30\%$).