

---

# Machine Learning Homework 5 Report

陳以瑄 (ID: 313551047)

---

## 1 Gaussian Proccess

### 1.1 Code

#### 1.1.1 Task 1

##### 1) Prepare data

Loading the training set from 'input.data' into the NumPy arrays 'X\_train' and 'Y\_train', and generating 1000 evenly spaced data points over the range [-60, 60] for the test set.

```
1 def LoadData(path):
2     X = []
3     Y = []
4     with open(path, 'r') as file:
5         for data in file.readlines():
6             xy = data.split()
7             X.append(float(xy[0]))
8             Y.append(float(xy[1]))
9     X = np.array(X).reshape(-1, 1)
10    Y = np.array(Y).reshape(-1, 1)
11    return X,Y
12
13 X_train,Y_train = LoadData("./data/input.data")
14 train_num = len(X_train)
15
16 test_num = 1000
17 X_test = np.linspace(-60, 60, test_num).reshape(-1, 1)
```

##### 2) Rational quadratic kernel

In this task, we use the rational quadratic kernel to compute the similarities between different data points. The formula is given as:

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

where  $\sigma$ , the signal variance, determines the overall amplitude of the kernel;  $l$ , the length scale, controls how quickly correlations between data points decay; and  $\alpha$ , the shape parameter, adjusts the flexibility.

For the initial parameter set, I set  $\sigma=1$ ,  $l=1$ , and  $\alpha=1$ .

```
1 sigma = 1
2 length = 1
3 alpha = 1

1 def RQK(xi,xj,sigma, l,a):
2     return (sigma**2)*(1+(xi-xj)**2/(2*a*l**2))**(-a)
```

### 3) Covariance matrix

The formula of covariance matrix of data is:

$$C(x_n, x_m) = k(x_n, x_m) + \beta^{-1} \delta_{n,m}$$

```
1 def Cov(X,beta, sigma, length, alpha):
2     N = len(X)
3     C = np.zeros((N,N))
4     for n in range(N):
5         for m in range(N):
6             C[n][m] = RQK(X[n],X[m], sigma, length, alpha)
7             if n==m:
8                 C[n][m]+=1/beta
9     return C
```

### 4) Gaussian process

After calculating the training set covariance matrix (C), we can make predictions. For a new testing data point, we compute the kernel function values between the test point and each training data point. Using these values, we calculate the mean and variance of the predictive distribution using the following formulas:

$$\mu(x^*) = k(x, x^*)^T C^{-1} y$$

$$\sigma(x^*) = k(x^*, x^*) + \beta^{-1} - k(x, x^*)^T C^{-1} k(x, x^*)$$

```
1 def Predict(X_train,Y_train,X_test, train_num, test_num,beta, sigma, length, alpha,C):
2     mean = np.zeros(test_num)
3     var = np.zeros(test_num)
4
5     for m in range (test_num):
6         # Calculate k(x,x*)
7         Km = np.zeros((train_num,1))
8         for n in range(train_num):
9             Km[n][0] = RQK(X_train[n],X_test[m], sigma, length, alpha)
10
11         # calculate mean and var
12         # mean = k(x,x*)T C-1 y
13         mean[m] = mul(mul(Km.T, inv(C)),Y_train)
14         # var = k(x*,x*) + B-1 - k(x,x*)T C-1 k(x,x*)
15         K = RQK(X_test[n],X_test[n], sigma, length, alpha)
16         var[m] = K + 1/beta - mul(mul(Km.T, inv(C)),Km)
17     return mean, var
```

### 5) Visualization

```
1 def Viz(X,Y,mean,var, test_num, sigma, l, a):
2     fig = plt.figure()
3     x_range = np.linspace(-60, 60, test_num)
4     lower_bound = mean - 1.96 * np.sqrt(var)
5     upper_bound = mean + 1.96 * np.sqrt(var)
6
7     plt.figure(figsize=(10, 6))
8     plt.scatter(X, Y, color="blue", label="Training Data Points", alpha=0.7)
9     plt.plot(x_range, mean, color="red", label="Prediction Mean", linewidth=2)
10    plt.fill_between(x_range, lower_bound, upper_bound, color="red", alpha=0.3,
11                    label="95% Confidence Interval")
12
13    plt.xlim(-60, 60)
14    plt.title(f"sigma: {sigma:.4f}, length scale: {l:.4f}, alpha: {a:.4f}", fontsize=14)
15    plt.legend(loc='upper left', bbox_to_anchor=(1, 1))
16    plt.show()
```

### 1.1.2 Task 2

#### 1) Negative marginal log-likelihood

The negative marginal log-likelihood is calculated using the following formulas:

$$likelihood = \frac{1}{|2\pi|^{\frac{N}{2}}} \frac{1}{|C_{\theta}|^{\frac{1}{2}}} e^{-\frac{1}{2}y^T C_{\theta}^{-1}y}$$

$$loglik = -\frac{N}{2}\ln(2\pi) - \frac{1}{2}\ln(C_{\theta}) - \frac{1}{2}y^T C_{\theta}^{-1}y$$

$$negloglik = \frac{N}{2}\ln(2\pi) + \frac{1}{2}\ln(|C_{\theta}|) + \frac{1}{2}y^T C_{\theta}^{-1}y$$

```
1 def NegLogLik(theta, X, Y, beta, N):
2     theta = theta.reshape(len(theta), 1)
3     sigma = theta[0]
4     length = theta[1]
5     alpha = theta[2]
6
7     C_theta = Cov(X, beta, sigma, length, alpha)
8
9     NLL = N * np.log(2 * np.pi) / 2
10    NLL += np.log(np.linalg.det(C_theta)) / 2
11    NLL += mul(mul(Y.T, inv(C_theta)), Y) / 2
12
13    return float(NLL[0])
```

#### 2) Optimize the parameter

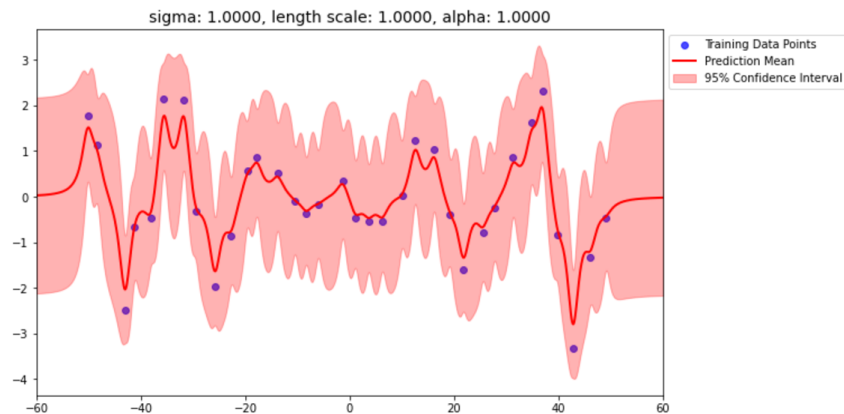
I use `scipy.optimize.minimize` to find the optimal parameters ( $\sigma$ ,  $l$ ,  $\alpha$ ) that minimize the negative marginal log-likelihood value. The initial guesses for the parameters are  $\sigma=1$ ,  $l=1$ ,  $\alpha=1$ . The function `NegLogLik` is used as the objective function, with additional arguments.

```
1 from scipy.optimize import minimize
2 Opt = minimize(NegLogLik, [sigma, length, alpha], args = (X_train, Y_train, beta, train_num))
3 sigma = Opt.x[0]
4 length = Opt.x[1]
5 alpha = Opt.x[2]
```

## 1.2 Experiments

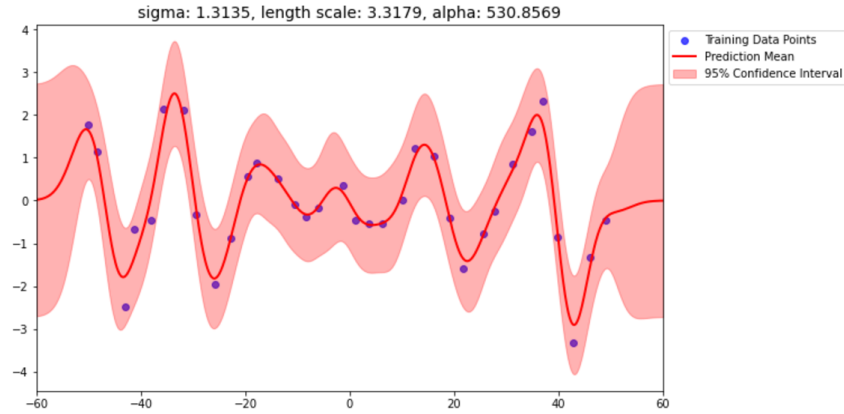
### 1.2.1 Task 1

The result of task 1, with  $\sigma=1$ ,  $l=1$ ,  $\alpha=1$ :



### 1.2.2 Task 2

The optimal parameter I found are:  $\sigma=1.3135$ ,  $l=3.3179$ ,  $\alpha=530.8569$ .



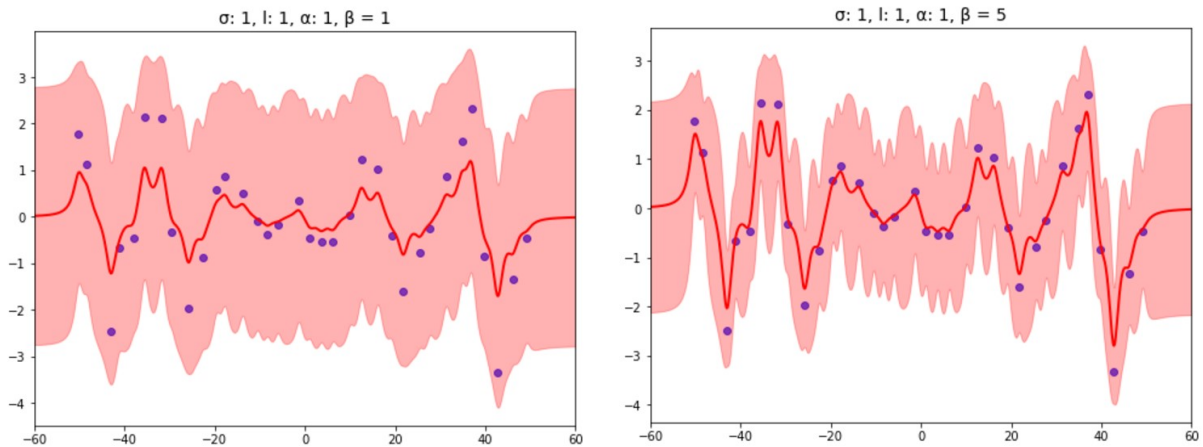
## 1.3 Observation and Discussion

### 1) Comparing the result of two tasks

The initial model has wider confidence intervals, especially in regions with sparse training data, meaning that it did not fully capture the underlying data structure. The optimized model has smaller and smoother confidence intervals, with a prediction mean that fits better with the training data points.

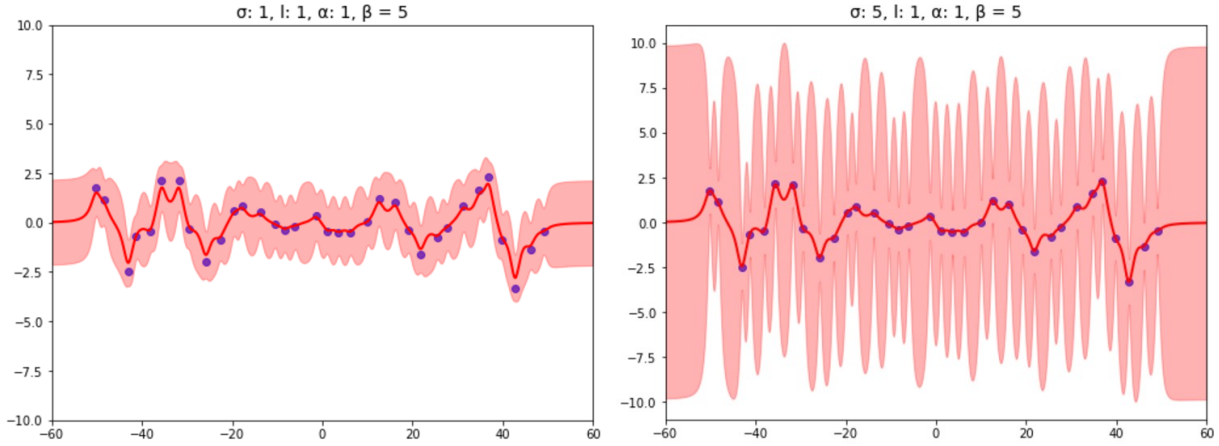
### 2) Impact of $\beta$ :

From the figure below, we can see that when  $\beta$  is lower, which means the variance of the random noise is larger, the prediction mean is farther from the training data points, and the confidence interval is wider.



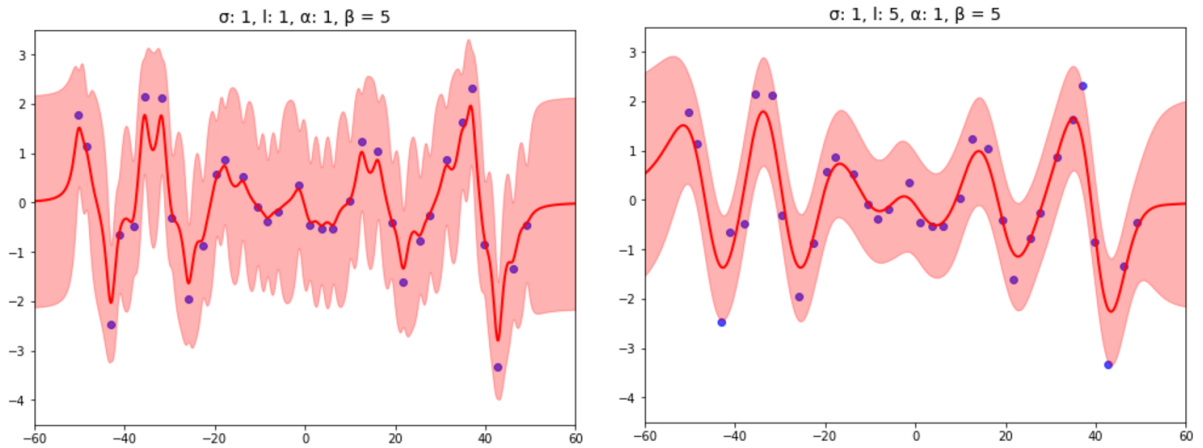
### 3) Impact of $\sigma$ :

From the figure below, we can see that when  $\sigma$  is higher, the confidence interval is wider. This is because, with a higher  $\sigma$ , the kernel function values become larger for all data points, meaning that the influence of each training point extends further. As a result, the model becomes more uncertain about the predictions, leading to greater variability and a wider confidence interval.



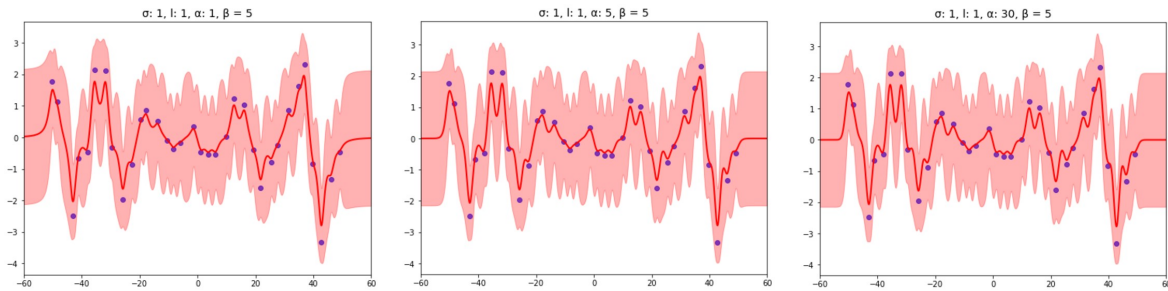
### 4) Impact of $l$ :

From the figure below, we can see that when the length scale is higher, the prediction mean and the confidence interval become smoother. This is because a larger length scale means that the kernel function assigns greater weights to distant data points, causing correlations between data points to decay more slowly. As a result, each training point has a broader influence on the predictions, leading to a smoother prediction mean and a more consistent, smoother confidence interval.



## 5) Impact of $\alpha$ :

From the figure below, we can see that the value of alpha seems to have no impact on the result.



## 2 SVM on MNIST

### 2.1 Code

#### 2.1.1 Task 1

##### 1) Prepare data

Loading the training and testing datasets from the given CSV files into the corresponding NumPy arrays.

```
1 def LoadData(X_path, Y_path):
2     X = np.loadtxt(X_path, delimiter=',')
3     Y = np.loadtxt(Y_path, delimiter=',')
4     return X, Y

1 X_train, Y_train = LoadData("./data/X_train.csv", "./data/Y_train.csv")
2 X_test, Y_test = LoadData("./data/X_test.csv", "./data/Y_test.csv")
```

##### 2) Training on different kernel and testing

First, I use libsvm's svm\_train function to train the SVM on the training set. The -t flag specifies the kernel type, where 0 corresponds to linear kernel, 1 to polynomial kernel, and 2 to RBF kernel. Then, I use the svm\_predict function to make predictions on the test set and obtain the accuracy.

```
1 # task 1 different kernel
2 from libsvm.svmutil import *
3 kernels = {'linear': 0, 'polynomial': 1, 'RBF': 2}
4
5 for key, value in kernels.items():
6     print(f"kernel: {key}", end = ", ")
7     model = svm_train(Y_train, X_train, f"-t {value}")
8     pred = svm_predict(Y_test, X_test, model)
```

## 2.1.2 Task 2

### 1) Grid search to find optimal parameters

Different kernels require different arguments, as shown in the table below:

Argument	Usage	Default	Grid Search	Linear	Polynomial	RBF
-c	Cost for C_SVC	1	[0.001,0.01,0.1,1,10]	v	v	v
-g	gamma of kernel	1/(28*28)	[0.001,0.1,1]		v	v
-d	degree of kernel	3	[2,3,4]		v	
-r	coefficient 0 of kernel	0	[2,4,10]		v	
-v	n-fold cross validation		fixed at 5	v	v	v

I applied both the default settings and the grid search values specific to each kernel. For each configuration, I performed 5-fold cross-validation to evaluate the performance.

```
1 def GridSearch(X,Y):
2     D = [2, 3, 4]
3     Gamma = [1e-2, 1e-1, 1]
4     Coef0 = [2, 4, 10]
5     C = [0.001, 0.01, 0.1, 1, 10]
6     v = 5
7
8     for key, kernel in kernels.items():
9         print(f"\nkernel: {key}")
10        print("Default:", end = " ")
11        res = svm_train(Y, X, f"-t {kernel} -v {v}")
12        if kernel==0:
13            for c in C:
14                arg = f"-t {kernel} -v {v} -c {c}"
15                print(f"-c {c:<5.3f}", end = " ")
16                res = svm_train(Y, X, arg)
17        elif kernel==1:
18            for c in C:
19                for d in D:
20                    for gamma in Gamma:
21                        for coef0 in Coef0:
22                            arg = f"-t {kernel} -v {v} -c {c} -d {d} -g {gamma} -r {coef0} "
23                            print(f"-c {c:<5.3f} -d {d:<1d} -g {gamma:<3.2f} -r {coef0:<2d} ", end = " ")
24                            res = svm_train(Y, X, arg)
25        else:
26            for c in C:
27                for gamma in Gamma:
28                    arg = f"-t {kernel} -v {v} -c {c} -g {gamma}"
29                    print(f"-c {c:<5.3f} -g {gamma:<3.2f}", end = " ")
30                    res = svm_train(Y, X, arg)
```

### 2) Performance of optimal parameters on testing set

I used the optimal parameters obtained from the grid search to train the best model for each kernel type. After training, I made predictions on the test set to evaluate the model's performance.

```
1 kernels = {'linear': 0, 'polynomial': 1, 'RBF': 2}
2 best_opt = {'linear': "-c 0.01",
3             'polynomial': "-c 0.1 -d 2 -g 1 -r 10",
4             'RBF': "-c 10 -g 0.01"}
5
6 for key, value in kernels.items():
7     print(f"kernel: {key}", end = ", ")
8     print(f"best opt: -t {value} {best_opt[key]}")
9     model = svm_train(Y_train, X_train, f"-t {value} "+best_opt[key])
10    pred = svm_predict(Y_test, X_test, model)
```

### 2.1.3 Task 3

#### 1) Kernels

The formula of linear kernels and RBF kernel is:

Kernel	$K_{ij}$	Matrix K
Linear	$K_{ij} = x_i \cdot x_j$	$K = XX^T$
RBF	$K_{ij} = \exp^{-\gamma \ x_i - x_j\ ^2}$	$K = e^{-\gamma D}, \text{ where } D = \ x_i\ ^2 + \ x_j\ ^2 - 2(x_i \cdot x_j)$

```
1 def Linear_K(X1,X2):  
2     return mul(X1,X2.T)
```

```
1 def RBF_K(X1,X2, gamma):  
2     square_X1 = np.sum(X1 ** 2, axis=1).reshape(-1, 1)  
3     square_X2 = np.sum(X2 ** 2, axis=1)  
4     d = square_X1 - 2 * mul(X1, X2.T) + square_X2  
5     k = np.exp(-gamma * d)  
6     return k
```

#### 2) Training on precompute kernel

I chose the optimal gamma value of 0.01 found in Task 2 and calculated the kernel values for both the training and testing sets by combining Linear and RBF kernels. Since precomputed kernels require the first column of the training matrix to be the ID of each instance, I used `np.column_stack` to append sequential indices as the first column to both `K_train` and `K_test`. Then, I used the precomputed kernel for training the SVM model (flag -t 4) and applied the model to the testing set to evaluate its performance.

```
1 gamma = 0.01  
2 K_train = Linear_K(X_train, X_train) + RBF_K(X_train, X_train, gamma)  
3 K_test = Linear_K(X_train,X_test).T + RBF_K(X_train,X_test, gamma).T  
4  
5 # add index for every Ki  
6 train_idx = np.arange(1,K_train.shape[0]+1)  
7 K_train = np.column_stack((train_idx, K_train))  
8  
9 test_idx = np.arange(1,K_test.shape[0]+1)  
10 K_test = np.column_stack((test_idx, K_test))  
  
1 # -t 4: precomputed kernel  
2 model = svm_train(Y_train, K_train, "-t 4")  
3 pred = svm_predict(Y_test, K_test, model)
```



## 2.2 Experiments

### 2.2.1 Task 1

Kernel	Accuracy
Linear	95.08%
Polynomial	34.68%
RBF	95.32%

### 2.2.2 Task 2

Kernel	Best parameters	Validation Accuracy	Testing Accuracy
Linear	-c 0.01	97.12%	95.96%
Polynomial	-c 0.1 -d 2 -g 1 -r 10	98.22%	97.84%
RBF	-c 10 -g 0.01	98.14%	98.20%

### 2.2.3 Task 3

Kernel	Parameter	Accuracy
Linear + RBF	gamma = 0.01	95.32%

## 2.3 Observation and Discussion

### 1) The degree in polynomial kernel

The optimal degree for the polynomial kernel is 2, which significantly improves performance compared to the default degree of 3 ( $97.84\% > 34.68\%$ ). This indicates that a higher degree makes the model overly complex, leading to overfitting. Reducing the degree to 2 simplifies the model and enhances its ability to generalize.

### 2) Best Kernel

The RBF kernel achieves the best results, suggesting that the dataset exhibits complex, non-linear relationships that are effectively captured by the RBF kernel.

### 3) Combining Kernels

Combining the linear and RBF kernels does not offer any significant improvement over using the RBF kernel alone with optimal parameters. This suggests that the RBF kernel already captures the critical patterns in the data, making the linear component redundant.