

2 Programming

In this assignment, we construct several SVM models with different kernels and slack variables to classify the Iris dataset.

The basic form of SVM is given below:

$$\min_{w,b} \frac{1}{2} \|w\|^2$$
$$s. t. \ 1 - y_i(w^T x_i + b) \leq 0, \forall i$$

Where w is the coefficient of different features, b is the intercept of the hyperplane, x_i and y_i are the features and labels of the Iris data.

2.1 (SVM without slack variables)

1. The optimization problem

We first get the dual problem of the original problem stated above.

The dual Lagrange function is:

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 + \sum_i^m \alpha_i (1 - y_i(w^T x_i + b))$$

The primal and dual optimal solutions should satisfy KKT conditions:

Stationarity:

$$\frac{\partial \mathcal{L}}{\partial w} = 0 \Rightarrow w = \sum_i^m \alpha_i y_i x_i$$
$$\frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow \sum_i^m \alpha_i y_i = 0$$

Feasibility:

$$\alpha_i \geq 0, 1 - y_i(w^T x_i + b) \leq 0, \forall i$$

Complementary slackness:

$$\alpha_i (1 - y_i(w^T x_i + b)) = 0, \forall i$$

Then, the dual problem can be derived by substituting all the stationary conditions into the primal problem, finally we get:

$$\max_{\alpha} \sum_i^m \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j,$$
$$s. t. \sum_i^m \alpha_i y_i = 0, \alpha_i \geq 0, \forall i$$

Also, b is given by:

$$y_j(w^T x_j + b) = 1, \forall j \in S$$

$$\Rightarrow y_j \left(\sum_i^m \alpha_i y_i x_i^T x_j + b \right) = 1, \forall j \in S$$

Since $y_j^2 = 1$,

$$\sum_i^m \alpha_i y_i x_i^T x_j + b = y_j, \forall j \in S$$

$$\Rightarrow b = \frac{1}{|S|} \sum_{j \in S} y_j - \sum_i^m \alpha_i y_i x_i^T x_j$$

2. Data processing and results

I use OneVsRestClassifier and sklearn.svm.SVC to do the following problems.

Since sklearn doesn't provide strict separation, we use $C = 1e5$, kernel = 'linear' to estimate the attributes and calculate errors. The result is shown below.

```
*SVM_linear - 记事本
文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)
training error: 0.041666666666666664
testing_error: 0.0
w_of_setosa: -0.04575352,0.52216766,-1.00294058,-0.46406882
b_of_setosa: 1.44746413
support_vector_indices_of_setosa: 78,13,31
w_of_versicolor: -0.75160959,-3.4187652,2.06714366,-4.63634689
b_of_versicolor: 11.31356887
support_vector_indices_of_versicolor: 1,2,3,14,15,20,28,31,32,81,82,83,84,86,88,89,
91,92,93,95,96,98,99,100,103,104,107,112,116,117,119,41,43,44,45,46,47,50,52,54,
55,56,57,58,59,62,64,65,66,68,69,71,73,74,75,76,77,78,79
w_of_virginica: -4.26389247,-6.19330415,8.64141632,12.56275266
b_of_virginica: -19.19066652
support_vector_indices_of_virginica: 50,52,57,63,97,99,103,108
```

When determining which class is linearly separable, we first calculate the train loss by 1-training score by the sklearn, then if the train loss is 0, it is linearly separable.

The statistics are shown below:

```
Label 0 linearly separable: True
train_loss for label 0: 0.0
Label 1 linearly separable: False
train_loss for label 1: 0.21666666666666667
Label 2 linearly separable: False
train_loss for label 2: 0.016666666666666672
```

In conclusion, only label 0 (setosa) is linearly separable in the dataset.

2.2 (SVM with slack variables)

1. The optimization problem

For SVM with slack variables, we can simply modify it by adding a penalty term.

The dual Lagrange function is:

$$\mathcal{L}(w, b, \alpha, \xi, \mu) = \frac{1}{2} \|w\|^2 + C \sum_i^m \xi_i + \sum_i^m [\alpha_i (1 - \xi_i - y_i (w^T x_i + b)) - \mu_i \xi_i]$$

The primal and dual optimal solutions should satisfy KKT conditions:

Stationarity:

$$\frac{\partial \mathcal{L}}{\partial w} = 0 \Rightarrow w = \sum_i^m \alpha_i y_i x_i$$

$$\frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow \sum_i^m \alpha_i y_i = 0$$

$$\frac{\partial \mathcal{L}}{\partial \xi_i} = 0 \Rightarrow \alpha_i = C - \mu_i, \forall i$$

Feasibility:

$$\alpha_i \geq 0, 1 - \xi_i - y_i(w^T x_i + b) \leq 0, \xi_i \geq 0, \mu_i \geq 0, \forall i$$

Complementary slackness:

$$\alpha_i(1 - \xi_i - y_i(w^T x_i + b)) = 0, \mu_i \xi_i = 0, \forall i$$

Then, the dual problem can be derived by substituting all the stationary conditions into the primal problem, finally we get:

$$\begin{aligned} \max_{\alpha} \quad & \sum_i^m \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j, \\ \text{s. t.} \quad & \sum_i^m \alpha_i y_i = 0, 0 \leq \alpha_i \leq C, \forall i \end{aligned}$$

Also, b is given by: $M = \{i | 0 < \alpha_i < C\}$

$$y_j(w^T x_i + b) = 1, \forall j \in M$$

$$\Rightarrow y_j \left(\sum_i^m \alpha_i y_i x_i^T x_j + b \right) = 1, \forall j \in M$$

Since $y_j^2 = 1$,

$$\sum_i^m \alpha_i y_i x_i^T x_j + b = y_j, \forall j \in M$$

$$\Rightarrow b = \frac{1}{|M|} \sum_{j \in M} y_j - \sum_i^m \alpha_j y_i x_i^T x_j$$

For the slack variables' calculation, $\max(0, 1 - y_i(w^T x_i + b))$ can return the value of the slack variables, where y_i is the label $\{-1, +1\}$, $w^T x_i + b$ is the result calculated by the input vectors x_i and the estimated parameters w and b .

2. Data processing and results

Below are two sample outputs $C = 0.3$ and $C = 0.7$:

2.3 (SVM with kernels and slack variables)

The derivation of the dual problem of the kernel is similar to 2.2:

The related kernel functions are:

$$k(x, x_i) = (1 + \frac{x^T x_i}{\sigma^2})^p, p > 0$$

Radical Basis Function (RBF) kernel:

$$k(x, x_i) = \exp \left\{ -\frac{\|x - x_i\|^2}{2\sigma^2} \right\}$$

Sigmoidal kernel:

$$k(x, x_i) = \frac{1}{1 + \exp \left(-\frac{x^T x_i + b}{\sigma^2} \right)}$$

2. Data processing and results

(a) In this scenario, we set $C = 1$, kernel = 'poly', degree = 2, gamma = 1, and get the following output:

```
training error: 0.025
testing_error: 0.0
b_of_setosa: 1.22094132
support_vector_indices_of_setosa: 78,13,31
b_of_versicolor: 4.33667006
support_vector_indices_of_versicolor: 14,31,89,93,96,97,99,103,108,48,50,52,57,58,63,64
b_of_virginica: -10.42876523
support_vector_indices_of_virginica: 50,52,57,63,96,97,99,103,108
```

(b) In this scenario, we set $C = 1$, kernel = 'poly', degree = 3, gamma = 1, and get the following output:

```
training error: 0.008333333333333333
testing_error: 0.0
b_of_setosa: 1.13434963
support_vector_indices_of_setosa: 78,13,31
b_of_versicolor: 1.54426028
support_vector_indices_of_versicolor: 31,89,97,99,101,103,108,119,50,52,57,63,70
b_of_virginica: -6.11788922
support_vector_indices_of_virginica: 50,52,57,63,89,103,108,119
```

(c) In this scenario, we set $C = 1$, kernel = 'rbf', gamma = 0.5, and get the following output:

```
training error: 0.03333333333333333
testing_error: 0.03333333333333333
b_of_setosa: -0.33592638
support_vector_indices_of_setosa: 42,45,78,84,87,88,89,101,104,106,4,5,12,14,31
b_of_versicolor: -0.41536232
support_vector_indices_of_versicolor: 4,5,13,14,31,80,88,89,91,93,96,97,99,101,103,108,116,119,40,43,46,48,50,52,56,57,58,63,64,65,66,78
b_of_virginica: -0.30449079
support_vector_indices_of_virginica: 3,4,5,12,14,31,40,43,46,48,50,52,56,57,58,63,64,65,66,80,87,88,89,91,93,96,97,99,101,103,104,108,111,116,119
```

(d) In this scenario, we set $C = 1$, kernel = 'sigmoid', gamma = 'auto' (gamma = 0.25), and get the following output: