

Experimental Report on Method and Evaluation of SVM Implemented by SMO

Yuxiang Zhang

I. METHODS

In this experiment, we try to evaluate the results performed by different models. The selected models contains SVM using SMO algorithm, which is implemented by us, and two well-developed toolboxes: sklearn and libsvm. This section includes two parts, in the first part, we will briefly introduce the formulation of SVM that implemented by us, such as choosing the kernel trick. In the second part, we will present the SMO using in SVM.

A. Formulation of SVM

For the convenience of the following description, the matrix of input sample features is called \mathbf{X} , the matrix of sample labels is called \mathbf{Y} , and a single sample with features is called \mathbf{S} .

Our very original goal is to find a classifier function $f(\mathbf{S}) = \mathbf{X}^T \beta + \beta_0$ to realize binary classification task. In order to avoid the possible two-type inseparable problem, the relaxation variable C is introduced to find a soft margin hyperplane. Then, the original task of SVM is transfer into minimize $E(\beta)$ subject to $\forall \varepsilon_i \geq 0, y_i(\mathbf{X}_i^T \beta + \beta_0) \geq 1 - \varepsilon_i$. By using Lagrange multipliers, we get the dual problem: maximize $E(\alpha)$ subject to $\sum_{i=1}^m \alpha_i y_i = 0, 0 \leq \alpha_i \leq C$. Then, we use SMO algorithm to solve the dual problem, getting the α_i, β and β_0 . Finally, we can gain the classification decision function.

$$E(\beta) = \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N \varepsilon_i. \quad (1)$$

$$E(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{X}_i^T \mathbf{X}_j. \quad (2)$$

Additionally, we used three kernel functions containing linear kernel, polynomial kernel and radial basis kernel, to convert input sample features to Kernel matrix \mathbf{K} , which could help simplify the calculation of the inner product after the sample is mapped to high-dimensional space. (\mathbf{K}_i represents the i th column of \mathbf{K} , \mathbf{K}_{ij} represents the i th row first j column of \mathbf{K}). Using Equation(3), Equation(4) and Equation(5) to generate the Kernel matrix for linear, polynomial kernel and radio basis kernel, respectively.

$$\mathbf{K}_i = \mathbf{X} \mathbf{S}^T. \quad (3)$$

$$\mathbf{K}_{ij} = \mathbf{K}_{ij}^k \quad (s.t. \mathbf{K}_i = \mathbf{X} \mathbf{S}^T). \quad (4)$$

$$\mathbf{K}_{ij} = \exp\left(-\frac{\mathbf{K}_{ij}}{\sigma^2}\right) \quad (s.t. \mathbf{K}_i = (\mathbf{X}_i - \mathbf{S})(\mathbf{X}_i - \mathbf{S})^T). \quad (5)$$

Then, the final solution of SVM can be written as:

$$f(\mathbf{S}) = \sum_{i=1}^N \alpha_i y_i \mathbf{K}(\mathbf{X}, \mathbf{S}) + \beta_0 \quad (6)$$

B. Algorithmic Description of SMO

In this section, we will introduce the SMO algorithm used in SVM. Firstly, we define a function: $smo_pre(\mathbf{X}, \mathbf{Y}, C, toler, i)$ (the detail is in Algorithm1). By using the smo_pre function in

Algorithm 1 $smo_pre(\mathbf{X}, \mathbf{Y}, C, toler, i)$

Input : $\mathbf{X}, \mathbf{Y}, C$, Tolerance: $toler$, Maximal training times: $maxIter$, Old threshold: β_0 , Old lagrange multipliers for solution: α , Number of samples: m .

```

1: Calculate  $E_i = f(X_i) - Y_i$  using Equation(6)
2: if  $(Y_i E_i < -toler \ \& \ \alpha_i < C)$  or  $(Y_i E_i > toler \ \& \ \alpha_i > 0)$  then
3:   Select  $j (j \neq i)$  that  $\alpha_j$  does not fulfill the Karush-Kuhn-Tucker conditions with a maximal  $E_j$ .
4:   Calculate  $E_i = f(X_i) - Y_i$  using Equation(6).
5:   Save old  $\alpha_i$  and  $\alpha_j$ ;  $i$  and  $j$ ,  $s = Y_i \times Y_j$ .
6:   if  $Y_i = Y_j$  then
7:      $L = \max(0, \alpha_i + \alpha_j - C)$ ;  $H = \min(C, \alpha_i + \alpha_j)$ .
8:   else
9:      $L = \max(0, \alpha_j - \alpha_i)$ ;  $H = \min(C, \alpha_j - \alpha_i + C)$ .
10:  end if
11:  if  $L = H$  then
12:    return 0
13:  end if
14:  Calculate Kernel matrix using Equation(3), (4) or (5).
15:   $\eta = \mathbf{K}_{ii} + \mathbf{K}_{jj} - 2 \times \mathbf{K}_{ij}$ 
16:  if  $\eta < 0$  then
17:    return 0
18:  end if
19:   $\alpha_{jn} = \alpha_j + Y_j \times \frac{(E_i - E_j)}{\eta}$ .
20:  Renew  $\alpha_j$  using  $\min(\alpha_j, H)$  and  $\alpha_j = \max(\alpha_j, L)$ .
21:  if  $|\alpha_{jn} - old \ \alpha_j| < 0.00001$  then
22:    return 0
23:  end if
24:  Renew  $\alpha_{in}$ , using  $\alpha_i + s \times (old \ \alpha_j - \alpha_{jn})$ ,  $E_i$  and  $E_j$ .
25:  Calculate  $\beta_1, \beta_2$  using Equation(7) and(8), respectively.
26:  Renew  $\beta_0 = \frac{1}{2} \times (\beta_1 + \beta_2)$ .
27:  return 1
28: else
29:  return 0
30: end if
```

Algorithm2, we can get the well-trained lagrange multipliers and the β_0 .

$$\beta_1 = Y_i \mathbf{K}_{ii}(\alpha_i - \alpha_{in}) + Y_j \mathbf{K}_{ij}(\alpha_j - \alpha_{jn}) + \beta_0 - E_i \quad (7)$$

$$\beta_2 = Y_i \mathbf{K}_{ij}(\alpha_i - \alpha_{in}) + Y_j \mathbf{K}_{jj}(\alpha_j - \alpha_{jn}) + \beta_0 - E_j \quad (8)$$

II. EXPERIMENTAL SETTINGS

The training dataset gp96.csv and test dataset gp97.csv come from Series GSE3494, which express signature for breast cancer. We use the entire training dataset to train the model. Noticed that that the features of some samples in the data set are missing, we preprocess the data by deleting the sample with feature missing, gaining a training dataset with 222 samples. In order to improve the training efficiency of model, we set the maximum number of iterations and the model will stop training when the parameters change to a certain extent. After comparing the prediction results of different

Algorithm 2 SMO for SVM

Input : $\mathbf{X}, \mathbf{Y}, C, \text{toler}, \text{maxIter}, \beta_0, m$.

Output : New β_0 and New α

- 1: Create a $m \times 1$ matrix and initialize α_i, β_0 to 0.
 - 2: alphachange times: alphachange to 0
 - 3: while alpha changes or *entireSet* = 1 do
 - 4: Initialize alphachange to 0.
 - 5: $\forall i < m, i \in \mathbb{N}$
 - 6: $\text{alphachange} + \text{smo_pre}(\mathbf{X}, \mathbf{Y}, C, \text{toler}, i)$
 - 7: \forall Index of $\alpha = 0$
 - 8: $\text{alphachange} + \text{smo_pre}(\mathbf{X}, \mathbf{Y}, C, \text{toler}, \text{index})$
 - 9: end while
-

kernel functions, we choose the kernel function with better performance for subsequent experiment. Using gradient search to select the best hyperparameters, we compare the model with these hyperparameters with sklearn and libsvm. The environment we use is python 3.9, numpy 1.20.3, pandas 1.3.4, scikit-learn 0.21.2 and libsvm-official 3.25.0.

III. RESULTS & DISCUSSION

A. Performance of Different Model On Test Dataset

Firstly, we employ our implementation, sklearn and libsvm on test dataset, using the ROC curve, AUC value and Accuracy as evaluation indexes. We find that the accuracy of SMO SVM with linear or polynomial kernel is obviously lower than the radio basis kernel (linear 0.73, poly 0.82 and radio basis 0.93), so we will mainly discuss the result performed by radio basis kernel. Then, the hyperparamaters of SMO SVM are C and σ , we used gradient search to find the best pair with the accuracy of 0.93: $\{C=1, \sigma=3\}$ (see Fig.1). Secondly, we use the SVM

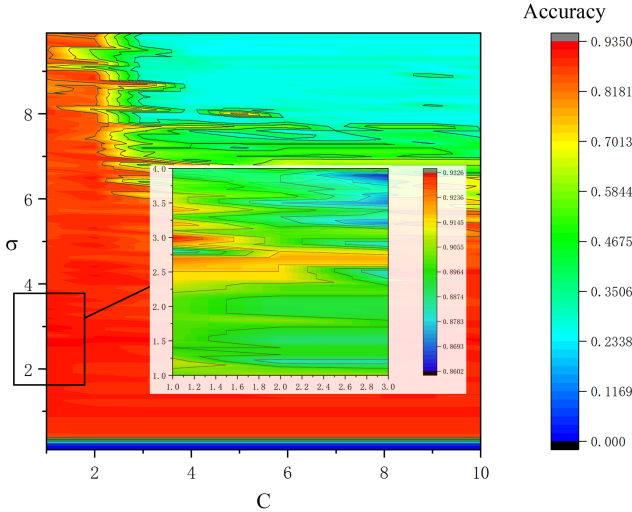


Fig. 1. The accuracy of SVM(SMO) using different C and σ .

model in two well-developed toolboxes, sklearn and libsvm, employing on the same train and test dataset. The accuracy from largest to smallest is sklearn with radio basis kernel, libsvm with radio basis kernel, SMO with radio basis kernel, sklearn with linear kernel and libsvm with linear kernel. We also drew the ROC-curve and compute the area under ROC-curve(AUC), the well-developed toolboxes both get a better performance comparing with our implementation(see Fig2).

Generally, we find that radio basis kernel is better than linear on the test dataset. Thus, we conclude that when the number

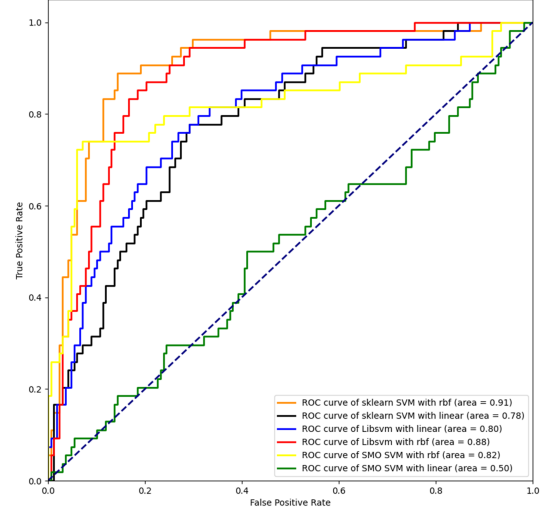


Fig. 2. ROC curve of different model perform on test dataset.

of samples is significantly larger than the number of features, a linear inseparable situation possibly occur, therefore, using a radio basis kernel may perform better at this circumstance.

B. Statistical Test of SMO SVM with Toolboxes

To compare whether there were significant differences between the predictive abilities of different models with radio basis kernel, we performed statistical tests for the two toolboxes and SMO SVM, respectively. Because of the use of the same dataset and the small number of models involved in the comparison, we used McNemar's test, which was tested with 95% confidence interval.

The null hypothesis H_0 is that "Two classifiers have similar error rate on the test dataset", while the alternative hypothesis H_1 is that "Two classifiers have different error rate on the test dataset".

	SMO T		SMO F	
Sklearn T	182	25	Libsvm T	171
Sklearn F	10	182	Libsvm F	10
				171

Fig. 3. 2×2 contingency table of McNemar's test.

Using the 2×2 contingency table showing in Fig.3, the p-value of test between sklearn and SMO is 0.018; between Libsvm is 0.00023, both of them is small than 0.05, so we have sufficient evidence to reject H_0 , the two well-developed toolboxes' predictive ability is significantly better than our implementation in 95% confidence interval.

REFERENCES

- [1] fengdu78, Support-vector-machine.ipynb, 2020, GitHub repository [Online]. Available: <https://github.com/fengdu78/lihang-code/blob/master/7.support-vector-machine>
- [2] Coding Notes, The Simplified SMO Algorithm, 2018, Codeproject [Online]. Available: <https://www.codeproject.com/Articles/1267445/An-Introduction-to-Support-Vector-Machine-SVM-and>