

1.

(1) F The  $L_2$  regularized logistic regression has  
only 1 global optimal solution

(2) F With penalty term  $\lambda \|w\|_2^2$ , the values  
of the weight should be distributed evenly

$\Rightarrow$  since there exist  $w'$  that can  
totally separate the data and achieve  
highest likely hood, multiplying a scalar  
of arbitrary large will not reduce  
the likelihood. Hence  $w=2w' \Rightarrow w$

4) F As  $\lambda$  increase,  $l(w, D_{train})$  decrease,  
hence  $l(w, D_{train})$  increase, meaning that  
the penalty cause the a worse performance  
in training set, hence a lower likely hood

5) F If  $\lambda$  is too large, & underfit  
tend to happen, hence  $l$  is lower

2. (1)

Let the problem with  $\xi_i \geq 0$  constraint  
be problem A, without the constraint  
be problem B. The goal is to prove  
A and B are equivalent

$B \rightarrow A$ : suppose  $w^*$  is the optimal solution  
 $A \rightarrow B$   
of problem A, so it satisfy:

①  $A \Rightarrow B$

Suppose there exist some case when  
the optimal solution of B is not that  
of A. the optimal value of problem B  
is less than that of problem A:

$$OPT(B) < OPT(A)$$

$$\frac{1}{2} \|w_B\|^2 + \frac{C}{2} \sum \xi_B < \frac{1}{2} \|w_A\|^2 + \frac{C}{2} \sum \xi_A;$$

for this to happen there must exist  $\xi_i > 0$  for some  
 $y^i(w_B^T x + b) \geq 1 - \xi_i, i=1, \dots, m$

supp let set S be the subset of  
 $\{1, \dots, m\}$  such that  $\xi_{Bi} > 0 \forall i \in S$ ,  
then  $y^i(w_B^T x + b) \geq 1 - \xi_{Bi} \geq 1 - \cancel{\xi_{Bi}} - \xi_{Bi}$

$$\text{where } \cancel{\xi_{Bi}} = 0$$

then

OPT.

$$\begin{aligned} OPT(B) &= \frac{1}{2} \|w_B\|^2 + \frac{C}{2} \sum \xi_B = \frac{1}{2} \|w_B\|^2 + \frac{C}{2} \sum_{i \in S} \xi_{Bi} \\ &\geq \frac{1}{2} \|w_B\|^2 + \frac{C}{2} \sum_{i \in S} \xi_{Bi}^2 + \frac{C}{2} \sum_{i \in S} \xi_{Bi} \geq \sum_{i \in S} \xi_{Bi}^2 \\ &= \frac{1}{2} \|w_A\|^2 + \frac{C}{2} \sum_{i \in S} \xi_{Bi}^2 \end{aligned}$$

hence this is not the optimal solution  
which is a contradiction.  
hence  $OPT(S) \geq OPT(A)$

②  $B \Rightarrow A$

since B has a less strict constraint,  
all value objective value in A can  
be achieved in B

hence  $OPT(B) \leq OPT(A)$

③  $\Rightarrow OPT(B) = OPT(A)$  Q.E.D.

Let the problem with  $\xi_i \geq 0$  constraint  
be problem A, without the constraint  
is problem B, it is to prove that  
A and B are equivalent:

$A \Leftrightarrow B$

$$(2) \min \frac{1}{2} \|w\|^2 + \frac{C}{2} \sum_{i=1}^m \xi_i^2$$

s.t.  $-y_i(w^T x_i + b) + 1 - \xi_i \leq 0, i=1, \dots, m$

(a gradient)

~~$L(w, b, \xi, \alpha) = \frac{1}{2} \|w\|^2 + \frac{C}{2} \sum_{i=1}^m \xi_i^2 - \sum_{i=1}^m \alpha_i (-y_i(w^T x_i + b) + 1 - \xi_i)$~~

~~$\nabla L(w, b, \xi, \alpha) = \frac{1}{2} \|w\|^2 + \frac{C}{2} \sum_{i=1}^m \xi_i^2 + \sum_{i=1}^m \alpha_i (-y_i(w^T x_i + b) + 1 - \xi_i)$~~

$$(3) \frac{\partial L}{\partial w} = w + \sum_{i=1}^m \alpha_i y_i x_i = 0$$

$$W = \sum_{i=1}^m \alpha_i y_i x_i$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^m \alpha_i$$

$$\frac{\partial L}{\partial \xi_i} = \alpha_i$$

Hence  $\frac{\partial L}{\partial \xi_i} = \alpha_i$

$$\frac{\partial L}{\partial \xi_i} = 0 \Rightarrow \alpha_i = \frac{C}{2}$$

$$(4) g(u, v) = \min_w L(w, b, \xi, \alpha)$$

s.t.  $\alpha \geq 0$

$$W = \sum_{i=1}^m \alpha_i y_i x_i \quad \sum_{i=1}^m \alpha_i y_i = 0 \quad \xi_i = \frac{\alpha_i}{C}$$

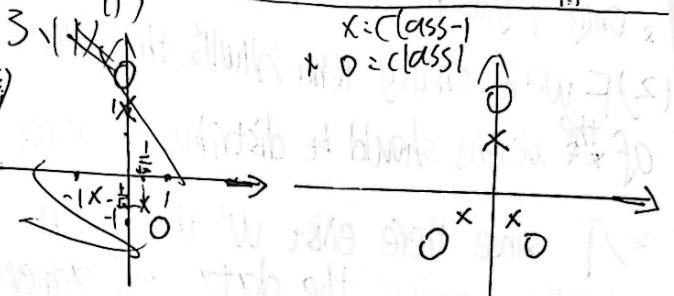
$$L = \frac{1}{2} \|w\|^2 + \frac{C}{2} \sum_{i=1}^m \xi_i^2 + \sum_{i=1}^m -y_i(w^T x_i + b) + \frac{1}{2} \sum_{i=1}^m \alpha_i$$

$$= \frac{1}{2} \|w\|^2 + \frac{C}{2} \sum_{i=1}^m \xi_i^2 - \sum_{i=1}^m \alpha_i y_i x_i^T w - b \sum_{i=1}^m \alpha_i y_i + \frac{1}{2} \sum_{i=1}^m \alpha_i$$

$$= \frac{1}{2} \|w\|^2 + \frac{C}{2} \left( \frac{1}{C} \right)^2 - W \sum_{i=1}^m \frac{1}{C} \alpha_i y_i x_i^T \left( \frac{1}{C} \alpha_i y_i x_i \right)$$

$$= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j (y_i^T x_i)(y_j^T x_j) - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j (y_i^T x_i)(y_j^T x_i) + \frac{1}{2} \sum_{i=1}^m \alpha_i^2 - \frac{1}{C} \sum_{i=1}^m \alpha_i + \frac{1}{2} \sum_{i=1}^m \alpha_i$$

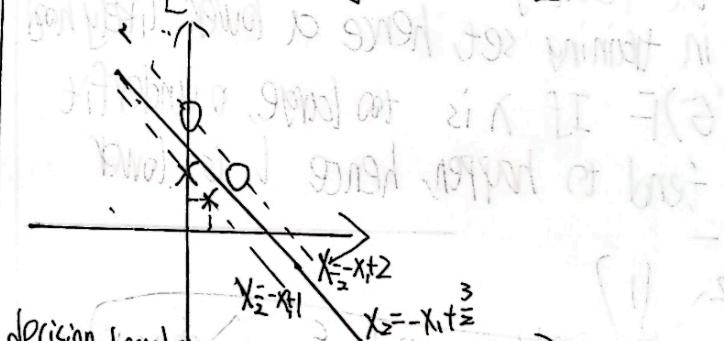
$$= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j - \frac{1}{2C} \sum_{i=1}^m \alpha_i^2 + \frac{m}{C}$$



If no kernel is applied, SVM classifier cannot split this dataset since it is linearly non-separable

(2)

$$\text{class -1: } \begin{bmatrix} (0, 1) \\ (\frac{1}{2}, \frac{1}{2}) \\ (\frac{1}{2}, \frac{1}{2}) \end{bmatrix} \quad \text{class 1: } \begin{bmatrix} (0, 2) \\ (1, 1) \\ (1, 1) \end{bmatrix}$$



decision boundary prediction:  $w^T x + b = 0 \Rightarrow x_2 - x_1 \geq \frac{3}{2} \Rightarrow x_1 + x_2 - \frac{3}{2} = 0$

$$w^T x + b = 0 \Rightarrow w_1 x_1 + w_2 x_2 + b = 0 \Rightarrow w_1 = 1, w_2 = -1, b = \frac{3}{2}$$

$$sgn(w^T x + b) = sgn(x_2 - x_1) = 1$$

$$sgn(w_1 x_1 + w_2 x_2 + b) = sgn(x_1 + x_2 - \frac{3}{2}) = 1$$

$$sgn(\frac{1}{\sqrt{2}}x_1 + \frac{1}{\sqrt{2}}x_2 - \frac{3}{2}) = 1$$

$$sgn(\frac{1}{\sqrt{2}}x_1 + \frac{1}{\sqrt{2}}x_2 - \frac{3}{2}) = -1$$

$$sgn(\frac{1}{\sqrt{2}}x_1 + \frac{1}{\sqrt{2}}x_2 - \frac{3}{2}) = -1$$

$$sgn(\frac{1}{\sqrt{2}}x_1 + \frac{1}{\sqrt{2}}x_2 - \frac{3}{2}) = -1$$

$$sgn(w^T x + b) = 1 \Rightarrow sgn(w_1 x_1 + w_2 x_2 + b) = 1$$

$$sgn(w_1 x_1 + w_2 x_2 + b) = 1 \Rightarrow sgn(w^T x + b) = 1$$

4

$$(1) H(\text{Pass}) = -\frac{4}{6} \log \frac{4}{6} - \frac{2}{6} \log \frac{2}{6} = 0.92$$

(2) ~~H(Pass|GPA)~~

$$P(T|L) = \frac{1}{2}$$

$$P(F|L) = \frac{1}{2}$$

$$P(T|M) = \frac{1}{2}$$

$$P(F|M) = \frac{1}{2}$$

$$P(T|H) = 1$$

$$P(F|H) = 0$$

$$P(T, L) = \frac{1}{6} \left| H(\text{Pass}|GPA) \right.$$

$$P(F, L) = \frac{1}{6} \left| -4 \times \frac{1}{6} \log \frac{1}{6} \right. = -4 \times \frac{1}{6} \log \frac{1}{2} - \frac{1}{3} \log \frac{1}{3}$$

$$P(T, M) = \frac{1}{6} \left| = -4 \times \frac{1}{6} \log \frac{1}{6} \right. = \frac{2}{3} = 0.67$$

$$P(F, M) = \frac{1}{6} \left| = \frac{2}{3} = 0.67 \right.$$

$$P(T, H) = \frac{1}{6} \left| \begin{array}{l} \\ \uparrow \end{array} \right.$$

$$P(F, H) = 0 \left| \begin{array}{l} \\ \uparrow \end{array} \right.$$

$$H(\text{Pass}|GPA) = \frac{1}{2} \times \frac{1}{6} \left( \frac{1}{2} \times \frac{1}{6} + \frac{1}{2} \times \frac{1}{6} + \frac{1}{2} \times \frac{1}{6} \right) + \frac{1}{2} \times \frac{1}{3} \left( 0 \times 0 \right) = \frac{2}{3} = 0.67$$

$$H(\text{Pass}|GPA) = \frac{1}{2} \times \frac{1}{6} \left( \frac{1}{2} \log \frac{1}{2} + \frac{1}{3} \log \frac{1}{3} + \frac{1}{6} \log \frac{1}{6} \right) = 0.18$$

(3)

$$P(T|T) = 1$$

$$P(T, T) = \frac{1}{2}$$

$$P(F, T) = 0$$

$$P(F|T) = 0$$

$$P(T, F) = \frac{1}{6}$$

$$P(F|F) = \frac{2}{3}$$

$$P(F, F) = \frac{1}{3}$$

$$H(\text{Pass}|\text{Studied}) = \frac{1}{2} \left( \frac{1}{2} \times \frac{1}{6} + \frac{1}{2} \times \frac{2}{3} \right) = 0.18$$

$$H(\text{Pass}|\text{Studied}) = \frac{1}{2} \log \frac{1}{2} + \frac{1}{6} \log \frac{1}{3} + \frac{1}{3} \log \frac{2}{3}$$

$$= 0.067 \quad 0.47$$

(4)



~~studied = True?~~



# Report-Assignment2

## Overview

In this assignment, I solve the wine analysis problem. Specifically, given its 13 features, a wine should be classified one of the three classes. To conduct such a multi-classes classification task, support vector machine(SVM) is utilized.

## Support Vector Machine

### Hard Margin SVM

The basic idea of SVM is to maximize the minimum distance between data and the decision boundary(the so-called margin). Hence, the optimization problem can be formulated as:

$$\max_{w,b} \min_i \frac{y_i(w^T x_i + b)}{\|w\|}$$

where  $w, b$  are the parameter of the linear classifier, and  $(x_i, y_i)$  are the supervised training data pair.

Since scaling  $(w,b)$  by some factor does not affect the objective function, the solution space of  $(w,b)$  can be reduced to a subspace where

$$y_i(w^T \hat{x} + b) = 1$$

where  $\hat{x}$  is the data sample that is closest to the decision boundary. Hence the optimization problem can be converted to

$$\begin{aligned} & \min_{w,b} \frac{1}{2} \|w\|^2 \\ & \text{s.t. } y_i(w^T x_i + b) \geq 1, \forall i \end{aligned}$$

## SVM With Slack Variables

For hard margin SVM, it is supposed that the data is linear separable. However, in real case, this is often impossible. Hence, slack variables are needed to relax this constraint. Specifically, each data sample is allowed to cross the margin even the decision boundary, but the number of these samples should not be too high. The problem can be formulated as

$$\begin{aligned} \min_{w,b} & \frac{1}{2} \|w\|^2 + C \sum_i^m \xi_i \\ \text{s.t. } & y_i(w^T x_i + b) \geq 1 - \xi_i, \forall i \\ & \xi \geq 0 \end{aligned}$$

## Solving SVM With Dual

By using KKT conditions, the dual problem of the original optimization of SVM with slack variable can be converted to

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

w can be obtained by

$$w = \sum_i^m \alpha_i y_i x_i$$

b can be obtained by

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

## SVM With Kernel

Since in many cases, the training samples are not linearly separable, it is necessary to project data into a dimension that can be linearly splitted. Hence kernel is needed. the simple dot product  $x_i x_j$  in the objective function is replaced by more complicated and flexible kernels  $k(x_i, x_j)$ . In this task, different kernels are tested.

## Training and Testing Setting

In this assignment, I use 80% of the data as training set, with a total number of 142, and the rest 20% of the data is used for testing, with 36 samples.

## Linear SVM Without Slack Variables

In this task, the kernel of SVM is the simplest linear model, and no slack variables are used. The training error rate is 0, meaning that the model well split the training data without error. However the testing error rate is 0.0278. Some samples are misclassified due to overfitting. For more information such as the support vector, weight and bias, please refer to A2\_119010156.ipynb.

## Linear SVM With Slack Variables

C	Train error	Test error
0.1	0.0141	0.0278
0.2	0.0211	0.0278
0.3	0.0211	0.0000
0.3	0.0211	0.0000
0.5	0.0211	0.0000
0.6	0.0211	0.0000
0.7	0.0211	0.0000
0.8	0.0070	0.0000

C	Train error	Test error
0.9	0.0211	0.0000
1.0	0.0211	0.0000

As is shown in the table, when the constraint put on the slack variables are small, the classifier tend to output a worse test error, meaning an underfit. As the constraint on slack variables grow stronger, the test error decrease and eventually fall to 0, which is suitable. For more information please refer to A2\_119010156.ipynb.

## SVM With Kernels

kernel	Train error	Test error
2nd-poly	0.0000	0.0278
3nd-poly	0.0000	0.0278
RBF	0.0070	0.0000
Sigmoidal	0.5634	0.4444

From the table, it can be concluded that RBF is the most suitable kernel, while Sigmoidal fail to map the data space into a linear separable space.