

EE5904/ME5404 Part II

Project 1  
SVM for Classification of Spam Email  
Messages

---

REPORT DUE ON 23 APRIL 2021

Thushara Sandakalum  
sandakalum@u.nus.edu



EE5904/ME5404 Part II

Project 1  
SVM for Classification of Spam Email  
Messages

---

REPORT DUE ON 23 APRIL 2021

Thushara Sandakalum  
[sandakalum@u.nus.edu](mailto:sandakalum@u.nus.edu)



# Outline

---

Project description

Recap

Task 1 : Train

Task 2 : Test

Task 3 : Evaluate

Important Notes



# Project Description

---

## Project Goal

- Implement a SVM to classify spam or not a spam for the **Spam Email Data Set**
- Spam Email Data Set
  - **4601 samples** of email metadata taken from **UC Irvine Machine Learning Repository**
  - **57 features** per sample
  - **Label : +1 (spam), -1 (non-spam)**
  - <http://archive.ics.uci.edu/ml/datasets/spambase>

```
0.00000 0.01043 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
0.00000 0.00000 0.01043 0.01043 0.02105 0.00000 0.00000 0.00000
0.00000 0.03166 0.06332 0.00000 0.02105 0.00000 0.00000 0.00000
0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
0.00000 0.00000 0.00000 0.00196 0.00000 0.00000 0.02601 0.12811
0.98827
```

48 continuous real [0,100] attributes of type word\_freq\_WORD  
= percentage of words in the e-mail that match WORD,  
i.e.  $100 * (\text{number of times the WORD appears in the e-mail}) /$   
total number of words in e-mail. A "word" in this case is any  
string of alphanumeric characters bounded by non-alphanumeric  
characters or end-of-string.

6 continuous real [0,100] attributes of type char\_freq\_CHAR  
= percentage of characters in the e-mail that match CHAR,  
i.e.  $100 * (\text{number of CHAR occurrences}) / \text{total characters in e-mail}$









# Project Description

---

## Project Goal

- Implement a SVM to classify spam or not a spam for the **Spam Email Data Set**
- Spam Email Data Set
  - **4601 samples** of email metadata taken from **UC Irvine Machine Learning Repository**
  - **57 features** per sample
  - **Label : +1 (spam), -1 (non-spam)**
  - <http://archive.ics.uci.edu/ml/datasets/spambase>
- Dataset divided into **3 subsets**
  - Training set
  - Test set
  - Evaluation set (Not provided)

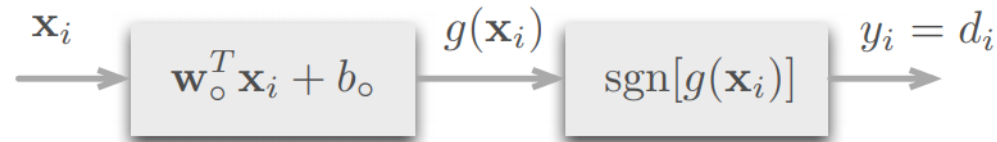
Name ▲	Value
 eval_data	57x600 double
 eval_label	600x1 double
 test_data	57x1536 double
 test_label	1536x1 double
 train_data	57x2000 double
 train_label	2000x1 double



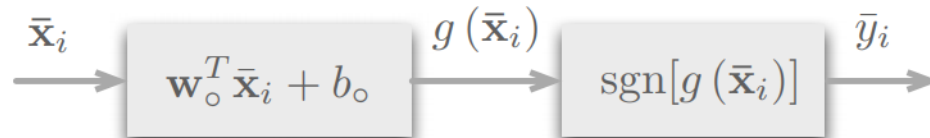
# Project Description

Slide 59

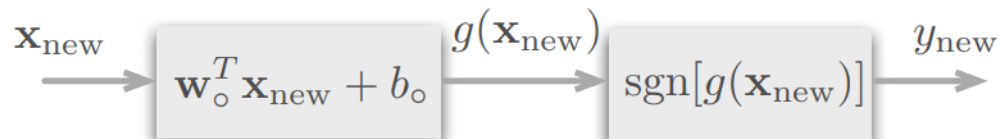
Train → **Construction:** For a given training set  $S = \{(\mathbf{x}_1, d_1), \dots, (\mathbf{x}_N, d_N)\}$ , find optimal hyperplane  $(\mathbf{w}_o, b_o)$  such that, for all  $i \in \{1, 2, \dots, N\}$ ,



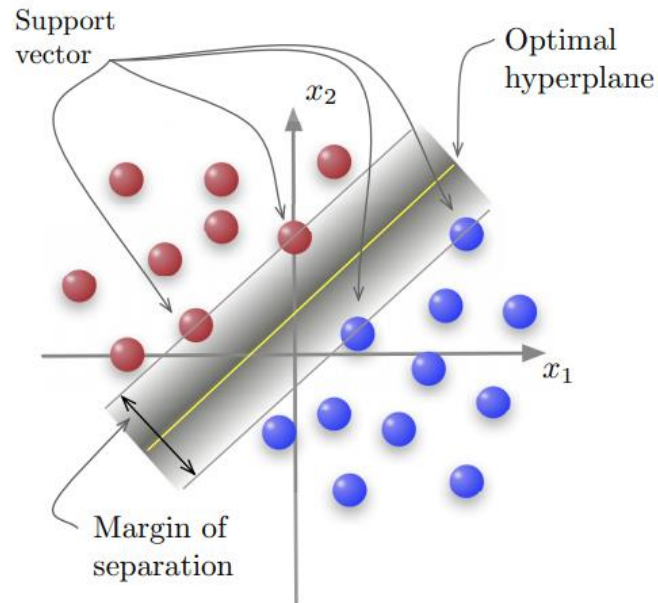
Test → **Testing:** For a given test set  $\bar{S} = \{(\bar{\mathbf{x}}_1, \bar{d}_1), \dots, (\bar{\mathbf{x}}_{\bar{N}}, \bar{d}_{\bar{N}})\}$ , compute output  $\bar{y}_i$  of SVM (with  $\mathbf{w}_o$  and  $b_o$ ) for all  $i \in \{1, 2, \dots, \bar{N}\}$ , and compare it against the known  $\bar{d}_i$  to evaluate performance of SVM



Evaluate → **Application:** Given a SVM with hyperplane  $(\mathbf{w}_o, b_o)$ , classify a data point  $\mathbf{x}_{\text{new}}$  that is not in  $\Sigma = S \cup \bar{S}$ :



# Recap – Hard Margin



Discriminant function

$$g(\mathbf{x}) = \mathbf{w}_o^T \mathbf{x} + b_o$$

Support vector:  $\mathbf{x}_i$  that satisfies

$$g(\mathbf{x}_i) = \pm 1$$

Primal problem

Given data set :  $S = \{(\mathbf{x}_i, d_i)\}, i = 1, 2, \dots, N$

Find :  $\mathbf{w}$  and  $b$

Minimizing :  $f(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$

Subject to :  $d_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

Known parameters:  $\mathbf{x}_i, d_i$

Unknown variables:  $\mathbf{w}, b$

Solve

Optimum hyperplane



# Recap – Hard Margin

Finding optimal hyperplane (primal problem) **Slide 73**

Given data set :  $S = \{(\mathbf{x}_i, d_i)\}$

Find :  $\mathbf{w}$  and  $b$

Minimizing :  $f(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$

Subject to :  $d_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

Finding optimal hyperplane (dual problem) **Slide 79**

Given :  $S = \{(\mathbf{x}_i, d_i)\}$

Find : Lagrange multipliers  $\{\alpha_i\}$

Maximizing :  $Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j \mathbf{x}_i^T \mathbf{x}_j$

Linear kernel

Subject to : (1)  $\sum_{i=1}^N \alpha_i d_i = 0$

(2)  $\alpha_i \geq 0$

} Karush-Kuhn-Tucker conditions

For data point  $\mathbf{x}_i$  that is a support vector  $\alpha_{o,i} \neq 0$

Alternative  
formulation using  
method of  
**Lagrange  
multipliers**

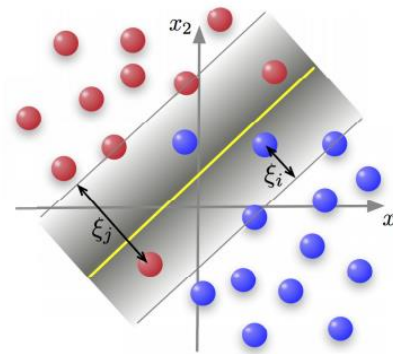
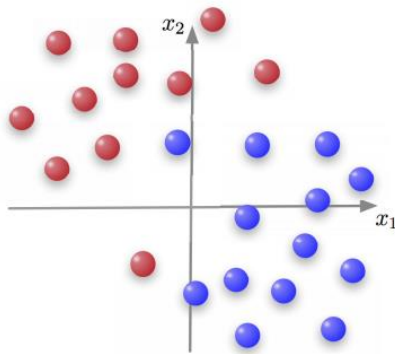




# Recap – Soft Margin

## Dealing with non-separable patterns:

1. Find optimal hyperplane to minimize classification error **Slide 94**



New function to be minimized

$$f(\mathbf{w}, \xi) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i$$

Dual problem (with soft margin) **Slide 100**

Find :  $\alpha_i$

Maximize :  $Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i c_j \mathbf{x}_i^T \mathbf{x}_j$

Subject to :  $\sum_{i=1}^N \alpha_i d_i = 0$  and  $0 \leq \alpha_i \leq C$

Linear kernel

**Slide 96**

- Value of  $C > 0$  reflects cost of violating constraints
  - A large  $C$  generally leads to smaller margin but also fewer misclassification of training data
  - A small  $C$  generally leads to larger margin but more misclassification of training data
- As a design parameter, value of  $C$  is set by user

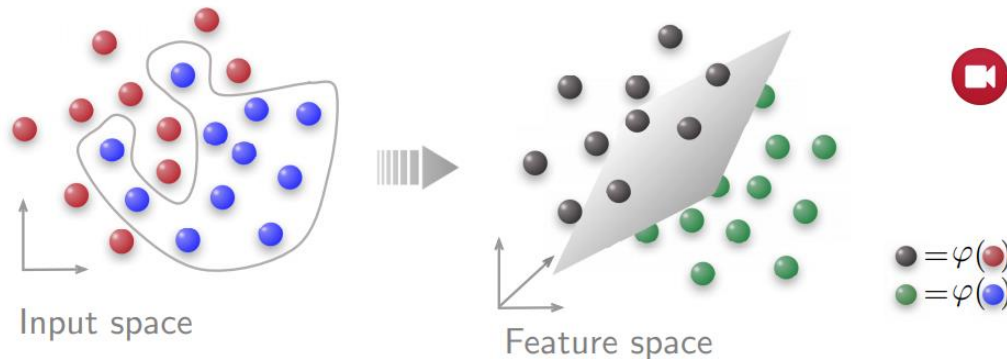
Soft Margin



# Recap – Soft Margin

Dealing with non-separable patterns:

2. Transform data into higher dimension space for separation [Slide 94](#)



Dual problem with soft margin and transformation [Slide 115](#)

Find :  $\alpha_i$

Maximize :  $Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j \varphi^T(\mathbf{x}_i) \varphi(\mathbf{x}_j)$  ← Nonlinear kernel

Subject to :  $\sum_{i=1}^N \alpha_i d_i = 0, 0 \leq \alpha_i \leq C$  ← Soft Margin



# Task 1 - Data

---

## Training set – 2000 samples

- Given – 'train.mat'

- Features – (57 x 2000)
- Label (2000 x 1)

- Features of a sample

```
0.00000 0.01043 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
0.00000 0.00000 0.01043 0.01043 0.02105 0.00000 0.00000 0.00000
0.00000 0.03166 0.06332 0.00000 0.02105 0.00000 0.00000 0.00000
0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
0.00000 0.00000 0.00000 0.00196 0.00000 0.00000 0.02601 0.12811
0.98827
```

- Label : +1 (spam), -1 (non-spam)



# Task 1 – Training set

---

Import the training set (i.e. train.mat)

- train\_data (57 x 2000)
- train\_label (2000 x 1)

Preprocess the 'data' (Various methods can be used including **Sample scaling** and **Standardization** [**CHOOSE ONE METHOD**])<sup>a, b</sup>

- **Scale** the data – Rescale the individual sample  $x$  such that  $\|x\| = 1$
- **Standardize** the data – Transform each feature by removing the mean value of each feature and then dividing by each feature's standard deviation

Please ensure the 'label' is mapped into the set of  $\{-1, +1\}$

<sup>a</sup> <https://scikit-learn.org/stable/modules/preprocessing.html>

<sup>b</sup> [https://en.wikipedia.org/wiki/Feature\\_scaling](https://en.wikipedia.org/wiki/Feature_scaling)



# Task 1 – Kernels

---

Hard-margin SVM with the linear kernel

$$K(x_1, x_2) = x_1^T x_2$$

Hard-margin SVM with a polynomial kernel

$$K(x_1, x_2) = (x_1^T x_2 + 1)^p$$

Soft-margin SVM with a polynomial kernel

$$K(x_1, x_2) = (x_1^T x_2 + 1)^p$$



# Task 1 – Hard and Soft Margins

---

Hard margin  $0 \leq \alpha_i$

- $C = +\infty$  (In theory)

$$0 \leq \alpha_i \leq C$$

- $C = \text{Large value (In practice e.g. } 10^6)$

Soft margin  $0 \leq \alpha_i \leq C$

- $C = 0.1, 0.6, 1.1, 2.1$



# Task 1 – Calculate $\alpha_i$

---

How to calculate  $\alpha_i$

$$\begin{aligned} \text{Maximize : } Q(\alpha) &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j K(\mathbf{x}_i, \mathbf{x}_j) \\ \text{Subject to : } \sum_{i=1}^N \alpha_i d_i &= 0, \quad 0 \leq \alpha_i \leq C \end{aligned}$$

Use **quadprog** function (Quadratic programming)

## Description

Solver for quadratic objective functions with linear constraints.

quadprog finds a minimum for a problem specified by

$$\min_x \frac{1}{2} x^T H x + f^T x \text{ such that } \begin{cases} A \cdot x \leq b, \\ Aeq \cdot x = beq, \\ lb \leq x \leq ub. \end{cases}$$

$H$ ,  $A$ , and  $Aeq$  are matrices, and  $f$ ,  $b$ ,  $beq$ ,  $lb$ ,  $ub$ , and  $x$  are vectors.

You can pass  $f$ ,  $lb$ , and  $ub$  as vectors or matrices; see [Matrix Arguments](#).

`x = quadprog(H,f,A,b,Aeq,beq,lb,ub,x0,options)` solves the preceding problem using the optimization options specified in `options`. Use `optimoptions` to create options. If you do not want to give an initial point, set `x0 = []`.



# Task 1 – Calculate $\alpha_i$ - quadprog

Maximize

$$Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j K(\mathbf{x}_i, \mathbf{x}_j)$$

Subject to :  $\sum_{i=1}^N \alpha_i d_i = 0, 0 \leq \alpha_i \leq C$

## Description

Solver for quadratic objective functions with linear constraints.

quadprog finds a minimum for a problem specified by

$$\min_x \frac{1}{2} x^T H x + f^T x \text{ such that } \begin{cases} A \cdot x \leq b, \\ Aeq \cdot x = beq, \\ lb \leq x \leq ub. \end{cases}$$

Convert the problem from 'Max' to 'Min'

- Max  $Q(\alpha) \rightarrow$  Min -  $Q(\alpha)$

If  $f$  is to be maximized instead, such a maximization problem can be expressed as a minimization problem by the transformation

$$\max_{\mathbf{w}} f(\mathbf{w}) = - \min_{\mathbf{w}} [-f(\mathbf{w})]$$

Slide 62





# Task 1 – Calculate $\alpha_i$ - quadprog

Maximize :

$$Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j K(\mathbf{x}_i, \mathbf{x}_j)$$

Subject to :

$$\sum_{i=1}^N \alpha_i d_i = 0, \quad 0 \leq \alpha_i \leq C$$

## Description

Solver for quadratic objective functions with linear constraints.

quadprog finds a minimum for a problem specified by

$$\min_x \frac{1}{2} x^T H x + f^T x \text{ such that } \begin{cases} A \cdot x \leq b, \\ Aeq \cdot x = beq, \\ lb \leq x \leq ub. \end{cases}$$

Not used

$$\min_x \frac{1}{2} x^T H x + f^T x \quad \left\{ \begin{array}{l} H_{ij} = d_i d_j K(x_i, x_j) \\ f = (-1, -1, \dots, -1)^T \end{array} \right.$$

$$A \cdot x \leq b, \quad \left\{ \begin{array}{l} A = [] \\ b = [] \end{array} \right.$$

$$Aeq \cdot x = beq, \quad \left\{ \begin{array}{l} Aeq = (d_1, d_2, \dots, d_N) \\ beq = 0 \end{array} \right.$$

$$lb \leq x \leq ub. \quad \left\{ \begin{array}{l} lb = (0, 0, \dots, 0)^T \\ ub = (C, C, \dots, C)^T \end{array} \right.$$

`x = quadprog(H,f,A,b,Aeq,beq,lb,ub,x0,options)`



# Task 1 – Calculate $\alpha_i$ - quadprog

`x = quadprog(H,f,A,b,Aeq,beq,lb,ub,x0,options)`

Hard-margin SVM with the linear kernel

$$K(x_1, x_2) = x_1^T x_2$$

For illustration only

$$\begin{aligned} \min_x \frac{1}{2} x^T H x + f^T x & \quad \left\{ \begin{array}{l} H(i,j) = d_i d_j x_i^T x_j \\ f = -\text{ones}(2000,1) \end{array} \right. \\ A \cdot x \leq b, & \quad \left\{ \begin{array}{l} A = [] \\ b = [] \end{array} \right. \\ Aeq \cdot x = beq, & \quad \left\{ \begin{array}{l} Aeq = \text{train\_label}' \\ beq = 0 \end{array} \right. \\ lb \leq x \leq ub. & \quad \left\{ \begin{array}{l} lb = \text{zeros}(2000,1) \\ ub = \text{ones}(2000,1) * C \end{array} \right. \end{aligned}$$

Hard margin  $0 \leq \alpha_i$   
C =  $+\infty$  (In theory)  
C = Large value (In practice e.g.  $10^6$ )

`x0 = []`    `options = optimset('LargeScale','off','MaxIter',1000)`



# Task 1 – Calculate $\alpha_i$ - quadprog

`x = quadprog(H,f,A,b,Aeq,beq,lb,ub,x0,options)`

Hard-margin SVM with a polynomial kernel

$$K(x_1, x_2) = (x_1^T x_2 + 1)^p$$

For illustration only

$$\min_x \frac{1}{2} x^T H x + f^T x \quad \left\{ \begin{array}{l} H(i,j) = d_i d_j (x_1^T x_2 + 1)^p \\ f = -\text{ones}(2000,1) \end{array} \right.$$

$$A \cdot x \leq b, \quad \left\{ \begin{array}{l} A = [] \\ b = [] \end{array} \right.$$

$$Aeq \cdot x = beq, \quad \left\{ \begin{array}{l} Aeq = \text{train\_label}' \\ beq = 0 \end{array} \right.$$

$$lb \leq x \leq ub. \quad \left\{ \begin{array}{l} lb = \text{zeros}(2000,1) \\ ub = \text{ones}(2000,1) * C \end{array} \right.$$

$$x0 = [] \quad \text{options} = \text{optimset('LargeScale','off','MaxIter',1000)}$$

Hard margin  $0 \leq \alpha_i$   
 $C = +\infty$  (In theory)  
 $C = \text{Large value}$  (In practice e.g.  $10^6$ )



# Task 1 – Calculate $\alpha_i$ - quadprog

`x = quadprog(H,f,A,b,Aeq,beq,lb,ub,x0,options)`

Soft-margin SVM with a polynomial kernel

$$K(x_1, x_2) = (x_1^T x_2 + 1)^p$$

For illustration only

$$\min_x \frac{1}{2} x^T H x + f^T x \quad \left\{ \begin{array}{l} H(i,j) = d_i d_j (x_1^T x_2 + 1)^p \\ f = -\text{ones}(2000,1) \end{array} \right.$$

$$A \cdot x \leq b, \quad \left\{ \begin{array}{l} A = [] \\ b = [] \end{array} \right.$$

$$Aeq \cdot x = beq, \quad \left\{ \begin{array}{l} Aeq = \text{train\_label}' \\ beq = 0 \end{array} \right.$$

$$lb \leq x \leq ub. \quad \left\{ \begin{array}{l} lb = \text{zeros}(2000,1) \\ ub = \text{ones}(2000,1) * C \end{array} \right.$$

Soft margin  $0 \leq \alpha_i \leq C$   
 $C = 0.1, 0.6, 1.1, 2.1$

`x0 = []`      `options = optimset('LargeScale','off','MaxIter',1000)`



# Task 1 – Select support vectors

---

Based on KKT conditions

- For a support vector,  $\alpha_i \neq 0$  (In theory,  $\alpha_i > 0$  )
- However, in practice,  $\alpha_i > \text{threshold}$
- How to decide?
  - Choose an **appropriate threshold** (e.g.  $1e-4$ ) to determine the corresponding  $\alpha_i$  to the support vectors



# Task 1 – Discriminant function $g(\mathbf{x})$

## Hard Margin SVM with Linear Kernel

Maximizing :  $Q(\boldsymbol{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j \mathbf{x}_i^T \mathbf{x}_j$

Subject to : (1)  $\sum_{i=1}^N \alpha_i d_i = 0$   
(2)  $\alpha_i \geq 0$

Discriminant function

$$g(\mathbf{x}) = \mathbf{w}_o^T \mathbf{x} + b_o$$

After  $\alpha_{o,i}$  is obtained, we can calculate  $\mathbf{w}_o$  and  $b_o$  as follows: **Slide 85**

$$\mathbf{w}_o = \sum_{i=1}^N \alpha_{o,i} d_i \mathbf{x}_i, \quad b_o = \frac{1}{d^{(s)}} - \mathbf{w}_o^T \mathbf{x}^{(s)}$$

where  $\mathbf{x}^{(s)}$  is a support vector with label  $d^{(s)}$



# Task 1 – Discriminant function $g(\mathbf{x})$

## Soft Margin SVM with Linear Kernel

$$\text{Maximize : } Q(\boldsymbol{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\text{Subject to : } \sum_{i=1}^N \alpha_i d_i = 0 \quad \text{and} \quad 0 \leq \alpha_i \leq C$$

Discriminant function

$$g(\mathbf{x}) = \mathbf{w}_o^T \mathbf{x} + b_o$$

After  $\alpha_{o,i}$  is obtained, we can calculate  $\mathbf{w}_o$  as follows: **Slide 104,105**

$$\mathbf{w}_o = \sum_{i=1}^N \alpha_{o,i} d_i \mathbf{x}_i$$

After  $\mathbf{w}_o$  is obtained, we can calculate  $b_o$  as follows: ② Take  $b_o$  as the average of all such  $b_{o,i}$

① For each example  $\mathbf{x}_i$  with  $0 < \alpha_i \leq C$ ,

$$b_{o,i} = \frac{1}{d_i} - \mathbf{w}_o^T \mathbf{x}_i$$

$$b_o = \frac{\sum_{i=1}^m b_{o,i}}{m}$$

where  $m$  is the total number of  $\mathbf{x}_i$  with  $0 < \alpha_i \leq C$ .



# Task 1 – Discriminant function $g(\mathbf{x})$

## Soft Margin SVM with Nonlinear Kernel

$$\text{Maximize : } Q(\boldsymbol{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j \boldsymbol{\varphi}^T(\mathbf{x}_i) \boldsymbol{\varphi}(\mathbf{x}_j)$$

$$\text{Subject to : } \sum_{i=1}^N \alpha_i d_i = 0, \quad 0 \leq \alpha_i \leq C$$

Discriminant function

$$g(\mathbf{x}) = \sum_{i=1}^N \alpha_{o,i} d_i K(\mathbf{x}, \mathbf{x}_i) - b_o$$

Determine  $b_o$  in **Slide 123**

$$g(\mathbf{x}) = \sum_{i=1}^N \alpha_{o,i} d_i K(\mathbf{x}, \mathbf{x}_i) + b_o$$

using the fact that for a support vector  $\mathbf{x}^{(s)}$

$$g(\mathbf{x}^{(s)}) = \pm 1 = d^{(s)}$$

Take  $b_o$  as the average of all such  $b_{o,i}$

$$b_o = \frac{\sum_{i=1}^m b_{o,i}}{m}$$

where  $m$  is the total number of  $\mathbf{x}_i$  with  $0 < \alpha_i \leq C$ .





# Task 1 – Summary

Given a training set  $S = \{(\mathbf{x}_i, d_i)\}, i = 1, \dots, N$

Slide 123

1 Find a suitable kernel

Choose expression  
then check Mercer's  
condition

4 Determine  $b_o$  in

$$g(\mathbf{x}) = \sum_{i=1}^N \alpha_{o,i} d_i K(\mathbf{x}, \mathbf{x}_i) + b_o$$

Kernel

Soft Margin

2 Choose a value for  $C$

3 Solve for  $\alpha_{o,i}$

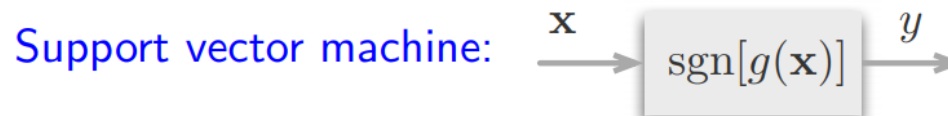
using the fact that for a support  
vector  $\mathbf{x}^{(s)}$

$$g(\mathbf{x}^{(s)}) = \pm 1 = d^{(s)}$$

Quadratic  
programming

Maximize :  $Q(\boldsymbol{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j K(\mathbf{x}_i, \mathbf{x}_j)$

Subject to :  $\sum_{i=1}^N \alpha_i d_i = 0, 0 \leq \alpha_i \leq C$



# Task 2 - Data

---

## Test set – 1536 samples

- Given – 'test.mat'

- Features – (57 x 1536)
- Label (1536 x 1)

- Features of a sample

```
0.00000 0.01043 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
0.00000 0.00000 0.01043 0.01043 0.02105 0.00000 0.00000 0.00000
0.00000 0.03166 0.06332 0.00000 0.02105 0.00000 0.00000 0.00000
0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000 0.00000
0.00000 0.00000 0.00000 0.00196 0.00000 0.00000 0.02601 0.12811
0.98827
```

- Label : +1 (spam), -1 (non-spam)



# Task 2 – Test set

---

Import the test set (i.e. test.mat)

- test\_data (57 x 1536)
- test\_label (1536 x 1)

Preprocess the 'data' (Various methods can be used including **Sample scaling** and **Standardization** [**CHOOSE ONE USE for TRAINING**])<sup>a, b</sup>

- **Scale** the data – Rescale the individual sample  $x$  such that  $\|x\| = 1$
- **Standardize** the data – Transform each feature in the same manner with the training data. Use the mean and variance of each feature from your training set.

Please ensure the 'label' is mapped into the set of  $\{-1, +1\}$

<sup>a</sup> <https://scikit-learn.org/stable/modules/preprocessing.html>

<sup>b</sup> [https://en.wikipedia.org/wiki/Feature\\_scaling](https://en.wikipedia.org/wiki/Feature_scaling)



# Task 2 – Test set

---

Discriminant function

$$g(\mathbf{x}) = \mathbf{w}_o^T \boldsymbol{\varphi}(\mathbf{x}) + b_o = \sum_{i=1}^N \alpha_{o,i} d_i \underbrace{\boldsymbol{\varphi}^T(\mathbf{x}_i) \boldsymbol{\varphi}(\mathbf{x})}_{K(\mathbf{x}_i, \mathbf{x})} + b_o$$

To classify a new data point  $\mathbf{x}_{\text{new}}$

$$d_{\text{new}} = \text{sgn} [g(\mathbf{x}_{\text{new}})]$$

For illustrations only

$$g(x_{\text{test}}) = \sum_{i=1}^N \alpha_{o,i} d_i K(x_i, x_{\text{test}}) + b_o$$

If  $g(x_{\text{test}}) > 0$   
     $x_{\text{test\_label}} = +1$   
Else  
     $x_{\text{test\_label}} = -1$



# Task 2 – Test set

---

Discriminant function

$$g(\mathbf{x}) = \mathbf{w}_o^T \boldsymbol{\varphi}(\mathbf{x}) + b_o = \sum_{i=1}^N \alpha_{o,i} d_i \underbrace{\boldsymbol{\varphi}^T(\mathbf{x}_i) \boldsymbol{\varphi}(\mathbf{x})}_{K(\mathbf{x}_i, \mathbf{x})} + b_o$$

To classify a new data point  $\mathbf{x}_{\text{new}}$

$$d_{\text{new}} = \text{sgn} [g(\mathbf{x}_{\text{new}})]$$

Type of SVM	Training accuracy				Test accuracy			
Hard margin with Linear kernel	?				?			
Hard margin with polynomial kernel	$p=2$ ?	$p=3$ ?	$p=4$ ?	$p=5$ ?	$p=2$ ?	$p=3$ ?	$p=4$ ?	$p=5$ ?
Soft margin with polynomial kernel	$C=0.1$	$C=0.6$	$C=1.1$	$C=2.1$	$C=0.1$	$C=0.6$	$C=1.1$	$C=2.1$
$p=1$	?	?	?	?	?	?	?	?
$p=2$	?	?	?	?	?	?	?	?
$p=3$	?	?	?	?	?	?	?	?
$p=4$	?	?	?	?	?	?	?	?
$p=5$	?	?	?	?	?	?	?	?

If  $g(\mathbf{x}_{\text{test}}) > 0$   
 $\mathbf{x}_{\text{test\_label}} = +1$   
 Else  
 $\mathbf{x}_{\text{test\_label}} = -1$



# Task 3 - Data

---

Evaluation set – 600 samples

- **Not Given** – 'eval.mat'
  - eval\_data (57 x 600)
  - eval\_label (600 x 1)



# Task 3 - Evaluation

---

Design your own SVM

- Hard margin or Soft margin?
- Linear or Polynomial kernel?
- What are the values for  $p$  and  $C$ ?



To produce  
the best  
performance

To classify the 600 samples in the evaluation set

Not Given – 'eval.mat'

eval\_data (57 x 600)

eval\_label (600 x 1)

**Output :** A column vector (600 x 1) named 'eval\_predicted'



# Task 3 - Evaluation

---

Hardcode the discriminant function  $g(x)$  in the file for evaluation

- If necessary, store the required variables in a separate \*.mat file and load at the beginning of the code

Prepare the code so that it could handle the evaluation dataset and able to preprocess the evaluation dataset

- Note: the `eval_data` is a (57 x 600) matrix

Your code should generate a column vector (600 x 1) named `'eval_predicted'`





# Task 3 - Evaluation

---

Please name your m-file for Task 3 as 'svm\_main.m'

**Do Not clear any variable** in the 'svm\_main' script

Before submitting your code, please ensure that the **code runs without errors** by testing it with a dummy data set (Dummy dataset can be created using the training set and test set)



# Important Notes

---

Preprocess your data – Choose one method

- Sample scaling/ Mean normalization/ standardization/ Rescaling ...

Use the training set statistics to preprocess the other data sets

Check Mercer Condition for Kernel suitability



# Important Notes

---

## Procedure to build SVM

- Preprocess data
- Choose a suitable kernel
  - Linear/ Nonlinear ?
- Choose C
  - Hard margin/ Soft margin
    - Hard margin  $0 \leq \alpha_i$ 
      - $C = +\infty$  (In theory)
      - $C = \text{Large value}$  (In practice e.g.  $10^6$ )
- Solve for  $\alpha_i$ 
  - Quadratic programming
- Support vector selection
  - Choose an appropriate threshold (e.g.  $1e-4$ ) to determine the corresponding  $\alpha_i$  to the support vectors
- Determine the discriminant function  $g(x)$



# Important Notes : Submission

---

Submit all your codes that you have implemented for the entire project

**Make sure your codes run without error**

All codes should be executable with the given datasets in the workspace without any additional inputs



# Important Notes : Submission

## Report

- Details on Implementation
- Completed Table 1
- Discuss the results and their implications
  - Admissibility of the kernels
  - Existence of optimal hyperplanes
  - Comments on results (with supporting arguments)

Type of SVM	Training accuracy				Test accuracy			
Hard margin with Linear kernel	?				?			
Hard margin with polynomial kernel	$p = 2$	$p = 3$	$p = 4$	$p = 5$	$p = 2$	$p = 3$	$p = 4$	$p = 5$
	?	?	?	?	?	?	?	?
Soft margin with polynomial kernel	$C = 0.1$	$C = 0.6$	$C = 1.1$	$C = 2.1$	$C = 0.1$	$C = 0.6$	$C = 1.1$	$C = 2.1$
	$p = 1$	?	?	?	?	?	?	?
	$p = 2$	?	?	?	?	?	?	?
	$p = 3$	?	?	?	?	?	?	?
	$p = 4$	?	?	?	?	?	?	?
	$p = 5$	?	?	?	?	?	?	?

TABLE I: Results of SVM classification.



# Important Notes : Submission



svm\_main.m



Other m-files, if needed \*.mat files



Results.mat



A0123456B\_SVM.pdf



Example student No = A0123456B



Folder name = A0123456B



Non password protected

Due date 23 April 2021



---

sandakalum@u.nus.edu

