

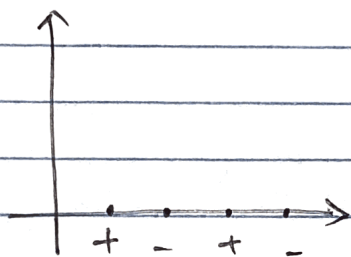
$$1. H = \{ \text{sgn}(ax^2 + bx + c); a, b, c \in \mathbb{R} \}$$

where  $\text{sgn}(ax^2 + bx + c) = 1$  if  $ax^2 + bx + c > 0$

$\text{sgn}(ax^2 + bx + c) = 0$  otherwise.

Since quadratic function could change sign only twice, it cannot shatter into alternating 4 VC dimension. The VC dimension of  $H$  is 3.

Eg:



This cannot be shattered.

$$2. K_{\beta}(x, z) = (1 + \beta x \cdot z)^3$$

$$= 1 + 3\beta x_1 z_1 + 3\beta x_2 z_2 + 3\beta^2 x_1^2 z_1^2 + 6\beta^2 x_1 x_2 z_1 z_2 + 3\beta^2 x_2^2 z_2^2 + \beta^3 x_1^3 z_1^3 + 3\beta^3 x_1 x_2^2 z_1 z_2^2 + 3\beta^3 x_1^2 x_2 z_1^2 z_2 + \beta^3 x_2^3 z_2^3$$

$$K(x, z) = \phi(x)^T \phi(z)$$

$$\therefore \phi_{\beta}(\vec{x}) = \begin{bmatrix} 1 \\ \sqrt{3} x_1 \sqrt{\beta} \\ \sqrt{3} x_1^2 \sqrt{\beta} \\ \sqrt{6} x_1 x_2 \beta \\ x_1^3 \sqrt{\beta^3} \\ \sqrt{3} x_1^2 x_2 \sqrt{\beta^3} \\ \sqrt{3} x_1 x_2^2 \sqrt{\beta^3} \\ \vdots \end{bmatrix}$$

The difference between  $K_{\beta}$  and  $K$  is that  $K_{\beta}$  has the parameter  $\beta$  to scale each term of  $\phi_{\beta}$ . Therefore,  $\beta$  plays the role of adjusting the magnitude of the kernel function.

$$3. a). \eta_1 = (1, 1)^T, \quad \eta_2 = (1, 0)^T$$

$$\therefore w_1 + w_2 \geq 1, \quad -w_1 \geq 1$$

$$1 - w_1, -w_2 \leq 0, \quad 1 + w_1 \leq 0$$

$$L(w, \alpha) = \frac{1}{2}(w_1^2 + w_2^2) + \alpha_1(1 - w_1 - w_2) + \alpha_2(1 + w_1)$$

$$d^* = \max_{\alpha} \min_w L(w, \alpha)$$

$$\frac{\partial L}{\partial w_1} = w_1 - \alpha_1 + \alpha_2 = 0,$$

$$\frac{\partial L}{\partial w_2} = w_2 - \alpha_1 = 0$$

$$\therefore w_1 = \alpha_1 - \alpha_2$$

$$w_2 = \alpha_1$$

Sub  $w_1, w_2$  into  $L(w, \alpha)$ .

$$\begin{aligned} L(w, \alpha) &= \frac{1}{2}((\alpha_1 - \alpha_2)^2 + \alpha_1^2) + \alpha_1(1 - \alpha_1 + \alpha_2 - \alpha_1) + \alpha_2(1 + \alpha_1 - \alpha_2) \\ &= \frac{1}{2}(2\alpha_1^2 - 2\alpha_1\alpha_2 + \alpha_2^2) + (\alpha_1 - \alpha_1^2 + \alpha_1\alpha_2) + (\alpha_2 + \alpha_1\alpha_2 - \alpha_2^2) \end{aligned}$$

$$\frac{\partial L}{\partial \alpha_1} = 2\alpha_1 - \alpha_2 + 1 - 4\alpha_1 + \alpha_2 + \alpha_2 = -2\alpha_1 + \alpha_2 + 1 = 0$$

$$\frac{\partial L}{\partial \alpha_2} = -\alpha_1 + \alpha_2 + \alpha_1 + 1 + \alpha_1 - 2\alpha_2 = \alpha_1 - \alpha_2 + 1 = 0$$

$$\therefore \alpha_1 = \frac{\alpha_2 + 1}{2}$$

$$\alpha_2 = \alpha_1 + 1$$

$$\therefore \alpha_1 = \frac{\alpha_1 + 2}{2}$$

$$\therefore \alpha_1 = 2, \quad \alpha_2 = 3.$$

$$\text{Since } w_1 = \alpha_1 - \alpha_2$$

$$w_2 = \alpha_1$$

$$\therefore w^* = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$$

b). With offset  $b$  be non-zero, we get:

$$\begin{aligned} w_1 + w_2 + b &\geq 1 &\Rightarrow 1 - w_1 - w_2 - b &\leq 0 \\ -(w_1 + b) &\geq 1 &1 + w_1 + b &\leq 0 \end{aligned}$$

$$\therefore L(w, \alpha, b) = \frac{1}{2}(w_1^2 + w_2^2) + \alpha_1(1 - w_1 - w_2 - b) + \alpha_2(1 + w_1 + b)$$

$$\frac{\partial L}{\partial w_1} = w_1 - \alpha_1 + \alpha_2 = 0$$

$$\frac{\partial L}{\partial w_2} = w_2 - \alpha_1 = 0$$

$$\therefore w_1 = \alpha_1 - \alpha_2$$

$$w_2 = \alpha_1$$

$$\Rightarrow L(w, \alpha, b) = \frac{1}{2}(2\alpha_1^2 - 2\alpha_1\alpha_2 + \alpha_2^2) + (\alpha_1 - 2\alpha_1^2 + \alpha_1\alpha_2 - \alpha_1 b) + (\alpha_2 + \alpha_1\alpha_2 - \alpha_2^2 + \alpha_2 b)$$

$$\frac{\partial L}{\partial \alpha_1} = -2\alpha_1 + \alpha_2 + 1 - b = 0$$

$$\frac{\partial L}{\partial \alpha_2} = -\alpha_2 + \alpha_1 + 1 + b = 0$$

$$\therefore \alpha_1 = \frac{\alpha_2 + 1 - b}{2}$$

$$\alpha_2 = \alpha_1 + 1 + b$$

$$\alpha_1 = \frac{\alpha_1 + 1 + b + 1 - b}{2}, \quad \alpha_1 = 2, \quad \alpha_2 = \alpha_1 = 2$$

$$\therefore b = -1, \quad w_1 = 0, \quad w_2 = 2$$

$$\therefore (w^*, b^*) = \left( \begin{bmatrix} 0 \\ 2 \end{bmatrix}, -1 \right), \quad r = \frac{1}{2}$$

$$\text{with } b = 0, \quad r = \frac{1}{\sqrt{2}}$$

Therefore, we could have larger margin within decision boundary with offset.

#### 4. 1

- d). The dimensionality of the feature matrix is:  
dimensionality of the all the tweet: (630, 1140930)  
dimensionality of training data: (560, 1014160)  
dimensionality of test data: (70, 126770)

#### 4.2

- b). It is beneficial to maintain class proportions across folds since the amount of data for each model needs to be reasonable, in order to make the model learn a proper decision boundary. If not doing this, in some extreme case, the model may even not be aware of one of the class for classification.

<b>C</b>	<b>accuracy</b>	<b>F1-score</b>	<b>AUROC</b>
<b>10<sup>-3</sup></b>	0.7089	0.8297	0.5
<b>10<sup>-2</sup></b>	0.7107	0.8306	0.5031
<b>10<sup>-1</sup></b>	0.8060	0.8755	0.7188
<b>10<sup>0</sup></b>	0.8146	0.8749	0.7531
<b>10<sup>1</sup></b>	0.8182	0.8766	0.7592
<b>10<sup>2</sup></b>	0.8182	0.8766	0.7592
<b>Best C</b>	100	100	100

Parameter C is used to decide the boundary, and it could adjust larger slack variables. Larger value C would give a relevantly small margin to avoid misclassification. Smaller value of C would instead give a large margin which may include more misclassification. In our case, the larger C leads to much better performance for our model.

#### 4.3

- a). The best C value is 100.

- c).

<b>Performance Metric</b>	<b>Linear-Kernel SVM score</b>
<b>Accuracy</b>	0.7429
<b>F1_score</b>	0.4375
<b>AUROC</b>	0.6259