

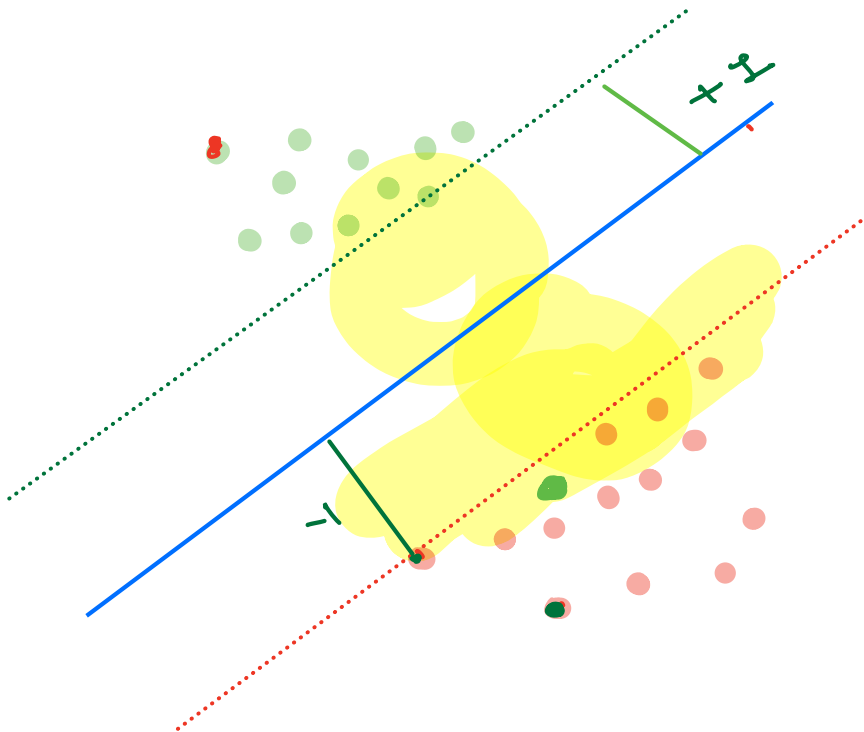
Recap

Hard Margin Classifier

$$\arg\max(\frac{2}{\|w\|})$$

$$\text{s.t. } \forall i: 1 \rightarrow N$$

$$y_i \times w^T x_i + b \geq 1$$



* Support Vectors are point which are responsible for Margin

* No misclassification allowed

Soft Margin Classifier

$$\min_{w, b} \frac{\|w\|}{2} + \frac{C}{N} \sum_{i=1}^N \xi_i$$

$$s.t. (w^T x + b) y_i \geq 1 - \xi_i \quad \forall i: 1 \rightarrow N$$

* Optimal values of w and $b \Rightarrow w^*, b^*$

* Primal Form of SVM

* Allows misclassification for a generalized solution

* In optimization theory

Primal - Dual Equivalence

For every optimization problem there exists a Dual form that gives same results

Dual form of Soft SVM

Kernel



$$\max_{\alpha_i} \left(\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \times \alpha_j \times y_i \times y_j \times x_i \cdot x_j \right)$$

$$\text{Set. } 0 \leq \alpha_i \leq C$$

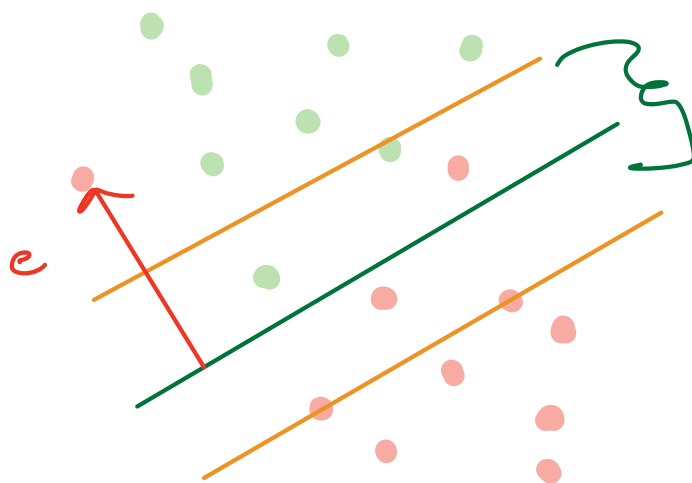
$$\sum_{i=1}^n \alpha_i y_i = 0$$

→ α_i 's are your parameter of Model
i.e. instead of finding w^* and b^*
Here we are going to find
 α_i^*

→ We are maximizing the dual form

→ The input samples always will be in pairs

	x_i	y_i
x_1	_____	—
x_2	_____	—
x_3	_____	—
x_n	_____	—



Support Vectors :

- 1) Points which are within Margin
- 2) Points which are on the Margin line
- 3) Points which are Mis-classified

Prediction in dual Form

$$f(x_q) = \sum_{i=1}^n \alpha_i y_i x_i^T \cdot x_q$$

↑
new query point)

learned during training

↑
Feature
of
new
Data
Point

labels and features
of

Support Vectors

(Training Data)

Let's say we have 10k points

⇒ $\alpha_i \Rightarrow$ for Non Support Vectors
will be 0.

⇒ $\alpha_i > 0$ for Support Vectors

10 k datapoint \rightarrow 1-2 %

100 - 200 Support Vectors

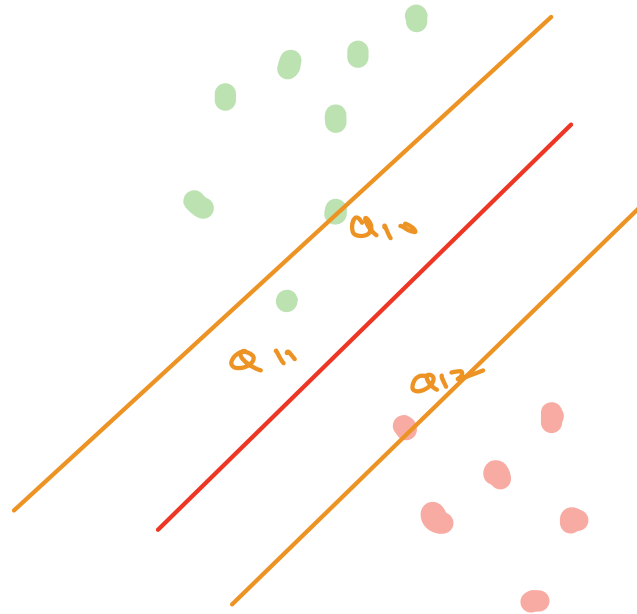
Let's assume we have 100 d.p.
as support vector

$\alpha_i > 0 \rightarrow$ 100 non-zero

$\alpha_i = 0 \rightarrow$ 9900 zero

Ex:

$\alpha_{10} \rightarrow [x_{10}, y_{10}]$
 $\alpha_{11} \rightarrow [x_{11}, y_{11}]$
 $\alpha_{12} \rightarrow [x_{12}, y_{12}]$



$$\hat{y}_q = \sum_{i=1}^n \alpha_i y_i x_i^T \cdot x_q$$

x_i	y_i
x_1	y_1
x_2	y_2
x_{10}	y_{10}
x_{11}	y_{11}

x_n y_n

$$\hat{y} = \alpha_{10} y_{10} x_{10}^T \cdot x_q + \alpha_{11} y_{11} x_{11}^T \cdot x_q + \alpha_{12} y_{12} x_{12}^T \cdot x_q$$

Summary

▷ α_i for support vector

▷ x_i, y_i for support vector

$$\max_{\alpha_i} \left(\sum_{i=1}^M \alpha_i - \frac{1}{2} \sum_{i=1}^M \sum_{j=1}^M \alpha_i \times \alpha_j \times y_i \times y_j \times x_i^T \cdot x_j \right)$$

s.t. $0 \leq \alpha_i \leq C$
 $\sum_{i=1}^M \alpha_i y_i = 0$

$K(x_i, x_j)$

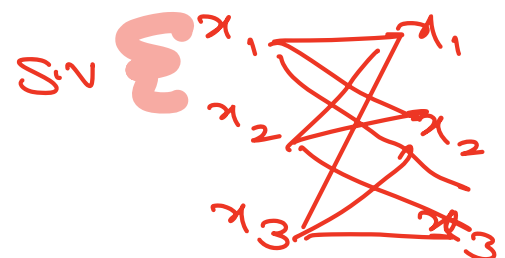
The input samples always will be in pairs

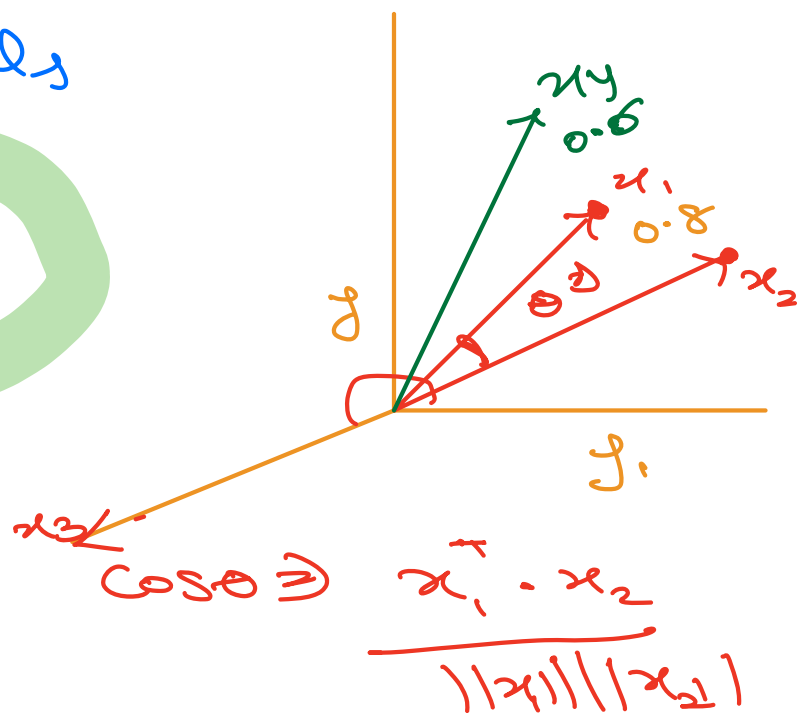
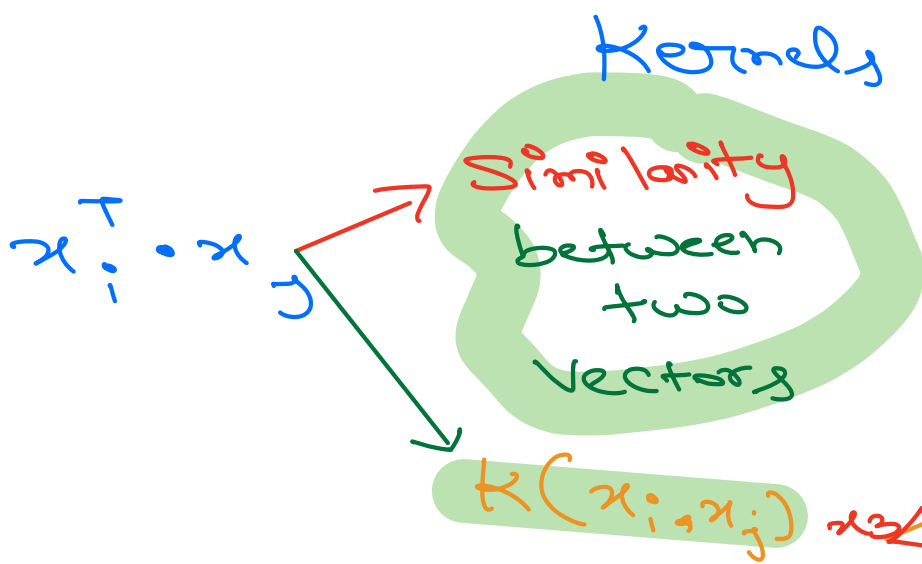
x_i and x_j are all the possible pairs

$$x_1 \cdot x_2 \Rightarrow \alpha > 0$$

$$x_3 \cdot x_2 \Rightarrow \alpha = 0$$

$$x_3 \cdot x_1 \Rightarrow \alpha = 0$$

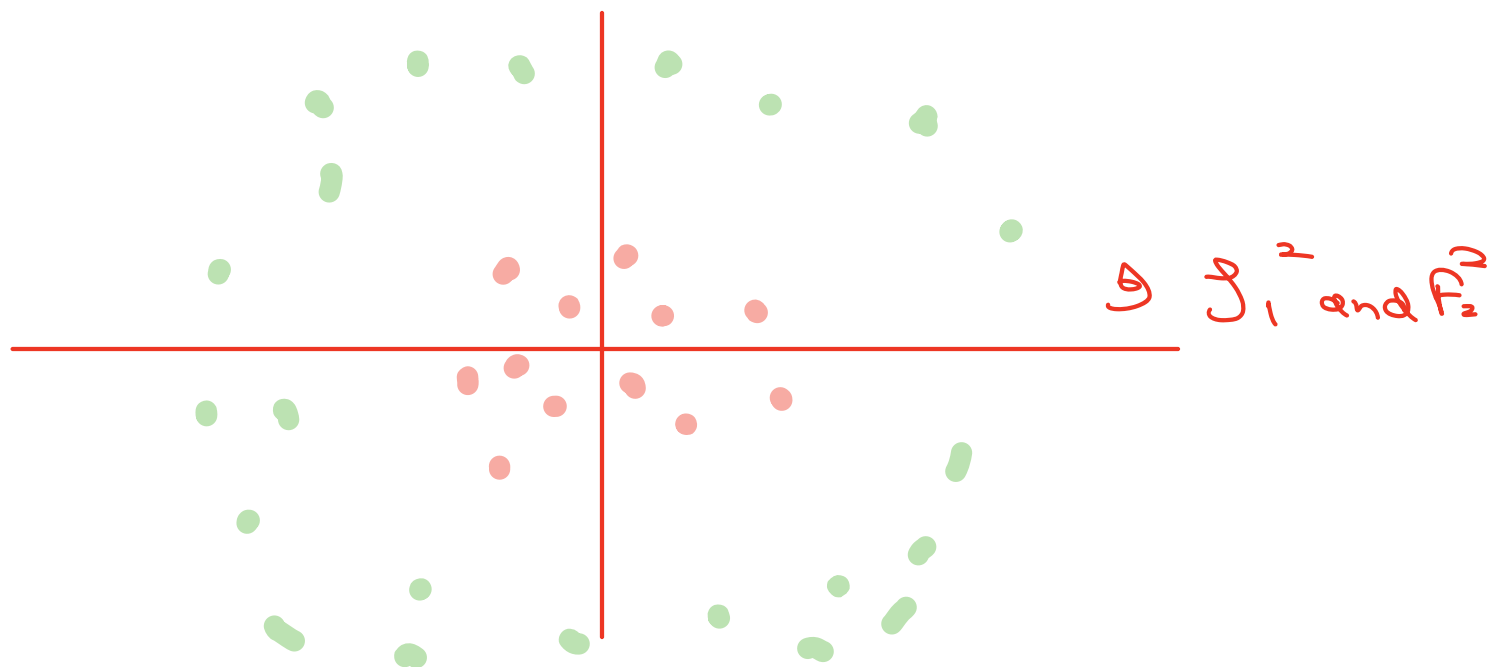




$$\theta \Rightarrow \cos^{-1} \frac{x_1^T \cdot x_2}{\|x_1\| \|x_2\|}$$

$$\theta_{x_1, x_2} \Rightarrow 0.8$$

$$\theta_{x_2, x_3} \Rightarrow -1$$





Linear Model
(Logistic)
or
SVM

Manual feature Transformation

Polynomial kernel

$$k(x_1, x_2) \ni (x_1^T x_2 + c)^n \leftarrow \begin{array}{l} \text{degree} \\ \text{Constant} \\ \Rightarrow 1 \end{array}$$

Quadratic kernel $\ni (x_1^T \cdot x_2 + 1)^2$

Cubic kernel $\ni (x_1^T \cdot x_2 + 1)^3$

Example: Quadratic kernel

$$x_1 \ni [x_{11}, x_{12}]$$

$$x_2 \ni [x_{21}, x_{22}]$$

$$\begin{array}{l} x_1 \ni [x_{11}, x_{12}] \\ x_2 \ni [x_{21}, x_{22}] \end{array}$$

$$K(x_1, x_2) \ni (x_1^T \cdot x_2 + 1)^2$$

$$\ni \left(1 + [x_{11}, x_{12}] \begin{bmatrix} x_{21} \\ x_{22} \end{bmatrix}\right)^2$$

\Downarrow

$$\ni (1 + x_{11}x_{21} + x_{12}x_{22})^2$$

$$(a+b+c)^2 \ni a^2 + b^2 + c^2 + 2ab + 2bc + 2ca$$

$$\ni 1 + x_{11}^2 + x_{12}^2 + 2x_{11}x_{21} + 2x_{12}x_{22} + 2x_{11}x_{21}x_{12}x_{22}$$

\downarrow

$$x_1' \ni [1, x_{11}, x_{12}, \sqrt{2}x_{11}, \sqrt{2}x_{12}, \sqrt{2}x_{11}x_{12}]$$

$$x_2' \ni [1, x_{21}, x_{22}, \sqrt{2}x_{21}, \sqrt{2}x_{22}, \sqrt{2}x_{21}x_{22}]$$

$$K(x_1, x_2) \ni (x_1'^T \cdot x_2')$$

$(1 + x_1^T \cdot x_2)^2$

6 dim vector

* Kernel trick performs implicit
feature transformation

$$\mathbb{R}^d \longrightarrow K_d \longrightarrow \mathbb{C}^d$$

\downarrow
 $d=2$

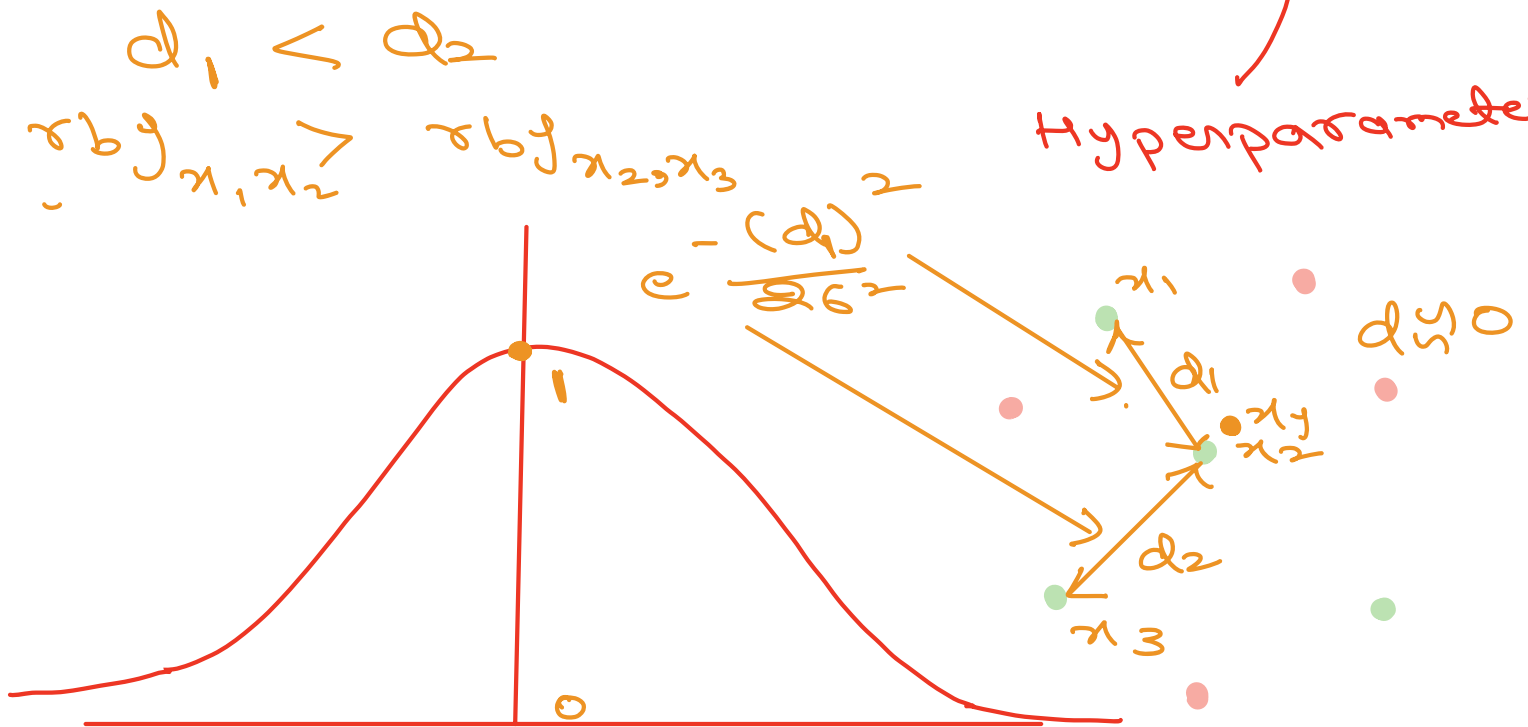
Solving the kernel function is
equivalent to finding hyperplane
in higher dimension

1) you need to Tune d
or
you can use RBF

RBF (Radial Basis Function)

$$K_{\text{rbf}}(x_1, x_2) = e^{-\frac{\overbrace{\|x_1 - x_2\|^2}^{\text{Euclidean Dist}}}{2\sigma^2}}$$

Hyperparameter



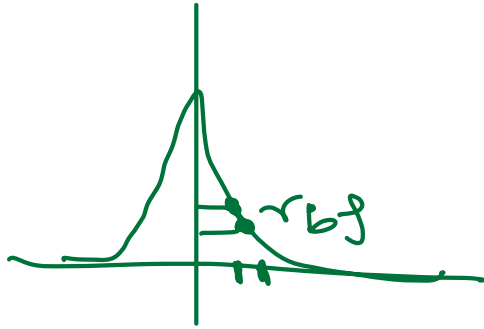
as σ increases the spread of distribution also increases

$\text{rbf}(x_2, x_3) \rightarrow d_{50}$
 $K \rightarrow e^{-0} \rightarrow 1$

(Point are highly Similar)

$\sigma \Rightarrow \text{low}$

σ , τ_{bf} value
will not change
much

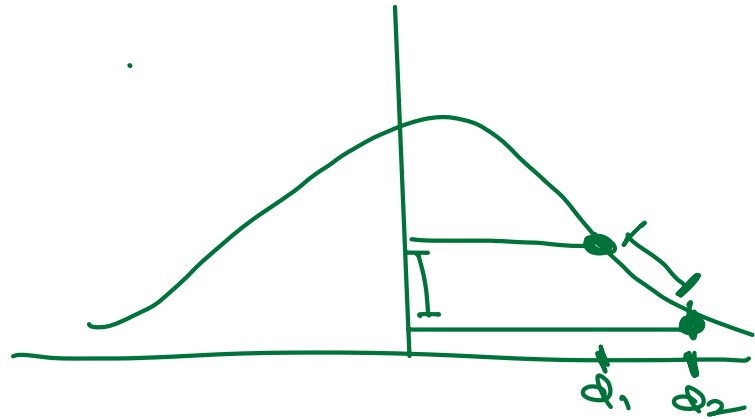


$x_1 \rightarrow x$

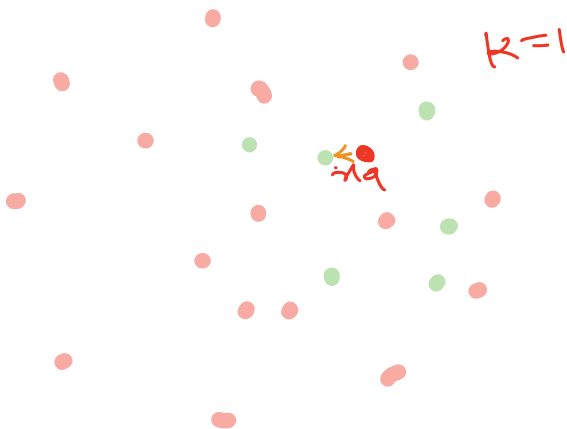
$x \rightarrow x_2$

$\sigma \Rightarrow \text{High}$

Similarity for
Closer points will
appear high



$k = ?$
 σ



Prediction

$$f(x_q) = \sum_{i=1}^M \alpha_i y_i x_i^T \cdot x_q$$



$$f(x_q) = \sum_{i=1}^M \alpha_i y_i K(x_i^T, x_q)$$



$\alpha_i > 0$ only for Support Vectors

$O(\text{no of Support Vector})$
↑
 $O(n \times d)$ \xrightarrow{x} \xrightarrow{d}

SVM will be much faster for predictions than KNN

Effect of Outlier

1) If outlier are Non Support Vectors then No problem (Non kernel SVM)

2) Kernel SVM since we are calculating distance/similarity, they will have higher impact of outliers

Training Time of SVM

$O(n \times n)$ \Rightarrow Very High

\Rightarrow High Training Time Complexity

