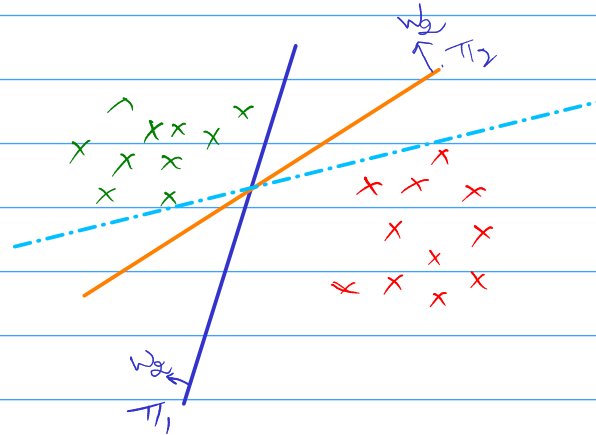


Support Vector Machine: (SVM)

Both Classification and Regression

Geometric Intuition:

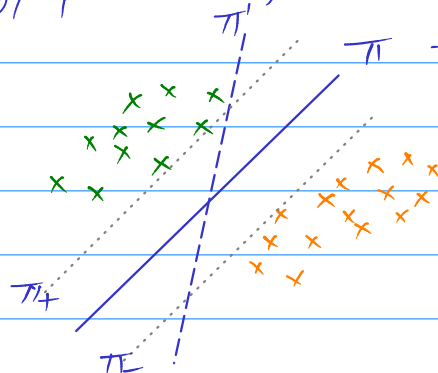
- Separates +ve pts from -ve points as widely as possible.



Many planes / Hyperplanes separate the two classes.

points that are close to hyperplane have less probability of belonging to a class that hyperplane classifies that pts. that are further to the hyperplane.

π_+ , π_- touch the 1st point in their respective class.



π - Separates points as farthest as possible than π'

π_- margin maximising hyperplane.

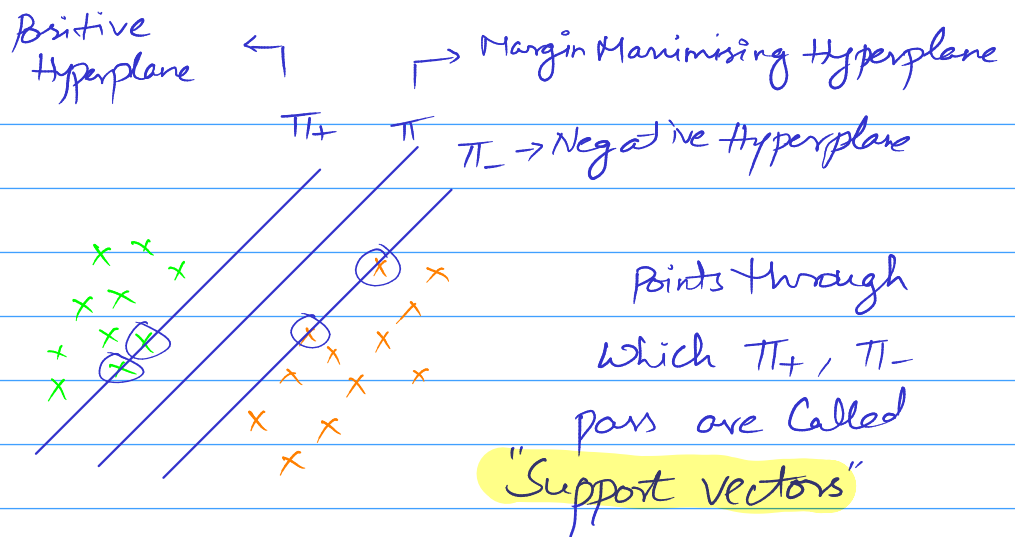
Let π_+ be a plane \parallel to π
 π_- be a plane \parallel to π

$\text{dist}(\pi_+, \pi_-) = d$, is constant
↳ Margin

SVM: Tries to minimise the margin

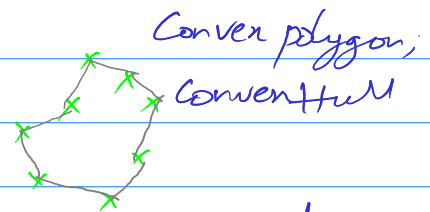
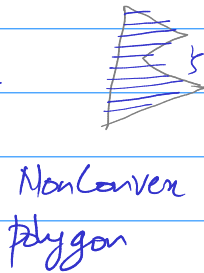
Margin \uparrow ; Generalization accuracy \uparrow
↳ Accuracy on unseen data; New data

Support Vector:



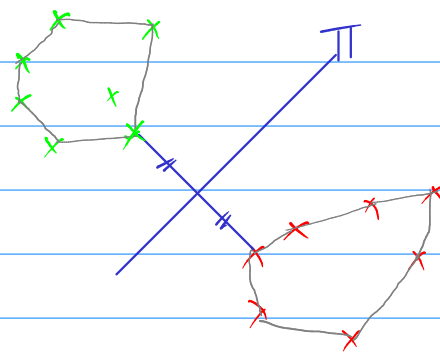
Alternative Geometric Intuition of SVM:

- ① Draw a Convex Hull for +ve points & -ve points separately.



Smallest polygon so that all points lie inside or on the polygon.

- ② Draw shortest line connecting both the hulls.
- ③ Bisect the line; the Hyperplane Bisecting the line will be the Margin minimising plane



Mathematical formulation of SVM:

π : Margin maximising Hyperplane

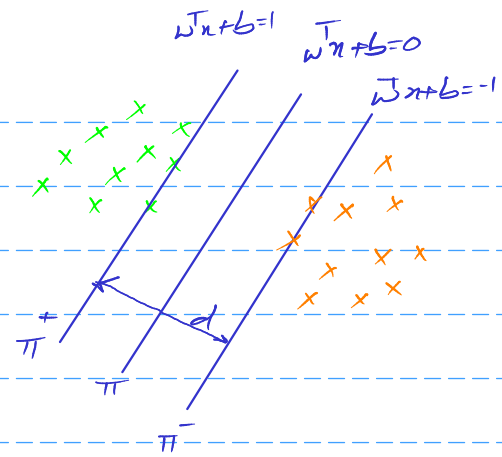
Let,

$$\pi : w^T x + b = 0$$

if,

$$\pi^+ : w^T x + b = 1$$

$$\pi^- : w^T x + b = -1$$



Margin $d = \frac{2}{\|w\|} \rightarrow L_2 \text{ Norm}$

$$(w^*, b^*) = \underset{w, b}{\operatorname{argmax}} \frac{2}{\|w\|} = \text{Margin maximising}$$

Such that

$$y_i(w^T x + b) \geq 1 \quad \forall x_i$$

$\left\{ \begin{array}{l} \vdots \\ n \text{ constraints for all } n \text{ points} \end{array} \right.$

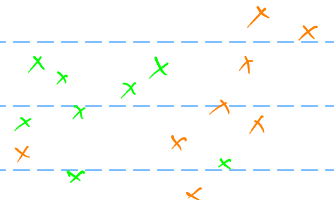
$$w^*, b^* = \underset{w, b}{\operatorname{argmax}} \frac{2}{\|w\|}$$

$$\text{such that } \forall i, y_i(w^T x_i + b) \geq 1$$

Hard Margin SVM

$\left\{ \begin{array}{l} \text{constrained} \\ \text{optimisation} \\ \text{problem} \end{array} \right.$ Assumes data is strictly linearly separable.

What if data is not linearly separable?



Almost linearly separable data \rightarrow

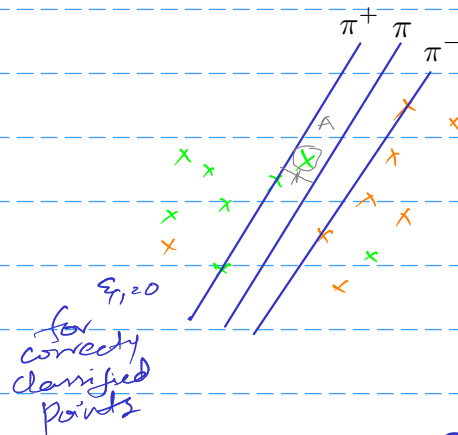
In such a case there will be no w^*, b^* that satisfies optimisation problem.

So Above optimisation problem is called Hard Margin SVM

$$\pi : w^T x + b = 0$$

$$\pi^+ : w^T x + b = +1$$

$$\pi^- : w^T x + b = -1$$



1. Create a new variable

ζ_i (zeta i)

$$\zeta_i = 0, \text{ if } y_i(w^T x + b) \geq 1$$

↳ Correctly classified

$\zeta_i > 0$ it is equal to

some units of distance

away from the hyperplane in incorrect direction.

For example point A,

$$y_i(w^T x + b) = 0.5$$

$$= 1 - \underbrace{(0.5)}_{\zeta_i}$$

0.5 Distance units

away from hyperplane π

Softmargin SVM: Regularization

$$(w^*, b^*) = \underset{w, b}{\operatorname{argmin}} \frac{\|w\|}{2} + C \cdot \frac{1}{n} \sum_i \zeta_i$$

\uparrow Hyperparameter
 \downarrow 1/(margin)

Loss (Hinge)

or
Avg distance of misclassified

points from correct π

such that

$$\begin{cases} y_i(w^T x_i + b) \geq 1 - \zeta_i \quad \forall i \\ \zeta_i \geq 0 \end{cases}$$

$$\min_{w, b} C \cdot \text{loss} + \text{regularization}$$

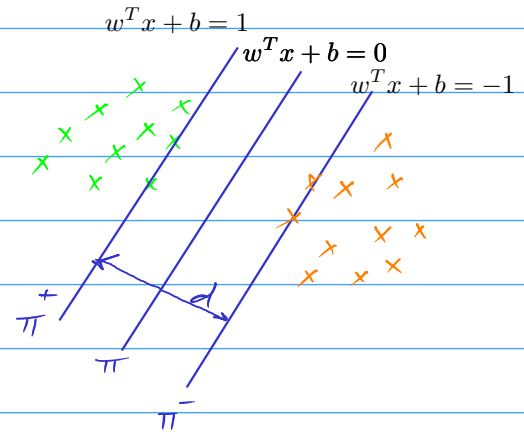
$C \uparrow$; Tendency to make mistakes on train data \downarrow (or) Overfit
High variance

$C \downarrow$ underfit; High bias

Why +1 and -1 for π^+ and π^- :

* w could be any vector and w does not have to be a unit vector
 $\|w\| \neq 1$

$$w^*, b^* = \underset{w, b}{\operatorname{argmax}} \frac{2}{\|w\|}$$



$$\textcircled{1} \pi^+ :- w^T x + b = k \quad k > 0$$

$$\pi^- :- w^T x + b = -k \quad k \text{ is a constant}$$

$$\text{margin} = \frac{2k}{\|w\|}$$

$$\underset{w, b}{\operatorname{argmax}} \frac{2}{\|w\|} = \underset{w, b}{\operatorname{argmax}} \frac{2k}{\|w\|}$$

Some optimisation problem
 as k is a constant

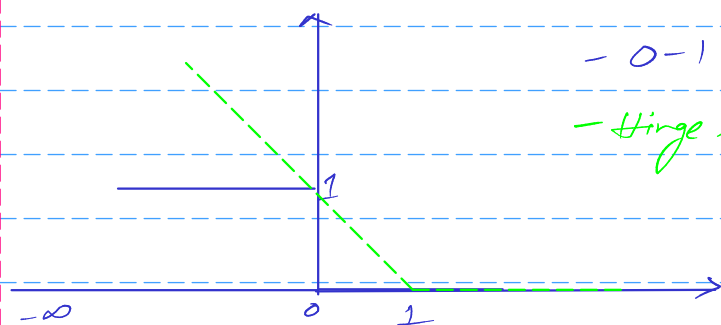
$$\textcircled{2} w^T x + b = k$$

$$\left(\frac{w}{k}\right)^T x + \frac{b}{k} = 1$$

you can have any value of k
 but it does not occur in
 optimisation problem.

This is possible only because
 w is not a unit vector and
 $\frac{2}{\|w\|}$ is not a constant

Loss minimization: Hinge - loss:



- 0-1 loss
- Hinge loss

Logistic Regression: Logistic Loss + Regularization

Linear Regression: Linear Loss + Regularization

SVM: Hinge Loss + Regularization

$$z_i = y_i (w^T x + b)$$

$z_i \geq 1$; hinge loss = 0

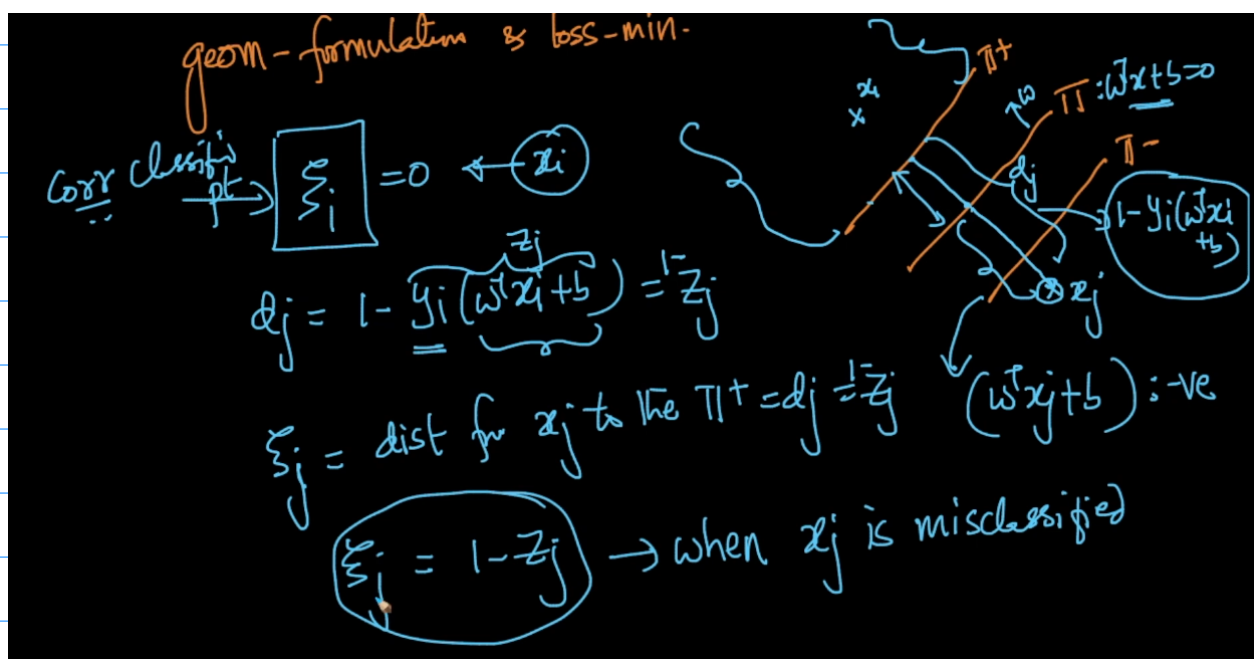
$z_i < 1$; hinge loss = $1 - z_i$

Case 1: $z_i \geq 1$, $1 - z_i$ is -ve $\Rightarrow \max(\quad) = 0$

Case 2: $z_i < 1$, $1 - z_i > 0 \Rightarrow \max(\quad) = 1 - z_i$

(or) $\max(0, 1 - z_i)$

Geometric formulation & loss minimization:



Soft SVM:

$C \uparrow \Rightarrow \text{overfit}$

$C \downarrow \Rightarrow \text{underfit}$

Primal of SVM

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

such that $(1 - y_i(w^T x_i + b)) \geq \xi_i$

$$\xi_i \geq 0$$

Loss min:-

$$\min_{w, b} \sum_{i=1}^n \max(0, 1 - y_i(w^T x_i + b)) + \lambda \|w\|^2$$

$\lambda \uparrow \Rightarrow \text{underfit}$

$\lambda \downarrow \Rightarrow \text{overfit}$

Equivalent

Dual form of SVM:

$$\max_{\alpha_i} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \underbrace{x_i^T x_j}_{\substack{\text{kernel} \\ \text{function}}} \quad \text{such that}$$

$0 \leq \alpha_i \leq C$

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

$\alpha_i = 0$ for non support vectors
 $\alpha_i > 0$ for support vectors

① for every $x_i \rightarrow \alpha_i$

$$x_i^T x_j = x_i \cdot x_j = \cos(\theta) \text{ if } \|x_i\| = \|x_j\| = 1$$

② x_i 's only occurs in form $x_i^T x_j$

③ $f(x_q) = \sum_{i=1}^n \alpha_i y_i x_i^T x_q + b \rightarrow \text{for a query point } x_q$

④ $\alpha_i > 0$ only for support vectors
 $\alpha_i = 0$ for non support vectors

SVM only points that matter are support vectors

since $\alpha_i = 0$ for non support vectors.

Kernel Trick: [Similarity matrix & similarity between pairs of points]

one kind of similarity is a kernel function

$K(x_i, x_j)$; K - kernel function

If no kernel function is used it is called \rightarrow Linear SVM

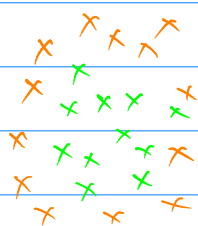
Else \rightarrow kernel SVM

Linear SVM: (x_i, x_j)

Linear SVM is logistic regression

Kernel SVM: $k(x_i, x_j)$

Except for Margin maximization



If dataset looks like this,

Linear SVM ✗ fail

Logistic Regression ✗ fail

Logistic Regression + Feature Transform ✓ Succeed

Kernel SVM ✓ Succeed
with right kernel

Kernel SVM also does similar like feature transforms

$$x_i \rightarrow x_i'$$

Polynomial Kernel:

Kernelization

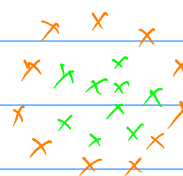
$$k(x_1, x_2) = (x_1^T x_2 + C)^d$$

Ex: $k(x_1, x_2) = (1 + x_1^T x_2)^2$; Quadratic kernel

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

$$x_1 = \langle x_{11}, x_{12} \rangle$$

$$x_2 = \langle x_{21}, x_{22} \rangle$$

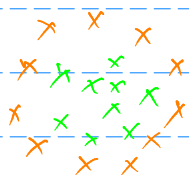


$$= 1 + x_1^2 x_2^2 + x_2^2 x_3^2 + 2x_1 x_2 + 2x_2 x_3 + 2x_1 x_2 x_3$$

$$= (x_1')^T (x_2') \quad \text{where,}$$

$$\begin{aligned} x_1' &= [1, x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2] \\ x_2' &= [1, x_2^2, x_3^2, \sqrt{2}x_2, \sqrt{2}x_3, \sqrt{2}x_2x_3] \end{aligned} \quad \left. \begin{array}{l} \text{From 2D to} \\ 6D \end{array} \right\}$$

Kernelization, $d \xrightarrow[\text{internally implicit transformation}]{\text{Feature transform}} d'$ where d' usually $> d$



For this dataset,
polynomial kernel with degree 2
will work

Kernelization vs Feature Transforms

↓	↓
Similarity b/w 2 points	Transformation
Implicit	Explicit
Defined on data points - pairs of data points	Defined on feature

The beauty of kernels is that we don't obtain the explicit representation of points in high-dimensions.

We simply obtain the similarity between the high dimensional points through the kernel function without the need to obtain the points explicitly.

High-dimensional data is not always bad and is sometimes preferred especially in the case of classification using hyperplanes as higher-dimensions allow for easy splitting of data-points using a hyper-plane than lower-dimensional data.

Higher dimensional data certainly poses problems with visualization, but if our goal is to classify data, why bother too much about visualization.

$$d = \|x_1 - x_2\|_2$$

x_1 \longleftrightarrow x_2

Radial Basis Function (RBF):

$$K_{RBF}(x_1, x_2) = \exp\left(\frac{-\|x_1 - x_2\|^2}{2\sigma^2}\right) \quad ; \quad \sigma - \text{hyper parameter}$$

$$= \exp\left(\frac{-d_{12}^2}{2\sigma^2}\right)$$

$$\frac{1}{e^{(d/2\sigma)^2}}$$

① $d_{12} \uparrow ; K(x_1, x_2) \downarrow$

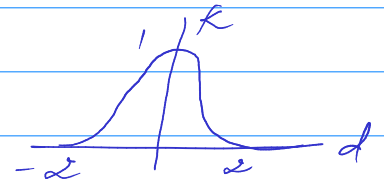
Like similarity

② Impact of σ :

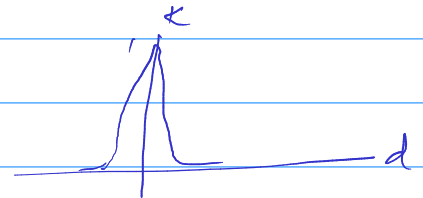
* $\sigma = 1 ; \quad d = 0$

$k = 1$

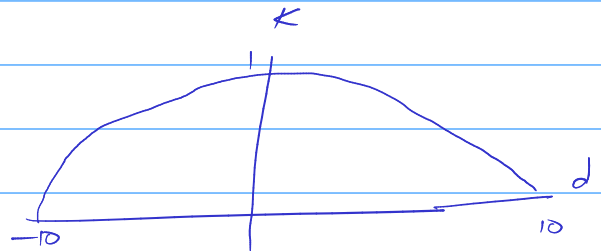
* $d \uparrow ; k \downarrow$ exponentially



* $\sigma = 0.1 ; \sigma^2 = 0.01$



* $\sigma = 10 ; \sigma^2 = 100$



KNN & RBF - Kernel :-

$\sigma \uparrow \implies k \uparrow$ in KNN ; Popular due to its similarity with KNN

RBF is a good approximation of KNN ; KNN requires high space complexity and time complexity.

Domain Specific Kernels:-

- String Kernels
- Genome Kernels

Feature Transformation is hard
↳ partially replaced by
Right Kernel

RBF is a General purpose Kernel

Train & RunTime Complexity:

Train: \rightarrow SGD

↳ specialized algorithm (dual) \rightarrow Sequential minimal optimization (SMO)

libsvm - library for training SVM

Training time: $O(n^2)$ for Kernel SVMs

SVM is typically not used when n is large

RunTime:

$$f(x_q) = \sum_{i=1}^n \sigma_i y_i k(x_i, x_q) + b$$

No. of support vectors $= K$

we have no
control over K

$$O(Kd)$$

nu-SVM: Control errors and support vectors: ν -SVM

Original formulation of SVM \rightarrow C-SVM $C \geq 0$

Alternative formulation of SVM:-

ν -SVM:- $0 \leq \nu \leq 1$

$\nu \geq$ fraction of errors

$\nu \leq$ fraction of support vectors

If $\nu = 0.01 \Rightarrow$ % of errors = 1%

Support vectors \geq 1% of n (n -datapoints)

$n = 100,000$

Support vectors $k \geq 1000$

SVM Regression:-

$$\min_{w, b} \frac{1}{2} \|w\|^2$$

$$f(x) = w^T x_i + b = \hat{y}_i$$

such that $y_i - (w^T x_i + b) \leq \epsilon$

$$(w^T x_i + b) - y_i \leq \epsilon$$

$$\epsilon \geq 0$$

} Linear
SVM

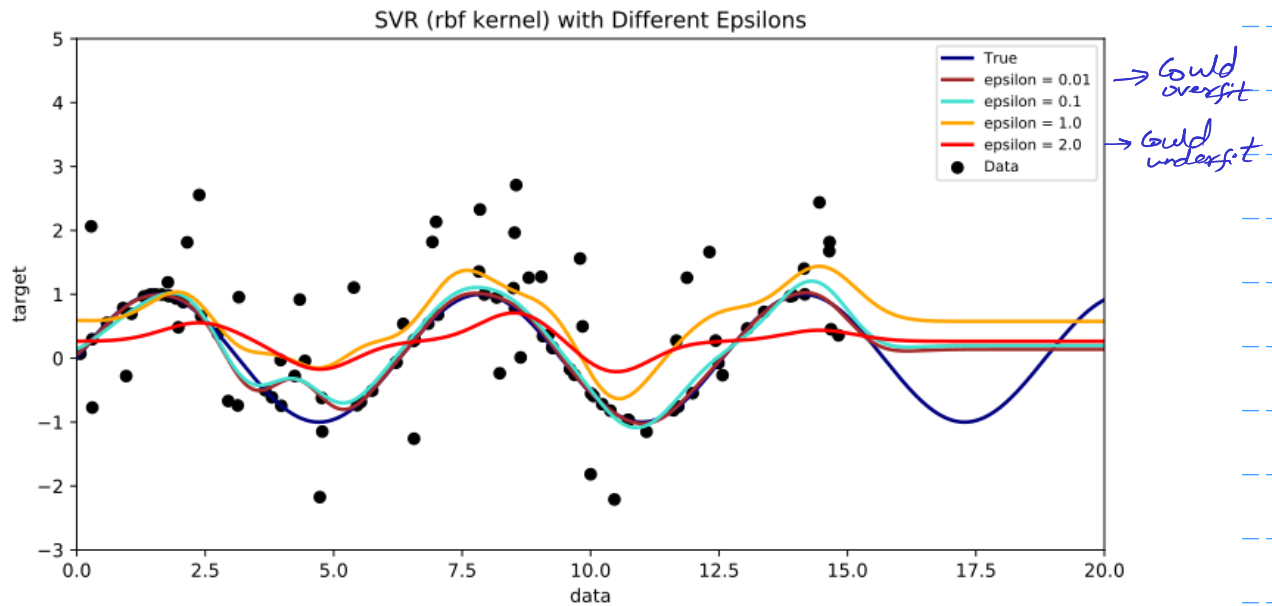
ϵ - Hyperparameter

(or) $|y_i - \hat{y}| \leq \epsilon$ and $\epsilon \geq 0$

$\epsilon \downarrow \Rightarrow$ Errors are low on Train data $\downarrow \Rightarrow$ overfitting \uparrow

$\epsilon \uparrow \Rightarrow$ Errors are Train data $\uparrow \Rightarrow$ underfitting \uparrow

RBF SVR Behaves similar to KNN:-



Cases: (SVM)

* Feature Engg & Feature Transform : Kernel design, finding right kernel

* Decision surface :- Linear SVM:- Hyperplane
Kernel SVM:- $d, x_i \rightarrow$ non linear surface
 \downarrow
 $d', x_i' \rightarrow$ Linear Surface

* Similarity function / distance fn
 \rightarrow kernel

Linear SVM is Interpretable

* Interpretability & Feature Importance \rightarrow Kernel SVM no interpretability
 \rightarrow Forward feature selection
 \rightarrow Can't get feature importance directly

* outliers:- - Very little impact (Because only support vectors matter)
- RBF with small σ will have outlier effect since it behaves like a KNN with small K

* Bias - Variance :- $C \uparrow \Rightarrow \text{overfit} \Rightarrow \text{high variance}$
 $C \downarrow \Rightarrow \text{underfit} \Rightarrow \text{high bias}$

For RBF SVM :- C, γ 2 hyperparameters \rightarrow Grid search

Large $D \uparrow \rightarrow$ Good for SVM

* Best Cases :- Right kernel

Worst Cases :- n is large \rightarrow Low latency not possible
 k is large