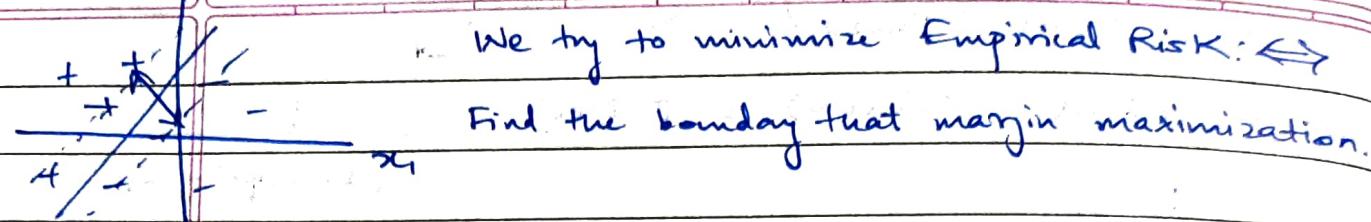


Week 7 SVM:

SVM for distribution free learning:

PAGE NO.	/ / /
DATE	/ / /



We try to minimize Empirical Risk: \Leftrightarrow

Find the boundary that margin maximization.

$t_i \in \{-1, +1\}$ for every classified point

$$t_i (w^T x_i + b) > 0 \quad (\text{Margin fold}) \quad \text{Subject to } \|w\| = 1$$

$$\text{Subject to: } \|w\| = 1 \Rightarrow$$

$$1. \underset{w, b}{\operatorname{Arg\,max}} \left[\frac{1}{\|w\|} \min_i [t_i (w^T x_i + b)] \right]$$

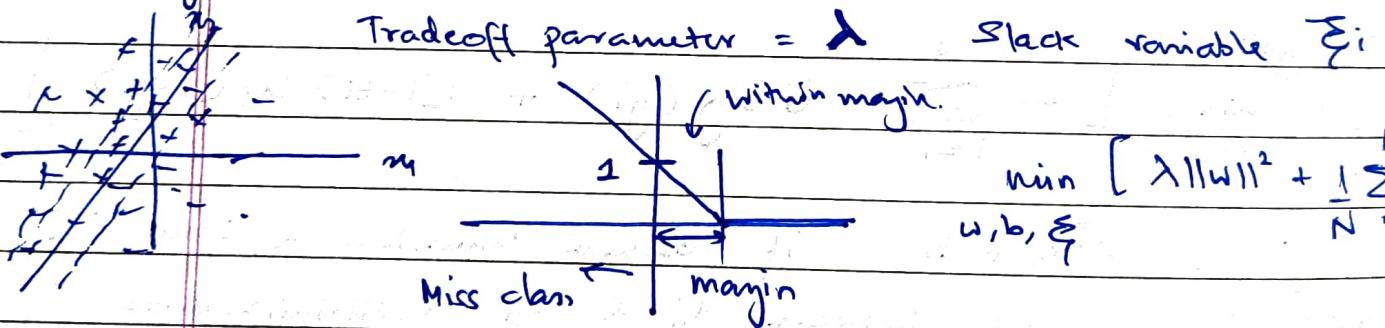
subject to $\|w\| = 1$ (when no normalization)

$$2. \underset{w, b}{\operatorname{arg\,min}} \|w\|^2 \text{ st } t_i (w^T x_i + b) \geq 1$$

Stretch the points from margin by increasing w .

Constrained Convex Problem: Quadratic Programming

Soft Margin for SVM for non-separable data:



Tradeoff parameter = λ Slack variable ξ_i

$$\min_{w, b, \xi} \left[\lambda \|w\|^2 + \frac{1}{N} \sum_{i=1}^N \xi_i \right]$$

$$\text{Such that } \forall i \quad t_i (w^T x_i + b) \geq 1 - \xi_i \quad (\xi_i \geq 0)$$

$$\min_{w, b} \left[\lambda \|w\|^2 + \frac{1}{N} \sum_{i=1}^N [1 - t_i (w^T x_i + b)] \right] \quad (\text{ReLU})$$

$$= \max(0, a) \quad \xrightarrow{\text{Convex fn}}$$

SVM - Also keep loss function

Such that margin is one & shrink w st margin maximized

(Robustness \Leftrightarrow Maximizing the margin.)

Different loss - in different properties

Perceptron loss fn: Not convex relaxation of step fn.
(Not always step function)

6.1.) Kernelized SVMs:-

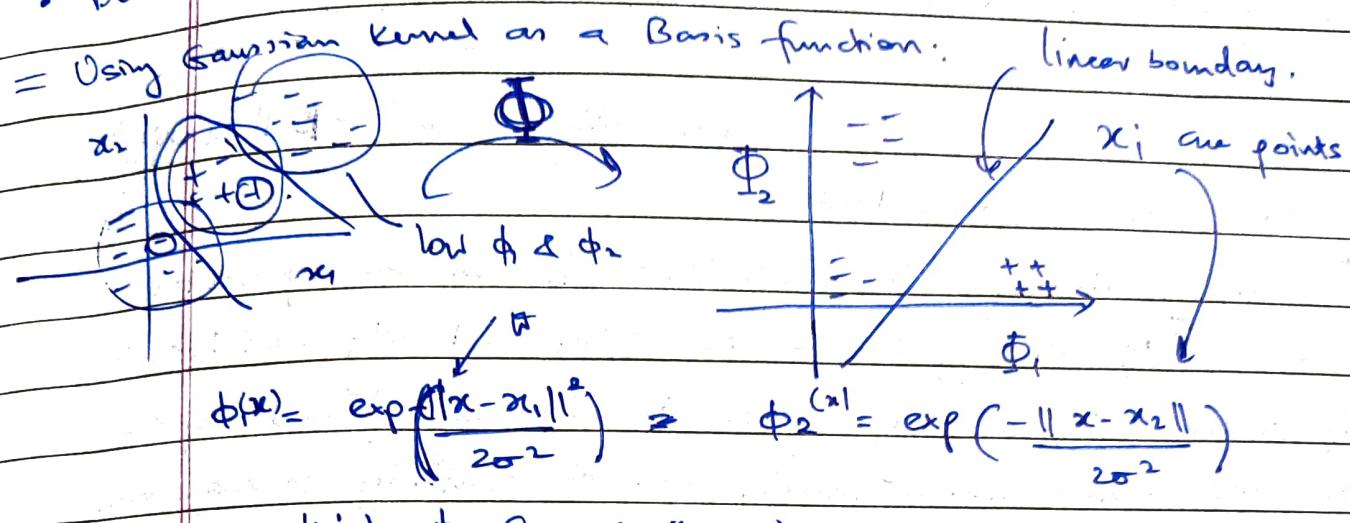
- Derive dual form of SVM
- Derive SV regression objective

Show raw data replaced with kernel

PAGE NO.

DATE

///



2 data points :) x where we want to compute feature
 x_i is the referral data points

Why use a Kernel?

- Using right kernel $\Phi(\cdot)$ can increase the dimension from x to that of $\Phi(x)$
- More importantly, transformed data can be linearly separable.

Eg.

$$\Phi(x) = \begin{bmatrix} x_1 \\ x_2 \\ \sqrt{x_1^2 + x_2^2} \end{bmatrix}$$

We may even need to transform into ∞ dimensions to have a linearly separable boundary.

$$K(x_i, x_j) = K(x_j, x_i) \quad [\text{Symmetric}]$$

High value when $x_i = x_j$ and low when x_i & x_j are not alike at all.

Also, we want K to form a positive semi definite.

(Non negative eigen-values.) Reproducing Kernel or Mercer Kernel in Hilbert space

So, we have: (Mercel Kernel)

$$K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$

PAGE No.

DATE

/ /

Example: $K(x_i, x_j) = \exp\left(\frac{-\|x_i - x_j\|^2}{2\sigma^2}\right)$

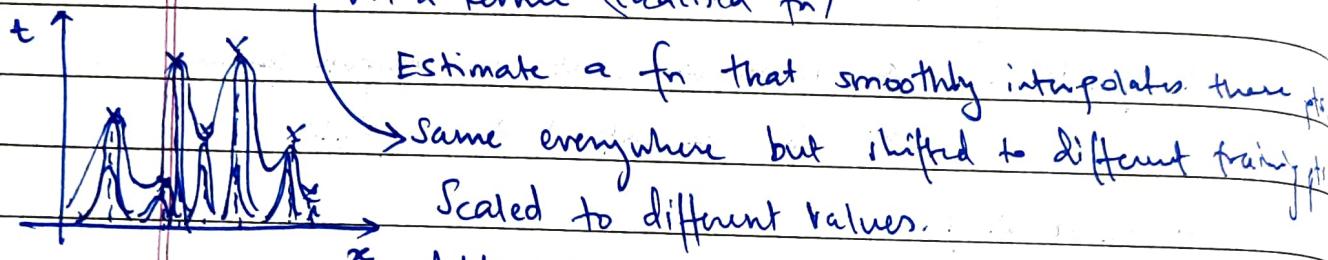
or equivalently $\exp(-K\|x_i - x_j\|^2)$

This gives a Gram Matrix.

* Change equation of SVM such that it can be represented as $x_i^T x_j$ as opposed to x [We want only DOT PRODUCT]

If we have this, we may replace x with any ϕ and replace the dot product with a kernel. Lot of flexibility to design non linear SVM.

Fit a kernel (localized fn)



Adding the linearly scaled functions we get a smooth function

$$f(x) = \sum_{i=1}^n a_i K(x, x_i) \quad [K \text{ is the gaussian kernel}]$$

If regularly spaced we may have a triangle kernel etc

Different kernels give different smoothness effects.

We use a sparse set of points \Rightarrow take out all non SV points.

Change form, introduce dot products and replace kernel with dot product

Change primal form into Dual form. So that it has $x_i^T x_j$

Then we replace dot product with a ϕ to get output of Kernel to be a scalar

$$K(x_i, x_j) = \phi(x_i)^T \phi(x_j) \quad \text{Has the semi definite gram matrix.}$$

6.2.1 Original SVM into dual form:-

$$L(w, b) = \frac{1}{2} \|w\|^2 - \sum_i a_i [t_i (w^T x_i + b)]$$

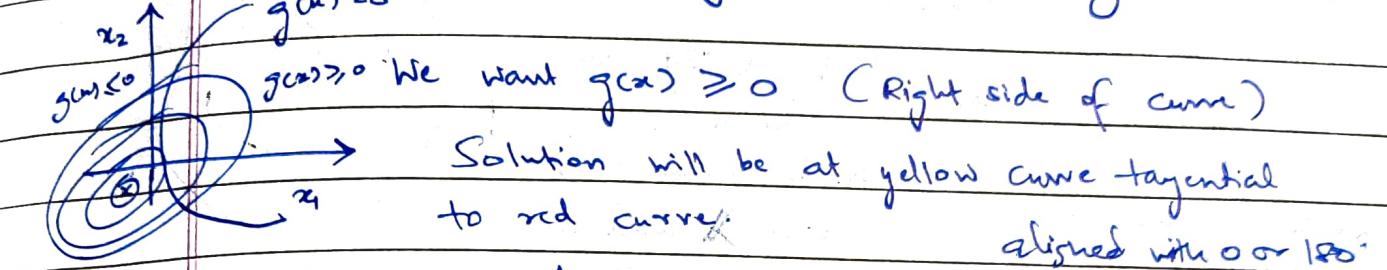
PAGE NO. _____
DATE _____

\uparrow shrink weights \uparrow lagrange multipliers \uparrow Margin.

[lagrange form of constraint opt over w] Subject to constraint $a_i \geq 0 \forall i$

$$\min L(\alpha, \lambda) = f(x) + \lambda g(x); \lambda \geq 0$$

① This function encodes minimizing $f(x)$ subject to $g(x) \geq 0$ ②



$$\text{Minimize and set } \nabla = 0 \quad \nabla f + \lambda \nabla g = 0$$

Thus lagrange multiplier gives constrained optimized solution

We want to replace x to convert into dual form replace with $x^T x$

$$\min L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_i a_i [t_i (w^T x_i + b) - 1] \quad a_i \geq 0$$

$$\text{Set derivatives } \frac{\partial L}{\partial w} = 0 \Rightarrow w \text{ value} = \sum_i a_i t_i x_i$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow 0 = \sum a_i t_i$$

Plug these in to eliminate w & b

$$L(\alpha) = \sum_i a_i - \frac{1}{2} \sum_i \sum_j a_i t_i t_j x_i^T x_j \quad K(x_i, x_j)$$

This on minimization gives values of α

They y Put values of α

$$y = \sum_i a_i t_i x_i^T x_i + b$$

Test training (If $\alpha_i = 0$ then more about support vectors)

$$k(x, x_i)$$

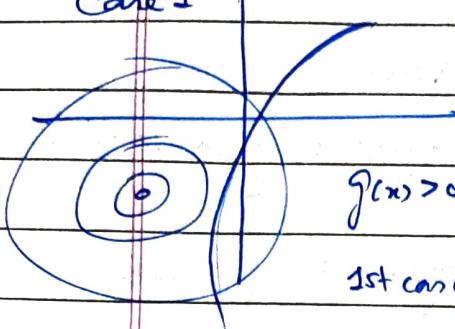
6.2.2 Solving the dual Lagrangian $L(a)$

PAGE No.	/ /
DATE	/ /

$$\min L(x) = f(x) + \lambda g(x)$$

$$\rightarrow g(x) \geq 0 \quad \lambda \geq 0 \quad \text{and} \quad \lambda g(x) = 0 \quad [\text{either } g(x) = 0 \text{ or } \lambda = 0]$$

Case 1



Ideal solution

will be tangential

$$g(x) > 0 \quad \text{to} \quad g(x) = 0$$

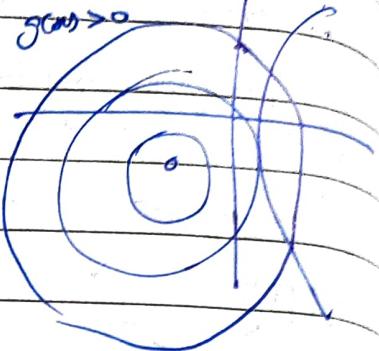
1st case ∇f & ∇g are aligned.

In this case $g(x)$

$$= 0 \quad \text{is used to}$$

cancel out.

Case 2



$g(x) \neq 0$ does not matter

$$\Rightarrow \lambda = 0$$

Case 1

Case 2

$$g(x) = 0$$

$$\lambda = 0$$

$$\lambda g(x) = 0 \quad \lambda g(x) = 0$$

$$\text{Then, } L(a) = \sum_i a_i - \frac{1}{2} \sum_i \sum_j a_i a_j t_i t_j x_i^T x_j$$

Using the observations from earlier, we have:

$$a_i \geq 0 \quad \text{and} \quad t_i y(x_i) - 1 \geq 0$$

$a_i [t_i y(x_i) - 1] = 0$ are the conditions

Using this, what we get for b:

$$t_i \left[\sum_j a_j t_j x_i^T x_j + b \right] = 1$$

$$\therefore b = \frac{1}{N} \sum_{i \in S} [t_i - \sum_j a_j t_j x_i^T x_j]$$

KKT conditions (Karush-Kuhn-Tucker condition)

- Points where $a_i = 0 \notin S$ (Not SV)

$t_i (y(x_i) - 1)$ need not be 0

Do not define decision boundary - can be moved around without changing the decision boundary.

- If $a_n > 0$ then $t_n(y(x_n) - 1) = 0$ (They are SV)
- They define the decision boundary
- These points need to be retained on for classification of test data

PAGE NO.	
DATE	/ /

Visualizing Gaussian Kernel:

$$x_i^T x_j = \exp(-\kappa \|x_i - x_j\|^2)$$

- Support vectors

0. Gaussian Kernel from these points

+ + + + Decision boundary linear in ϕ space

Other useful Kernels: $-(x_i^T x_j + c)^d$ - Polynomial

Cosine = $x_i^T x_j / \|x_i\| \|x_j\|$ and linear kernels etc.

6.3.) What if the classes are not separable?

(Allow some missclass^n)

- An error function that gives an error on misclassifn needs to be modified
- Introduce slack variable $\xi_i \geq 0$ for each point n st:
 - $\xi_i = 0$ if the point is on correct side
 - $|t_n - y(x_n)|$ for points outside the margin

For this formulation :-

$$L(w, b, a, y) = \frac{1}{2} \|w\|^2 - \sum_i a_i [t_i y(x_i) - 1 + \xi_i] + C \sum_i \xi_i - \sum_i y_i \xi_i$$

Now we will have 6 KKT condns:

- | | | |
|----------------|-----------------------------------|--------------------------------------|
| • $a_i \geq 0$ | • $t_i y(x_i) - 1 + \xi_i \geq 0$ | • $a_i (t_i y(x_i) - 1 + \xi_i) = 0$ |
| • $y_i \geq 0$ | • $\xi_i \geq 0$ | • $y_i \xi_i = 0$ |

Now we eliminate w and have everything term of x_i & x_j dots

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_i a_i t_i x_i \quad \frac{\partial L}{\partial b} = 0 \Rightarrow \sum_i a_i t_i = 0$$

$$\frac{\partial L}{\partial \xi_i} = 0 \Rightarrow a_i = C - y_i \quad (\text{Eliminate } \xi, w \text{ & } b)$$

Consequently get dual form of soft SVM

$$\hat{L}(a) = \sum_i a_i - \frac{1}{2} \sum_i \sum_j a_i a_j b_i b_j \frac{x_i^T x_j}{c}$$

PAGE NO.	
DATE	/ /

Condition: $0 \leq a_i \leq C$ & $\sum a_i b_i = 0$

Boxed into a and C . Therefore 3 types of a_i .

- 1) $a_i = 0$ Not support vectors, don't matter
- 2) $\hookrightarrow a_i \leq 1$ correct points but inside margin
- 2.) $a_i = C$; $\epsilon_i \geq 1 \Rightarrow$ Incorrectly classified points

Analyzing further & Support Vector Regression: \rightarrow Hyperparameter

\therefore We finally minimize $\hat{L}(a)$ where $(C > 0)$ controls tradeoff between slack variable penalty & the margin.

\rightarrow Note that this is sensitive to outliers as the penalty is linearly proportional to distance from the margin.

- Large $C \rightarrow$ More emphasis on reducing ϵ_i , few misclassified
- ~~Small $C \rightarrow$~~ More SV \Rightarrow Smoother boundary. (GWS)

6.4) SVM for Regression:

- Box Constraints: Four types of points:
-
- $a_n > 0$ ($\epsilon_n > 0$) $\bullet a_n = 0$ (Inside margin & correctly class)
 - $a_n < C$ $\bullet a_n < C, \epsilon_n = 0$ (On the margin)
 - $a_n = C$, $\epsilon_n \leq 1$ ($\epsilon_n < 1$) (Across margin, correctly class)
 - $a_n = C$, $\epsilon_n > 1$ ($\epsilon_n > 1$) (Across margin, misclassified)

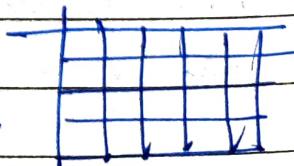
Influence of point across margin is capped by C . [Capped]

Hyperparameters:

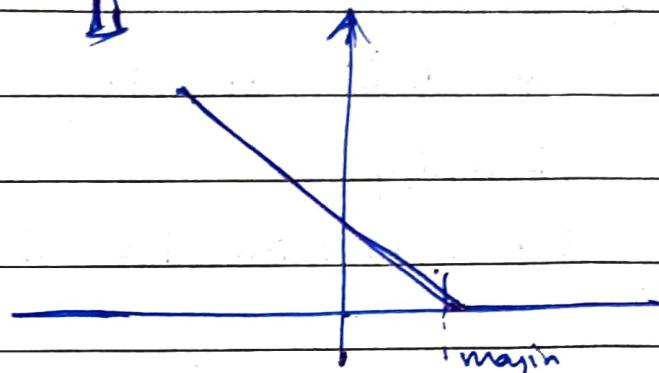
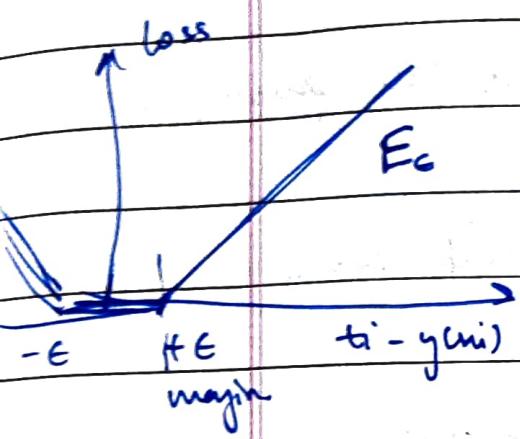
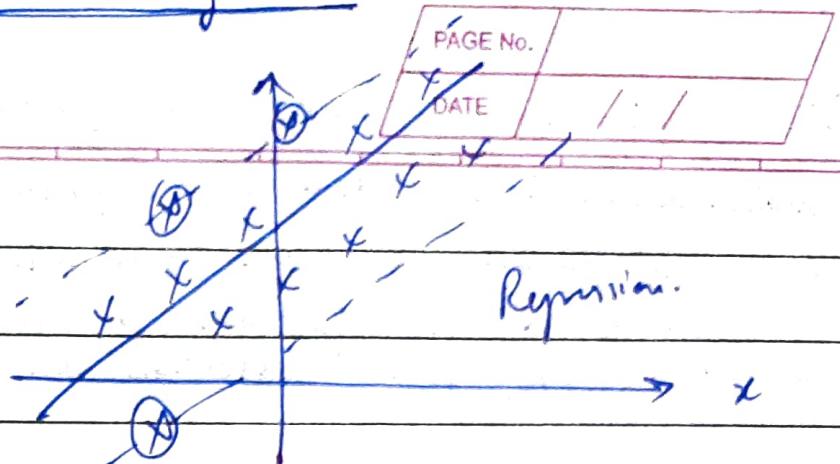
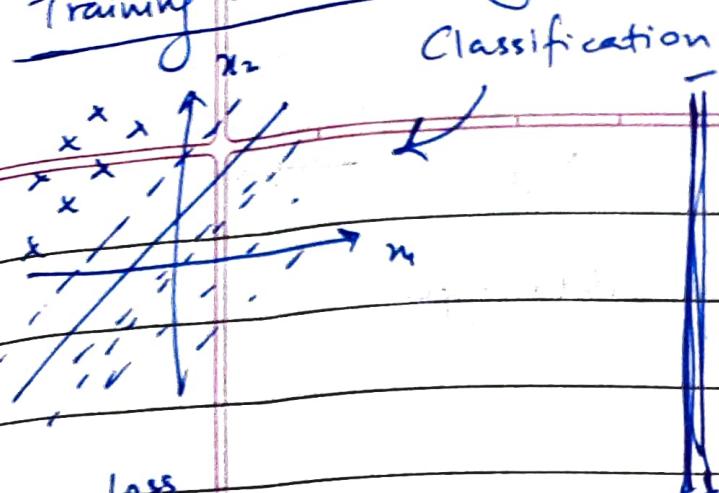
C , & Kernel specific: K , Φ matrix

And then we choose the best K

Combination of C, K



Training data sparsity in regression using SVM:



Angular linear reg.

$$\text{Mean Square Error} = \frac{1}{2} \sum_i (y_i - t_i)^2 + \frac{\lambda}{2} \|w\|^2$$

$$\text{SVR} = C \sum_i E_\xi (y_i - t_i) + \frac{1}{2} \|w\|^2$$

$$E_\xi = \max(0, |y_i - t_i| - \epsilon)$$

Tors slack variables $t_i \leq y_i + \sum \xi_i \quad | \quad t_i \geq y_i - \sum \bar{\xi}_i$

$$C \sum_i (\xi_i + \bar{\xi}_i) + \frac{1}{2} \|w\|^2$$