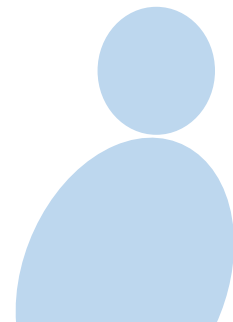


Elective Module:

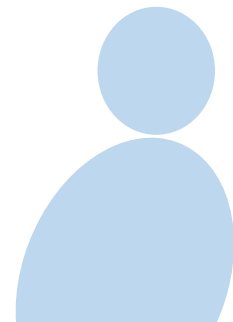
**Advanced ML
Techniques**



Chapter 11

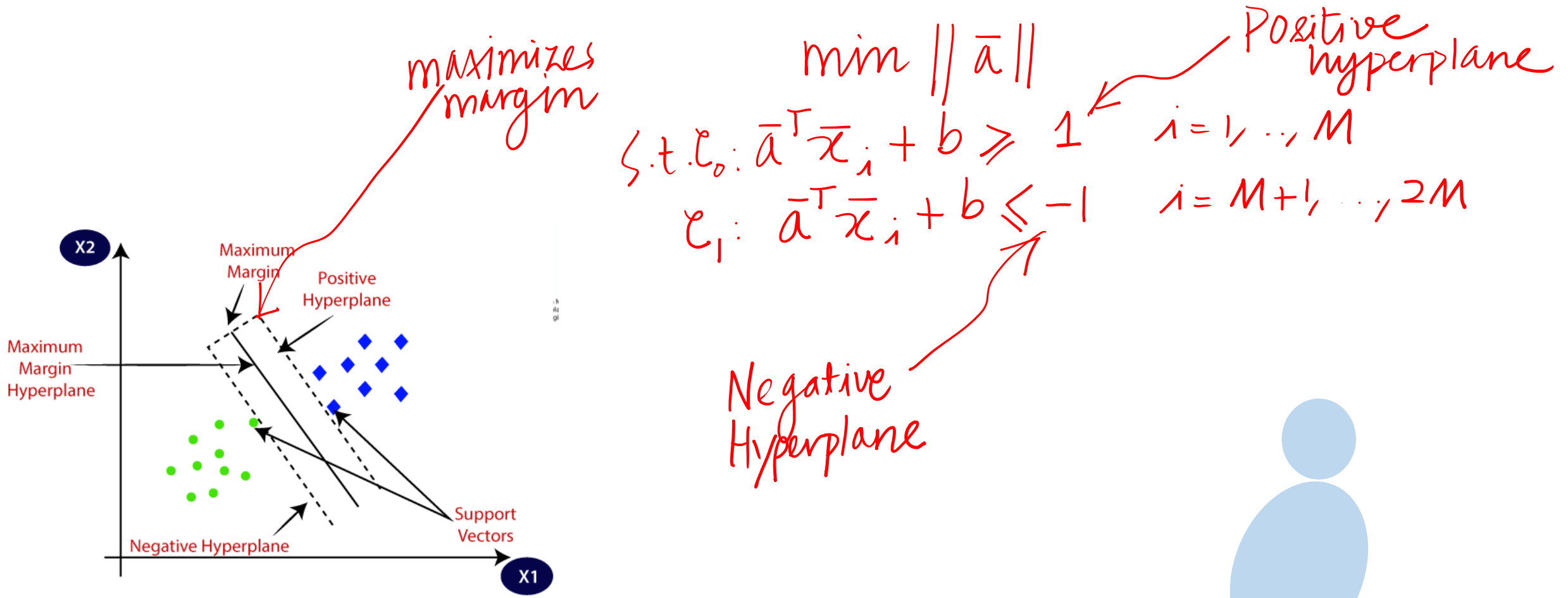
Support
Vector Machine

Dual SVM and Kernel SVM



SVM Classifier

- Recall, the problem to determine classifier with maximum margin is



SVM Classifier

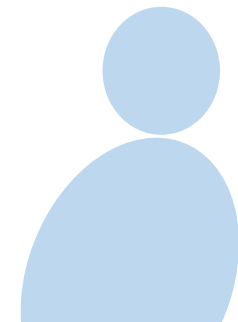
- Recall, the problem to determine classifier with maximum margin is

$$\min \|\bar{\mathbf{a}}\|_2$$

convex
optimization

$$\mathcal{C}_0: \bar{\mathbf{a}}^T \bar{\mathbf{x}}_i + b \geq 1, 1 \leq i \leq M$$

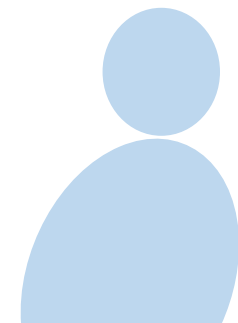
$$\mathcal{C}_1: \bar{\mathbf{a}}^T \bar{\mathbf{x}}_i + b \leq -1, M + 1 \leq i \leq 2M$$



SVM Response

- Let the response be defined as

$$\begin{aligned} \mathcal{C}_0 : y_i &= 1, \quad i = 1, \dots, M \\ \mathcal{C}_1 : y_i &= -1, \quad i = M+1, \dots, 2M \end{aligned}$$

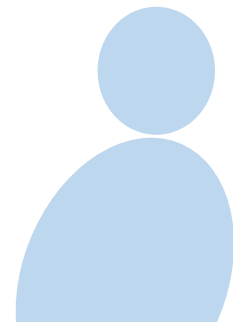


SVM Response

- Let the response be defined as

$$\mathcal{C}_0: y_i = 1, 1 \leq i \leq M$$

$$\mathcal{C}_1: y_i = -1, M + 1 \leq i \leq 2M$$



SVM Constraints

- The constraints can be expressed as

$$C_0: 1(\bar{a}^T \bar{x}_i + b) \geq 1$$

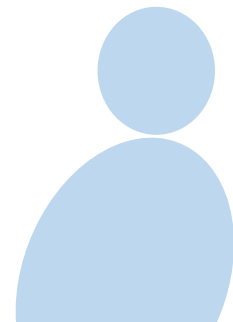
$$\Rightarrow y_i(\bar{a}^T \bar{x}_i + b) \geq 1$$

$$C_1: \bar{a}^T \bar{x}_i + b \leq -1$$

$$\Rightarrow -1(\bar{a}^T \bar{x}_i + b) \geq 1$$

$$\Rightarrow y_i(\bar{a}^T \bar{x}_i + b) \geq 1$$

constraints become
same for both
classes!

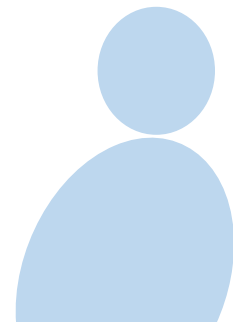


SVM Constraints

- The constraints can be expressed as

$$\mathcal{C}_0: y_i(\bar{\mathbf{a}}^T \bar{\mathbf{x}}_i + b) \geq 1, 1 \leq i \leq M$$

$$\mathcal{C}_1: y_i(\bar{\mathbf{a}}^T \bar{\mathbf{x}}_i + b) \geq 1, M + 1 \leq i \leq 2M$$



SVM Constraints

- The constraints can be combined as

$i = 1, 2, \dots, 2M$. For all i

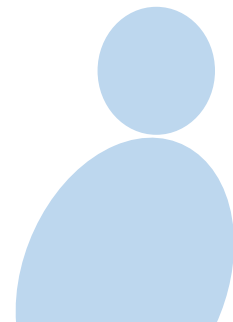
$$y_i (\bar{a}^T \bar{x}_i + b_i) \geq 1$$

$$\Rightarrow -y_i (\bar{a}^T \bar{x}_i + b_i) \leq -1$$

$$\Rightarrow -\left(y_i (\bar{a}^T \bar{x}_i + b_i) - 1\right) \leq 0$$

Valid for all i

Final version
of the
constraint



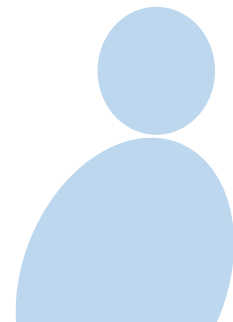
SVM Constraints

- The constraints can be combined as

$$y_i(\bar{\mathbf{a}}^T \bar{\mathbf{x}}_i + b) \geq 1, 1 \leq i \leq 2M$$

$$\Rightarrow -y_i(\bar{\mathbf{a}}^T \bar{\mathbf{x}}_i + b) \leq -1$$

$$\Rightarrow -(y_i(\bar{\mathbf{a}}^T \bar{\mathbf{x}}_i + b) - 1) \leq 0$$



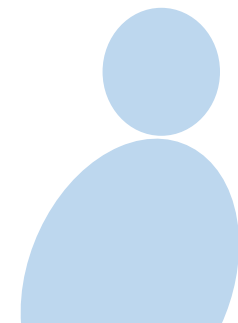
SVM Classifier

- The SVM classifier problem can be recast as

$$\begin{aligned} \min \|\bar{a}\| &\equiv \min \frac{1}{2} \|\bar{a}\|^2 \\ \text{s.t. } &-(y_i(\bar{a}^T \bar{x}_i + b) - 1) \leq 0 \\ &\text{for all } i \end{aligned}$$

CONVEX

Equivalent
SVM optimization
problem



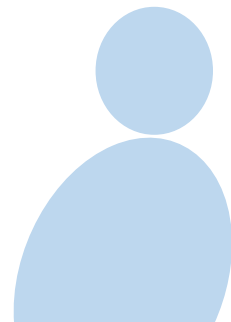
SVM Classifier

- The SVM classifier problem can be recast as

$$\min \frac{1}{2} \|\bar{\mathbf{a}}\|^2$$

subject to

$$-(y_i(\bar{\mathbf{a}}^T \bar{\mathbf{x}}_i + b) - 1) \leq 0, 1 \leq i \leq 2M$$



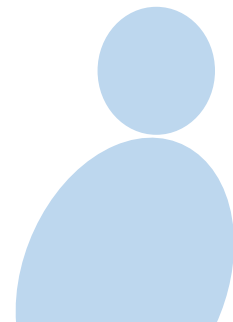
Lagrangian

- The Lagrangian for this problem is

constrained
optimization problem

$$\frac{1}{2} \|\bar{a}\|^2 + \sum_{i=1}^{2M} \lambda_i \left(-\left(y_i (\bar{a}^T \bar{x}_i + b) - 1 \right) \right) \quad \left. \vphantom{\sum_{i=1}^{2M}} \right\} \text{Lagrangian Function}$$

$$= \frac{1}{2} \|\bar{a}\|^2 - \sum_{i=1}^{2M} \lambda_i \left(y_i (\bar{a}^T \bar{x}_i + b) - 1 \right)$$



Lagrangian

- The **Lagrangian** for this problem is

$$\frac{1}{2} \|\bar{\mathbf{a}}\|^2 - \sum_{i=1}^{2M} \lambda_i (y_i (\bar{\mathbf{a}}^T \bar{\mathbf{x}}_i + b) - 1)$$

Handwritten notes in red:

- $\frac{1}{2} \cdot \bar{\mathbf{a}}^T \bar{\mathbf{a}} \quad \nabla \left(\frac{1}{2} \cdot \bar{\mathbf{a}}^T \bar{\mathbf{a}} \right) = \bar{\mathbf{a}}$
- $\lambda_i y_i \bar{\mathbf{a}}^T \bar{\mathbf{x}}_i \xrightarrow{\nabla} \lambda_i y_i \bar{\mathbf{x}}_i$
- $\nabla = 0$

Lagrangian

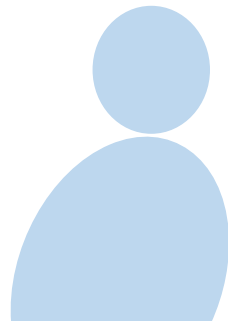
KKT
Karush Kuhn Tucker

- Setting gradient wrto $\bar{\mathbf{a}}$ to zero

$$\nabla_{\bar{\mathbf{a}}} \left(\frac{1}{2} \bar{\mathbf{a}}^T \bar{\mathbf{a}} - \sum_i \lambda_i (y_i (\bar{\mathbf{a}}^T \bar{\mathbf{x}}_i + b) - 1) \right) = 0$$

$$\Rightarrow \bar{\mathbf{a}} - \sum_i \lambda_i y_i \bar{\mathbf{x}}_i = 0$$

$$\Rightarrow \bar{\mathbf{a}} = \sum_i \lambda_i y_i \bar{\mathbf{x}}_i$$

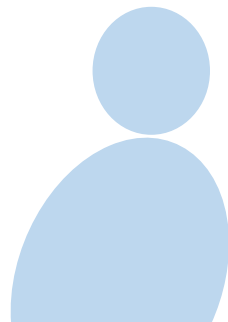


Lagrangian

- Setting gradient wrto $\bar{\mathbf{a}}$ to zero

$$\nabla_{\bar{\mathbf{a}}} \left(\frac{1}{2} \|\bar{\mathbf{a}}\|^2 - \sum_{i=1}^{2M} \lambda_i (y_i (\bar{\mathbf{a}}^T \bar{\mathbf{x}}_i + b) - 1) \right) = 0$$

$$\Rightarrow \bar{\mathbf{a}} - \sum_{i=1}^{2M} \lambda_i y_i \bar{\mathbf{x}}_i = 0 \Rightarrow \bar{\mathbf{a}} = \sum_{i=1}^{2M} \lambda_i y_i \bar{\mathbf{x}}_i$$



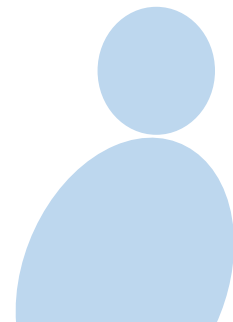
Lagrangian

$$\bar{\mathbf{a}} = \sum_{i=1}^{2M} \lambda_i y_i \bar{\mathbf{x}}_i$$

Handwritten annotations in red:

- A bracket above the summation index $i=1$ to $2M$ is labeled α_i .
- The term $\lambda_i y_i$ is circled.
- The term $\bar{\mathbf{x}}_i$ is circled.
- An arrow points from the circled $\bar{\mathbf{x}}_i$ to the expression $\sum_i \alpha_i \bar{x}_i$.

- Thus, $\bar{\mathbf{a}}$ can be expressed as linear combination of $\bar{\mathbf{x}}_i$.



Lagrangian

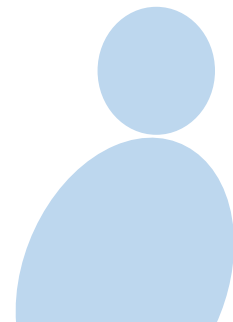
$$\bar{\mathbf{a}} = \sum_{i=1}^{2M} \lambda_i y_i \bar{\mathbf{x}}_i$$

$$\alpha_i \neq 0$$

- The points for which $\lambda_i \neq 0$ are termed the **support vectors**.

Support vectors.

Linear combination
of support vectors
since $\alpha_i \neq 0$.



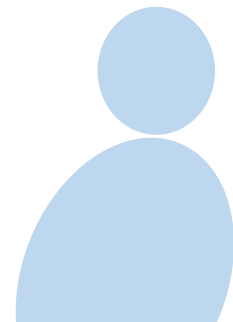
Lagrangian

- Setting gradient wrto b to zero

$$\nabla_b \left(\underbrace{\frac{1}{2} \tilde{a}^T \tilde{a}} - \sum_i \lambda_i \left(y_i (\tilde{a}^T \tilde{x}_i + b) - 1 \right) \right) = 0$$
$$\sum_i \lambda_i y_i \tilde{a}^T \tilde{x}_i + \underbrace{\lambda_i y_i b}_{\nabla = \lambda_i y_i}$$

$\nabla = 0$

$$\Rightarrow \sum_i \lambda_i y_i = 0$$

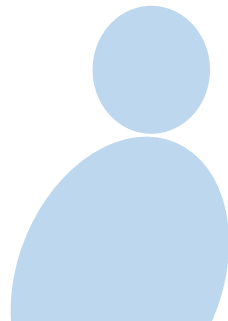


Lagrangian

- Setting gradient wrto b to zero

$$\nabla_b \left(\frac{1}{2} \|\bar{\mathbf{a}}\|^2 - \sum_{i=1}^{2M} \lambda_i (y_i (\bar{\mathbf{a}}^T \bar{\mathbf{x}}_i + b) - 1) \right) = 0$$

$$\Rightarrow \sum_{i=1}^{2M} \lambda_i y_i = 0$$



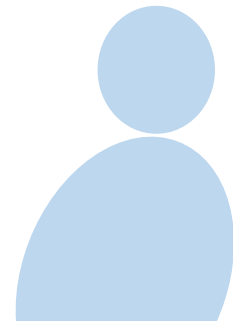
Lagrangian

$$\sum_i \lambda_i y_i \bar{x}_i$$

- The expression for \bar{a} can be substituted in the Lagrangian
- The resulting expression can be simplified as shown next

Original problem: Primal problem \rightarrow min

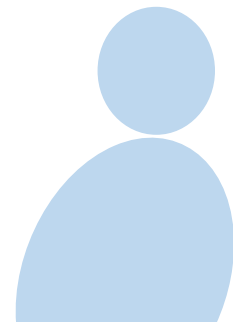
\rightarrow Dual objective: maximization



Lagrangian

- Consider the **Lagrangian** given as

$$\frac{1}{2} \bar{a}^T \bar{a} - \sum_i \lambda_i \left(y_i (\bar{a}^T \bar{x}_i + b) - 1 \right)$$



Lagrangian

- Consider the **Lagrangian** given as

$$\begin{aligned} & \frac{1}{2} \|\bar{\mathbf{a}}\|^2 - \sum_{i=1}^{2M} \lambda_i (y_i (\bar{\mathbf{a}}^T \bar{\mathbf{x}}_i + b) - 1) \\ &= \frac{1}{2} \bar{\mathbf{a}}^T \bar{\mathbf{a}} - \sum_{i=1}^{2M} \lambda_i (y_i (\bar{\mathbf{a}}^T \bar{\mathbf{x}}_i + b) - 1) \end{aligned}$$

$\sum_i \lambda_i y_i \bar{\mathbf{x}}_i$

Lagrangian

- Substitute $\bar{\mathbf{a}} = \sum_{i=1}^{2M} \lambda_i y_i \bar{\mathbf{x}}_i$

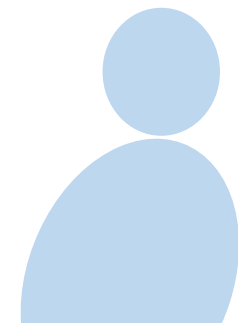
Lagrangian

$$= \frac{1}{2} \left(\sum_i \lambda_i y_i \bar{\mathbf{x}}_i \right)^T \left(\sum_j \lambda_j y_j \bar{\mathbf{x}}_j \right) - \sum_i \lambda_i \left(y_i \left(\left(\sum_j \lambda_j y_j \bar{\mathbf{x}}_j \right)^T \bar{\mathbf{x}}_i + b \right) - 1 \right)$$

Messy!

Lagrangian

- Substitute $\bar{\mathbf{a}} = \sum_{i=1}^{2M} \lambda_i y_i \bar{\mathbf{x}}_i$
$$= \frac{1}{2} \left(\sum_{i=1}^{2M} \lambda_i y_i \bar{\mathbf{x}}_i \right)^T \left(\sum_{i=1}^{2M} \lambda_j y_j \bar{\mathbf{x}}_j \right)$$
$$- \sum_{i=1}^{2M} \lambda_i \left(y_i \left(\left(\sum_{j=1}^{2M} \lambda_j y_j \bar{\mathbf{x}}_j \right)^T \bar{\mathbf{x}}_i + b \right) - 1 \right)$$

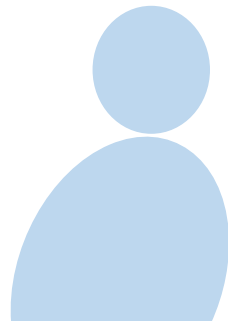


Lagrangian

$$= \sum_i \lambda_i - \frac{1}{2} \sum_{i,j=1}^{2M} \lambda_i \lambda_j y_i y_j \bar{x}_i^T \bar{x}_j - \underbrace{b \sum_i \lambda_i y_i}_0$$

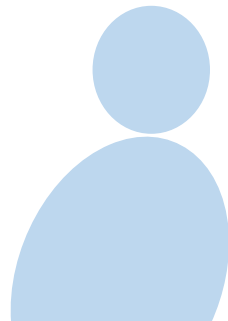
$$= \sum_i \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \bar{x}_i^T \bar{x}_j \} \max.$$

Dual Objective:



Lagrangian

$$\begin{aligned} &= \sum_{i=1}^{2M} \lambda_i - \frac{1}{2} \sum_{i,j=1}^{2M} \lambda_i \lambda_j y_i y_j \bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_j - b \sum_{i=1}^{2M} \lambda_i y_i \\ &= \sum_{i=1}^{2M} \lambda_i - \frac{1}{2} \sum_{i,j=1}^{2M} \lambda_i \lambda_j y_i y_j \bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_j \end{aligned}$$



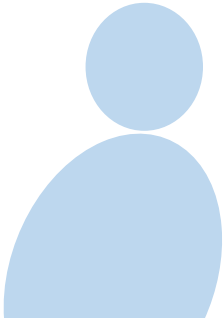
Dual SVM

- Therefore, the **dual problem** can be formulated as

Dual Problem

$$\begin{aligned} \max \cdot & \sum_i \lambda_i - \sum_i \sum_j \lambda_i \lambda_j y_i y_j \bar{x}_i^T \bar{x}_j \\ \text{s.t.} \quad & \lambda_i \geq 0 : \text{Lagrange multipliers} \geq 0 \\ & \sum_i \lambda_i y_i = 0 \end{aligned}$$

non-negative



Dual SVM

- Therefore, the **dual problem** can be formulated as

$$\max \sum_{i=1}^{2M} \lambda_i - \frac{1}{2} \sum_{i,j=1}^{2M} \lambda_i \lambda_j y_i y_j \bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_j$$

subject to $\lambda_i \geq 0$

$$\sum_{i=1}^{2M} \lambda_i y_i = 0$$

DUAL
PROBLEM
for SVM.

Dual Objective

Variables:

λ_i

Lagrange multipliers

Primal convex

Same solution!

Dual concave

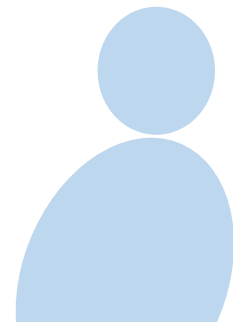
Duality Gap = 0

Dual SVM

- How to calculate b ?
- For any point for which $\lambda_i \neq 0$

$$y_i (\bar{a}^T \bar{x}_i + b) = 1 \quad \left. \vphantom{y_i (\bar{a}^T \bar{x}_i + b) = 1} \right\} \begin{array}{l} \text{complementary} \\ \text{slackness} \end{array}$$

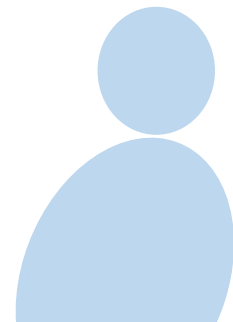
Solving this one can
determine b .



Dual SVM

- How to calculate b ?
- For any point for which $\lambda_i \neq 0$

$$y_i(\bar{\mathbf{a}}^T \bar{\mathbf{x}}_i + b) - 1 = 0$$

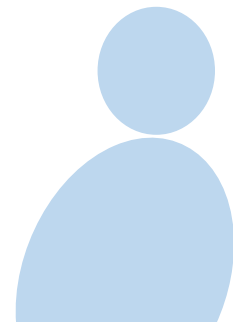


Dual SVM

- Note that the quantity $\bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_j$ denotes the inner product.
- This can be represented as

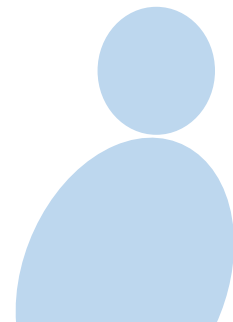
$$\langle \bar{\mathbf{x}}_i, \bar{\mathbf{x}}_j \rangle$$

inner product.



Dual SVM

- Note that the quantity $\bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_j$ denotes the inner product.
- This can be represented as $\langle \bar{\mathbf{x}}_i, \bar{\mathbf{x}}_j \rangle$



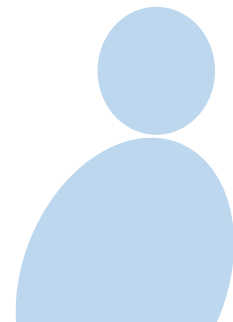
Dual SVM

- Using this notation, the **dual SVM** problem can be defined as

Depends only on inner products.

$$\max. \sum_i \lambda_i - \frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \langle \bar{x}_i, \bar{x}_j \rangle$$

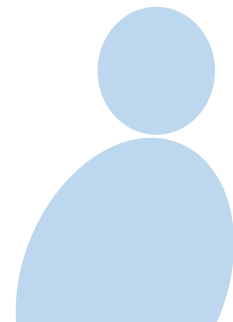
$$\text{subject to } \left. \begin{array}{l} \lambda_i \geq 0 \\ \sum_i \lambda_i y_i = 0 \end{array} \right\} \text{constraints.}$$



Dual SVM

- Using this notation, the **dual SVM** problem can be defined as

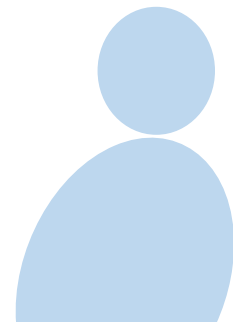
$$\begin{aligned} \max \quad & \sum_{i=1}^{2M} \lambda_i - \frac{1}{2} \sum_{i,j=1}^{2M} \lambda_i \lambda_j y_i y_j \langle \bar{\mathbf{x}}_i, \bar{\mathbf{x}}_j \rangle \\ \text{subject to } & \lambda_i \geq 0 \\ & \sum_{i=1}^{2M} \lambda_i y_i = 0 \end{aligned}$$



Dual SVM

- For any new point $\bar{\mathbf{x}}$, y can be calculated as

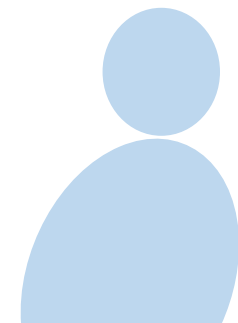
$$\begin{aligned}\bar{\mathbf{a}} &= \sum_i \lambda_i y_i \bar{\mathbf{x}}_i \\ \bar{\mathbf{a}}^T \bar{\mathbf{x}} + b &= \left(\sum_i \lambda_i y_i \bar{\mathbf{x}}_i \right)^T \bar{\mathbf{x}} + b \\ &= \sum_i \lambda_i y_i \underbrace{\langle \bar{\mathbf{x}}_i, \bar{\mathbf{x}} \rangle}_{\text{inner product}} + b\end{aligned}$$



Dual SVM

- For any new point $\bar{\mathbf{x}}$, y can be calculated as

$$\begin{aligned}\bar{\mathbf{a}}^T \bar{\mathbf{x}}_i + b &= \left(\sum_{i=1}^{2M} \lambda_i y_i \bar{\mathbf{x}}_i \right)^T \bar{\mathbf{x}} + b \\ &= \sum_{i=1}^{2M} \lambda_i y_i \langle \bar{\mathbf{x}}_i, \bar{\mathbf{x}} \rangle + b\end{aligned}$$



Kernel

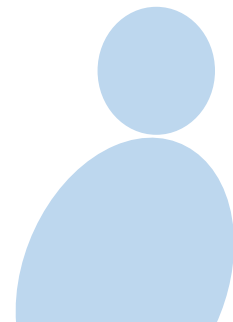
- One can now replace $\langle \bar{\mathbf{x}}_i, \bar{\mathbf{x}}_j \rangle$ by a kernel

$$K(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}_j) = \phi(\bar{\mathbf{x}}_i)^T \phi(\bar{\mathbf{x}}_j)$$

$\bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_j$: Linear


Feature mapping
Non-Linear

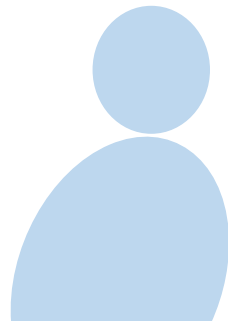
Kernel. Kernel SVM.



Kernel

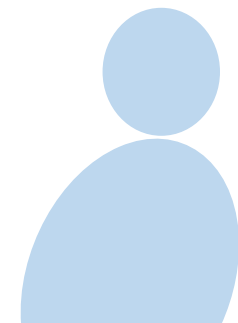
- One can now replace $\langle \bar{\mathbf{x}}_i, \bar{\mathbf{x}}_j \rangle$ by a kernel

$$K(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}_j) = \phi^T(\bar{\mathbf{x}}_i)\phi(\bar{\mathbf{x}}_j)$$




Feature mapping

- The quantity $\phi(\bar{\mathbf{x}}_i)$ is termed as a feature mapping.
- This can be used to model non-linear features.



Kernel SVM

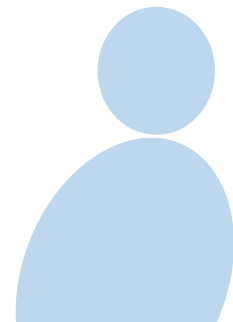

- Using this notation, the **kernel SVM** problem can be defined as

$$\min. \sum_i \lambda_i - \frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j K(\bar{x}_i, \bar{x}_j)$$

$$\text{s.t.} \quad \lambda_i \geq 0$$

$$\sum_i \lambda_i y_i = 0$$

Kernel.



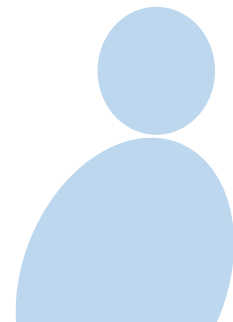
Kernel SVM

- Using this notation, the **kernel SVM** problem can be defined as *Kernel.*

$$\max \sum_{i=1}^{2M} \lambda_i - \frac{1}{2} \sum_{i,j=1}^{2M} \lambda_i \lambda_j y_i y_j K(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}_j)$$

subject to $\lambda_i \geq 0$

$$\sum_{i=1}^{2M} \lambda_i y_i = 0$$



Feature mapping

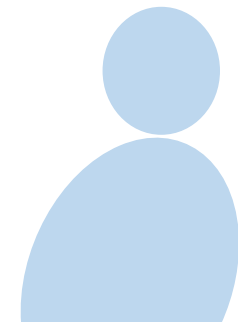
- For example, for $N = 3$,

3 dimensional Feature vectors.

$$\bar{x}_i = \begin{bmatrix} x_i(1) \\ x_i(2) \\ x_i(3) \end{bmatrix}$$

$$\phi(\bar{x}_i) = \begin{bmatrix} x_i(1) x_i(1) \\ x_i(1) x_i(2) \\ x_i(1) x_i(3) \\ x_i(2) x_i(1) \\ x_i(2) x_i(2) \\ x_i(2) x_i(3) \\ x_i(3) x_i(1) \\ x_i(3) x_i(2) \\ x_i(3) x_i(3) \end{bmatrix}$$

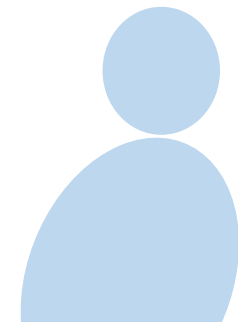
Non Linear
Feature map.



Feature mapping

- For example, for $N = 3$,

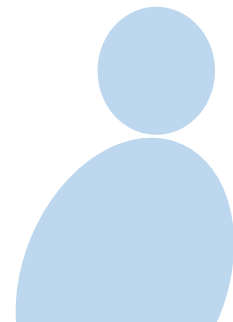
$$\phi(\bar{\mathbf{x}}_i) = \begin{bmatrix} x_i(1)x_i(1) \\ x_i(1)x_i(2) \\ x_i(1)x_i(3) \\ x_i(2)x_i(1) \\ x_i(2)x_i(2) \\ x_i(2)x_i(3) \\ x_i(3)x_i(1) \\ x_i(3)x_i(2) \\ x_i(3)x_i(3) \end{bmatrix}$$



Feature mapping

- For this kernel,

$$\begin{aligned} K(\bar{x}_i, \bar{x}_j) &= \phi^T(\bar{x}_i) \phi(\bar{x}_j) \\ &= (\bar{x}_i^T \bar{x}_j)^2 \} \text{Non-linear} \end{aligned}$$

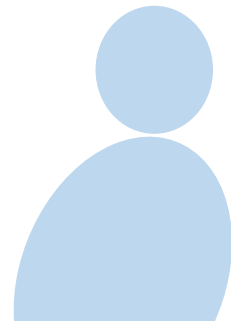


Feature mapping

- For this kernel,

$$\underbrace{K(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}_j)}_{\substack{\uparrow \\ \text{Need to know} \\ \text{only Kernel!}}} = \underbrace{\phi^T(\bar{\mathbf{x}}_i)\phi(\bar{\mathbf{x}}_j)}_{\text{circled}} = (\bar{\mathbf{x}}_i^T \bar{\mathbf{x}}_j)^2$$

Need to know
only Kernel!



Gaussian Kernel

- Another interesting kernel is the

Gaussian kernel defined as,

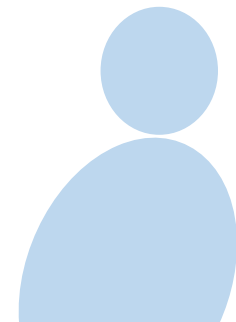
$$K(\bar{x}_i, \bar{x}_j) = \exp\left(-\frac{\|\bar{x}_i - \bar{x}_j\|^2}{2\sigma^2}\right)$$

Gaussian Kernel.

~~$\phi^T(\bar{x}_i) \phi(\bar{x}_j)$~~

Don't Need!

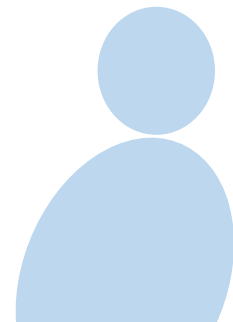
Directly
evaluate
Kernel!



Gaussian Kernel

- Another interesting kernel is the **Gaussian kernel** defined as,

$$K(\bar{\mathbf{x}}_i, \bar{\mathbf{x}}_j) = \exp\left(-\frac{\|\bar{\mathbf{x}}_i - \bar{\mathbf{x}}_j\|^2}{2\sigma^2}\right)$$

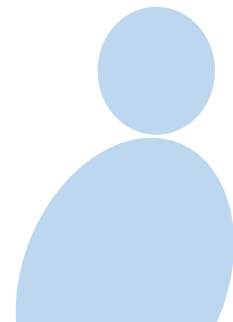


Gaussian Kernel — Good Performance!

- Example- Handwritten digit
recognition, from 16×16 images
- Gaussian kernel SVMs yield very
good performance!

Handwritten Digits
Recognition

76



Instructors may use this white area (14.5 cm / 25.4 cm) for the text.
Three options provided below for the font size.

Font: Avenir (Book), Size: 32, Colour: Dark Grey

Font: Avenir (Book), Size: 28, Colour: Dark Grey

Font: Avenir (Book), Size: 24, Colour: Dark Grey

Do not use the space below.

