# Support Vector Machines

## Somak Aditya        Sudeshna Sarkar
## CSE Department, IIT Kharagpur
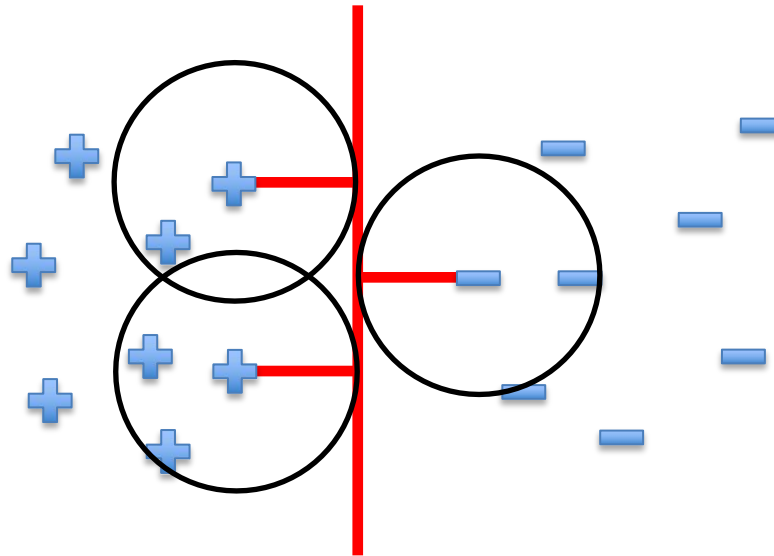## Aug 30, 2024
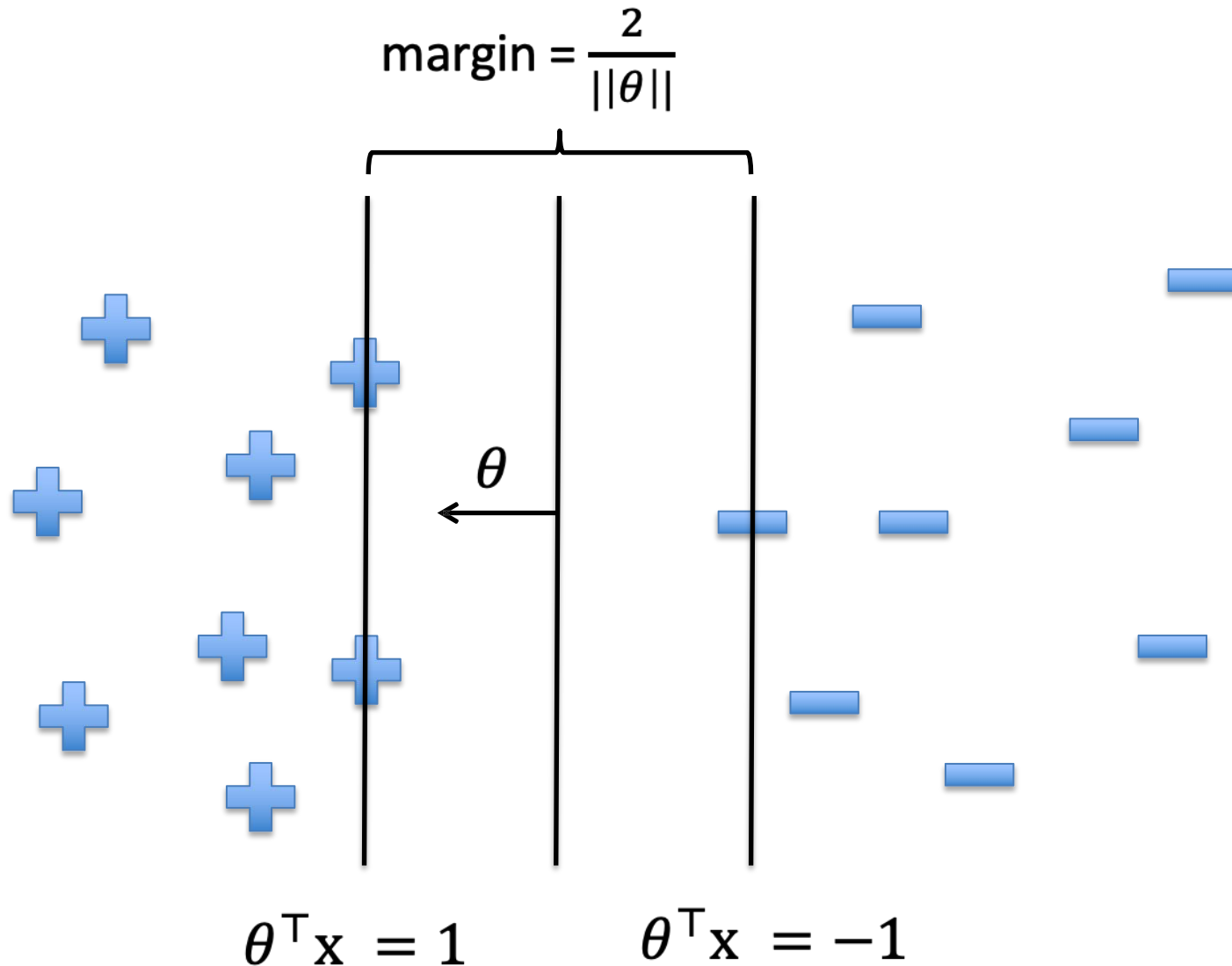
# Last Time: SVMs, Maximizing Margin

The SVM problem (assuming data is linearly separable):
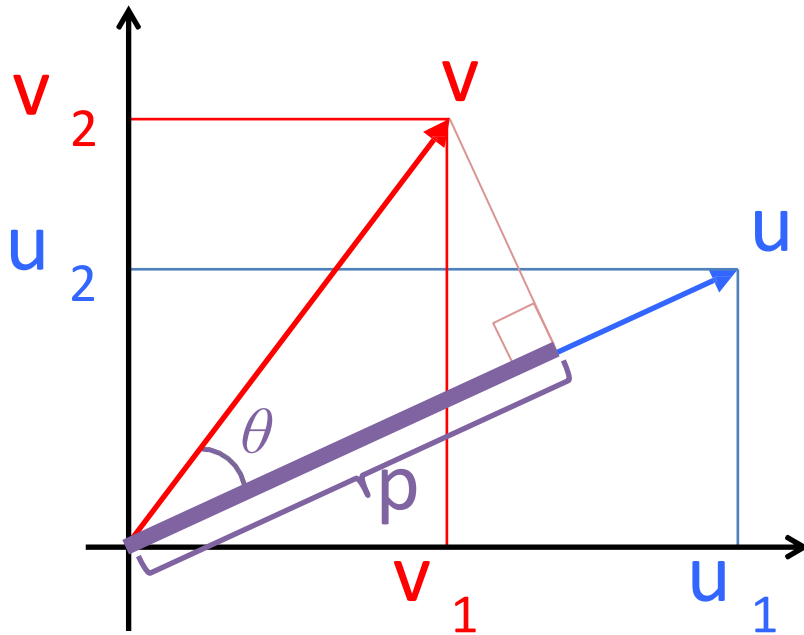
$$\min_\theta \frac{1}{2} \sum_{j=1}^{d} \theta_j^2$$

$$\text{s. t. } y_i(\theta^\top x_i) \geq 1 \forall i$$

# Maximum Margin Hyperplane

$$\text{margin} = \frac{2}{||\theta||}$$

$\theta$

$$\theta^\top x = 1 \qquad \theta^\top x = -1$$

# Vector Inner Product



$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \qquad v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$\left|\left|\mathbf{u}\right|\right|_2 = \text{length}(\mathbf{u}) \in \mathbb{R}$$

$$= \sqrt{u_1^2 + u_2^2}$$

$$\mathbf{u}^\mathsf{T}\mathbf{v} = \mathbf{v}^\mathsf{T}\mathbf{u}$$

$$= u_1 v_1 + u_2 v_2$$

$$= \left\|\mathbf{u}\right\|_2 \left\|\mathbf{v}\right\|_2 \cos\theta$$

$$= p\|\mathbf{u}\|_2 \quad \text{where } p = \|\mathbf{v}\|_2 \cos\theta$$
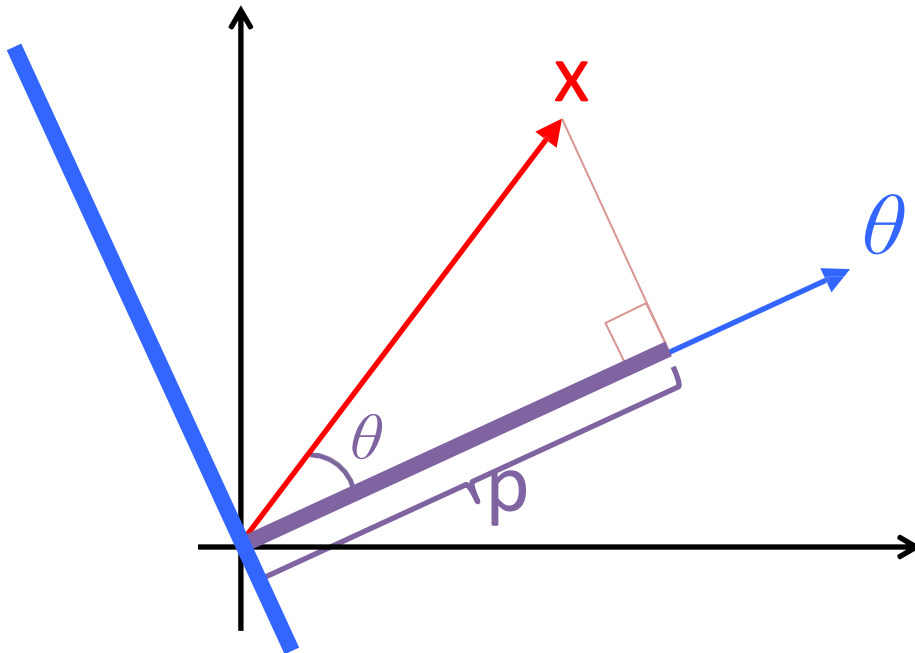
# Understanding the Hyperplane

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^{d} \theta_j^2$$

$$\text{s.t.} \ \theta^T x_i \geq 1, \quad \text{if } y_i = 1$$

$$\text{s.t.} \ \theta^T x_i \leq -1, \quad \text{if } y_i = -1$$

Assume $\theta_0 = 0$ so that the hyperplane is centered at the origin, and that d = 2



$$\theta^\top \mathbf{x} = \left|\left|\boldsymbol{\theta}\right|\right|_2 \left|\left|\mathbf{x}\right|\right|_2 \cos(\boldsymbol{\theta})$$
$$= p \left|\left|\boldsymbol{\theta}\right|\right|_2$$
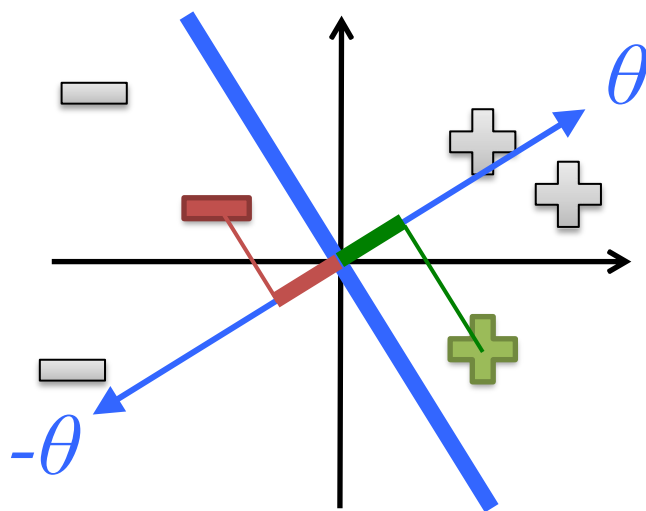
# Maximizing the Margin

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^{d} \theta_j^2$$

$$\text{s.t. } \theta^T x_i \geq 1, \quad \text{if } y_i = 1$$

$$\text{s.t. } \theta^T x_i \leq -1, \quad \text{if } y_i = -1$$

Assume $\theta_0 = 0$ so that the hyperplane is centered at the origin, and that d = 2

Let $p_i$ be the projection of $x_i$ onto the vector $\theta$



Since p is small, therefore $||\theta||_2$ must be large to have $p||\theta||_2 \geq 1$ (or $\leq$ -1)

Since p is larger, $||\theta||_2$ can be smaller and still satisfy $p||\theta||_2 \geq 1$ (or $\leq -1$)

# Support Vectors



$$\theta$$

$$\theta^\top x = 1 \qquad \theta^\top x = -1$$

# Size of the Margin

For the support vectors, we have. $p\left|\left|\boldsymbol{\theta}\right|\right|_2 = \pm 1$

p is the length of the projection of the SVs onto $\boldsymbol{\theta}$



Therefore,

$$p = \frac{1}{\left|\left|\boldsymbol{\theta}\right|\right|_2}$$

$$\text{Margin} = 2p = \frac{2}{\left|\left|\boldsymbol{\theta}\right|\right|_2}$$

# The SVM Dual Problem

The primal SVM problem was given as

$$J(\boldsymbol{\theta}) = \min_{\theta} \frac{1}{2} \sum_{j=1}^{d} \theta_j^2$$

$$\text{s.t.} \ y_i(\boldsymbol{\theta}^T x_i) \geq 1, \quad \text{if } \forall i$$

Can solve it more efficiently by taking the Lagrangian dual

- Duality is a common idea in optimization
- Transforms into a simpler optimization
- Key idea:    introduce slack variables $\alpha_i$ for each constraint
  - $\alpha_i$ indicates how important a particular constraint is to the solution

# The SVM Dual Problem

- The Lagrangian is given by

$$L(\boldsymbol{\theta}, \boldsymbol{\alpha}) = \frac{1}{2}\sum_{j=1}^{d} \theta_j^2 - \sum_{i=1}^{n} \alpha_i (y_i(\theta^T x_i) - 1)$$

$$\text{s.t. } \alpha_i \geq 0. \quad \forall i$$

- By definition this new formulation

$$\min_{\theta} \max_{\alpha} L(\boldsymbol{\theta}, \boldsymbol{\alpha}) \equiv \min_{\theta} J(\boldsymbol{\theta})$$

- We must minimize over $\theta$ and maximize over $\boldsymbol{\alpha}$
- At optimal solution, partials w.r.t $\theta$'s are 0

Solve by a bunch of algebra and calculus ... and we obtain ...

# Solving the Optimization Problem (Primal to Dual)

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^{d} \theta_j^2$$

$$\text{s. t. } \forall i \ y_i(\theta^T x_i) \geq 1$$

Quadratic programming with linear constraints

**Minimize**

$$L(\boldsymbol{\theta}, \boldsymbol{\alpha}) = \frac{1}{2} \sum_{j=1}^{d} \theta_j^2 - \sum_{i=1}^{n} \alpha_i (y_i(\theta^T x_i) - 1)$$

$$\text{s. t. } \forall i \ \alpha_i \geq 0$$

Lagrangian Function

10

# Solving the Optimization Problem

Minimize

$$L(\boldsymbol{\theta}, \boldsymbol{\alpha}) = \frac{1}{2}\sum_{j=1}^{d} \theta_j^2 - \sum_{i=1}^{n} \alpha_i(y_i(\theta^T x_i) - 1)$$

$$\text{s.t.} \,\forall i \,\, \alpha_i \geq 0$$

$$\frac{\partial L}{\partial \boldsymbol{\theta}} = 0 \Rightarrow \boldsymbol{\theta} - \sum_{i=1}^{n} \alpha_i y_i \mathbf{x}_i = 0$$

$$\Rightarrow \boldsymbol{\theta} = \sum_{i=1}^{n} \alpha_i y_i \mathbf{x}_i$$

The representer theorem: $\boldsymbol{\theta}$ as linear combination of training data

$$\frac{\partial L}{\partial \theta_0} = 0 \Rightarrow \sum_{i=1}^{n} \alpha_i y_i = 0$$

Where does $\theta_0$ come from?

# Solving the Optimization Problem

$$L = \frac{1}{2}\sum_{j=1}^{d}\theta_j^2 - \sum_{i=1}^{n}\alpha_i(y_i(\theta^T x_i) - 1)$$

If we substitute $\boldsymbol{\theta} = \sum_{i=1}^{n}\alpha_i y_i \mathbf{x}_i$ to L, we have

Details may be skipped

$$L = \frac{1}{2}\sum_{i=1}^{n}\alpha_i y_i \mathbf{x}_i^T \sum_{j=1}^{n}\alpha_j y_j \mathbf{x}_j - \sum_{i=1}^{n}\alpha_i(y_i(\theta^T x_i) - 1)$$

$$L = \frac{1}{2}\sum_{i=1}^{n}\alpha_i y_i \mathbf{x}_i^T \sum_{j=1}^{n}\alpha_j y_j \mathbf{x}_j - \sum_{i=1}^{n}\alpha_i\left(y_i\left(\sum_{j=1}^{n}\alpha_j y_j \mathbf{x}_j^T \mathbf{x}_i\right) - 1\right)$$

$$L = \frac{1}{2}\sum_{i=1}^{n}\sum_{j=1}^{n}\alpha_i\alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=1}^{n}\alpha_i - \sum_{i=1}^{n}\alpha_i y_i \sum_{j=1}^{n}\alpha_j y_j \mathbf{x}_j^T \mathbf{x}_i$$

$$L = -\frac{1}{2}\sum_{i=1}^{n}\sum_{j=1}^{n}\alpha_i\alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=1}^{n}\alpha_i$$

This is a function of $\alpha_i$

# The Dual Problem

The objective function is in terms of $\alpha_i$ only.
- It is known as the dual problem: if we know $\boldsymbol{\theta}$, we know all $\alpha_i$; if we know all $\alpha_i$, we know $\boldsymbol{\theta}$
- The original problem = primal problem
- The objective function of the dual problem needs to be maximized (comes out from the KKT theory)
- *Learn $d$ parameters for primal. $N$ parameters for dual. Efficient if $N \ll d$*

The dual problem is:

$$\textbf{max} \quad W(\alpha) = -\frac{1}{2}\sum_{i=1}^{n}\sum_{j=1}^{n}\alpha_i\alpha_j y_i y_j \mathbf{x}_i^T\mathbf{x}_j + \sum_{i=1}^{n}\alpha_i$$

Subject to $\quad \alpha_i \geq 0$, $\qquad \sum_{i=1}^{n}\alpha_i y_i = 0$

Properties of $\alpha_i$ when we introduce the Lagrange multipliers

The result when we differentiate the original Lagrangian w.rt. $\theta_0$

This is a quadratic programming (QP) problem. A global maximum of $\alpha_i$ can always be found.
$\boldsymbol{\theta}$ can be recovered by $\boldsymbol{\theta} = \sum_{i=1}^{n}\alpha_i y_i \mathbf{x}_i$

Rememeber the classifier becomes $= \boldsymbol{\theta}^T\mathbf{x} = (\sum_{i=1}^{n}\alpha_i y_i \mathbf{x_i})^T\mathbf{x}$

# One Slide Summary

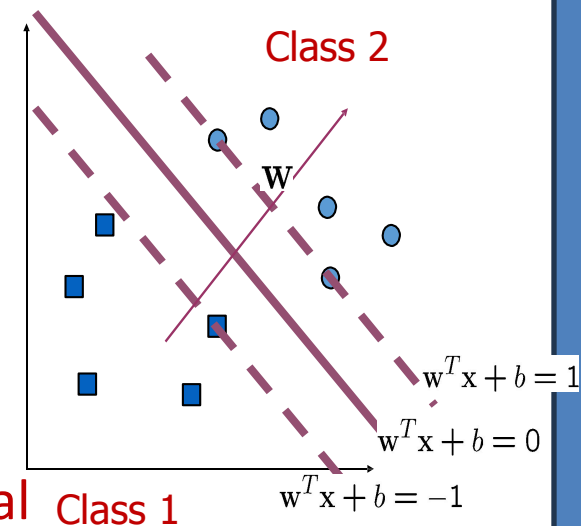Minimize $\frac{1}{2}\boldsymbol{\theta}^\top\boldsymbol{\theta}$

subject to $y_n(\boldsymbol{\theta}^\top x_n) \geq 1$, for $n = 1,2,\ldots N$

**Maximize Margin**
**Learn $\theta$ (hyperplane)**

Minimize                                                    **Primal**

$$L(\boldsymbol{\theta}, \boldsymbol{\alpha}) = \frac{1}{2}\boldsymbol{\theta}^\top\boldsymbol{\theta} - \sum_{i=1}^{n} \alpha_i(y_i(\boldsymbol{\theta}^T x_i) - 1)$$

$$\text{s.t.} \ \forall i \ \ \alpha_i \geq 0$$



Class 2

w

$w^T x + b = 1$

$w^T x + b = 0$

**Dual**   Class 1            $w^T x + b = -1$

Maximize $J(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2}\sum_{i=1}^{n}\sum_{j=1}^{n} \alpha_i\alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$

s.t. $\alpha_i \geq 0 \ \forall i$          (comes from Lagrangian Assumptions)

$\sum_i \alpha_i y_i = 0$          (comes from differentiating w.r.t $\theta_0$)

**Maximize Margin**
**Learn $\alpha$ (weight**
**of support vectors)**

# SVM Dual

$$\text{Maximize } J(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

$$\text{s. t. } \alpha_i \geq 0 \; \forall i$$

$$\sum_i \alpha_i y_i = 0$$

The decision function is given by

$$h(\mathbf{x}) = \text{sign}\left( \sum_{i \in \mathcal{SV}} \alpha_i y_i \langle \mathbf{x}, \mathbf{x}_i \rangle + b \right)$$

$$\text{where} \quad b = \frac{1}{|\mathcal{SV}|} \sum_{i \in \mathcal{SV}} \left( y_i - \sum_{j \in \mathcal{SV}} \alpha_j y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \right)$$

Interesting Twist:
Many $\alpha_i$'s are zero.
Only SVs have non-zero $\alpha_i$'s

11

# Understanding the Dual

$$\text{Maximize } J(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

$$\text{s. t. } \alpha_i \geq 0 \; \forall i$$

$$\sum_i \alpha_i y_i = 0$$

Balances between the weight of constraints for different classes

Constraint weights ($\alpha_i$'s) cannot be negative

# Understanding the Dual

$$\text{Maximize } J(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2}\sum_{i=1}^{n}\sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j \boxed{\langle \mathbf{x}_i, \mathbf{x}_j \rangle}$$

$$\text{s. t. } \alpha_i \geq 0 \ \forall i$$

$$\sum_i \alpha_i y_i = 0$$

Points with different labels increase the sum

Points with same label decrease the sum

Measures the similarity between points

Intuitively, we should be more careful around points near the margin

# Understanding the Dual

$$\text{Maximize } J(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j \boxed{\langle \mathbf{x}_i, \mathbf{x}_j \rangle}$$

$$\text{s. t. } \alpha_i \geq 0 \ \forall i$$
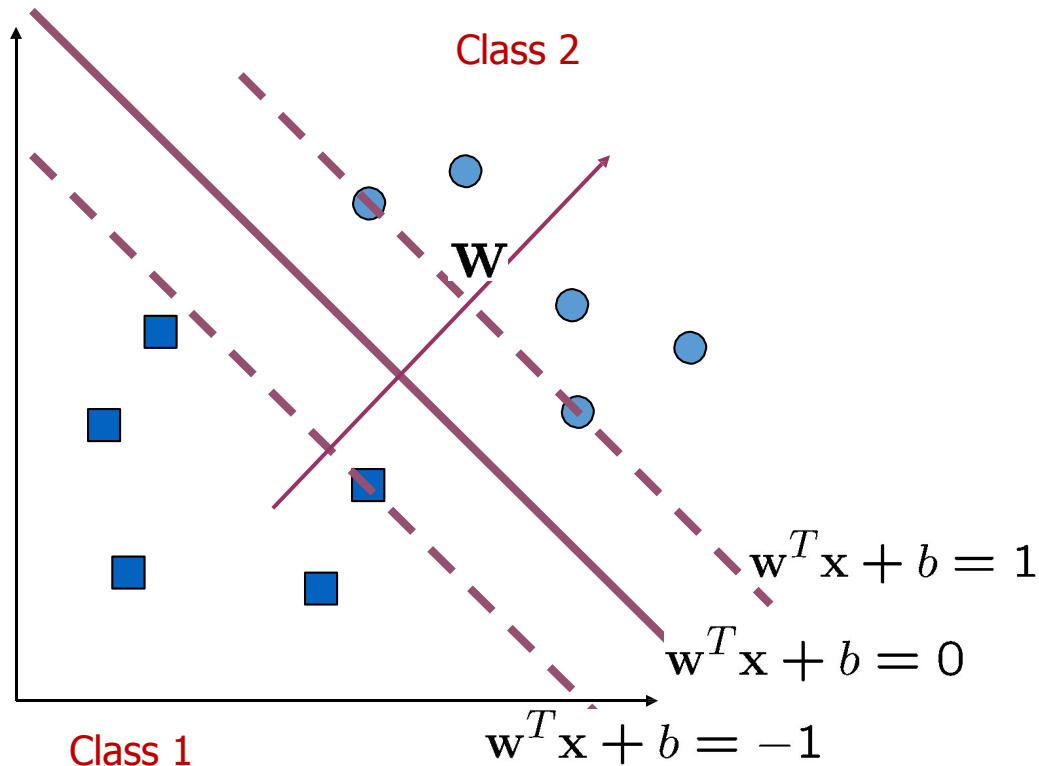
$$\sum_i \alpha_i y_i = 0$$

In the solution, either:

- $\alpha_i > 0$ and the constraint is tight $(y_i(\boldsymbol{\theta}^\top x_i) = 1)$
  - ➢ point is a support vector

- $\alpha_i = 0$
  - ➢ point is not a support vector

# Deploying the Solution

Given the optimal solution $\boldsymbol{\alpha} *$, optimal weights are

$$\boldsymbol{\theta}^* = \sum_{i \in SVs} \alpha_i^* y_i \mathbf{x}_i$$

# A Geometrical Interpretation

Class 2

Class 1

$$\mathbf{w}^T\mathbf{x} + b = 1$$

$$\mathbf{w}^T\mathbf{x} + b = 0$$

$$\mathbf{w}^T\mathbf{x} + b = -1$$

$\mathbf{W}$

# A Geometrical Interpretation



Class 2

$\mathbf{W}$

$\mathbf{w}^T\mathbf{x} + b = 1$

$\mathbf{w}^T\mathbf{x} + b = 0$

$\mathbf{w}^T\mathbf{x} + b = -1$

Class 1

# A Geometrical Interpretation



Class 2

$\alpha_8 = 0.6$  $\alpha_{10} = 0$

$\mathbf{W}$

$\alpha_7 = 0$

$\alpha_2 = 0$

$\alpha_5 = 0$

$\alpha_5 = 0$

$\alpha_1 = 0.8$

$\alpha_6 = 1.4$

$\mathbf{w}^T \mathbf{x} + b = 1$

$\alpha_5 = 0$  $\alpha_5 = 0$

$\mathbf{w}^T \mathbf{x} + b = 0$

Class 1

$\mathbf{w}^T \mathbf{x} + b = -1$

# Characteristics of the Solution

For testing with a new data $\mathbf{z}$
Compute

$$\boldsymbol{\theta}^T \mathbf{z} = \sum_{j=1}^{s} \alpha_{t_j} y_{t_j} \left( \mathbf{x}_{t_j}^T \mathbf{z} \right)$$

Classify $\mathbf{z}$ as class 1 if the sum is positive, and class 2 otherwise. Note $\boldsymbol{\theta}$ need not be formed explicitly.

Given the optimal solution $\boldsymbol{\alpha}^*$, optimal weights are

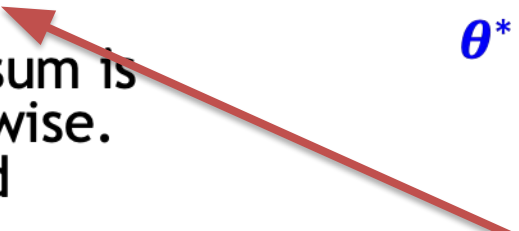$$\boldsymbol{\theta}^* = \sum_{i \in \text{SVs}} \alpha_i^* y_i \mathbf{x}_i$$

Note: The computation relies on a dot product between the test point and the support vectors

# What if Data Are Not Linearly Separable?

Cannot find $\boldsymbol{\theta}$ that satisfies. $y_i(\theta^T x_i) \geq 1 \ \forall i$

Introduce slack variables $\xi_i$

$$y_i(\theta^T x_i) \geq 1 - \xi_i \ \forall i$$

New Problem

$$\min_{\theta} \frac{1}{2}\sum_{j=1}^{d} \theta_j^2 + C \sum_{i} \xi_i$$

$$\text{s.t.} \ y_i(\theta^T x_i) \geq 1 - \xi_i, \quad \text{if } \forall i$$

# Strengths of SVMs

- Good generalization in theory
- Good generalization in practice
- Work well with few training instances
- Find globally best model
- Efficient algorithms
- Amenable to the kernel trick …