

EE5907: Pattern Recognition & EE5026: Machine Learning for Data Analytics

Dr. WANG Si

Email: si.wang@nus.edu.sg

Office: E1a-04-01

Outlines

- Pattern Representation Learning
 - Unsupervised Representation Learning (week 7)
 - Supervised Representation Learning (week 8)
 - Unified Framework: Graph Embedding (week 8)
- Pattern Recognition Models
 - Clustering (K-means, Agglomerative)(week 9)
 - Gaussian Mixture Model and Boosting (week 10)
 - Support Vector Machines (week 11)
- Deep Learning (week 12)
- Revision and Q&A (week 13)

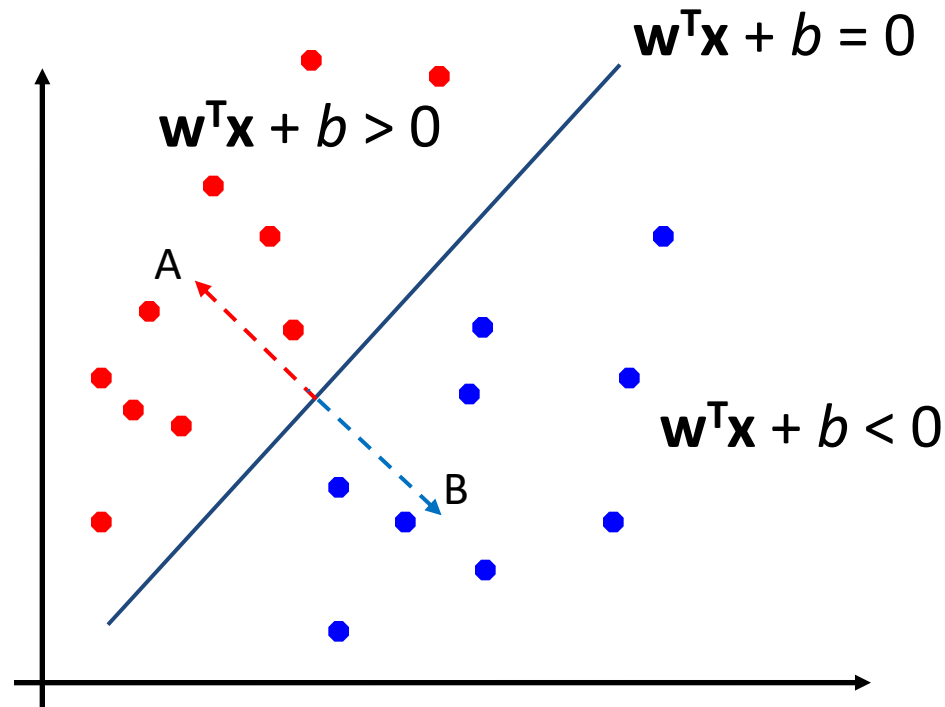
Support Vector Machine (SVM)

- We will discuss:
 - Hard-margin SVM (Primal and Dual formulations)
 - Soft-margin SVM (Primal and Dual formulations)
 - Kernel SVM
 - Multi-class SVM
- At the end of this lecture, you should be able to:
 - Understand different formulations for different SVMs
 - Elaborate the two tricks to deal with non linearly separable data
 - Apply SVM to the multi-class problem

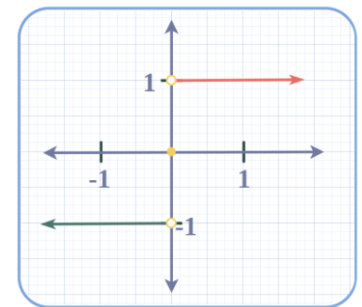
Recap: Linear Classifier

Class labels:

- +1
- -1



Graph of Signum Function



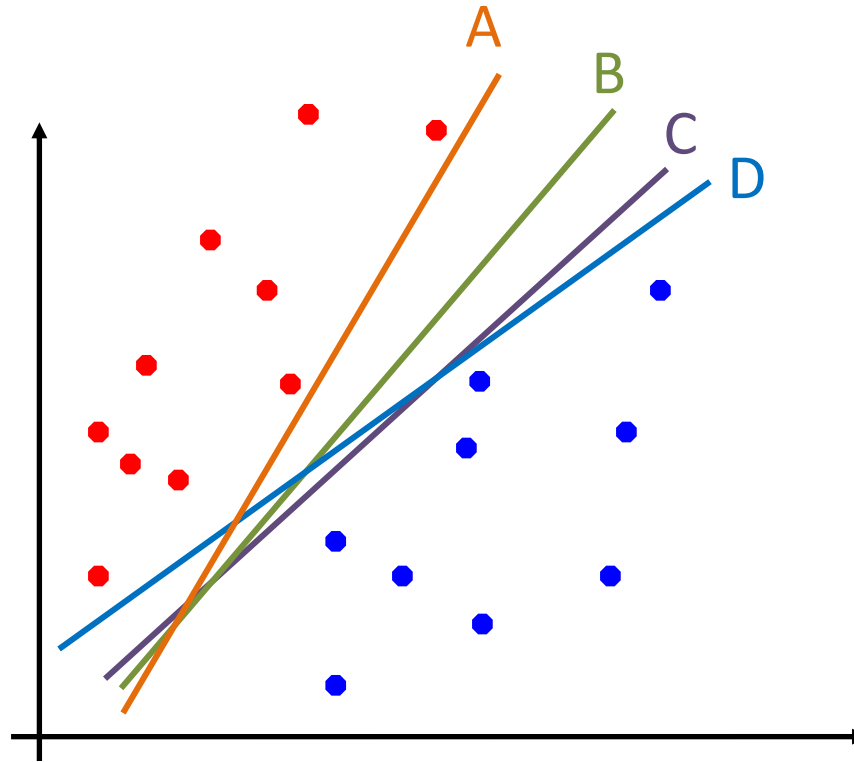
What is the direction of \mathbf{w} ?

$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

$$\text{sign}(z) = \begin{cases} -1, & z < 0 \\ 0, & z = 0 \\ +1, & z > 0 \end{cases}$$

Recap: Linear Classifier

- Which one is optimal?



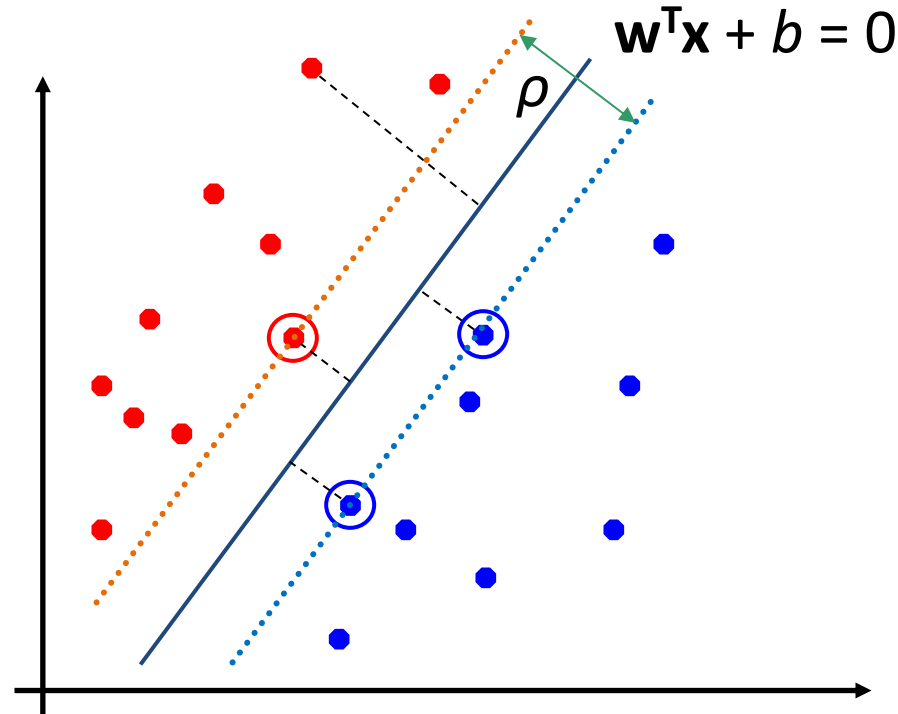
Support Vector Machine (SVM)

[Hard Margin]

Hard-Margin Support Vector Machine (SVM)

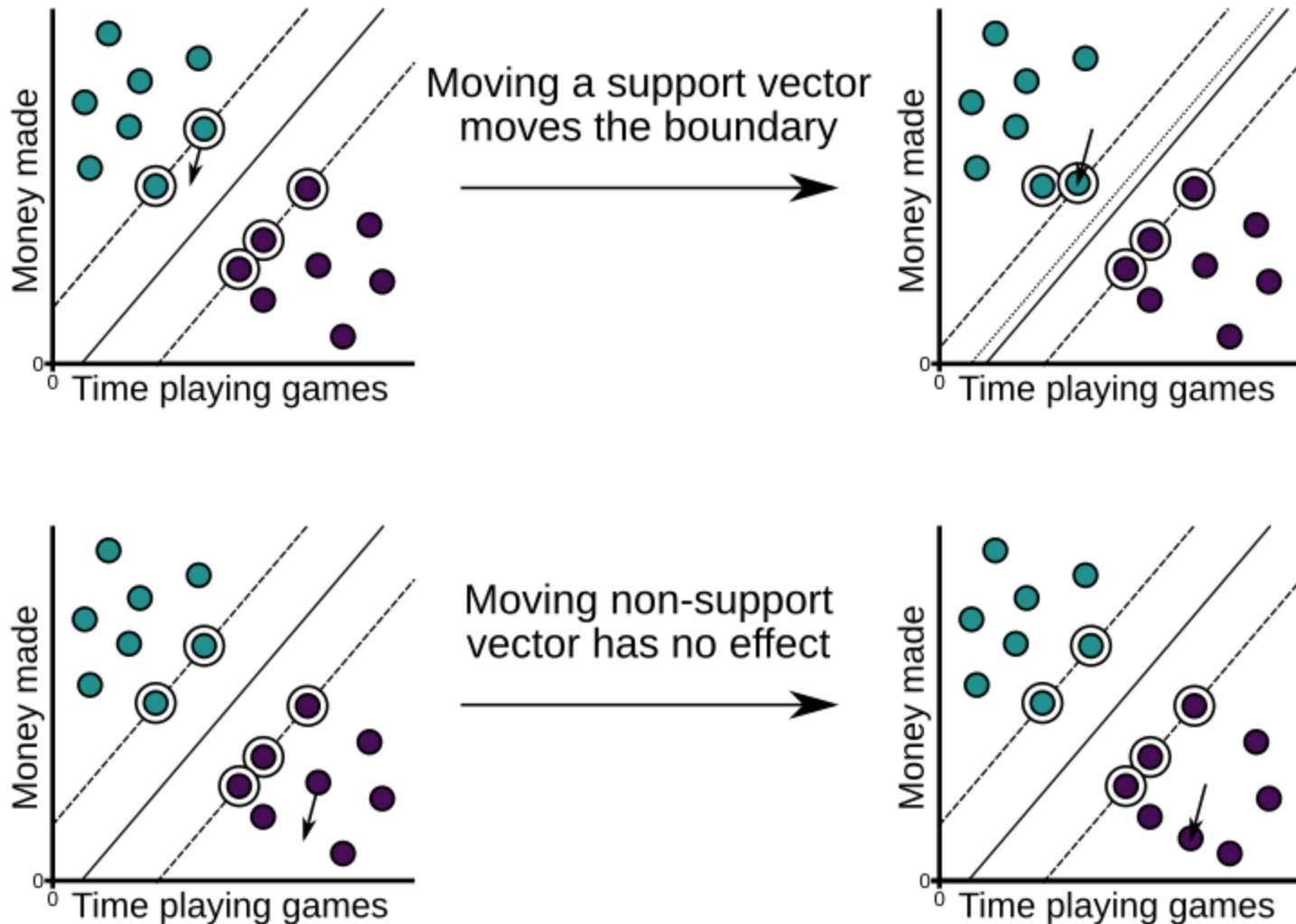
- Maximize the **margin** around the decision boundary

⊙ Support vector



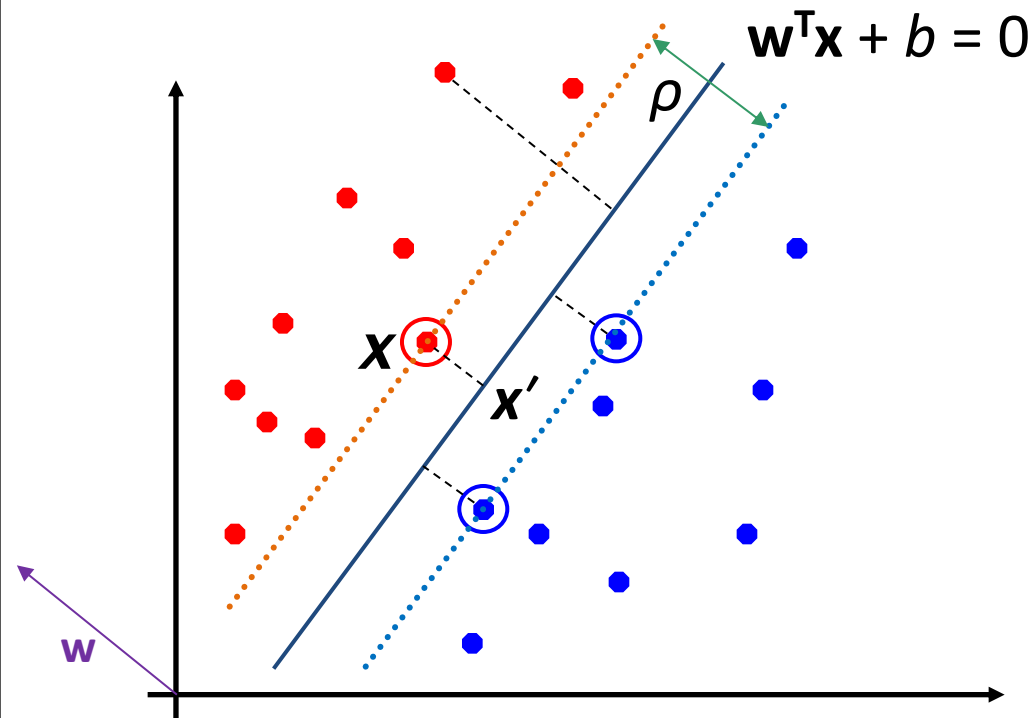
Hard-Margin Support Vector Machine (SVM)

- Move Support Vectors



Hard-Margin Support Vector Machine (SVM)

Equations for orange dashed line and blue dashed line?



Orange dashed line:

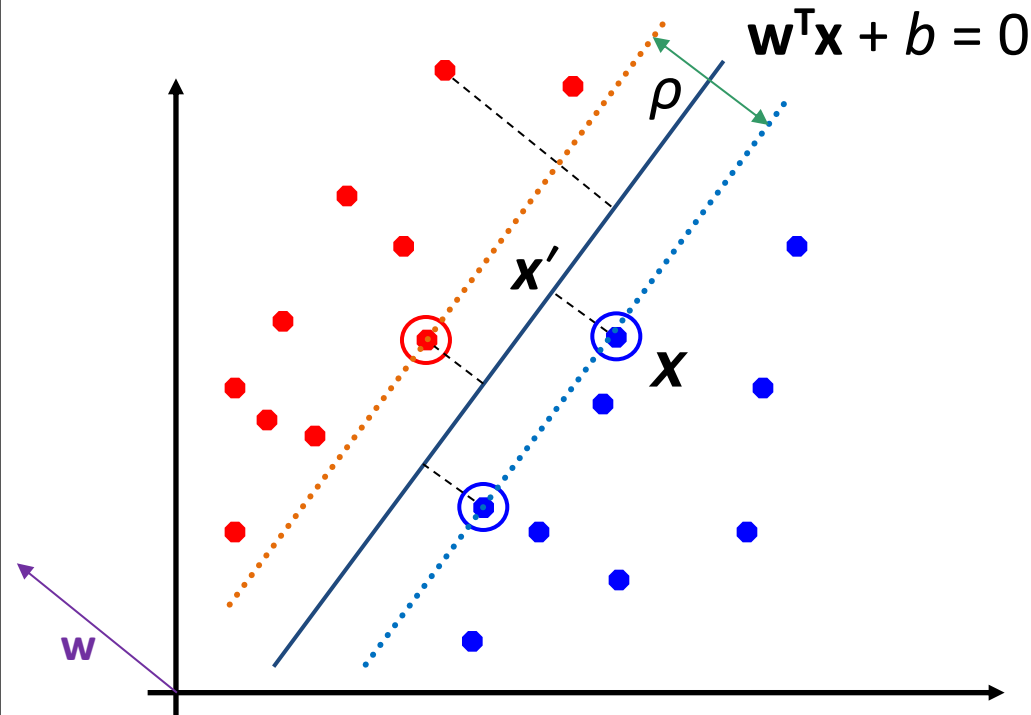
$$\begin{cases} w^T x' + b = 0 \\ \frac{x - x'}{\rho/2} = \frac{w}{\|w\|} \rightarrow x' = x - \frac{\rho w}{2\|w\|} \end{cases}$$

$$w^T \left(x - \frac{\rho w}{2\|w\|} \right) + b = 0$$

$$w^T x + b = \frac{\rho \|w\|}{2}$$

Hard-Margin Support Vector Machine (SVM)

Equations for orange dashed line and blue dashed line?



Blue dashed line:

$$\begin{cases} \mathbf{w}^T \mathbf{x}' + b = 0 \\ \frac{\mathbf{x} - \mathbf{x}'}{\rho/2} = -\frac{\mathbf{w}}{\|\mathbf{w}\|} \rightarrow \mathbf{x}' = \mathbf{x} + \frac{\rho \mathbf{w}}{2\|\mathbf{w}\|} \end{cases}$$

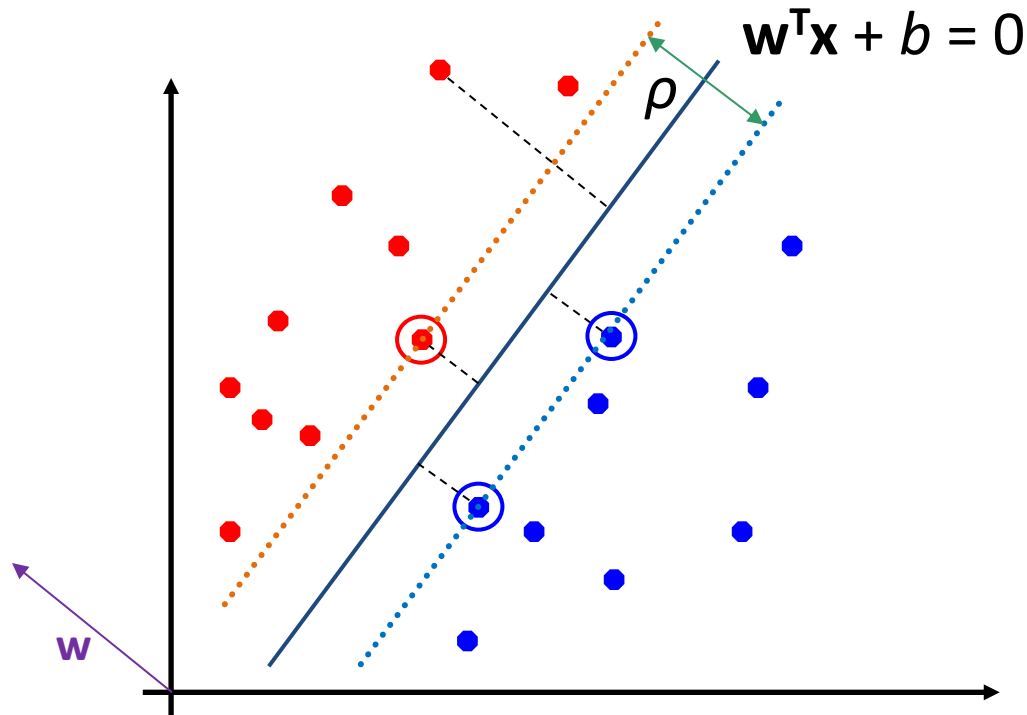
$$\mathbf{w}^T \left(\mathbf{x} + \frac{\rho \mathbf{w}}{2\|\mathbf{w}\|} \right) + b = 0$$

$$\mathbf{w}^T \mathbf{x} + b = -\frac{\rho \|\mathbf{w}\|}{2}$$

Hard-Margin Support Vector Machine (SVM)

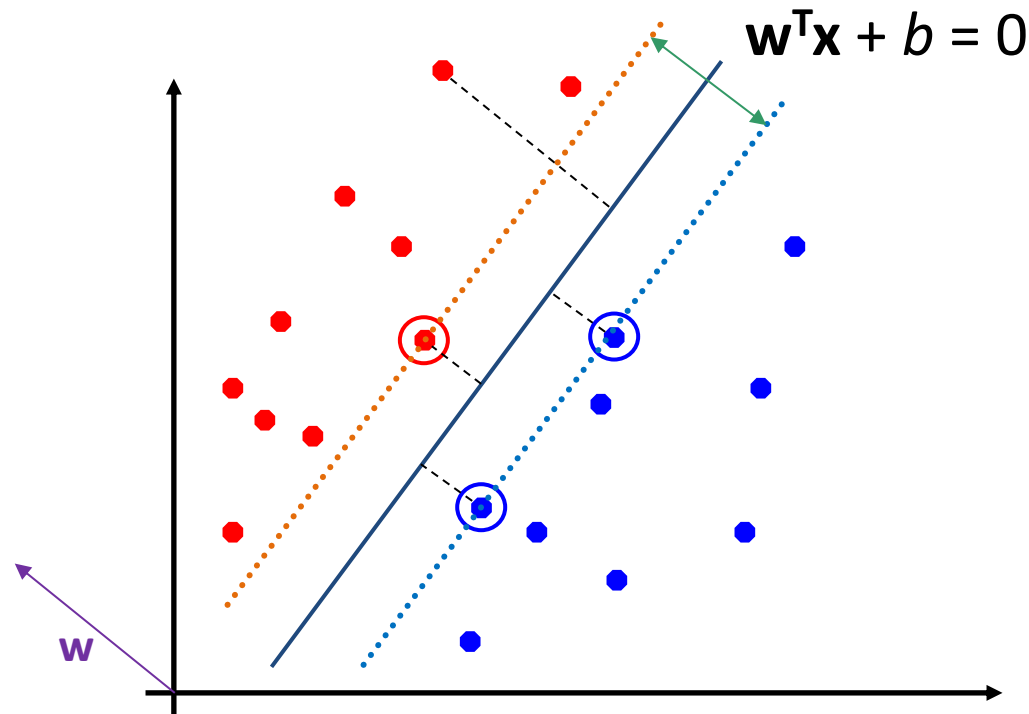
- Plus Plane: $\mathbf{w}^T \mathbf{x} + b = 1$
- Minus Plane: $\mathbf{w}^T \mathbf{x} + b = -1$

$$\Rightarrow \rho = \frac{2}{\|\mathbf{w}\|}$$



Hard-Margin Support Vector Machine (SVM)

- Objective: find \mathbf{w} , b such that
 - ✓ All positive training points ($\mathbf{x}_i, y_i = 1$) in the red zone $\Rightarrow \mathbf{w}^T \mathbf{x}_i + b \geq 1, \text{ if } y_i = 1$
 - ✓ All negative training points ($\mathbf{x}_i, y_i = -1$) in the blue zone $\Rightarrow \mathbf{w}^T \mathbf{x}_i + b \leq -1, \text{ if } y_i = -1$
 - ✓ The margin ρ is maximized $\Rightarrow \max \frac{2}{\|\mathbf{w}\|}$



Hard-Margin SVM Primal as Quadratic Program (QP)

- Objective: find \mathbf{w} , b such that

$$\left. \begin{array}{l} \mathbf{w}^T \mathbf{x}_i + b \geq 1, \text{ if } y_i = 1 \\ \mathbf{w}^T \mathbf{x}_i + b \leq -1, \text{ if } y_i = -1 \end{array} \right\} \Rightarrow y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

$$\max \frac{2}{\|\mathbf{w}\|} \Rightarrow \min \|\mathbf{w}\|, \text{ or } \min \|\mathbf{w}\|^2, \text{ or } \min \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

Reformulated as a quadratic optimization problem:

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

Subject to $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$, for all $(\mathbf{x}_i, y_i), i = 1 \dots N$

Quadratic objective function
Linear constraints

Quadratic
program

❖ QP solver libraries:
cvxpy, quadprog, OSQP,
scipy.optimize, etc.

Hard-Margin SVM Primal & Dual

Primal
Problem

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

Subject to $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$, for all $(\mathbf{x}_i, y_i), i = 1 \dots N$

$d+1$ variables

N constraints

Which one is
better?



Dual
Problem

$$\max_{\lambda} -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_i \lambda_i$$

s.t. $\lambda_i \geq 0, \sum_i \lambda_i y_i = 0, i = 1, \dots, N$

N variables

$N+1$ constraints

λ_i is the Lagrange Multiplier associated with each constraint in the primal problem

Hard-Margin SVM Primal -> Dual

Primal
Problem

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

Subject to $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$, for all $(\mathbf{x}_i, y_i), i = 1 \dots N$

■ Inequality Constraints:

Karush-Kuhn-Tucker (KKT) conditions

$$f(\mathbf{w}, b) = \frac{1}{2} \mathbf{w}^T \mathbf{w},$$

1) Define the Lagrangian function:

$$g_i(\mathbf{w}, b) = y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1$$

$$\mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) = f(\mathbf{w}, b) - \sum_i \lambda_i g_i(\mathbf{w}, b), \text{ where } \lambda_i \geq 0, \forall i$$

2) KKT conditions for optimality:

✓ Stationarity:

$$\left[\begin{array}{l} \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{0} \rightarrow \mathbf{w} - \sum_i \lambda_i y_i \mathbf{x}_i = \mathbf{0} \rightarrow \mathbf{w} = \sum_i \lambda_i y_i \mathbf{x}_i \\ \frac{\partial \mathcal{L}}{\partial b} = 0 \rightarrow \sum_i \lambda_i y_i = 0 \end{array} \right.$$

Hard-Margin SVM Primal -> Dual

■ Inequality Constraints:

Karush-Kuhn-Tucker (KKT) conditions

1) Define the Lagrangian function:

$$\mathcal{L}(\mathbf{w}, b, \lambda) = f(\mathbf{w}, b) - \sum_i \lambda_i g_i(\mathbf{w}, b), \text{ where } \lambda_i \geq 0, \forall i$$

$$f(\mathbf{w}, b) = \frac{1}{2} \mathbf{w}^T \mathbf{w},$$

$$g_i(\mathbf{w}, b) = y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 \geq 0$$

2) KKT conditions for optimality:

✓ Stationarity: $\mathbf{w}^* = \sum_i \lambda_i y_i \mathbf{x}_i, \sum_i \lambda_i y_i = 0$

✓ Primal Feasibility: $g_i(\mathbf{w}, b) = y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 \geq 0, \forall i$

✓ Dual Feasibility: $\lambda_i \geq 0, \forall i$

✓ Complementary Slackness: $\lambda_i g_i(\mathbf{w}, b) = 0, \forall i$



If $g_i(\mathbf{w}, b) > 0$ (inactive), $\lambda_i = 0$

If $g_i(\mathbf{w}, b) = 0$ (active), $\lambda_i > 0$

λ_i only plays a role
when a constraint is
exactly binding

Hard-Margin SVM Primal -> Dual

$$\begin{aligned}
 \min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b, \lambda) &= \min_{\mathbf{w}, b} f(\mathbf{w}, b) - \sum_i \lambda_i g_i(\mathbf{w}, b) \\
 &= \min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_i \lambda_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1) \\
 &= \min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_i \lambda_i y_i \mathbf{w}^T \mathbf{x}_i - \sum_i \lambda_i y_i b + \sum_i \lambda_i \\
 &= \frac{1}{2} \mathbf{w}^{*T} \mathbf{w}^* - \mathbf{w}^{*T} \mathbf{w}^* + \sum_i \lambda_i \\
 &= -\frac{1}{2} \sum_i \lambda_i y_i \mathbf{x}_i^T \sum_i \lambda_i y_i \mathbf{x}_i + \sum_i \lambda_i \\
 &= -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_i \lambda_i
 \end{aligned}$$

$$f(\mathbf{w}, b) = \frac{1}{2} \mathbf{w}^T \mathbf{w},$$

$$g_i(\mathbf{w}, b) = y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 \geq 0$$

$$\begin{cases} \mathbf{w}^* = \sum_i \lambda_i y_i \mathbf{x}_i \\ \sum_i \lambda_i y_i = 0 \end{cases}$$

Primal

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w}, \text{ s. t. } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

WHY?

Dual

$$\begin{aligned}
 \max_{\lambda} & -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_i \lambda_i \\
 \text{s.t. } & \lambda_i \geq 0, \sum_i \lambda_i y_i = 0
 \end{aligned}$$

Lagrange Duality

Prove

$$\min_{\mathbf{w}, b} f(\mathbf{w}, b), \text{ s.t. } g_i(\mathbf{w}, b) \geq 0$$

\equiv

$$\max_{\lambda} \min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b, \lambda) \text{ s.t. } \lambda_i \geq 0$$

$$\mathcal{L}(\mathbf{w}, b, \lambda) = f(\mathbf{w}, b) - \sum_i \lambda_i g_i(\mathbf{w}, b)$$



$$f(\mathbf{w}, b) \geq \mathcal{L}(\mathbf{w}, b, \lambda)$$

$$\left[\begin{array}{l} l(\lambda) = \min_{\mathbf{w}, b} \mathcal{L}(\mathbf{w}, b, \lambda), \text{ s.t. } \lambda_i \geq 0 \\ p^* = \min_{\mathbf{w}, b} f(\mathbf{w}, b), \text{ s.t. } g_i(\mathbf{w}, b) \geq 0 \end{array} \right]$$

$$p^* \geq l(\lambda)$$



$$p^* = \max_{\lambda} l(\lambda)$$

$$f(\mathbf{w}, b) = \frac{1}{2} \mathbf{w}^T \mathbf{w},$$
$$g_i(\mathbf{w}, b) = y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1$$

$$\left[\begin{array}{l} \mathbf{w}^* = \sum_i \lambda_i y_i \mathbf{x}_i \\ \sum_i \lambda_i y_i = 0 \end{array} \right]$$

Lagrange Duality

Primal

$$\min_x f(x), \text{ s.t. } g_i(x) \geq 0$$

=

Dual

$$\max_{\lambda} \min_x f(x) - \sum_i \lambda_i g_i(x), \text{ s.t. } \lambda_i \geq 0$$

$$\min_x f(x), \text{ s.t. } g_i(x) \leq 0$$

=

$$\max_x f(x), \text{ s.t. } g_i(x) \geq 0$$

=

$$\max_x f(x), \text{ s.t. } g_i(x) \leq 0$$

=

Lagrange Duality

Primal

Dual

$$\min_x f(x), \text{ s.t. } g_i(x) \geq 0$$

=

$$\max_{\lambda} \min_x f(x) - \sum_i \lambda_i g_i(x), \text{ s.t. } \lambda_i \geq 0$$

$$\min_x f(x), \text{ s.t. } g_i(x) \leq 0$$

=

$$\max_{\lambda} \min_x f(x) + \sum_i \lambda_i g_i(x), \text{ s.t. } \lambda_i \geq 0$$

$$\max_x f(x), \text{ s.t. } g_i(x) \geq 0$$

=

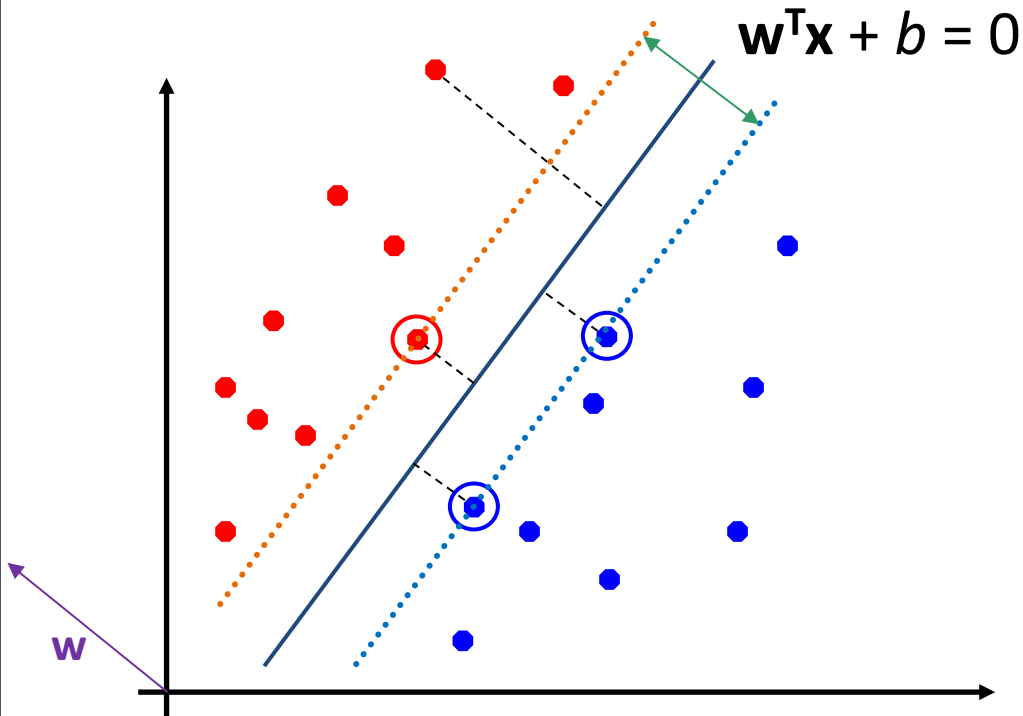
$$\min_{\lambda} \max_x f(x) + \sum_i \lambda_i g_i(x), \text{ s.t. } \lambda_i \geq 0$$

$$\max_x f(x), \text{ s.t. } g_i(x) \leq 0$$

=

$$\min_{\lambda} \max_x f(x) - \sum_i \lambda_i g_i(x), \text{ s.t. } \lambda_i \geq 0$$

Hard-Margin SVM: Decision Boundary



$$\begin{cases} \mathbf{w}^* = \sum_i \lambda_i y_i \mathbf{x}_i \\ \sum_i \lambda_i y_i = 0 \end{cases}$$

$$\mathbf{w}^{*T} \mathbf{x} + b^* = \sum_i \lambda_i y_i \mathbf{x}_i^T \mathbf{x} + b^* = 0$$

for any support vector \mathbf{x}_s ($\lambda_i > 0$),

$$y_s (\mathbf{w}^{*T} \mathbf{x}_s + b^*) = 1$$

↓

$$y_s^2 (\mathbf{w}^{*T} \mathbf{x}_s + b^*) = y_s$$

↓

$$b^* = y_s - \sum_i \lambda_i y_i \mathbf{x}_i^T \mathbf{x}_s$$

Quick Question

For SVM with the following objective function:

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

$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \forall i$$

or

$$\max_{\lambda} -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_i \lambda_i$$

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

Which of the following statements correctly describe *support vectors*?

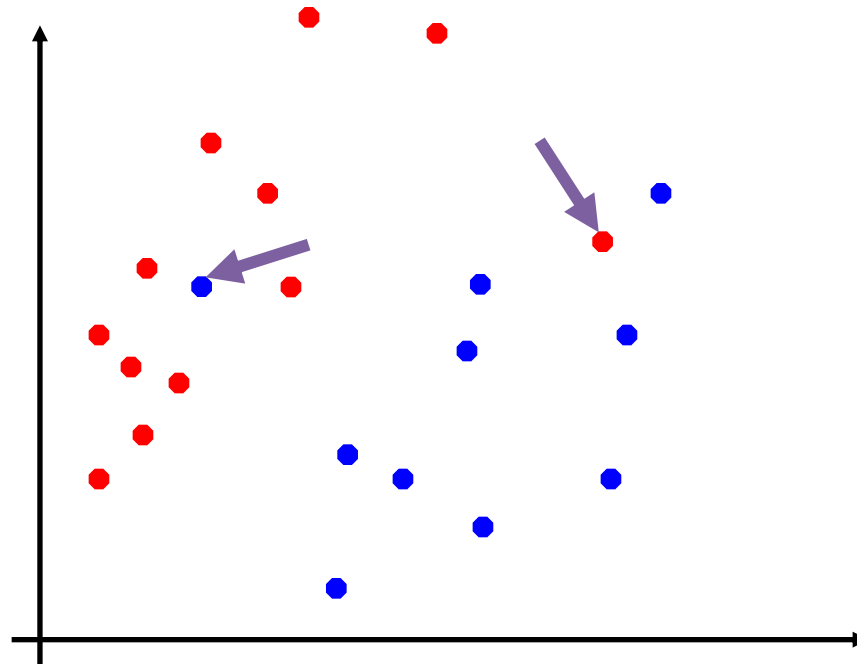
- A. Data points that lie exactly on the margin boundaries.
- B. Data points that have $\lambda_i > 0$.
- C. Data points that determine the decision boundary
- D. Data points inside the margin
- E. Data points that are misclassified

SVM: linearly non-separable data

- What if the data is not linearly separable?

Trick 1: Allow a few points to violate the margin (Soft-margin SVM)

Trick 2: Map data to a higher dimensional space using a kernel function, do linear classification there (Kernel SVM)



Support Vector Machine (SVM)

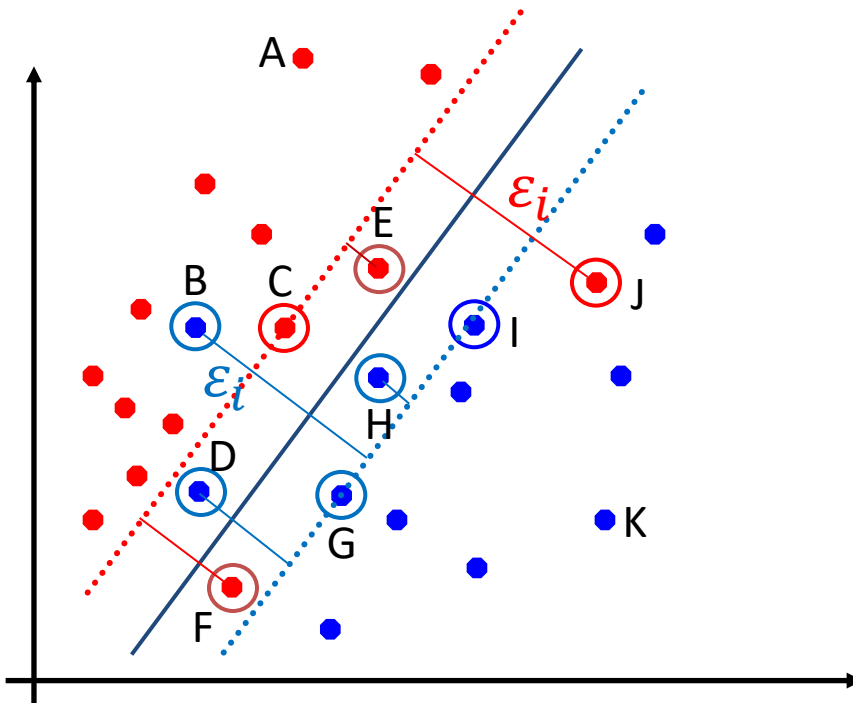
[Soft Margin]

Trick 1: **Soft-margin** SVM---Slack Variable

- Relax hard-margin constraints by introducing **slack variables** ε_i

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \varepsilon_i, \varepsilon_i \geq 0, \forall i$$

Measures how much a data point violates the margin



$$\text{if } \varepsilon_i = 0: y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

Correctly classified, on or outside the margin

$$\text{if } 0 < \varepsilon_i \leq 1: 0 \leq y_i(\mathbf{w}^T \mathbf{x}_i + b) < 1$$

Correctly classified, inside the margin

$$\text{if } \varepsilon_i > 1: y_i(\mathbf{w}^T \mathbf{x}_i + b) < 0$$

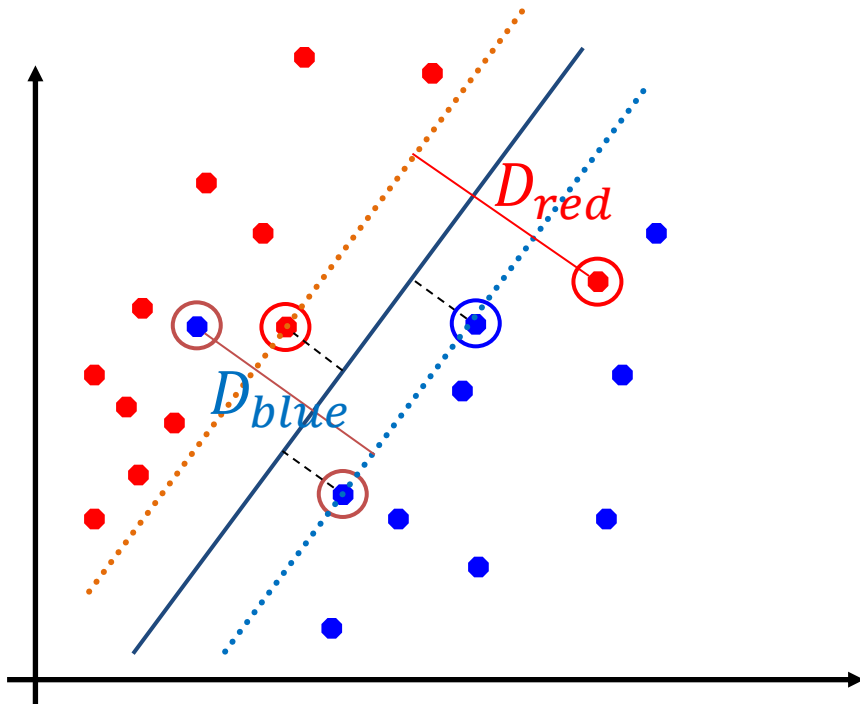
Misclassified, wrong side

ε_i is not the distance to the corresponding margin boundary!

Trick 1: **Soft-margin** SVM---Slack Variable

- **Slack variables** ε_i are not the distances to the corresponding margin boundary!

$D=0$ If the data points do not violate the margin



if $y_i = -1$ (blue):

$$\begin{cases} \mathbf{w}^T \mathbf{x}' + b = -1 \\ \frac{\mathbf{x}_i - \mathbf{x}'}{D_{blue}} = \frac{\mathbf{w}}{\|\mathbf{w}\|} \end{cases}$$

if $y_i = 1$ (red):

$$\begin{cases} \mathbf{w}^T \mathbf{x}' + b = 1 \\ \frac{\mathbf{x}_i - \mathbf{x}'}{D_{red}} = \frac{-\mathbf{w}}{\|\mathbf{w}\|} \end{cases}$$



$$\begin{cases} \mathbf{w}^T \mathbf{x}' + b = y_i \\ \frac{\mathbf{x}_i - \mathbf{x}'}{D} = \frac{-y_i \mathbf{w}}{\|\mathbf{w}\|} \rightarrow \mathbf{x}' = \mathbf{x}_i + \frac{D y_i \mathbf{w}}{\|\mathbf{w}\|} \end{cases}$$

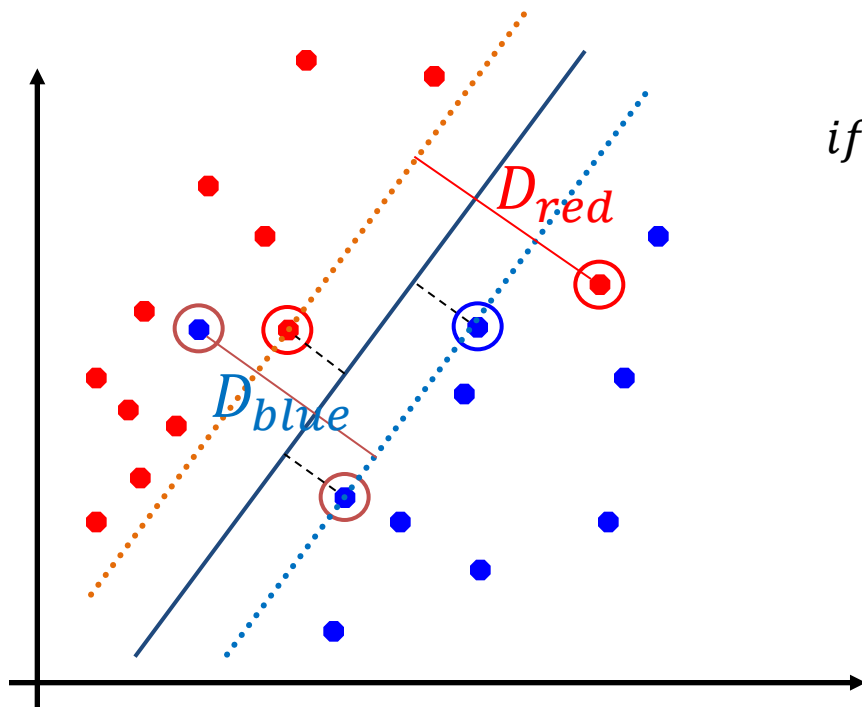
$$\mathbf{w}^T \left(\mathbf{x}_i + \frac{D y_i \mathbf{w}}{\|\mathbf{w}\|} \right) + b = y_i$$

$$D = \frac{y_i - \mathbf{w}^T \mathbf{x}_i - b}{y_i \|\mathbf{w}\|}$$

Trick 1: **Soft-margin** SVM---Slack Variable

- **Slack variables** ε_i are not the distances to the corresponding margin boundary!

$$D = \frac{y_i - \mathbf{w}^T \mathbf{x}_i - b}{y_i \|\mathbf{w}\|}$$



if $D = 0$: $\mathbf{w}^T \mathbf{x}_i + b = y_i$, or $y_i(\mathbf{w}^T \mathbf{x}_i + b) > 1$

Correctly classified, on or outside the margin

if $0 < D \leq 1$: $0 < \frac{y_i - \mathbf{w}^T \mathbf{x}_i - b}{y_i \|\mathbf{w}\|} \leq 1$

[if $y_i = -1$: $-1 < \mathbf{w}^T \mathbf{x}_i + b \leq \|\mathbf{w}\| - 1$
if $y_i = 1$: $1 - \|\mathbf{w}\| \leq \mathbf{w}^T \mathbf{x}_i + b < 1$

$1 - \|\mathbf{w}\| \leq y_i(\mathbf{w}^T \mathbf{x}_i + b) < 1$

Correctly classified inside the margin or misclassified

if $D > 1$: $y_i(\mathbf{w}^T \mathbf{x}_i + b) < 1 - \|\mathbf{w}\|$

Correctly classified inside the margin or misclassified

Trick 1: **Soft-margin** SVM---Slack Variable

- **Slack variables** ε_i are not the distances are not the distances to the corresponding margin boundary!

ε_i



$$D = \frac{y_i - \mathbf{w}^T \mathbf{x}_i - b}{y_i \|\mathbf{w}\|}$$

if $\varepsilon_i = 0$: $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

Correctly classified, on or outside the margin

if $0 < \varepsilon_i \leq 1$: $0 \leq y_i(\mathbf{w}^T \mathbf{x}_i + b) < 1$

Correctly classified, inside the margin

if $\varepsilon_i > 1$: $y_i(\mathbf{w}^T \mathbf{x}_i + b) < 0$

Misclassified, wrong side

if $D = 0$: $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

Correctly classified, on or outside the margin

if $0 < D \leq 1$: $1 - \|\mathbf{w}\| \leq y_i(\mathbf{w}^T \mathbf{x}_i + b) < 1$

Inside the margin or misclassified

if $D > 1$: $y_i(\mathbf{w}^T \mathbf{x}_i + b) < 1 - \|\mathbf{w}\|$

Inside the margin or misclassified

What is the relationship
between ε_i and D ?

Trick 1: **Soft-margin** SVM---Slack Variable

- **Slack variables** ε_i are not the distances!

$$\varepsilon_i \neq D = \frac{y_i - \mathbf{w}^T \mathbf{x}_i - b}{y_i \|\mathbf{w}\|}$$

$$\text{if } \varepsilon_i = 0: y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

Correctly classified, on or outside the margin

$$\text{if } 0 < \varepsilon_i \leq 1: 0 \leq y_i(\mathbf{w}^T \mathbf{x}_i + b) < 1$$

Correctly classified, inside the margin

$$\text{if } \varepsilon_i > 1: y_i(\mathbf{w}^T \mathbf{x}_i + b) < 0$$

Misclassified, wrong side

$$\text{if } D = 0: y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

Correctly classified, on or outside the margin

$$\text{if } 0 < D \leq 1: 1 - \|\mathbf{w}\| \leq y_i(\mathbf{w}^T \mathbf{x}_i + b) < 1$$

Inside the margin or misclassified

$$\text{if } D > 1: y_i(\mathbf{w}^T \mathbf{x}_i + b) < 1 - \|\mathbf{w}\|$$

Inside the margin or misclassified

$$\varepsilon_i = D, \text{ when } \|\mathbf{w}\| = 1$$

Trick 1: **Soft-margin** SVM Primal

Soft-margin Primal

$$\begin{aligned} \min_{\mathbf{w}, b, \boldsymbol{\varepsilon}} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \varepsilon_i \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \varepsilon_i, \varepsilon_i \geq 0, \forall i \end{aligned}$$

$d+1+N$ variables

$2N$ constraints

Hard-margin Primal

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \forall i \end{aligned}$$

$d+1$ variables

N constraints

C	ε_i	Tolerance for misclassification	Margin
Larger			
Smaller			

Trick 1: **Soft-margin** SVM Primal

Soft-margin Primal

$$\begin{aligned} \min_{\mathbf{w}, b, \boldsymbol{\varepsilon}} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \varepsilon_i \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \varepsilon_i, \varepsilon_i \geq 0, \forall i \end{aligned}$$

$d+1+N$ variables

$2N$ constraints

Hard-margin Primal

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \forall i \end{aligned}$$

$d+1$ variables

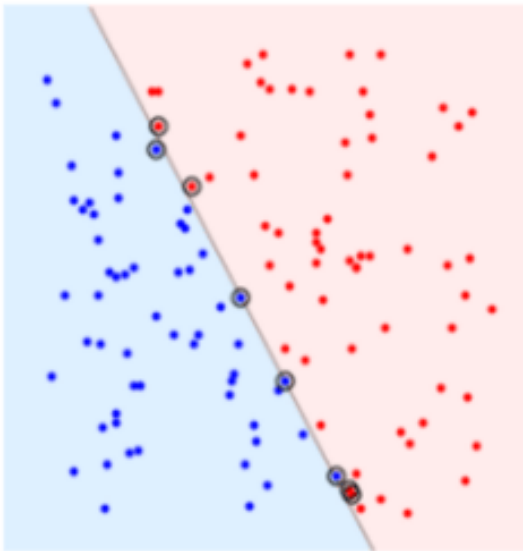
N constraints

C	ε_i	Tolerance for misclassification	Margin
Larger	Smaller	Less	Smaller
Smaller	Larger	More	Larger

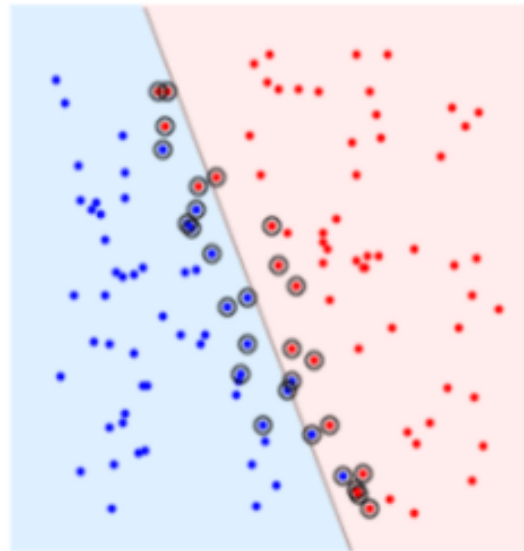
Trick 1: **Soft-margin** SVM Primal

Different values of C

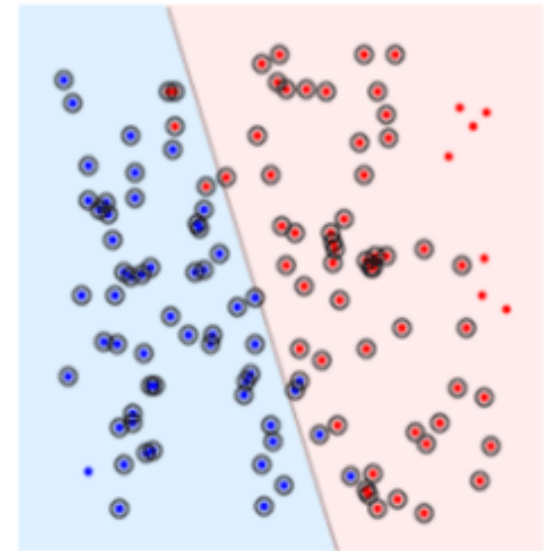
$C=1000$



$C=10$



$C=0.1$



The data points circled are the points lying within the margin region including the margin boundaries

Trick 1: **Soft-margin** SVM Primal->Dual

Soft-margin Primal

$$\begin{aligned} \min_{\mathbf{w}, b, \boldsymbol{\varepsilon}} & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \varepsilon_i \\ \text{s.t. } & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \varepsilon_i, \varepsilon_i \geq 0, \forall i \end{aligned}$$

■ Inequality Constraints:

Karush-Kuhn-Tucker (KKT) conditions

1) Define the Lagrangian function:

$$\mathcal{L}(\mathbf{w}, b, \boldsymbol{\varepsilon}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{w}, b, \boldsymbol{\varepsilon}) - \sum_i \lambda_i g_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) + \mu_i h_i(\mathbf{w}, b, \boldsymbol{\varepsilon}), \text{ where } \lambda_i, \mu_i \geq 0, \forall i$$

$$\begin{aligned} f(\mathbf{w}, b, \boldsymbol{\varepsilon}) &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \varepsilon_i, \\ g_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) &= y_i(\mathbf{w}^T \mathbf{x}_i + b) + \varepsilon_i - 1 \geq 0, \\ h_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) &= \varepsilon_i \geq 0 \end{aligned}$$

Trick 1: **Soft-margin** SVM Primal->Dual

- Inequality Constraints:

Karush-Kuhn-Tucker (KKT) conditions

1) Define the Lagrangian function:

$$\mathcal{L}(\mathbf{w}, b, \boldsymbol{\varepsilon}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{w}, b, \boldsymbol{\varepsilon}) - \sum_i \lambda_i g_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) + \mu_i h_i(\mathbf{w}, b, \boldsymbol{\varepsilon}), \text{ where } \lambda_i, \mu_i \geq 0, \forall i$$

$$f(\mathbf{w}, b, \boldsymbol{\varepsilon}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \varepsilon_i, \quad g_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = y_i(\mathbf{w}^T \mathbf{x}_i + b) + \varepsilon_i - 1 \geq 0, \quad h_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = \varepsilon_i \geq 0$$

2) KKT conditions for optimality:

✓ Stationarity:

$$\left[\begin{array}{l} \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{0} \rightarrow \mathbf{w} - \sum_i \lambda_i y_i \mathbf{x}_i = \mathbf{0} \rightarrow \mathbf{w} = \sum_i \lambda_i y_i \mathbf{x}_i \\ \frac{\partial \mathcal{L}}{\partial b} = 0 \rightarrow \sum_i \lambda_i y_i = 0 \\ \frac{\partial \mathcal{L}}{\partial \varepsilon_i} = 0 \rightarrow C - \lambda_i - \mu_i = 0 \rightarrow C = \lambda_i + \mu_i \end{array} \right.$$

$(\lambda_i, \mu_i \leq C)$

Trick 1: **Soft-margin** SVM Primal->Dual

■ Inequality Constraints:

Karush-Kuhn-Tucker (KKT) conditions

1) Define the Lagrangian function:

$$\mathcal{L}(\mathbf{w}, b, \boldsymbol{\varepsilon}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{w}, b, \boldsymbol{\varepsilon}) - \sum_i \lambda_i g_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) + \mu_i h_i(\mathbf{w}, b, \boldsymbol{\varepsilon}), \text{ where } \lambda_i, \mu_i \geq 0, \forall i$$

$$f(\mathbf{w}, b, \boldsymbol{\varepsilon}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \varepsilon_i, \quad g_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = y_i(\mathbf{w}^T \mathbf{x}_i + b) + \varepsilon_i - 1 \geq 0, \quad h_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = \varepsilon_i \geq 0$$

2) KKT conditions for optimality:

✓ Stationarity: $\mathbf{w}^* = \sum_i \lambda_i y_i \mathbf{x}_i$, $\sum_i \lambda_i y_i = 0$, $C = \lambda_i + \mu_i$ ($\lambda_i, \mu_i \leq C$)

✓ Primal Feasibility: $g_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = y_i(\mathbf{w}^T \mathbf{x}_i + b) + \varepsilon_i - 1 \geq 0$, $h_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = \varepsilon_i \geq 0$

✓ Dual Feasibility: $\lambda_i, \mu_i \geq 0, \forall i$

✓ Complementary Slackness: $\lambda_i g_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = 0$, $\mu_i h_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = 0, \forall i$



If $g_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) > 0$ (inactive), $\lambda_i = 0$

If $g_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = 0$ (active), $\lambda_i > 0$

If $h_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) > 0$ (inactive), $\mu_i = 0$

If $h_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = 0$ (active), $\mu_i > 0$

How to
interpret?

Trick 1: **Soft-margin** SVM Primal->Dual

■ Interpretation of Complementary Slackness

$$\lambda_i g_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = 0,$$



If $g_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) > 0$ (inactive), $\lambda_i = 0$

If $g_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = 0$ (active), $\lambda_i > 0$

$$\mu_i h_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = 0, \forall i$$



If $h_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) > 0$ (inactive), $\mu_i = 0$

If $h_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = 0$ (active), $\mu_i > 0$

$$C = \lambda_i + \mu_i$$

(a) When $\lambda_i = 0$, $\mu_i = C > 0$

$$g_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = y_i(\mathbf{w}^T \mathbf{x}_i + b) + \varepsilon_i - 1 > 0$$

$$h_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = \varepsilon_i = 0$$

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) > 1$$

Correctly classified and outside the margin

(b) When $\lambda_i > 0$,

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1 - \varepsilon_i$$

Case 1: $\varepsilon_i = 0, \mu_i > 0$

On the margin boundaries

Case 2: $0 < \varepsilon_i \leq 1, \mu_i = 0$

Correctly classified, inside the margin

Case 3: $\varepsilon_i > 1, \mu_i = 0$

Misclassified

Trick 1: **Soft-margin** SVM Primal->Dual

$$\min_{\mathbf{w}, b, \boldsymbol{\varepsilon}} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\varepsilon}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \min_{\mathbf{w}, b, \boldsymbol{\varepsilon}} f(\mathbf{w}, b, \boldsymbol{\varepsilon}) - \sum_i \lambda_i g_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) + \mu_i h_i(\mathbf{w}, b, \boldsymbol{\varepsilon}), \text{ where } \lambda_i, \mu_i \geq 0, \forall i$$

$$= \min_{\mathbf{w}, b, \boldsymbol{\varepsilon}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \varepsilon_i - \sum_i \lambda_i (y_i (\mathbf{w}^T \mathbf{x}_i + b) + \varepsilon_i - 1) + \mu_i \varepsilon_i$$

$$= \min_{\mathbf{w}, b, \boldsymbol{\varepsilon}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \varepsilon_i - \sum_i \lambda_i y_i \mathbf{w}^T \mathbf{x}_i - \sum_i \lambda_i y_i b - \sum_i \lambda_i \varepsilon_i + \sum_i \lambda_i - \sum_i \mu_i \varepsilon_i$$

$$= \frac{1}{2} \mathbf{w}^{*T} \mathbf{w}^* - \mathbf{w}^{*T} \mathbf{w}^* + \sum_i \lambda_i$$

$$= -\frac{1}{2} \sum_i \lambda_i y_i \mathbf{x}_i^T \sum_i \lambda_i y_i \mathbf{x}_i + \sum_i \lambda_i$$

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

$$f(\mathbf{w}, b, \boldsymbol{\varepsilon}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \varepsilon_i,$$

$$g_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = y_i (\mathbf{w}^T \mathbf{x}_i + b) + \varepsilon_i - 1 \geq 0,$$

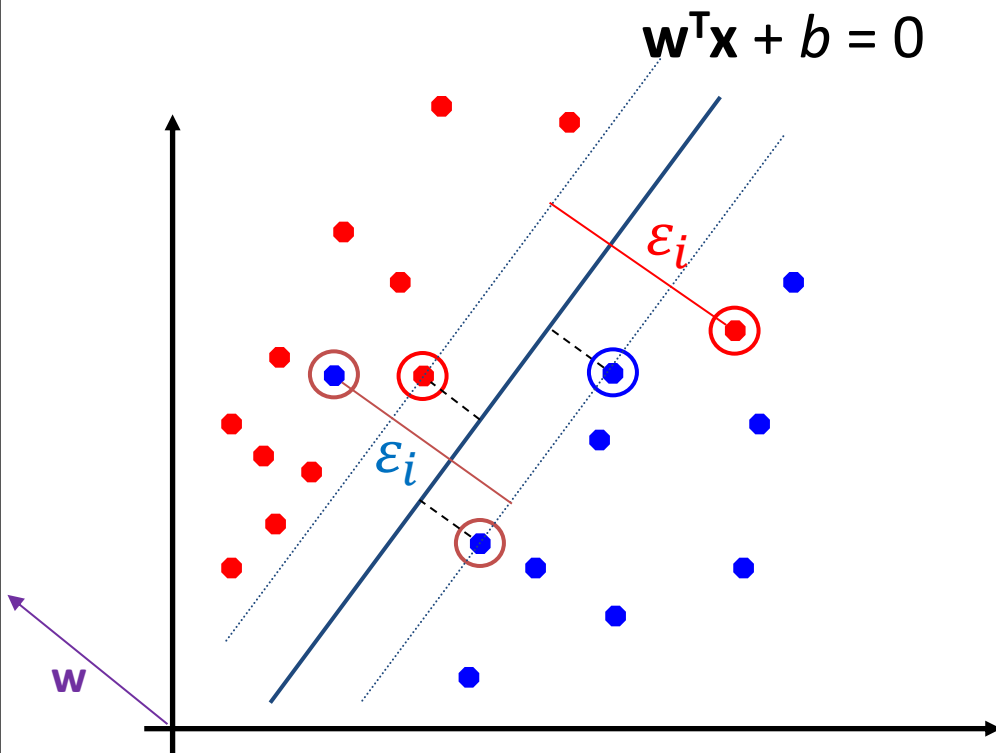
$$h_i(\mathbf{w}, b, \boldsymbol{\varepsilon}) = \varepsilon_i \geq 0$$

$$\left[\begin{array}{l} \mathbf{w}^* = \sum_i \lambda_i y_i \mathbf{x}_i \\ \sum_i \lambda_i y_i = 0 \\ C = \lambda_i + \mu_i \end{array} \right.$$

**Soft-margin
Dual**

$$\begin{array}{ll} \max_{\boldsymbol{\lambda}} & -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_i \lambda_i \\ \text{s.t.} & 0 \leq \lambda_i \leq C, \sum_i \lambda_i y_i = 0, \forall i \end{array}$$

Soft-margin SVM: Decision Boundary



$$\begin{cases} \mathbf{w}^* = \sum_i \lambda_i y_i \mathbf{x}_i \\ \sum_i \lambda_i y_i = 0 \end{cases}$$

$$\mathbf{w}^{*T} \mathbf{x} + b^* = \sum_i \lambda_i y_i \mathbf{x}_i^T \mathbf{x} + b^* = 0$$

for any support vector \mathbf{x}_s ($\lambda_i > 0$),

$$y_s(\mathbf{w}^{*T} \mathbf{x}_s + b^*) = 1 - \epsilon_s$$

↓

$$y_s^2(\mathbf{w}^{*T} \mathbf{x}_s + b^*) = y_s(1 - \epsilon_s)$$

↓

$$b^* = y_s(1 - \epsilon_s) - \sum_i \lambda_i y_i \mathbf{x}_i^T \mathbf{x}_s$$

Quick Question

For Soft-margin SVM with the following objective function:

$$\min_{\mathbf{w}, b, \varepsilon} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \varepsilon_i$$

or

$$\max_{\lambda} -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_i \lambda_i$$

$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \varepsilon_i, \varepsilon_i \geq 0, \forall i$$

$$\text{s.t. } 0 \leq \lambda_i \leq C, \sum_i \lambda_i y_i = 0, \forall i$$

Which of the following data points are considered support vectors?

- A. Data points that lie exactly on the margin boundaries.
- B. Data points that have $\lambda_i > 0$.
- C. Data points that determine the decision boundary
- D. Data points inside the margin
- E. Data points that are misclassified
- F. Data points that have $\varepsilon_i \geq 0$

Soft-margin VS Hard-margin SVMs

Hard-margin Primal:

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

$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \forall i$$

Optimize over $d+1$ variables
 N constraints

Hard-margin Dual:

$$\max_{\lambda} -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_i \lambda_i$$

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

Optimize over N variables
 $N+1$ constraints

Soft-margin Primal:

$$\min_{\mathbf{w}, b, \varepsilon} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \varepsilon_i$$

$$\text{s.t. } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \varepsilon_i, \varepsilon_i \geq 0, \forall i$$

Optimize over $d+1+N$ variables
 $2N$ constraints

Soft-margin Dual:

$$\max_{\lambda} -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_i \lambda_i$$

$$\text{s.t. } 0 \leq \lambda_i \leq C, \sum_i \lambda_i y_i = 0, \forall i$$

Optimize over N variables
 $(2)N+1$ constraints

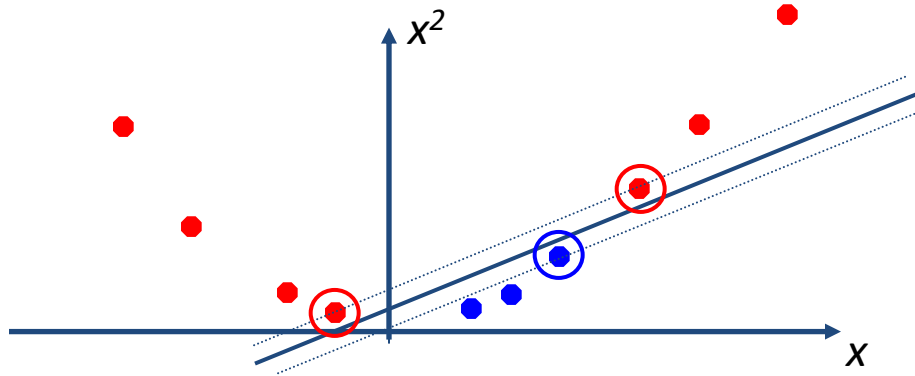
Kernel SVM

Trick 2: Kernel SVM

- Here's a linearly non-separable dataset in 1-D space



- We can map data x from 1-D to 2-D by $x \rightarrow (x, x^2)$

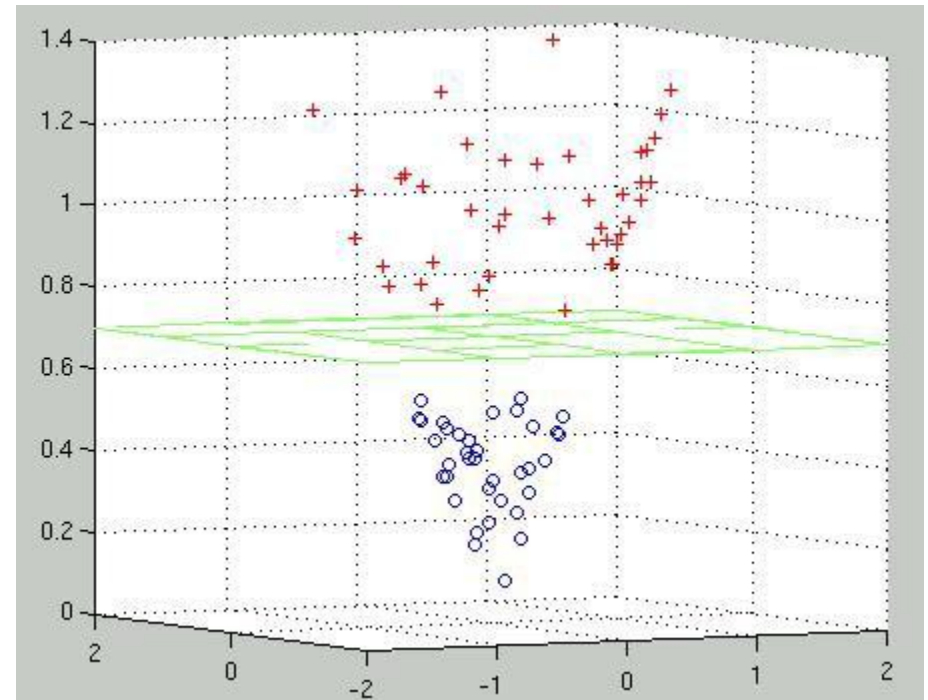
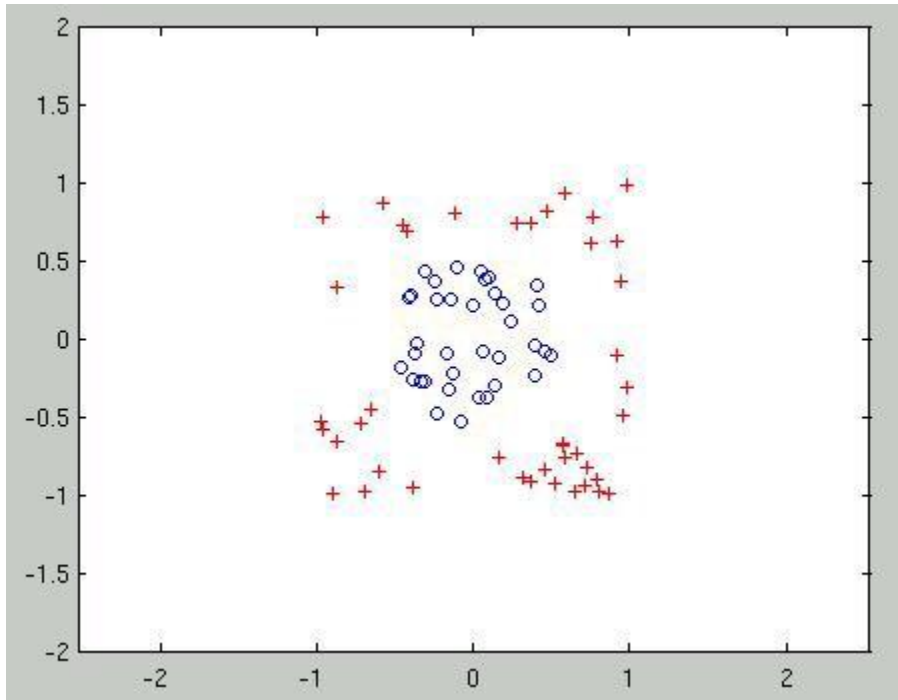


- ❖ Now the data is linearly separable in the new space
- ❖ We can run SVM in the new space
- ❖ The linear decision boundary in the new space corresponds to a non-linear decision boundary in the old space

Trick 2: Kernel SVM

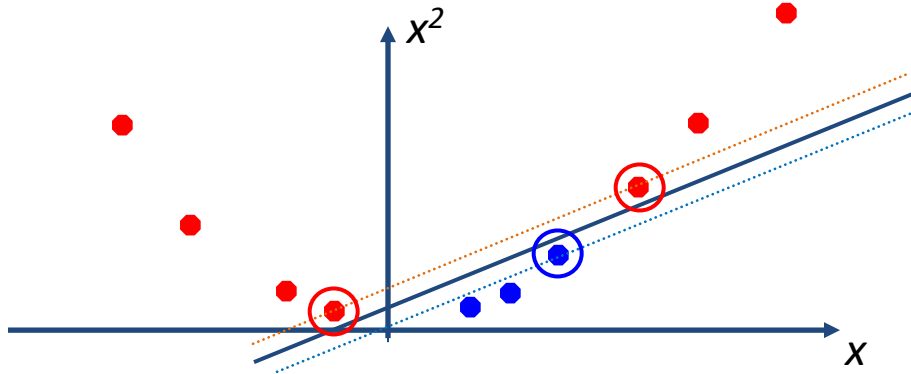
Another Example

$$(x_1, x_2) \Rightarrow (x_1, x_2, \sqrt{x_1^2 + x_2^2})$$



Trick 2: Kernel SVM

- **General idea:** map the original feature space $\mathbf{x} = (x_1, \dots, x_d)^T$ to some higher-dimensional feature space $\Phi(\mathbf{x})$ where the data is separable



Decision Boundary:

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

Plus Plane:

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

Minus Plane:

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

Trick 2: Kernel SVM

- Optimization Formulation after mapping

Hard-margin Primal:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ \text{s.t.} \quad & y_i (\mathbf{w}^T \Phi(\mathbf{x}_i) + b) \geq 1, \forall i \end{aligned}$$

Hard-margin Dual:

$$\begin{aligned} \max_{\lambda} \quad & -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j) + \sum_i \lambda_i \\ \text{s.t.} \quad & \lambda_i \geq 0, \sum_i \lambda_i y_i = 0, \forall i \end{aligned}$$

Soft-margin Primal:

$$\begin{aligned} \min_{\mathbf{w}, b, \varepsilon} \quad & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \varepsilon_i \\ \text{s.t.} \quad & y_i (\mathbf{w}^T \Phi(\mathbf{x}_i) + b) \geq 1 - \varepsilon_i, \varepsilon_i \geq 0, \forall i \end{aligned}$$

Soft-margin Dual:

$$\begin{aligned} \max_{\lambda} \quad & -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j) + \sum_i \lambda_i \\ \text{s.t.} \quad & 0 \leq \lambda_i \leq C, \sum_i \lambda_i y_i = 0, \forall i \end{aligned}$$

Trick 2: Kernel SVM

- Optimal decision boundary after mapping

Hard-margin SVM:

$$\mathbf{w}^{*T} \Phi(\mathbf{x}) + b^* = \sum_i \lambda_i y_i \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}) + b^* = 0$$

where $b^* = y_s - \sum_i \lambda_i y_i \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_s)$
for any support vector \mathbf{x}_s ($\lambda_i > 0$)

$$\begin{cases} \mathbf{w}^* = \sum_i \lambda_i y_i \Phi(\mathbf{x}_i) \\ \sum_i \lambda_i y_i = 0 \end{cases}$$

Soft-margin SVM:

$$\mathbf{w}^{*T} \Phi(\mathbf{x}) + b^* = \sum_i \lambda_i y_i \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}) + b^* = 0$$

where $b^* = y_s(1 - \varepsilon_s) - \sum_i \lambda_i y_i \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_s)$
for any support vector \mathbf{x}_s ($\lambda_i > 0$)

Trick 2: Kernel SVM

- The solution relies on the inner product $\Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j)$
- A kernel function K : $K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j)$

Hard-margin Dual:

$$\begin{aligned} \max_{\lambda} & -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) + \sum_i \lambda_i \\ \text{s.t. } & \lambda_i \geq 0, \sum_i \lambda_i y_i = 0, \forall i \end{aligned}$$

$$\mathbf{w}^{*T} \Phi(\mathbf{x}) + b^* = \sum_i \lambda_i y_i K(\mathbf{x}_i, \mathbf{x}) + b^* = 0$$

$$\text{where } b^* = y_s - \sum_i \lambda_i y_i K(\mathbf{x}_i, \mathbf{x}_s)$$

for any support vector \mathbf{x}_s ($\lambda_i > 0$)

Soft-margin Dual:

$$\begin{aligned} \max_{\lambda} & -\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) + \sum_i \lambda_i \\ \text{s.t. } & 0 \leq \lambda_i \leq C, \sum_i \lambda_i y_i = 0, \forall i \end{aligned}$$

$$\mathbf{w}^{*T} \Phi(\mathbf{x}) + b^* = \sum_i \lambda_i y_i K(\mathbf{x}_i, \mathbf{x}) + b^* = 0$$

$$\text{where } b^* = y_s(1 - \varepsilon_s) - \sum_i \lambda_i y_i K(\mathbf{x}_i, \mathbf{x}_s)$$

for any support vector \mathbf{x}_s ($\lambda_i > 0$)

Examples of Kernel Functions

- Linear: $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
Mapping Φ : $\mathbf{x} \rightarrow \Phi(\mathbf{x})$, where $\Phi(\mathbf{x})$ is \mathbf{x} itself
- Polynomial of power p : $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^p$
Mapping Φ : $\mathbf{x} \rightarrow \Phi(\mathbf{x})$, where $\Phi(\mathbf{x})$ has $\binom{d+p}{p} = \frac{(d+p)!}{p! \cdot d!}$ dimensions
- Gaussian (radial-basis function): $K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$
Mapping Φ : $\mathbf{x} \rightarrow \Phi(\mathbf{x})$, where $\Phi(\mathbf{x})$ is *infinite-dimensional*: every point is mapped to *a function* (a Gaussian)

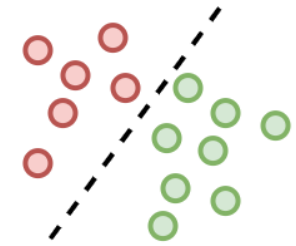


Right Kernel for SVM?

Right Kernel for SVM

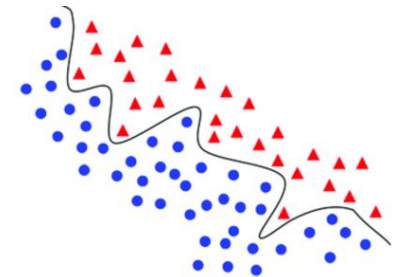
- Linear: $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$

When to Use? ✓ The data is linearly separable



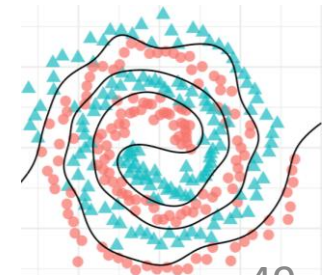
- Polynomial : $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^p$ (need to choose the right p)

When to Use? ✓ Nonlinear and polynomial-like decision boundary
✓ Capture feature interactions
✓ Low-dimensional data



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

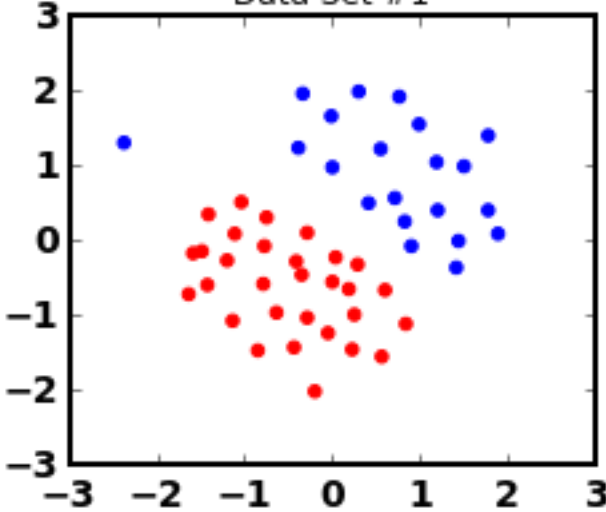
When to Use? ✓ Highly complex and nonlinear decision boundary
✓ Unknown underlying feature transformations
✓ Moderate-sized dataset



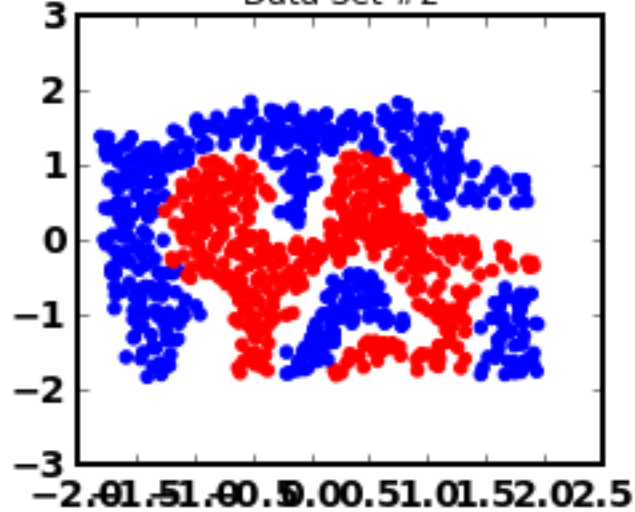
Right Kernel for SVM

Which Kernel to Use?

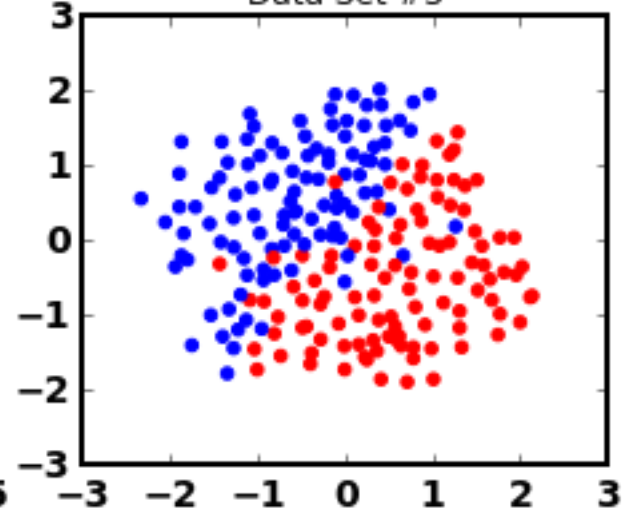
Data Set #1



Data Set #2



Data Set #3



What Functions Can be Kernels?

- Can we always find out a K (kernel) for a given Φ ?
- Given a K , can we always find out a Φ such that $K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j)$?

What Functions Can be Kernels?

- Mercer's Theorem

For $K(\mathbf{x}_i, \mathbf{x}_j)$ to be a valid kernel, it must satisfy:

- i. Symmetry: $K(\mathbf{x}_i, \mathbf{x}_j) = K(\mathbf{x}_j, \mathbf{x}_i)$
- ii. The kernel matrix \mathbf{G} : $\mathbf{v}^T \mathbf{G} \mathbf{v} \geq 0$, for any non-zero vector \mathbf{v}

$$\mathbf{G} = \begin{bmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & \cdots & K(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ K(\mathbf{x}_N, \mathbf{x}_1) & \cdots & K(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix} \quad (\text{Positive Semi-Definiteness})$$



What Functions Can be Kernels?

- Mercer's Theorem

For $K(\mathbf{x}_i, \mathbf{x}_j)$ to be a valid kernel, it must satisfy:

- Symmetry: $K(\mathbf{x}_i, \mathbf{x}_j) = K(\mathbf{x}_j, \mathbf{x}_i)$ $\Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j) = \Phi(\mathbf{x}_j)^T \Phi(\mathbf{x}_i)$
- The kernel matrix \mathbf{G} : $\mathbf{v}^T \mathbf{G} \mathbf{v} \geq 0$, for any non-zero vector \mathbf{v}

$$\mathbf{G} = \begin{bmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & \cdots & K(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ K(\mathbf{x}_N, \mathbf{x}_1) & \cdots & K(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix} \quad (\text{Positive Semi-Definiteness})$$

$$\begin{aligned} \mathbf{v}^T \mathbf{G} \mathbf{v} &= [v_1 \quad \cdots \quad v_N] \begin{bmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & \cdots & K(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ K(\mathbf{x}_N, \mathbf{x}_1) & \cdots & K(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix} \begin{bmatrix} v_1 \\ \vdots \\ v_N \end{bmatrix} \\ &= \left[\sum_{i=1}^N v_i K(\mathbf{x}_i, \mathbf{x}_1) \quad \cdots \quad \sum_{i=1}^N v_i K(\mathbf{x}_i, \mathbf{x}_N) \right] \begin{bmatrix} v_1 \\ \vdots \\ v_N \end{bmatrix} = \sum_{j=1}^N \sum_{i=1}^N v_i K(\mathbf{x}_i, \mathbf{x}_j) v_j \end{aligned}$$

What Functions Can be Kernels?

$$\mathbf{v}^T \mathbf{G} \mathbf{v} = \sum_{j=1}^N \sum_{i=1}^N v_i K(\mathbf{x}_i, \mathbf{x}_j) v_j$$

If K is a valid kernel, $K(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\Phi}(\mathbf{x}_i)^T \boldsymbol{\Phi}(\mathbf{x}_j)$

$$\mathbf{v}^T \mathbf{G} \mathbf{v} = \sum_{j=1}^N \sum_{i=1}^N v_i K(\mathbf{x}_i, \mathbf{x}_j) v_j = \sum_{j=1}^N \sum_{i=1}^N v_i \boldsymbol{\Phi}(\mathbf{x}_i)^T \boldsymbol{\Phi}(\mathbf{x}_j) v_j$$

Because of the Linearity of the inner product:

$$\left[\begin{aligned} (a\boldsymbol{\Phi}(\mathbf{x}_i))^T \boldsymbol{\Phi}(\mathbf{x}_j) &= a\boldsymbol{\Phi}(\mathbf{x}_i)^T \boldsymbol{\Phi}(\mathbf{x}_j) \\ (\boldsymbol{\Phi}(\mathbf{x}_i) + \boldsymbol{\Phi}(\mathbf{x}_k))^T \boldsymbol{\Phi}(\mathbf{x}_j) &= \boldsymbol{\Phi}(\mathbf{x}_i)^T \boldsymbol{\Phi}(\mathbf{x}_j) + \boldsymbol{\Phi}(\mathbf{x}_k)^T \boldsymbol{\Phi}(\mathbf{x}_j) \end{aligned} \right]$$

$$\mathbf{v}^T \mathbf{G} \mathbf{v} = \sum_{j=1}^N \sum_{i=1}^N v_i \boldsymbol{\Phi}(\mathbf{x}_i)^T \boldsymbol{\Phi}(\mathbf{x}_j) v_j = \sum_{j=1}^N \left(\sum_{i=1}^N v_i \boldsymbol{\Phi}(\mathbf{x}_i) \right)^T \boldsymbol{\Phi}(\mathbf{x}_j) v_j = \left(\sum_{i=1}^N v_i \boldsymbol{\Phi}(\mathbf{x}_i) \right)^T \left(\sum_{j=1}^N \boldsymbol{\Phi}(\mathbf{x}_j) v_j \right)$$

Because of the Positive Semi-definiteness of the inner product: $\boldsymbol{\Phi}(\mathbf{x}_i)^T \boldsymbol{\Phi}(\mathbf{x}_i) \geq 0$

$$\mathbf{v}^T \mathbf{G} \mathbf{v} = \left(\sum_{i=1}^N v_i \boldsymbol{\Phi}(\mathbf{x}_i) \right)^T \left(\sum_{j=1}^N \boldsymbol{\Phi}(\mathbf{x}_j) v_j \right) \geq 0$$

What Functions Can be Kernels?

- Can we always find out a K (kernel) for a given Φ ?

Yes, such K always satisfies Mercer's conditions

- Given a K , can we always find out a Φ such that $K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j)$?

No, only for K that satisfies Mercer's conditions

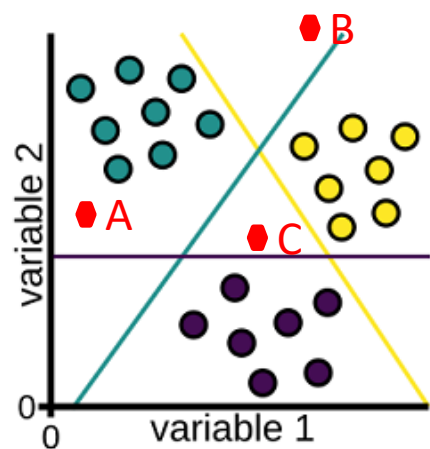
Multi-class SVM

SVM: more than two classes

- N-class problem

1. One-vs-All approach

Training: train a binary SVM classifier for each class



- Classifier 1: Class 1 vs. the rest (class 2-- N)
Label: +1 Label: -1
- Classifier 2: Class 2 vs. the rest (class 1, 3-- N)
Label: +1 Label: -1
- ...
- Classifier N : Class N vs. the rest (class 1--($N-1$))
Label: +1 Label: -1

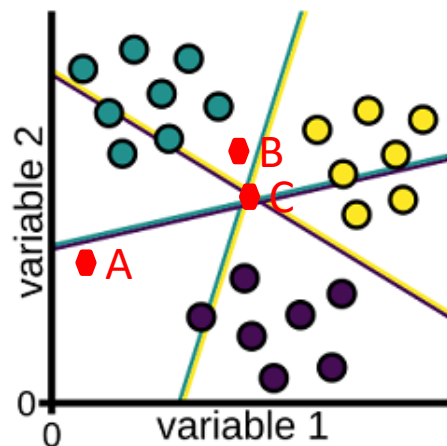
Prediction: choose the class whose classifier outputs the largest value

SVM: more than two classes

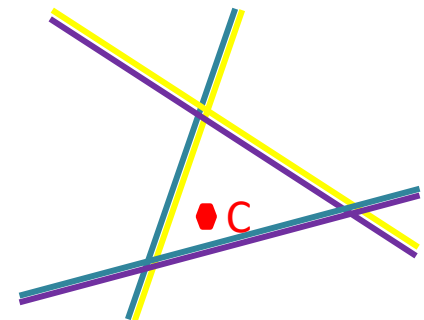
- N-class problem

2. One-vs-One approach

Training: train a binary SVM classifier for each possible pair of classes



- Classifier (1, 2): Class 1 vs. Class 2
- Classifier (1, 3): Class 1 vs. Class 3
- ...
- Classifier ($N-1$, N): Class $N-1$ vs. Class N



Prediction: Majority voting among all the classifiers

More References

- Christopher J. C. Burges
A Tutorial on Support Vector Machines for Pattern Recognition (1998).
- Kayleigh Calder, et al.
Quadratic programming
https://optimization.cbe.cornell.edu/index.php?title=Quadratic_programming
- S. Boyd and L. Vandenberghe
Convex Optimization, Chapter 5 (2004)
- John C. Platt
Sequential Minimal Optimization: A Fast Algorithm for Training Support Vector Machines (1998)