

CMSC 478

Machine Learning

KMA Solaiman
ksolaima@umbc.edu

(Adapted from Tommi Jaakkola, MIT CSAIL)

Today's topics

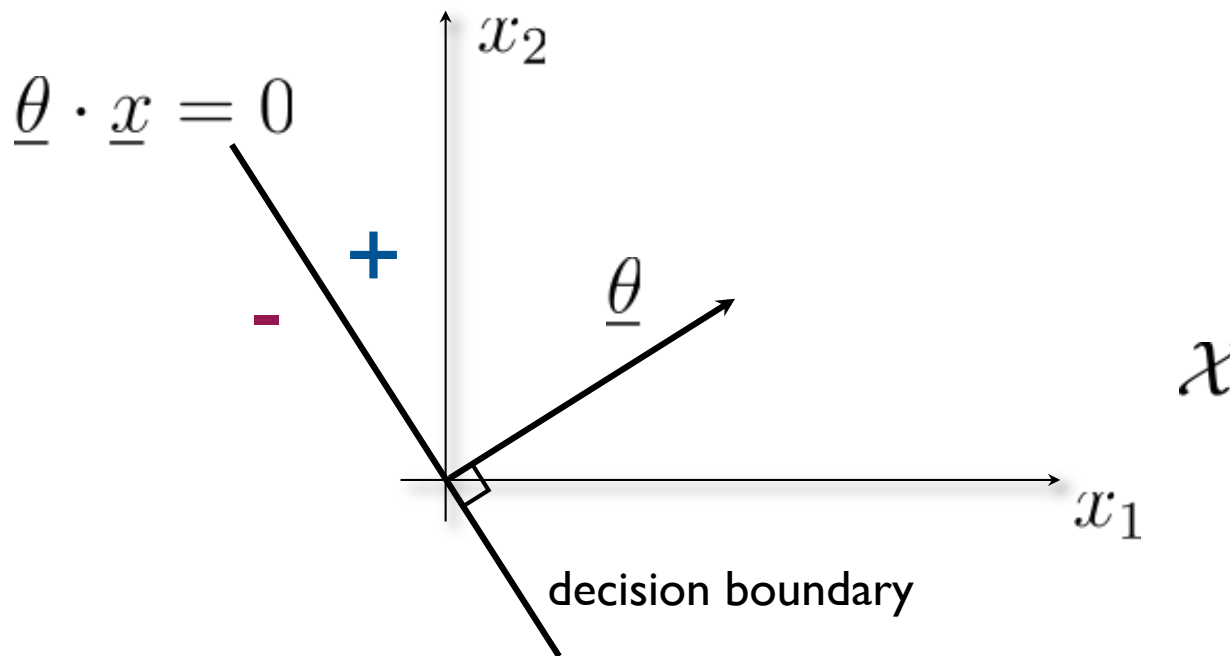
- Perceptron, convergence
 - the prediction game
 - mistakes, margin, and generalization
- Maximum margin classifier -- support vector machine
 - estimation, properties
 - allowing misclassified points

Recall: linear classifiers

- A linear classifier (through origin) with parameters $\underline{\theta}$ divides the space into positive and negative halves

$$\begin{aligned} f(\underline{x}; \underline{\theta}) &= \text{sign}(\underline{\theta} \cdot \underline{x}) = \text{sign}(\theta_1 x_1 + \dots + \theta_d x_d) \\ &= \begin{cases} +1, & \text{if } \underline{\theta} \cdot \underline{x} > 0 \\ -1, & \text{if } \underline{\theta} \cdot \underline{x} \leq 0 \end{cases} \end{aligned}$$

discriminant function



The perceptron algorithm

- A sequence of examples and labels

$$(\underline{x}_t, y_t), \quad t = 1, 2, \dots$$

- The perceptron algorithm applied to the sequence

Initialize: $\underline{\theta} = 0$

For $t = 1, 2, \dots$

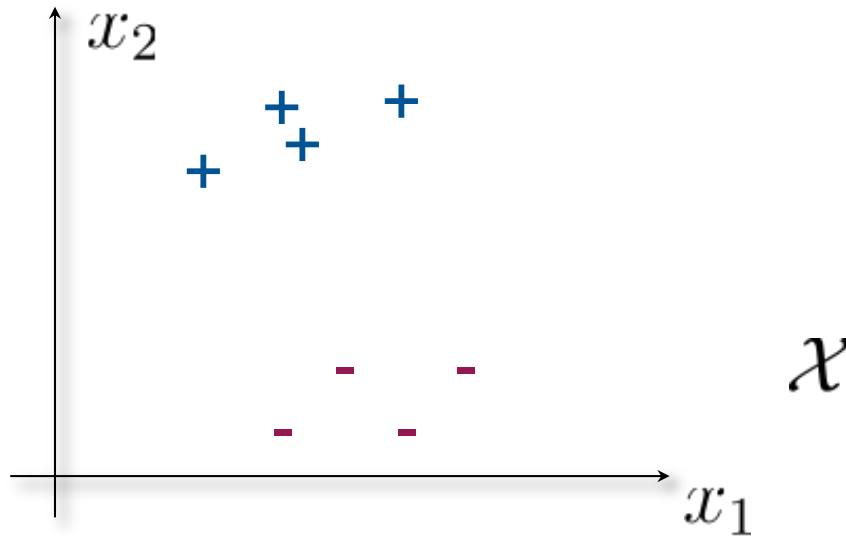
if $y_t(\underline{\theta} \cdot \underline{x}_t) \leq 0$ (mistake)

$$\underline{\theta} \leftarrow \underline{\theta} + y_t \underline{x}_t$$

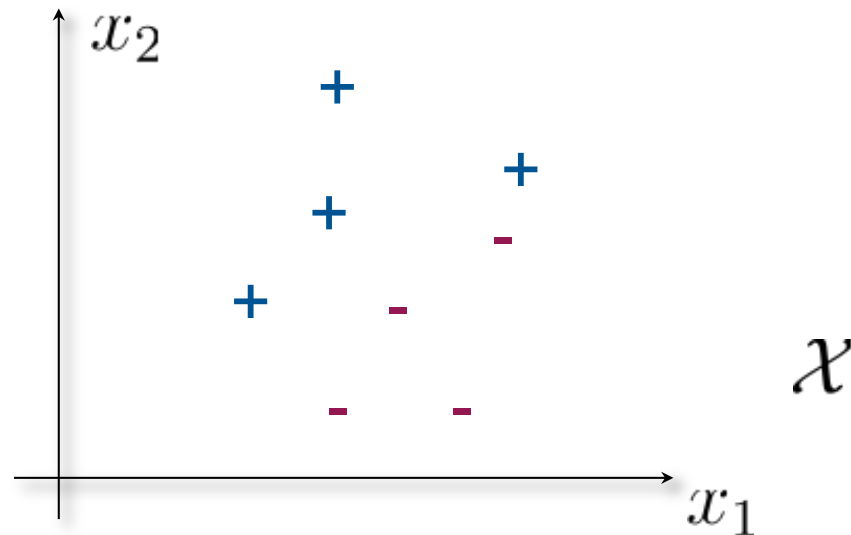
- We would like to bound the number of mistakes that the algorithm makes

Mistakes and margin

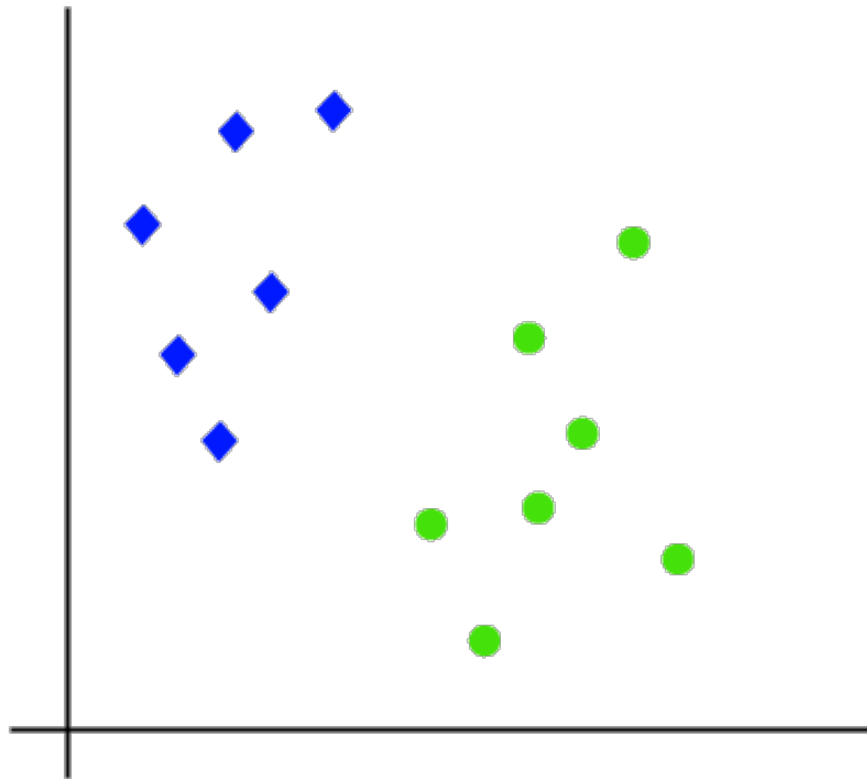
Easy problem
- large margin
- few mistakes

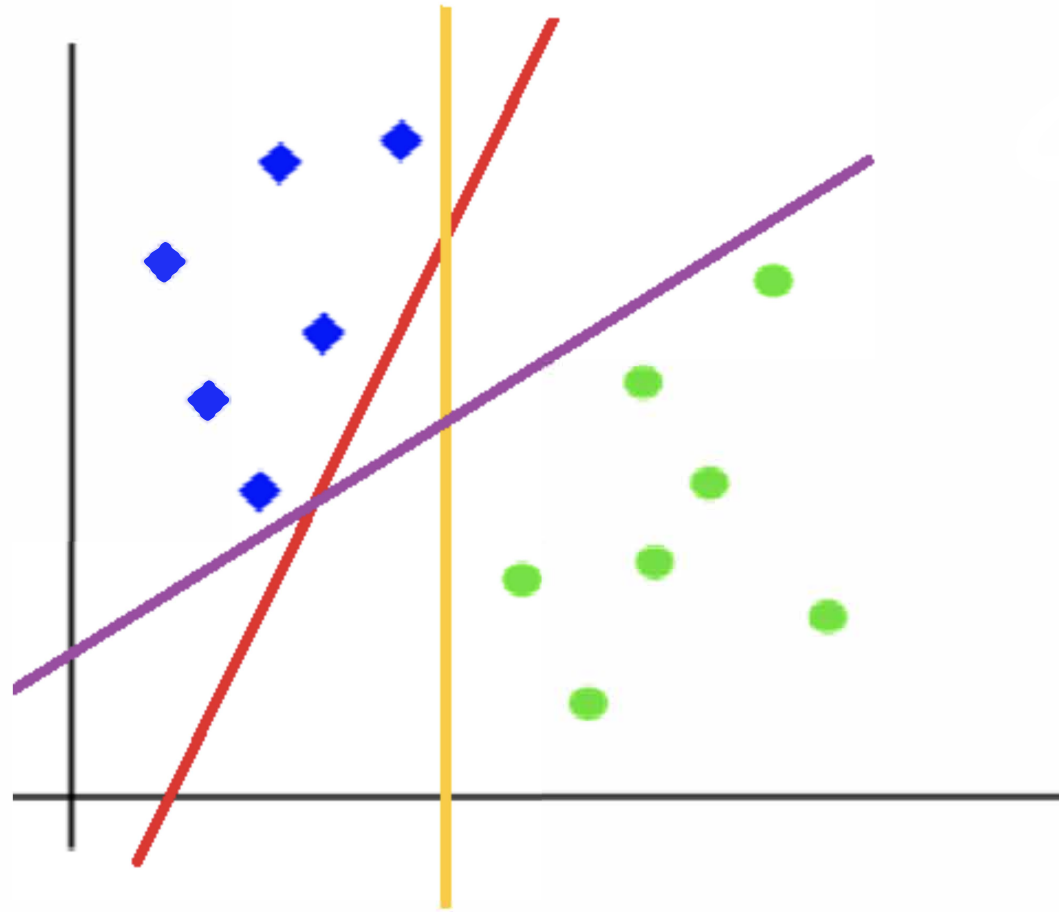


Harder problem
- small margin
- many mistakes

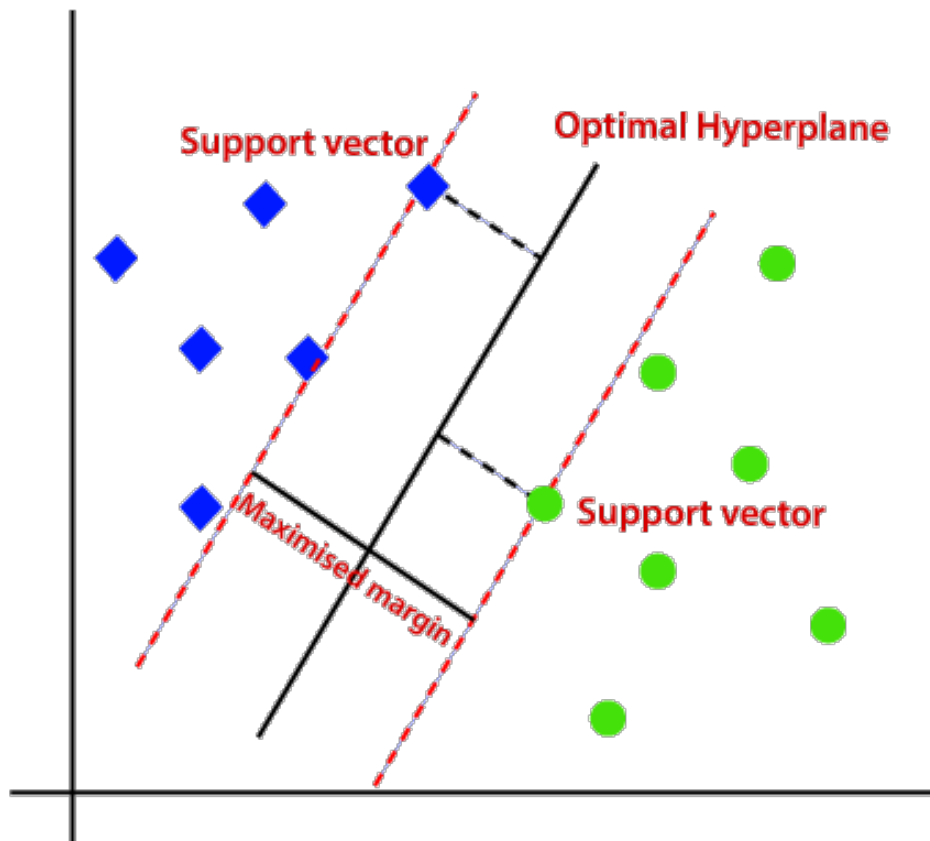


How to choose the margin?





Which
line is
best?



Margin in SVM

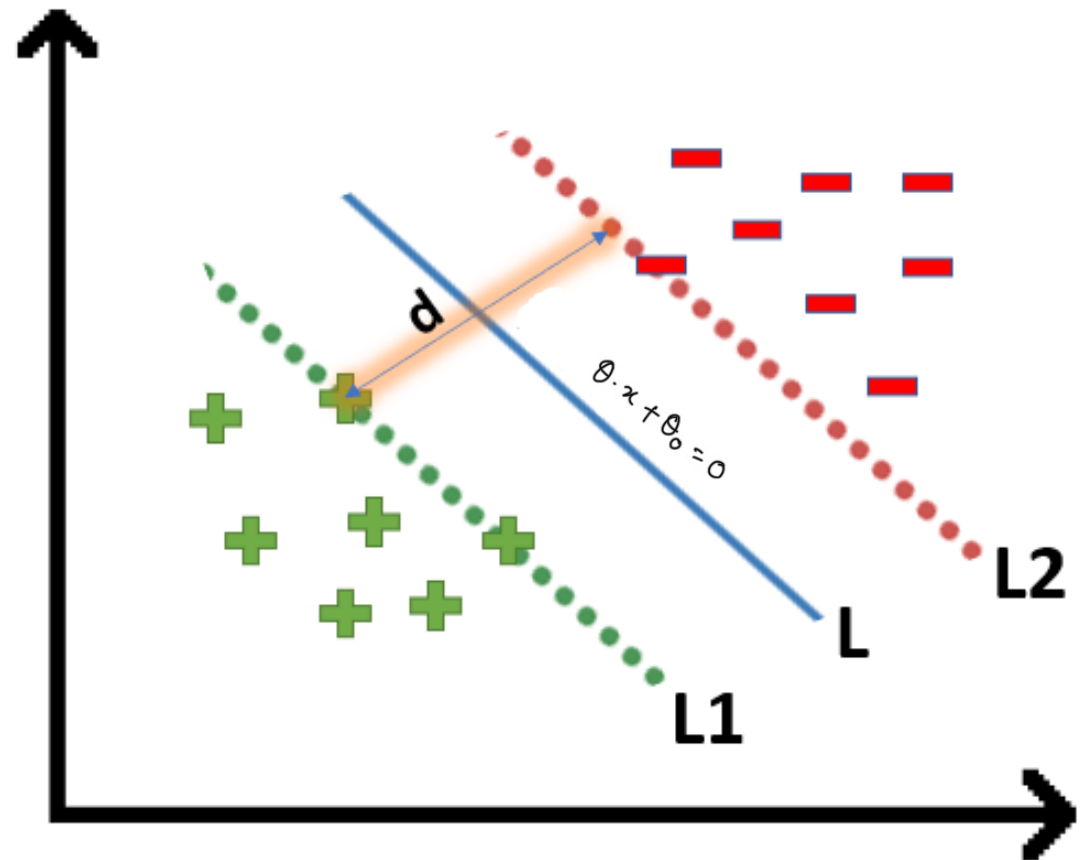
Without offset

$$y = \begin{cases} +1, & \text{if } \underline{\theta} \cdot \underline{x} > 0 \\ -1, & \text{if } \underline{\theta} \cdot \underline{x} \leq 0 \end{cases}$$

- $b = 0$
- Hyperplane through origin

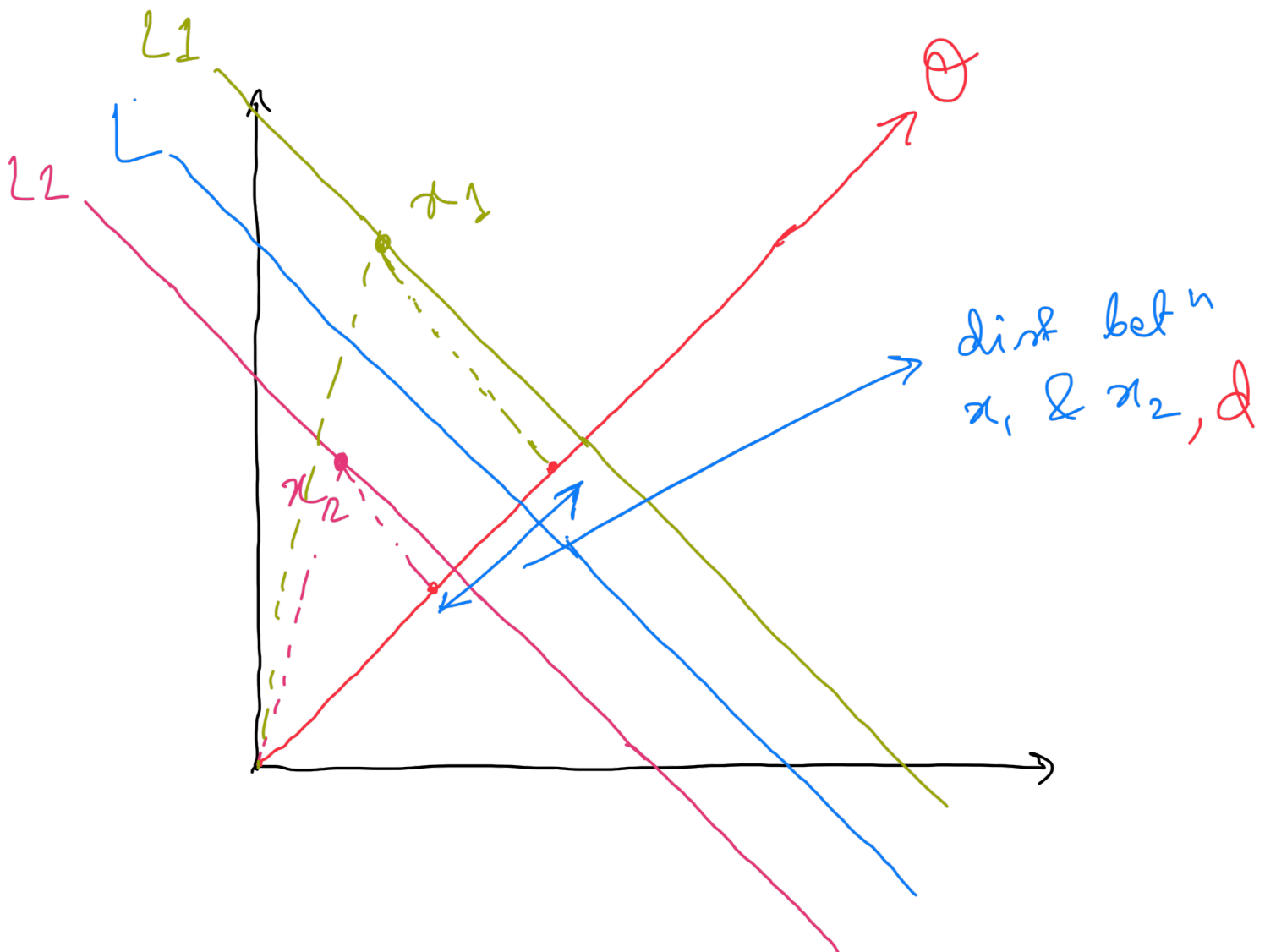
With offset

$$y = \begin{cases} +1, & \text{if } \underline{\theta} \cdot \underline{x} + \theta_0 > 0 \\ -1, & \text{if } \underline{\theta} \cdot \underline{x} + \theta_0 \leq 0 \end{cases}$$

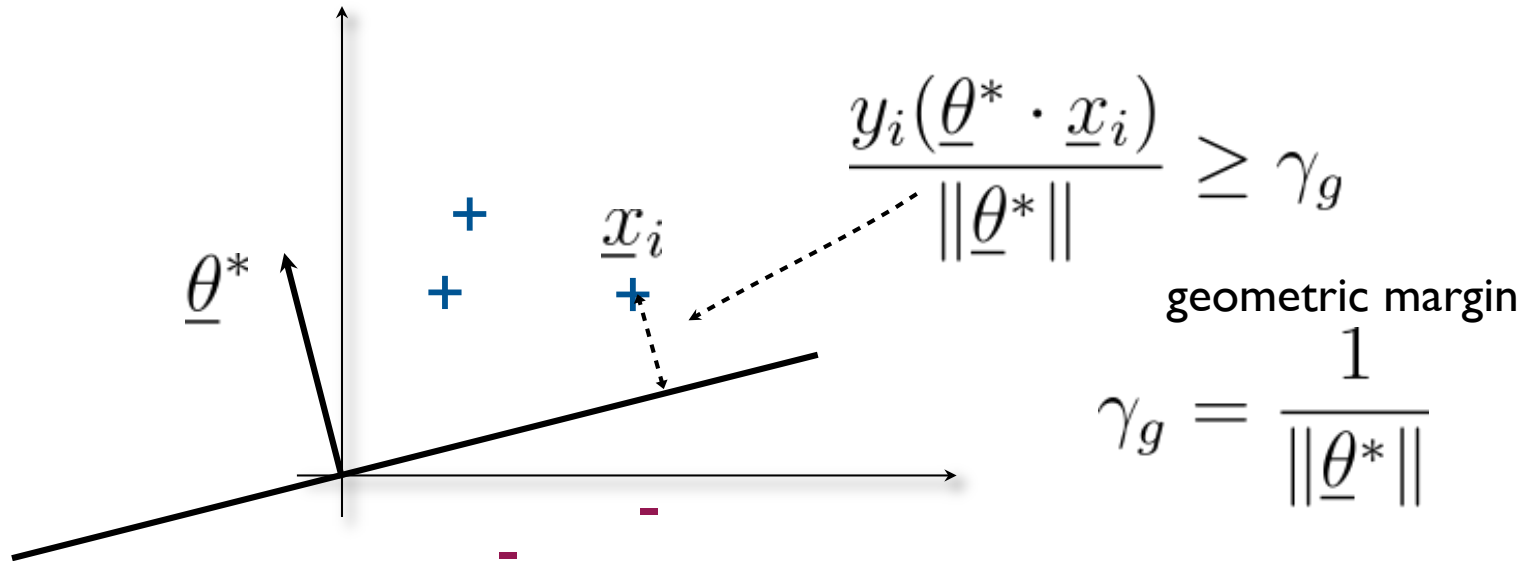


$$d = \frac{\bar{x}_2 \cdot \bar{\theta}}{\|\bar{\theta}\|} - \frac{\bar{x}_1 \cdot \bar{\theta}}{\|\bar{\theta}\|}$$

$$= \frac{\bar{x}_2 \cdot \bar{\theta} - \bar{x}_1 \cdot \bar{\theta}}{\|\bar{\theta}\|}$$



Maximum margin classifier

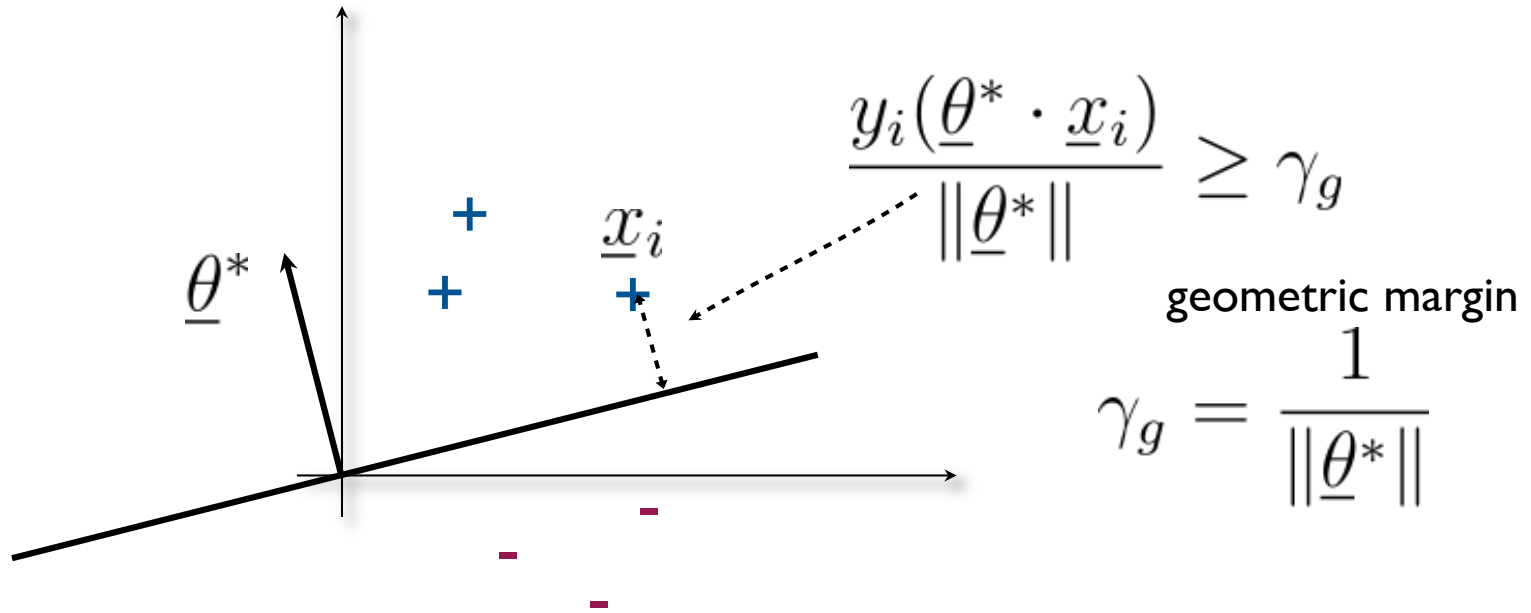


To find $\underline{\theta}^*$:

maximize $\frac{1}{\|\underline{\theta}\|}$ subject to

$$y_i(\underline{\theta} \cdot \underline{x}_i) \geq 1, \quad i = 1, \dots, n$$

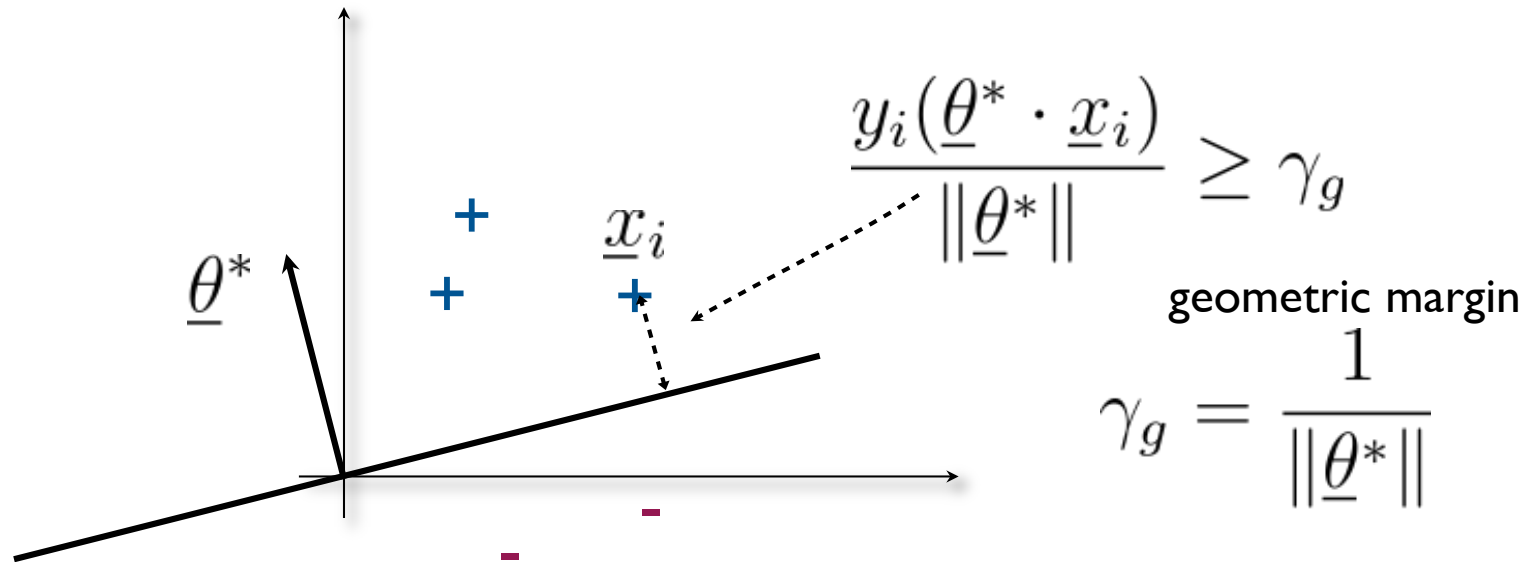
Maximum margin classifier



To find $\underline{\theta}^*$: minimize $\|\underline{\theta}\|$ subject to

$$y_i(\underline{\theta} \cdot \underline{x}_i) \geq 1, \quad i = 1, \dots, n$$

Support vector machine

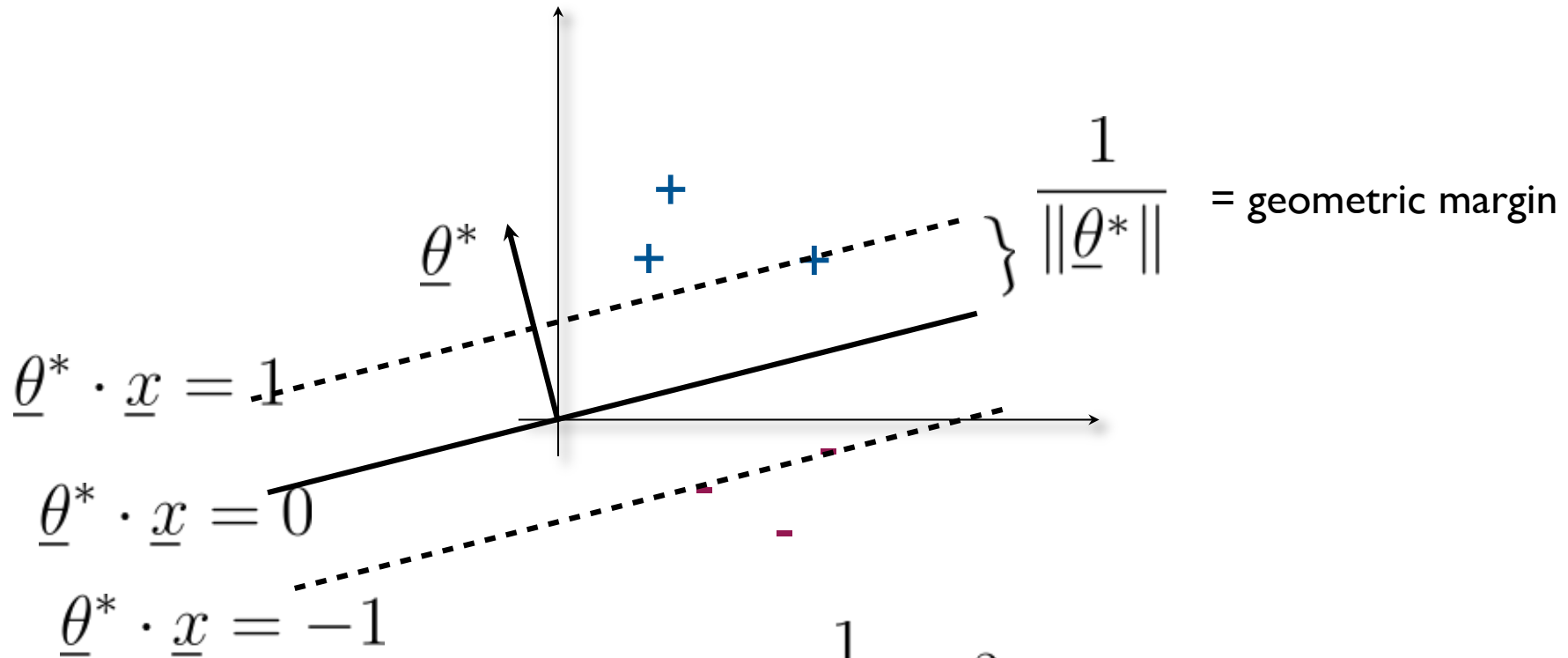


To find $\underline{\theta}^*$: minimize $\frac{1}{2} \|\underline{\theta}\|^2$ subject to

$$y_i(\underline{\theta} \cdot \underline{x}_i) \geq 1, \quad i = 1, \dots, n$$

- This is a quadratic programming problem (quadratic objective, linear constraints)
- The solution is unique, typically obtained in the dual

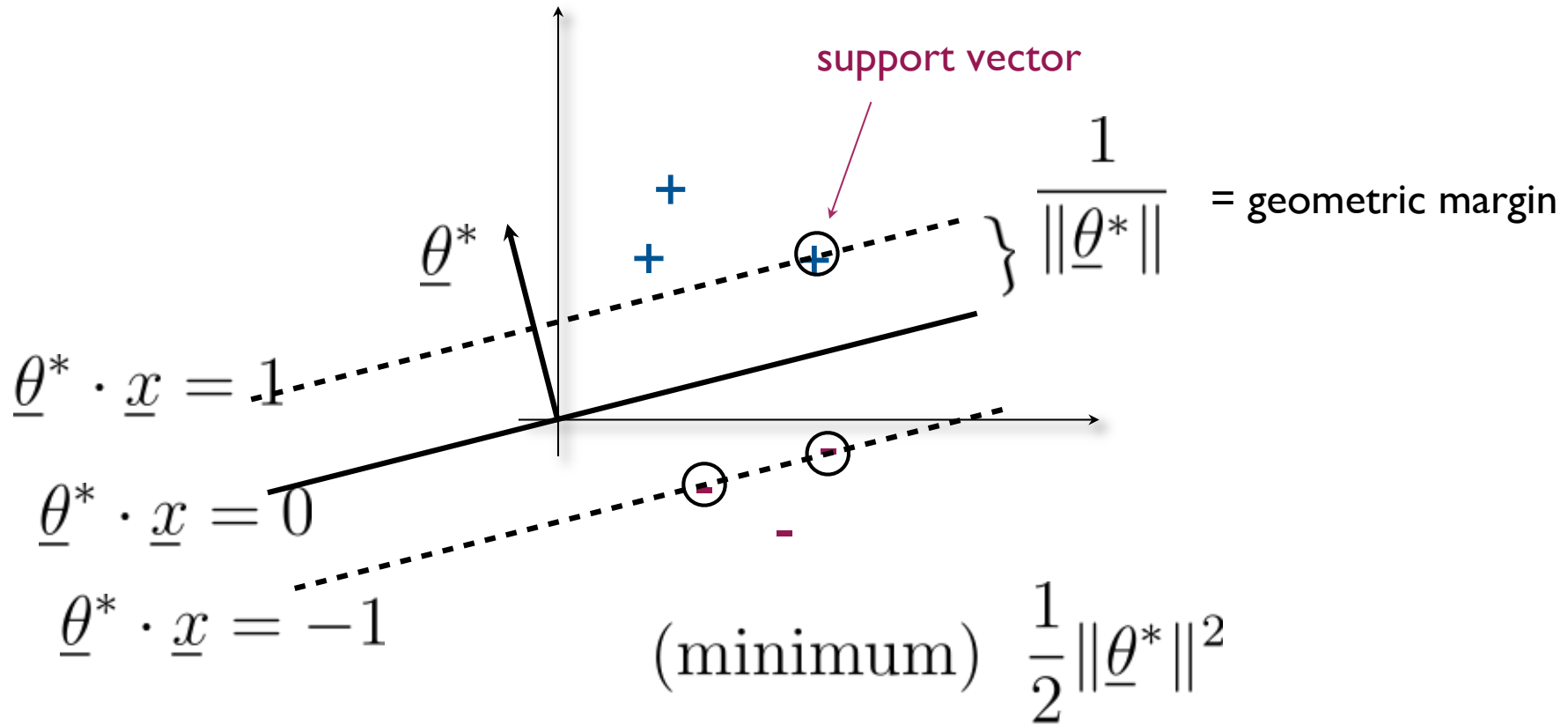
Support vector machine



To find $\underline{\theta}^*$:

minimize $\frac{1}{2} \|\underline{\theta}\|^2$ subject to
 $y_i(\underline{\theta} \cdot \underline{x}_i) \geq 1, \quad i = 1, \dots, n$

Support vector machine

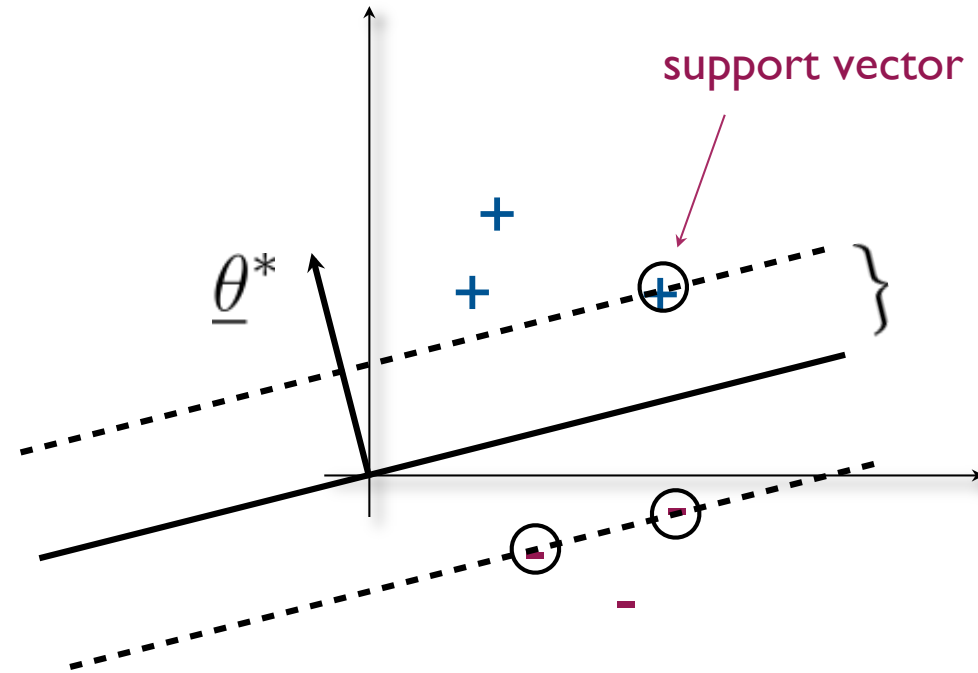


The solution is **sparse**

$$\begin{aligned}
 y_1(\underline{\theta}^* \cdot \underline{x}_1) &= 1 \\
 y_2(\underline{\theta}^* \cdot \underline{x}_2) &> 1 \\
 y_3(\underline{\theta}^* \cdot \underline{x}_3) &= 1 \\
 &\dots
 \end{aligned}$$

active constraints
 = support vectors

Is sparse solution good?



- We can simulate test performance by evaluating Leave-One-Out Cross-Validation error

$$\text{LOOCV}(\underline{\theta}^*) \leq \frac{\# \text{ of support vectors}}{n}$$

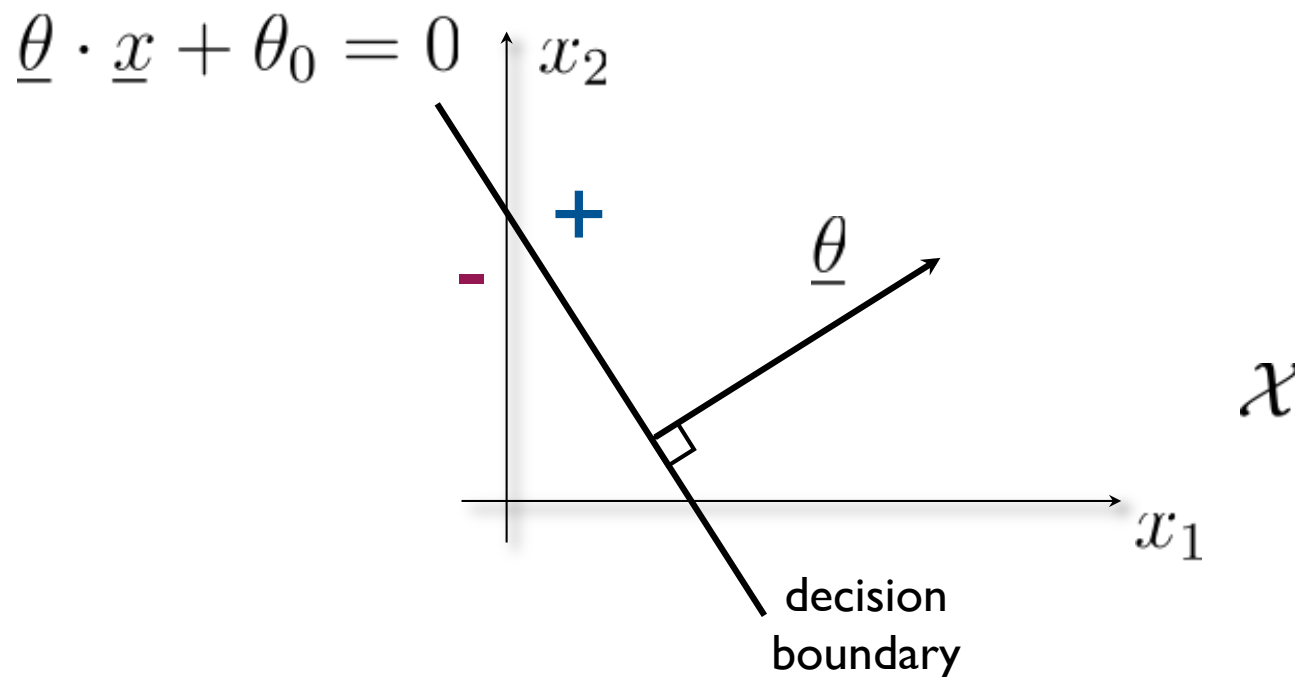
Intuitively:

if you remove the support vector from the training set, and you receive the support vector as a test point, then you would make a mistake

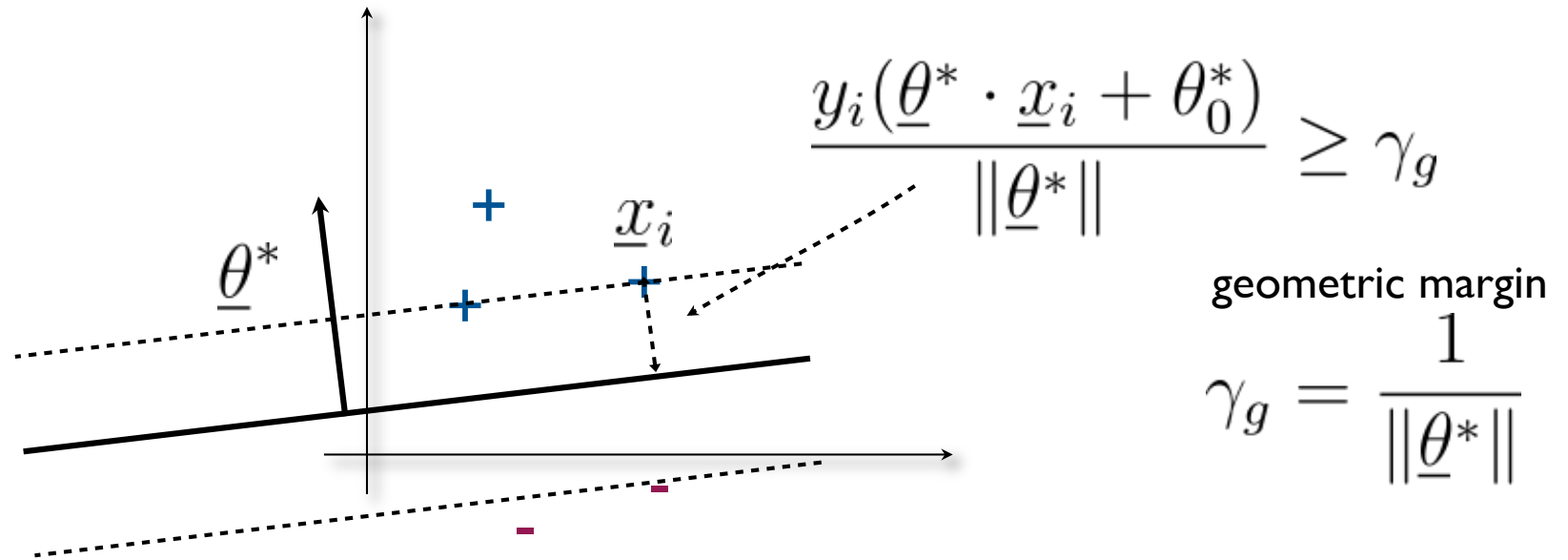
Linear classifiers (with offset)

- A linear classifier with parameters $(\underline{\theta}, \theta_0)$

$$\begin{aligned} f(\underline{x}; \underline{\theta}, \theta_0) &= \text{sign}(\underline{\theta} \cdot \underline{x} + \theta_0) \\ &= \begin{cases} +1, & \text{if } \underline{\theta} \cdot \underline{x} + \theta_0 > 0 \\ -1, & \text{if } \underline{\theta} \cdot \underline{x} + \theta_0 \leq 0 \end{cases} \end{aligned}$$



Support vector machine



$$\frac{y_i(\underline{\theta}^* \cdot \underline{x}_i + \theta_0^*)}{\|\underline{\theta}^*\|} \geq \gamma_g$$

geometric margin

$$\gamma_g = \frac{1}{\|\underline{\theta}^*\|}$$

To find $\underline{\theta}^*, \theta_0^*$:

minimize $\frac{1}{2} \|\underline{\theta}\|^2$ subject to

$$y_i(\underline{\theta} \cdot \underline{x}_i + \theta_0) \geq 1, \quad i = 1, \dots, n$$

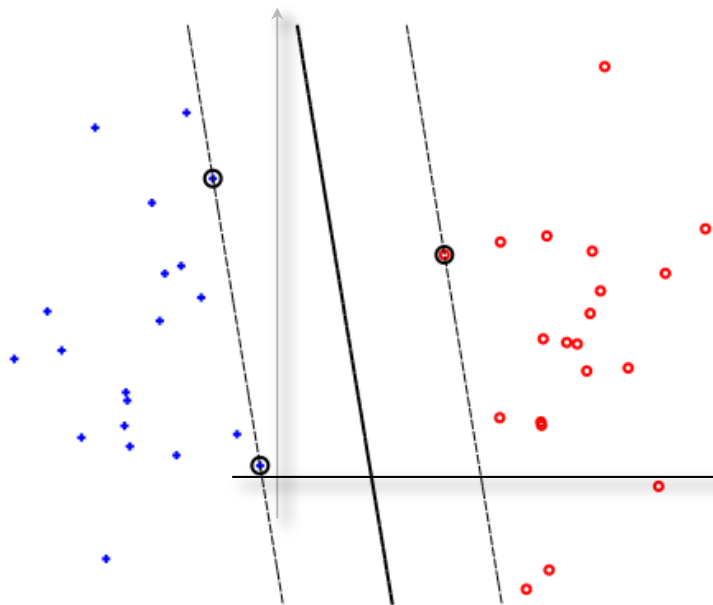
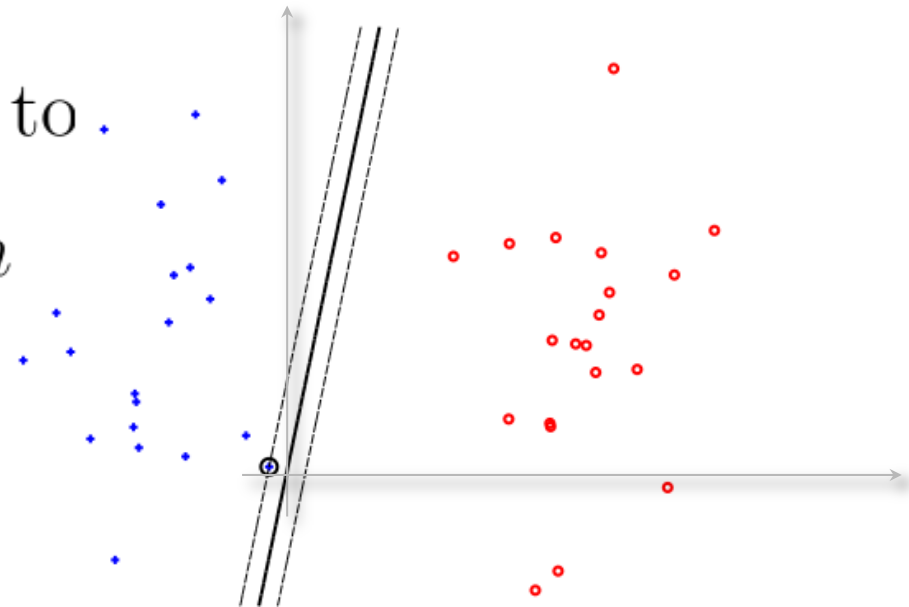
- Still a quadratic programming problem (quadratic objective, linear constraints)

The impact of offset

- Adding the offset parameter to the linear classifier can substantially increase the margin

minimize $\frac{1}{2} \|\underline{\theta}\|^2$ subject to

$$y_i(\underline{\theta} \cdot \underline{x}_i) \geq 1, \quad i = 1, \dots, n$$



minimize $\frac{1}{2} \|\underline{\theta}\|^2$ subject to

$$y_i(\underline{\theta} \cdot \underline{x}_i + \theta_0) \geq 1, \quad i = 1, \dots, n$$

Support vector machine

- Several desirable properties
 - maximizes the margin on the training set (\approx good generalization)
 - the solution is unique and sparse (\approx good generalization)
- But...
 - the solution is sensitive to outliers, labeling errors, as they may drastically change the resulting max-margin boundary
 - if the training set is not linearly separable, there's no solution!

Support vector machine

- Relaxed quadratic optimization problem

penalty for constraint violation

$$\text{minimize } \frac{1}{2} \|\underline{\theta}\|^2 + C \sum_{i=1}^n \xi_i \quad \text{subject to}$$

$$y_i(\underline{\theta} \cdot \underline{x}_i + \theta_0) \geq 1 - \xi_i, \quad i = 1, \dots, n$$

$$\xi_i \geq 0, \quad i = 1, \dots, n$$

slack variables
permit us to violate
some of the margin
constraints

Support vector machine

- Relaxed quadratic optimization problem

penalty for constraint violation

$$\text{minimize } \frac{1}{2} \|\underline{\theta}\|^2 + C \sum_{i=1}^n \xi_i \quad \text{subject to}$$

$$y_i(\underline{\theta} \cdot \underline{x}_i + \theta_0) \geq 1 - \xi_i, \quad i = 1, \dots, n$$

$$\xi_i \geq 0, \quad i = 1, \dots, n$$

large $C \Rightarrow$ few (if any) violations

small $C \Rightarrow$ many violations

slack variables
permit us to violate
some of the margin
constraints

Support vector machine

- Relaxed quadratic optimization problem

penalty for constraint violation

$$\text{minimize } \frac{1}{2} \|\underline{\theta}\|^2 + C \sum_{i=1}^n \xi_i \text{ subject to}$$

$$y_i(\underline{\theta} \cdot \underline{x}_i + \theta_0) \geq 1 - \xi_i, \quad i = 1, \dots, n$$

$$\xi_i \geq 0, \quad i = 1, \dots, n$$

large $C \Rightarrow$ few (if any) violations

small $C \Rightarrow$ many violations

slack variables
permit us to violate
some of the margin
constraints

we can still interpret the margin as $1/\|\underline{\theta}^*\|$

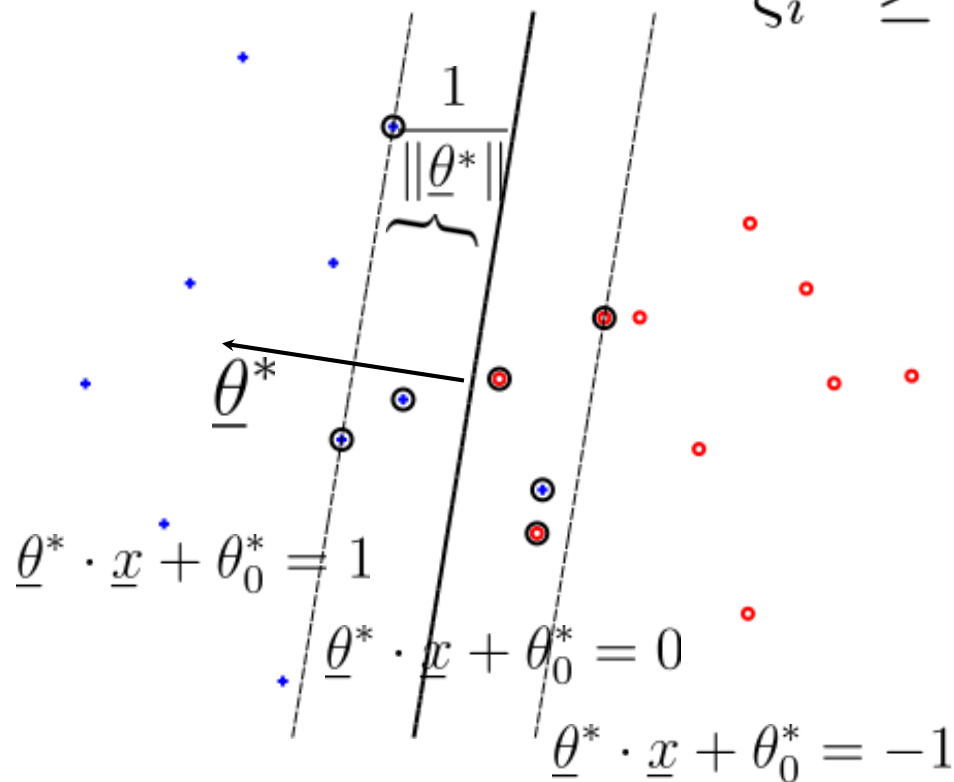
Support vector machine

- Relaxed quadratic optimization problem

$$\text{minimize } \frac{1}{2} \|\underline{\theta}\|^2 + C \sum_{i=1}^n \xi_i \quad \text{subject to}$$

$$y_i(\underline{\theta} \cdot \underline{x}_i + \theta_0) \geq 1 - \xi_i, \quad i = 1, \dots, n$$

$$\xi_i \geq 0, \quad i = 1, \dots, n$$



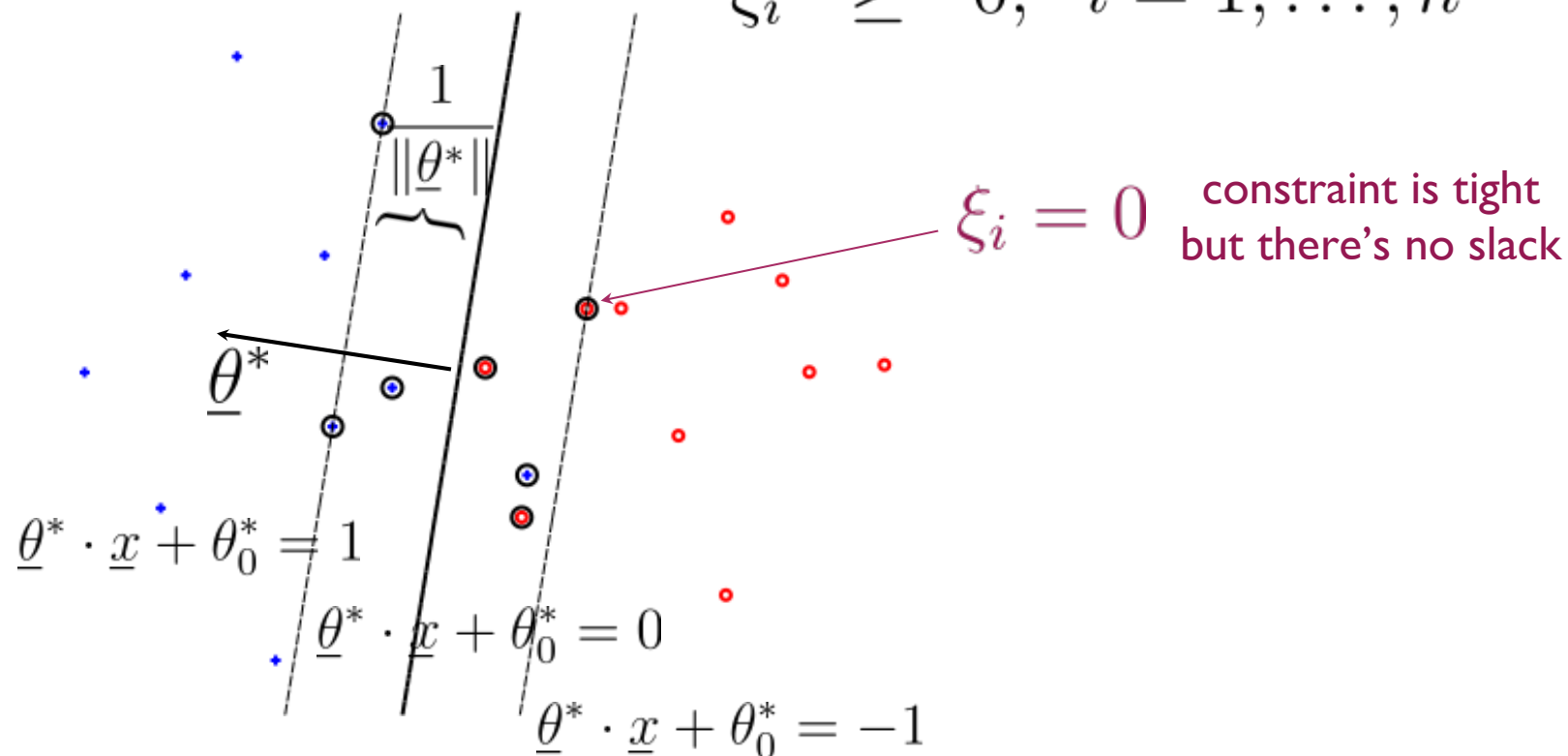
Support vectors and slack

- The solution now has three types of support vectors

$$\text{minimize } \frac{1}{2} \|\underline{\theta}\|^2 + C \sum_{i=1}^n \xi_i \quad \text{subject to}$$

$$y_i(\underline{\theta} \cdot \underline{x}_i + \theta_0) \geq 1 - \xi_i, \quad i = 1, \dots, n$$

$$\xi_i \geq 0, \quad i = 1, \dots, n$$



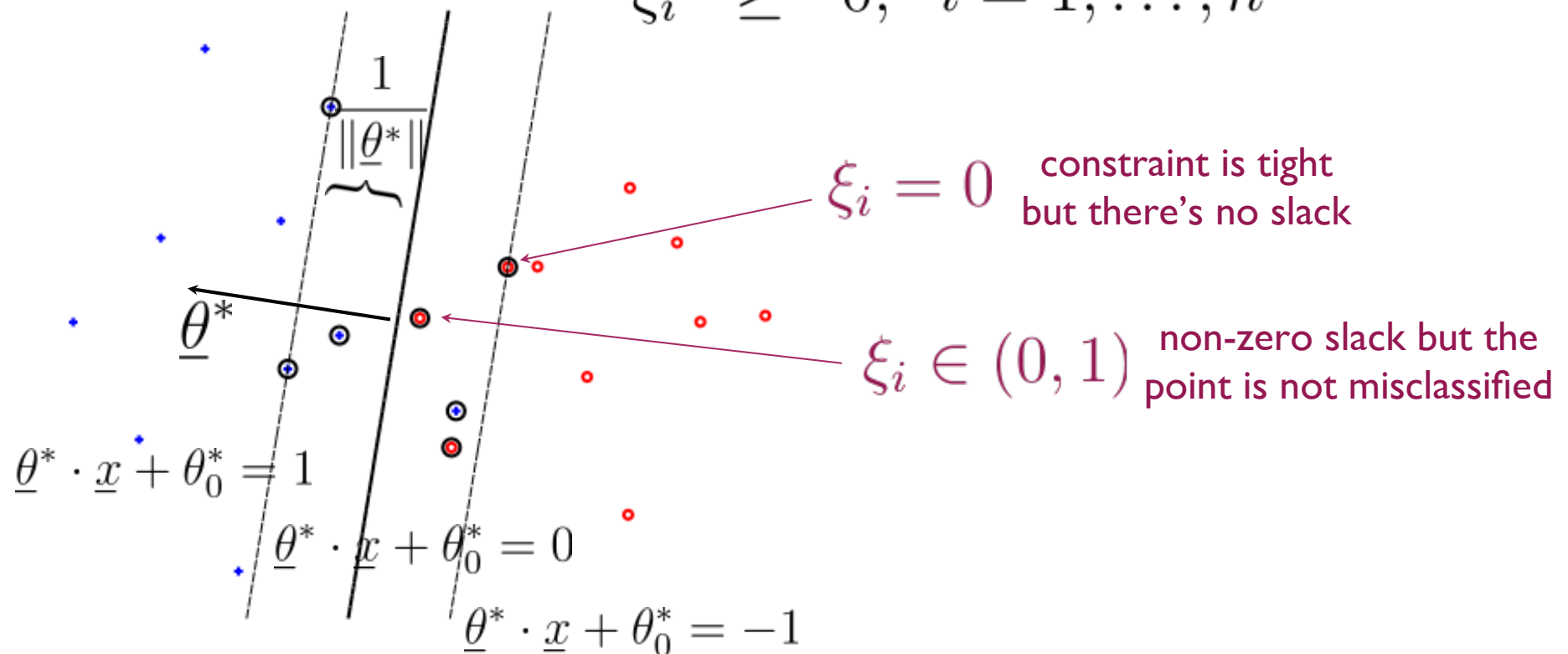
Support vectors and slack

- The solution now has three types of support vectors

$$\text{minimize } \frac{1}{2} \|\underline{\theta}\|^2 + C \sum_{i=1}^n \xi_i \quad \text{subject to}$$

$$y_i(\underline{\theta} \cdot \underline{x}_i + \theta_0) \geq 1 - \xi_i, \quad i = 1, \dots, n$$

$$\xi_i \geq 0, \quad i = 1, \dots, n$$



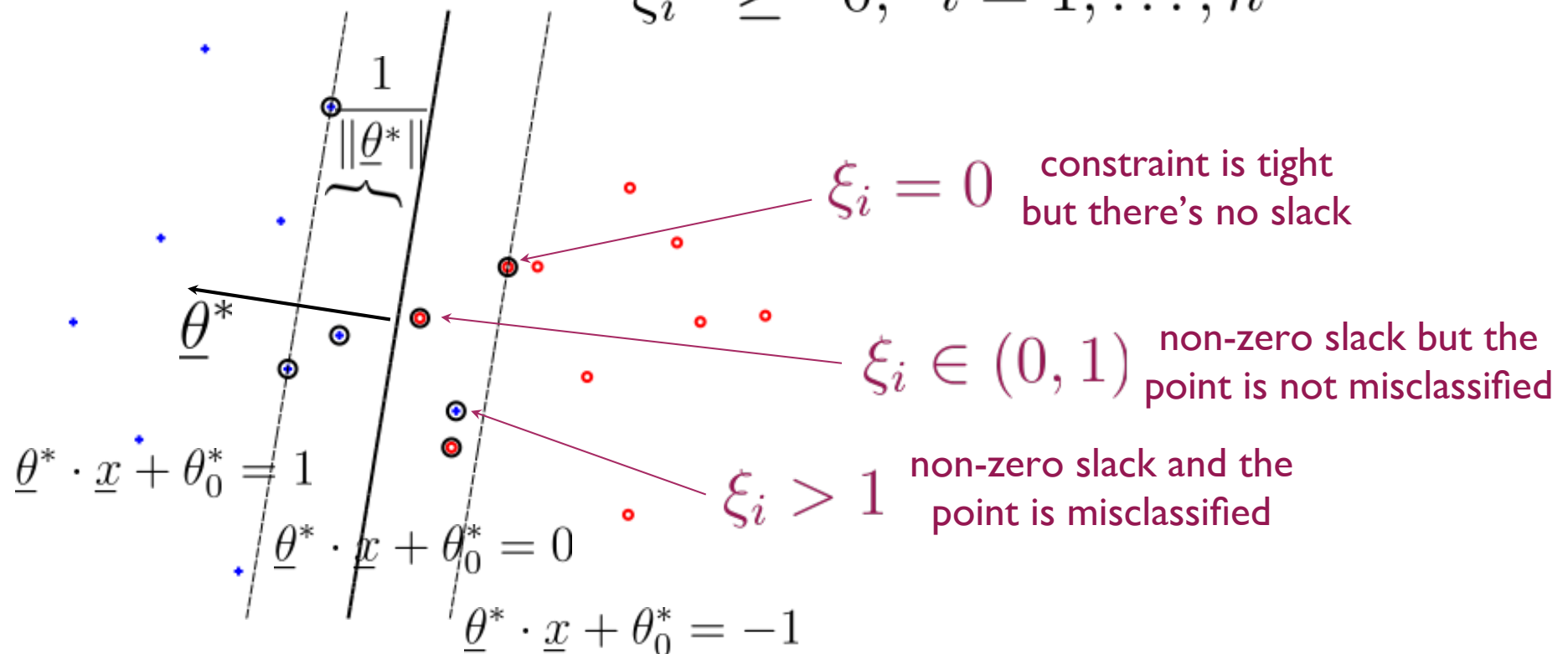
Support vectors and slack

- The solution now has three types of support vectors

$$\text{minimize } \frac{1}{2} \|\underline{\theta}\|^2 + C \sum_{i=1}^n \xi_i \quad \text{subject to}$$

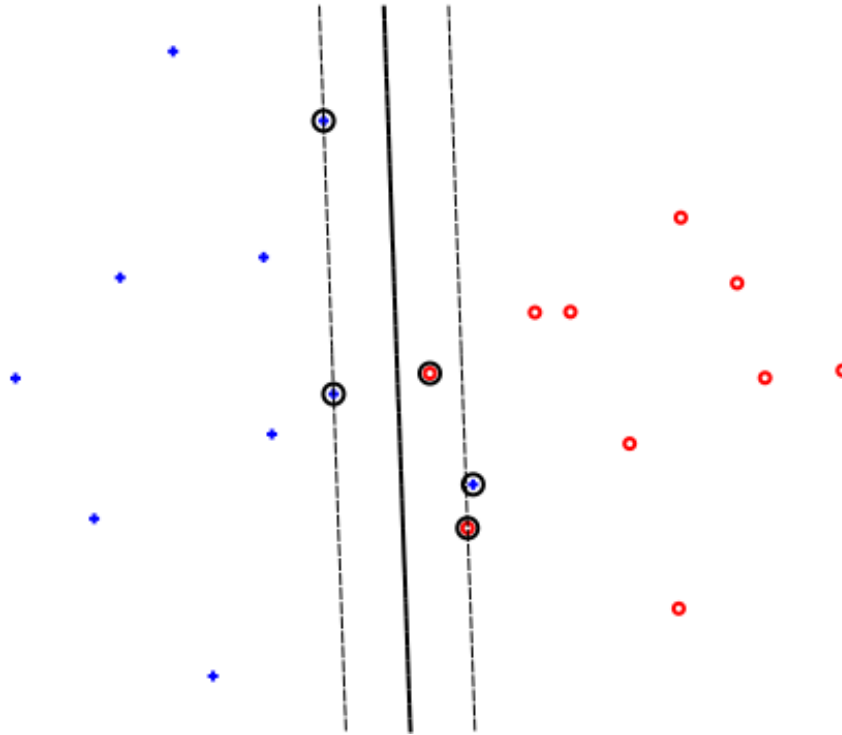
$$y_i(\underline{\theta} \cdot \underline{x}_i + \theta_0) \geq 1 - \xi_i, \quad i = 1, \dots, n$$

$$\xi_i \geq 0, \quad i = 1, \dots, n$$



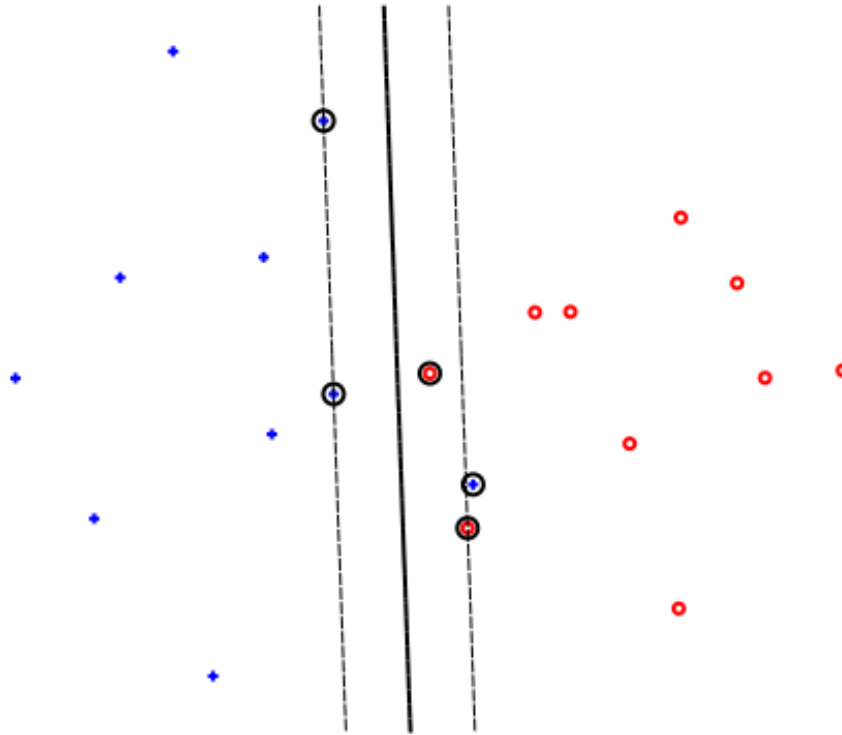
Examples

- $C=100$



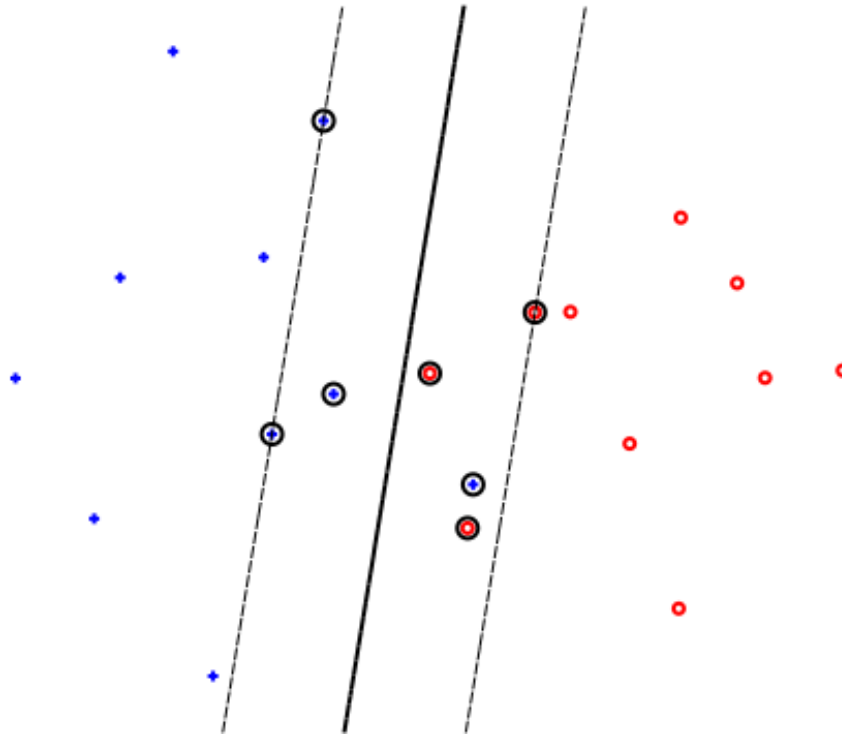
Examples

- $C=10$



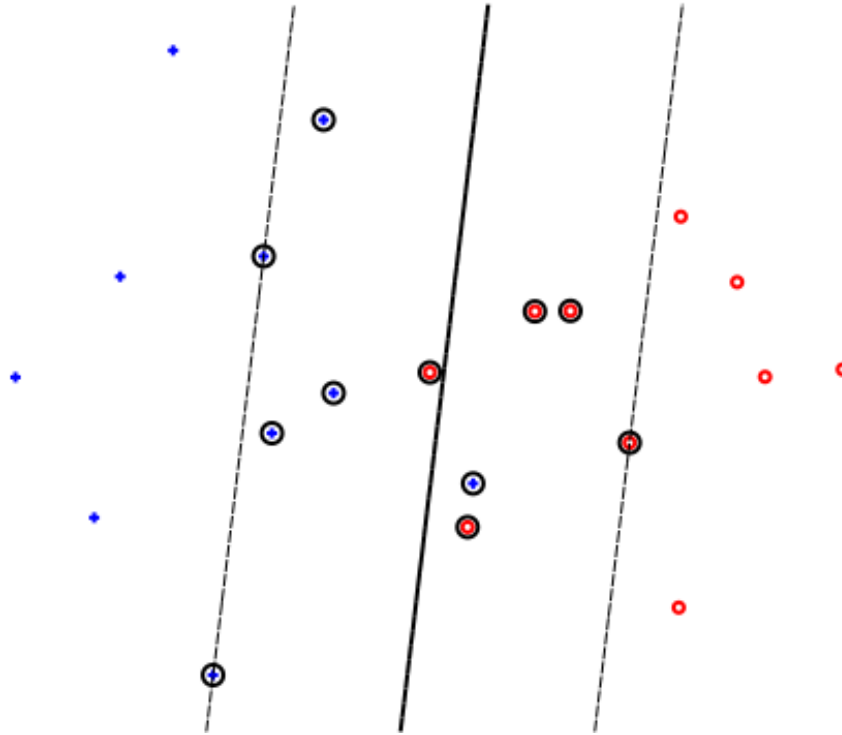
Examples

- $C=1$



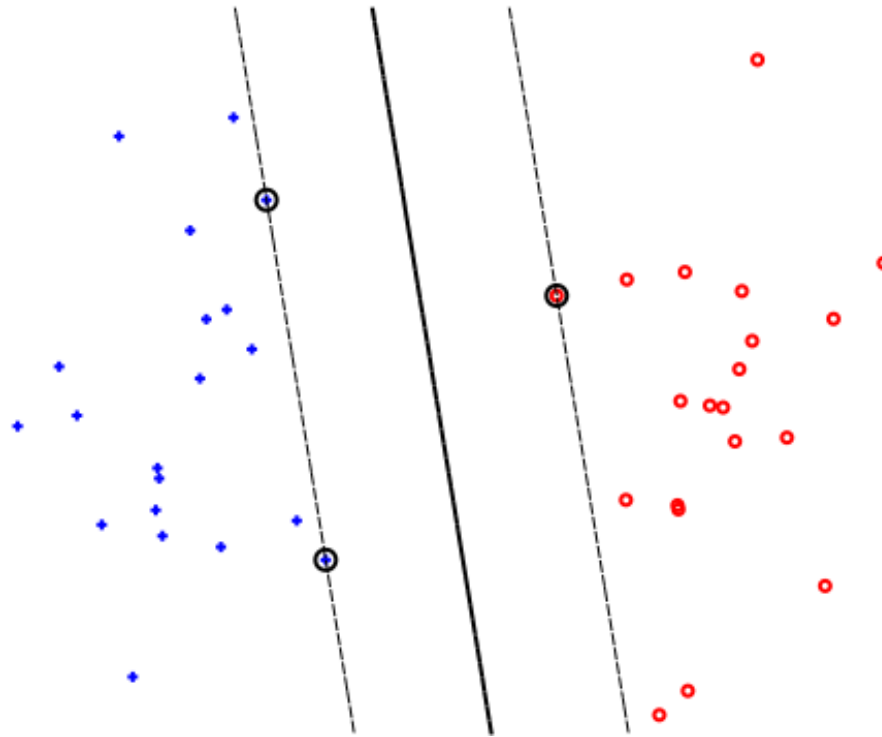
Examples

- $C=0.1$



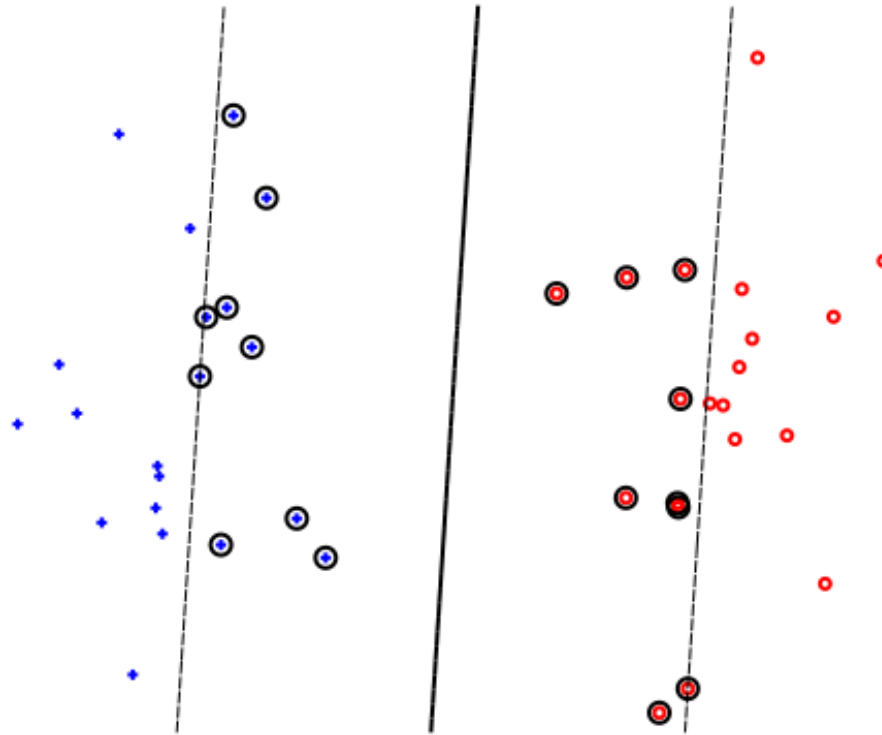
Examples

- C potentially affects the solution even in the separable case
- $C = I$



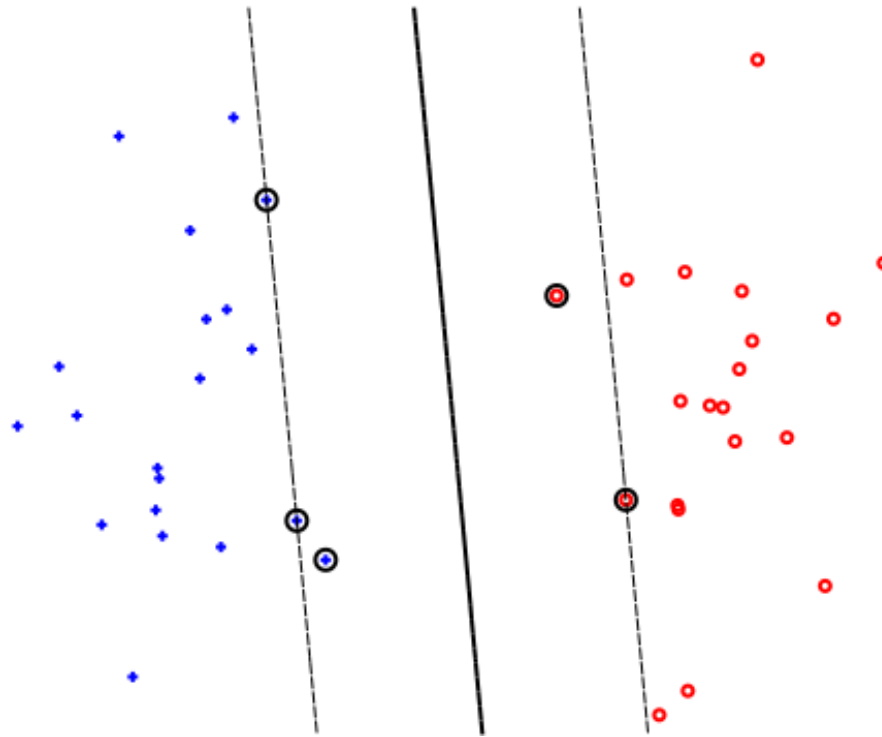
Examples

- C potentially affects the solution even in the separable case
- C = 0.01

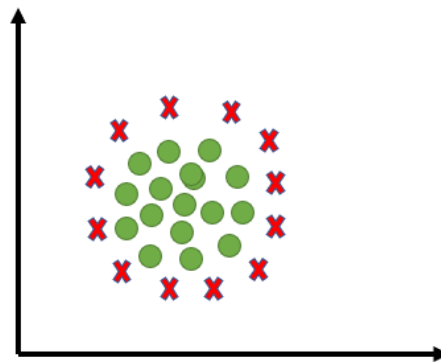


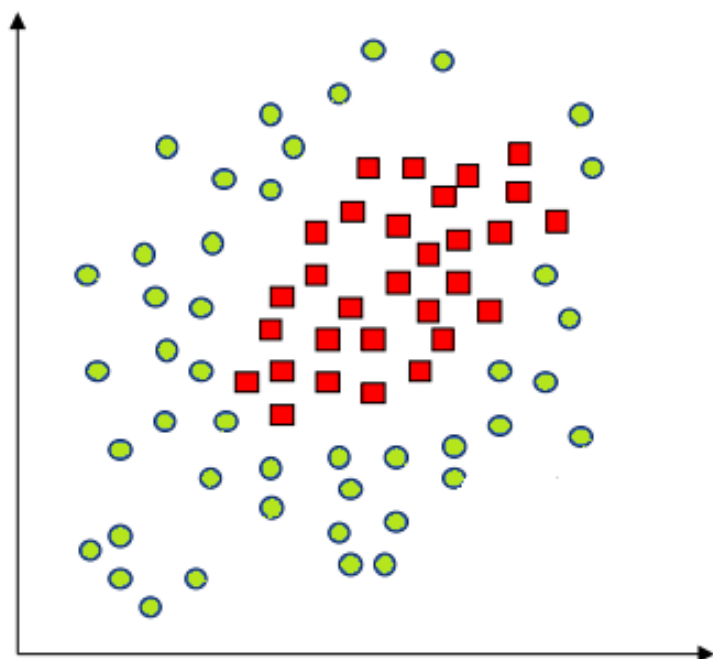
Examples

- C potentially affects the solution even in the separable case
- $C = 0.1$

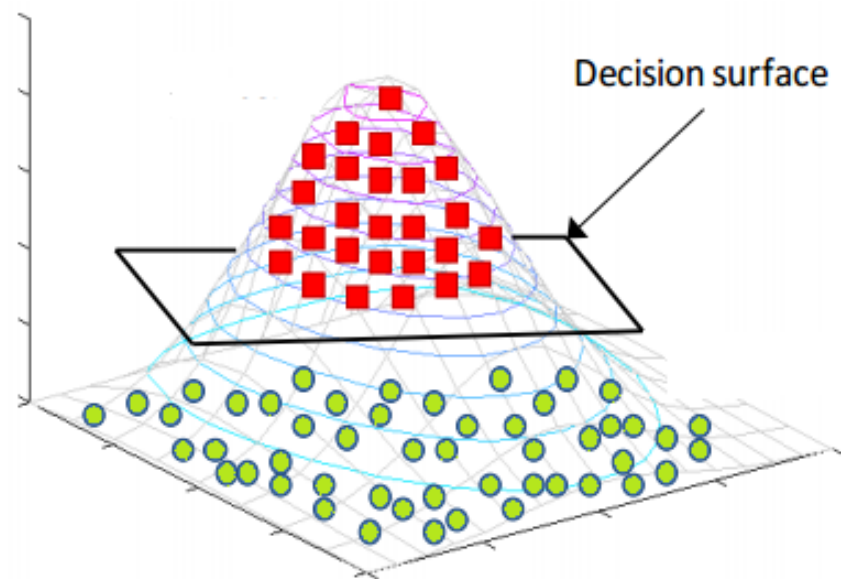


Non-linear dataset





kernel
→



Different Types of kernel

Polynomial

Sigmoid

RBF

$$\kappa(X1, X2) = (X1^T \cdot X2 + 1)^d$$

$$\kappa(x1, x2) = \tanh(\alpha x^T y + x)$$

$$\kappa(x1, x2) = e^{\frac{-\|x1 - x2\|^2}{2\sigma^2}}$$

Polynomial Kernel

- $K(X1, X2) = \phi(X1) \cdot \phi(X2)$

$$\begin{aligned} X1^T \cdot X2 &= \begin{bmatrix} X1 \\ X2 \end{bmatrix} \cdot [X1 \quad X2] \\ &= \begin{bmatrix} X1^2 & X1 \cdot X2 \\ X1 \cdot X2 & X2^2 \end{bmatrix} \end{aligned}$$