

CS578 Statistical Machine Learning Lecture 2

Jean Honorio
Purdue University

(based on slides by Tommi Jaakkola, MIT CSAIL)

Today's topics

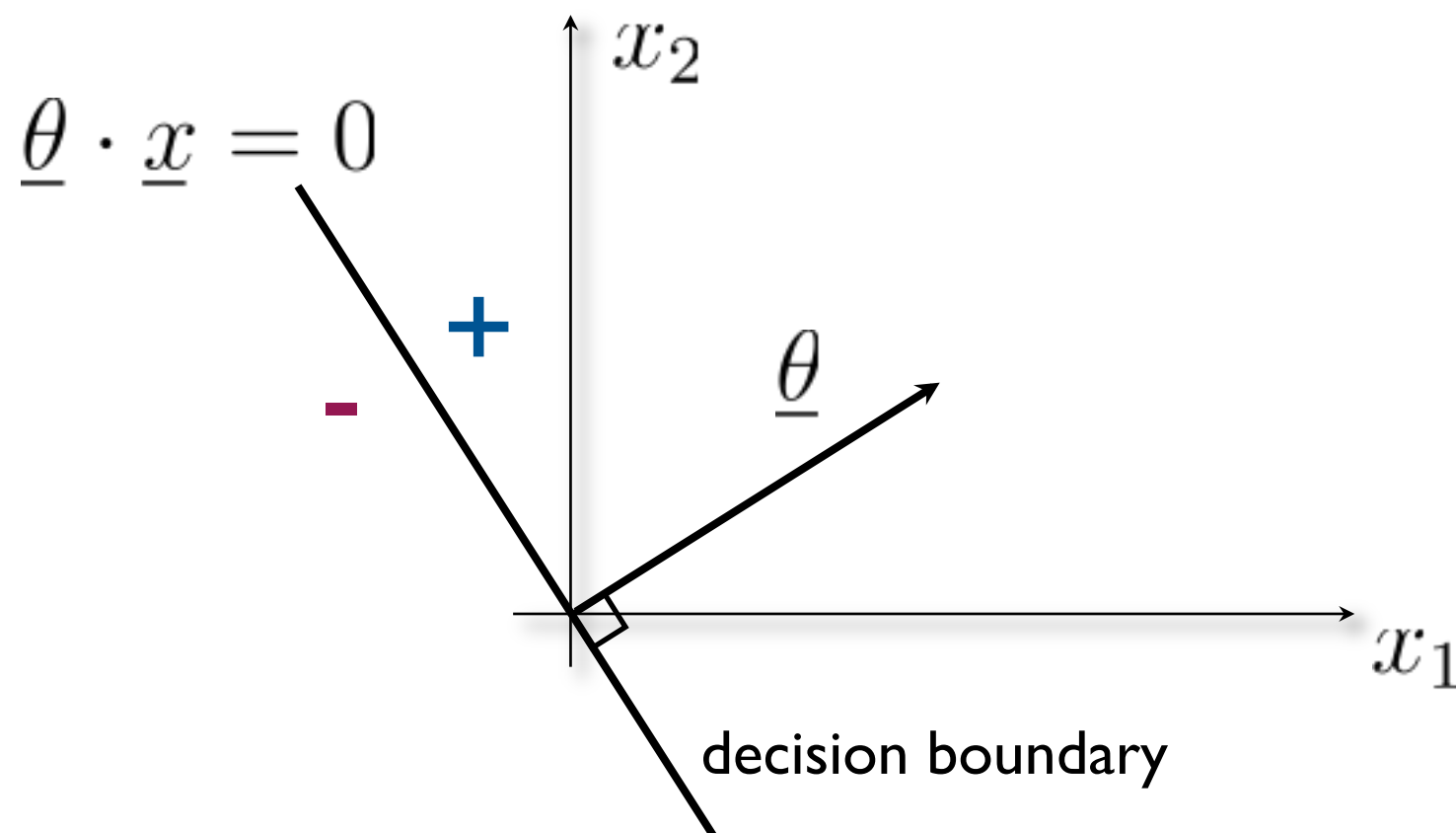
- Perceptron, convergence
 - the prediction game
 - mistakes, margin, and generalization
- Maximum margin classifier -- support vector machine
 - estimation, properties
 - allowing misclassified points

Recall: linear classifiers

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

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

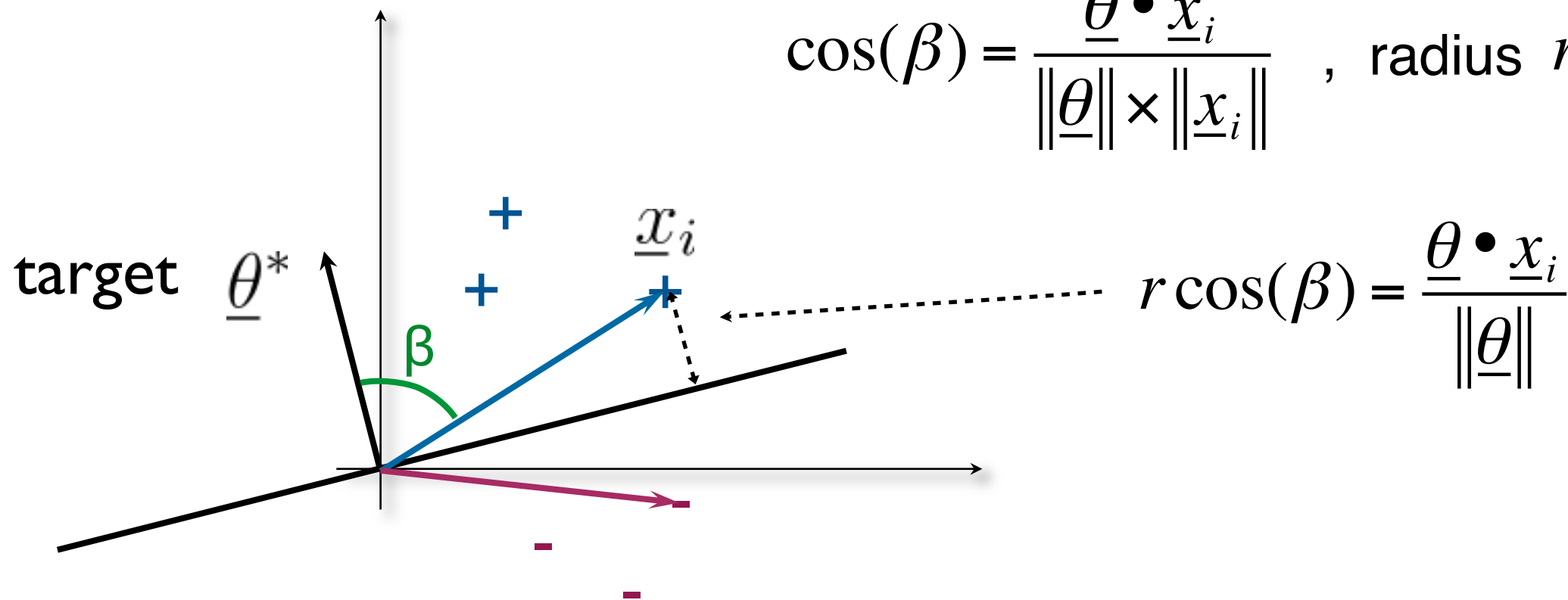
discriminant function



Recall: some linear algebra

- More details about slide 1-52. Consider a positive point:

$$\cos(\beta) = \frac{\underline{\theta} \bullet \underline{x}_i}{\|\underline{\theta}\| \times \|\underline{x}_i\|} \quad , \quad \text{radius } r = \|\underline{x}_i\|$$



- Positive point: $y_i=1$, $\beta < 90^\circ$, $\cos(\beta) > 0$, $\underline{\theta} \bullet \underline{x}_i > 0$
- Negative point: $y_i=-1$, $\beta > 90^\circ$, $\cos(\beta) < 0$, $\underline{\theta} \bullet \underline{x}_i < 0$

- General formula for positive and negative points: $\frac{y_i(\underline{\theta}^* \bullet \underline{x}_i)}{\|\underline{\theta}^*\|}$

Perceptron algorithm

- The perceptron algorithm considers each training point in turn, adjusting the parameters to correct any mistakes

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

Repeat until convergence:

for $t = 1, \dots, n$

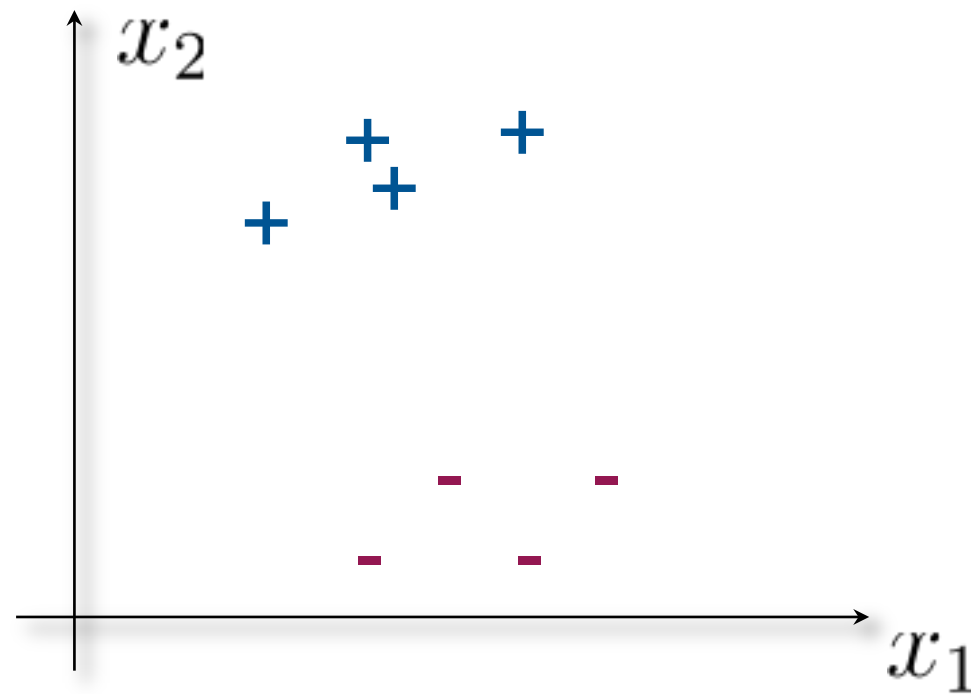
if $y_t(\underline{\theta} \cdot \underline{x}_t) \leq 0$ (mistake)

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

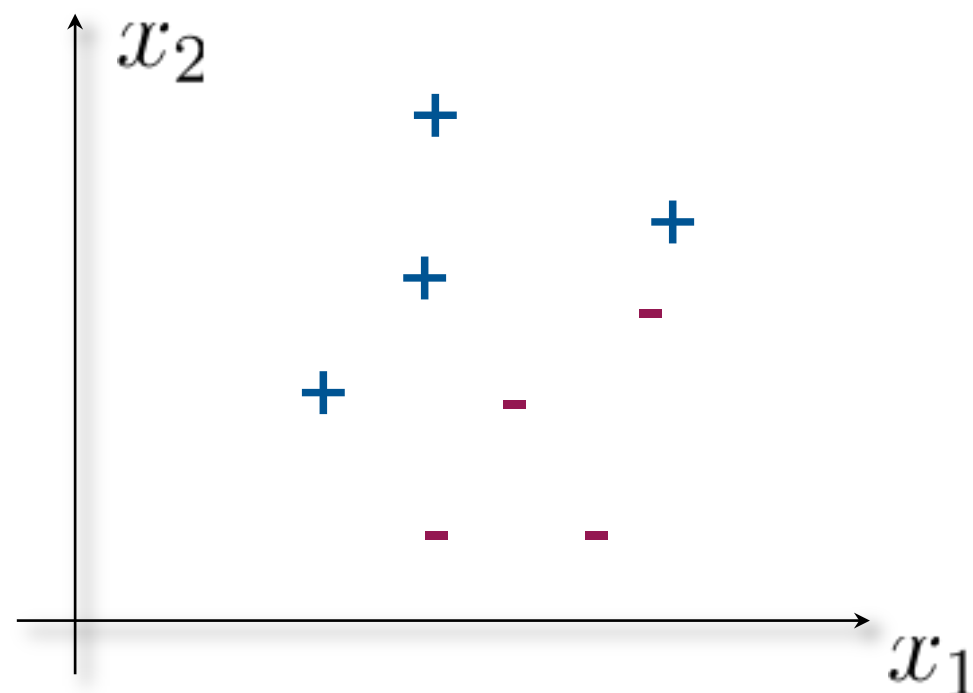
- We would like to bound the number of mistakes that the algorithm makes

Mistakes and margin

Easy problem
- large margin
- few mistakes

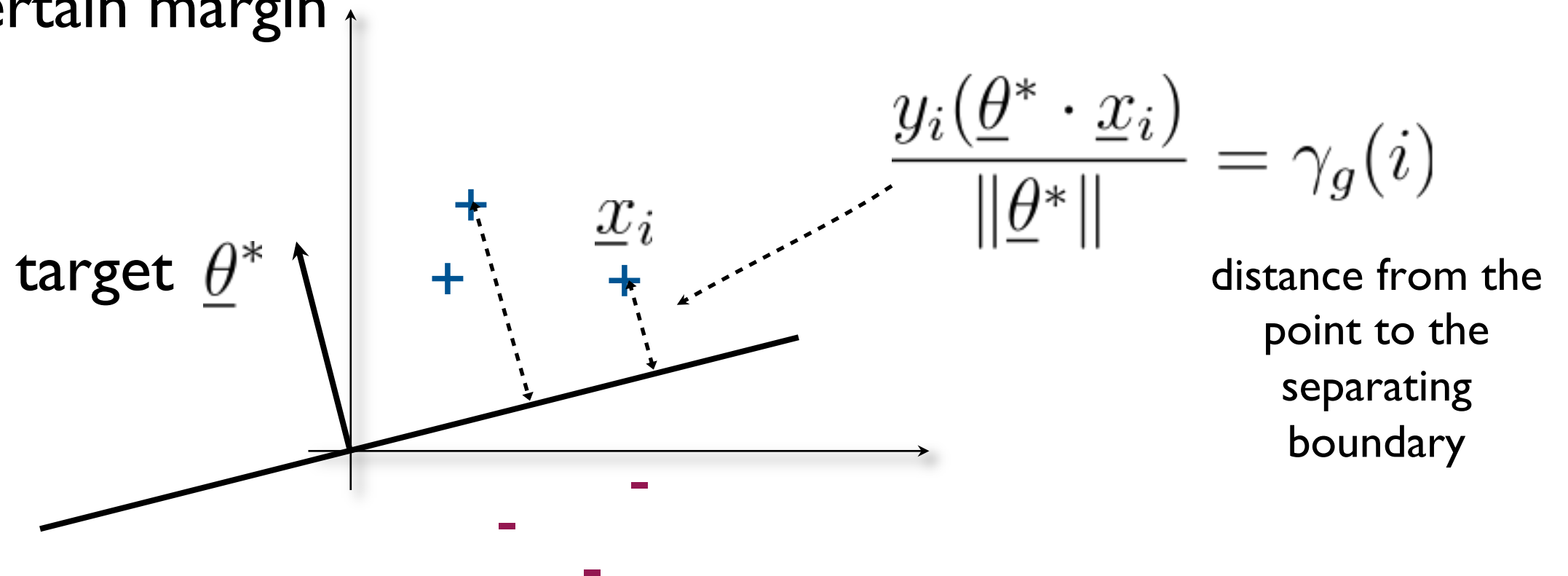


Harder problem
- small margin
- many mistakes



The target classifier

- We can quantify how hard the problem is by assuming that there exists a target classifier that achieves a certain margin



- The geometric margin γ_g is the closest distance to the separating boundary $\gamma_g = \min_i \gamma_g(i)$

Perceptron convergence theorem

- If there exists $\underline{\theta}^*$ such that

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

and $\|\underline{x}_i\| \leq R$ then the perceptron algorithm makes at most

$$\frac{R^2}{\gamma_g^2}$$

mistakes (on the training set).

- Key points
 - large geometric margin relative to the norm of the examples implies few mistakes
 - the result does not depend on the dimension d of the examples (the number of parameters)

Mistake guarantee: proof

- We show that after k updates (mistakes),

$$\frac{\underline{\theta}^{(k)} \cdot \underline{\theta}^*}{\|\underline{\theta}^{(k)}\|^2} \geq k\gamma_g \|\underline{\theta}^*\|$$
$$\|\underline{\theta}^{(k)}\|^2 \leq kR^2$$

Mistake guarantee: proof

- We show that after k updates (mistakes),

$$\begin{aligned}\underline{\theta}^{(k)} \cdot \underline{\theta}^* &\geq k\gamma_g \|\underline{\theta}^*\| \\ \|\underline{\theta}^{(k)}\|^2 &\leq kR^2\end{aligned}$$

- Let the k th mistake be on the i th example

$$\begin{aligned}\underline{\theta}^{(k)} \cdot \underline{\theta}^* &= [\underline{\theta}^{(k-1)} + y_i \underline{x}_i] \cdot \underline{\theta}^* \\ &= \underline{\theta}^{(k-1)} \cdot \underline{\theta}^* + \underbrace{y_i \underline{x}_i \cdot \underline{\theta}^*}_{\text{margin}} \\ &\geq \underline{\theta}^{(k-1)} \cdot \underline{\theta}^* + \gamma_g \|\underline{\theta}^*\|\end{aligned}$$

Note:

Since $\underline{\theta}^{(0)} = 0$ then $\underline{\theta}^{(k)} \cdot \underline{\theta}^* \geq k \gamma_g \|\underline{\theta}^*\|$

Mistake guarantee: proof

- We show that after k updates (mistakes),

$$\frac{\underline{\theta}^{(k)} \cdot \underline{\theta}^*}{\|\underline{\theta}^{(k)}\|^2} \geq k\gamma_g \|\underline{\theta}^*\|$$
$$\|\underline{\theta}^{(k)}\|^2 \leq kR^2$$

- Let the k th mistake be on the i th example

$$\begin{aligned}\|\underline{\theta}^{(k)}\|^2 &= \|\underline{\theta}^{(k-1)} + y_i \underline{x}_i\|^2 && \text{mistake: } \leq 0 \\ &= \|\underline{\theta}^{(k-1)}\|^2 + 2y_i \underline{\theta}^{(k-1)} \cdot \underline{x}_i + \|\underline{x}_i\|^2 \\ &\leq \|\underline{\theta}^{(k-1)}\|^2 + \|\underline{x}_i\|^2 \\ &\leq \|\underline{\theta}^{(k-1)}\|^2 + R^2\end{aligned}$$

Note:

For any two vectors, \underline{a} and \underline{b} :

$$\|\underline{a} + \underline{b}\|^2 = (\underline{a} + \underline{b}) \cdot (\underline{a} + \underline{b}) = \underline{a} \cdot \underline{a} + 2 \underline{a} \cdot \underline{b} + \underline{b} \cdot \underline{b} = \|\underline{a}\|^2 + 2 \underline{a} \cdot \underline{b} + \|\underline{b}\|^2$$

Since $\underline{\theta}^{(0)} = 0$ then $\|\underline{\theta}^{(k)}\|^2 \leq k R^2$

Mistake guarantee: proof

- We have shown that after k updates (mistakes),

$$\begin{aligned}\frac{\underline{\theta}^{(k)} \cdot \underline{\theta}^*}{\|\underline{\theta}^{(k)}\|^2} &\geq k\gamma_g \|\underline{\theta}^*\| \\ \|\underline{\theta}^{(k)}\|^2 &\leq kR^2\end{aligned}$$

- As a result,

$$1 \geq \frac{\overbrace{\underline{\theta}^{(k)} \cdot \underline{\theta}^*}^{\text{cosine}}}{\|\underline{\theta}^{(k)}\| \|\underline{\theta}^*\|} \geq \frac{k\gamma_g}{\sqrt{k}R}$$

$$\Rightarrow k \leq \frac{R^2}{\gamma_g^2}$$

Summary (perceptron)

- By analyzing the simple perceptron algorithm, we were able to relate the number of mistakes and geometric margin
- In cases where we are given a fixed set of training examples, and they are linearly separable, we can find and use the maximum margin classifier directly

Optimization

- Example problem:

minimize $(z-1)^2$

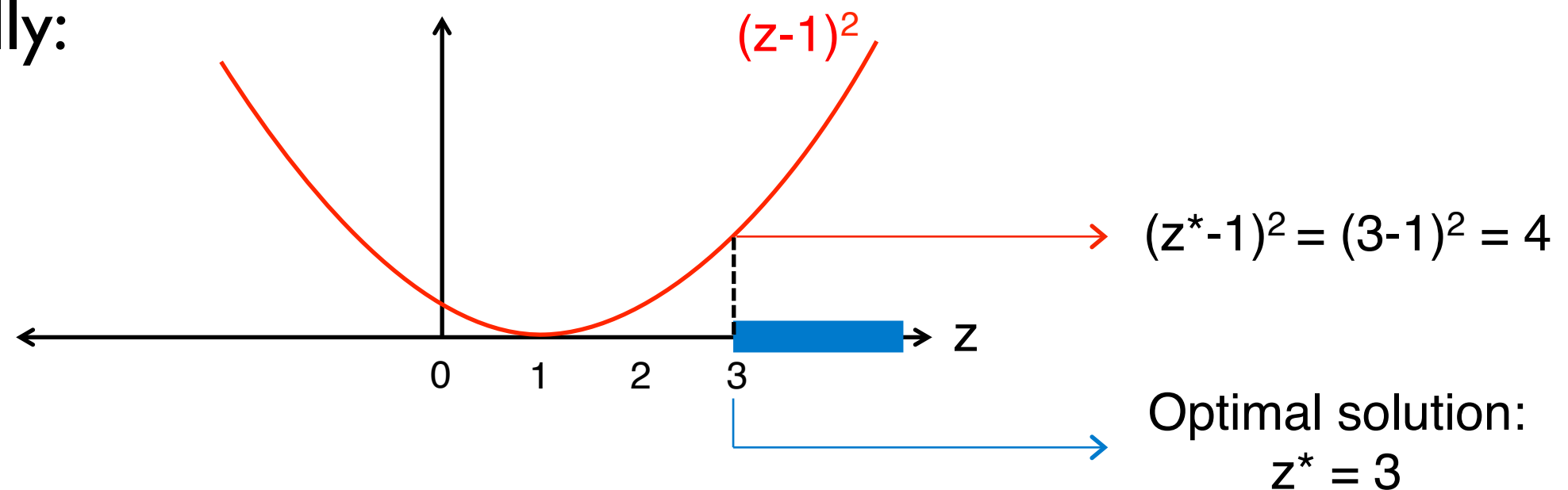
subject to $z \geq 3$

Objective function

Optimization variable z

Constraint

- Graphically:



- What if we have several optimization variables?

Optimization

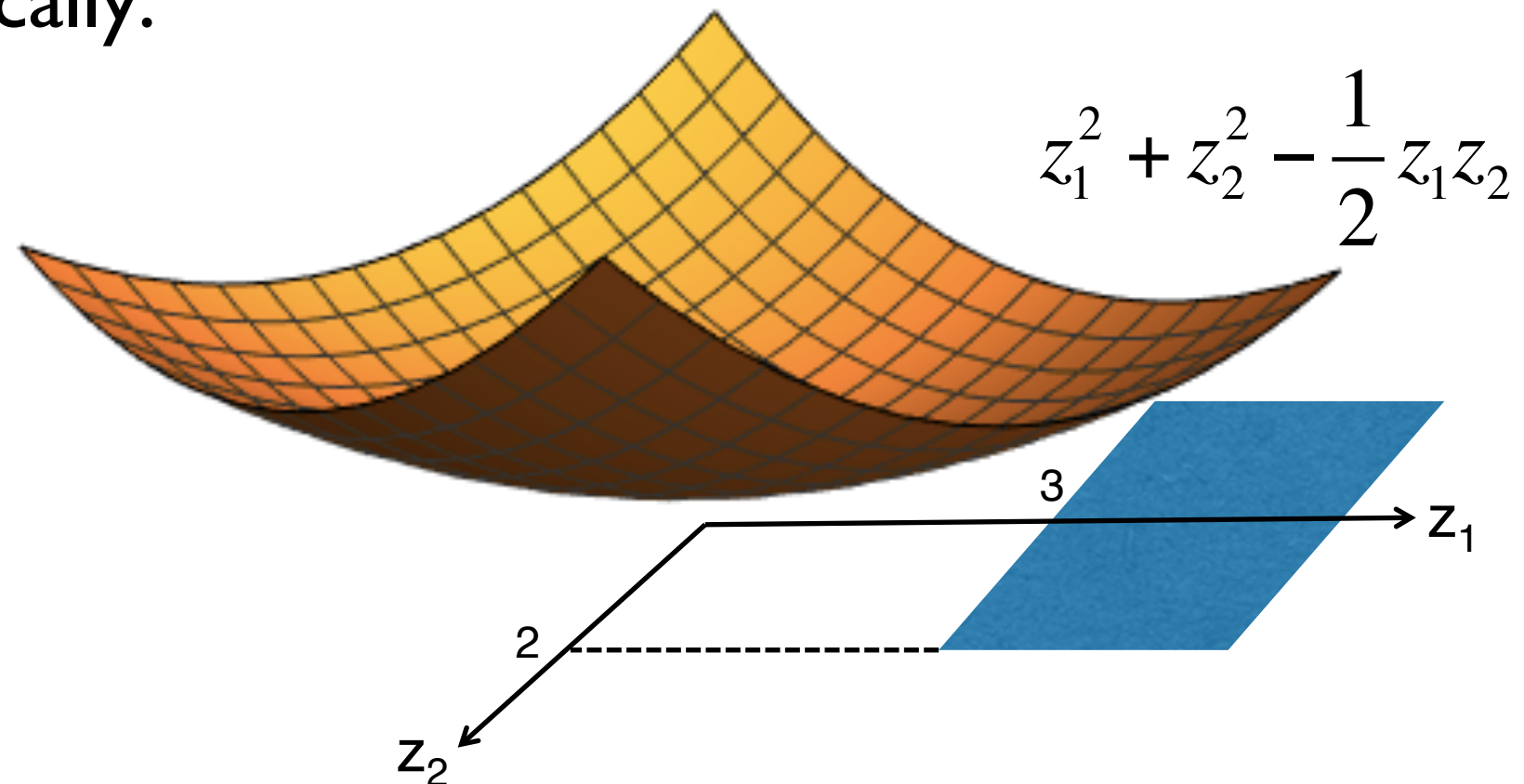
- Example problem:

minimize $z_1^2 + z_2^2 - \frac{1}{2}z_1z_2$ → Objective function

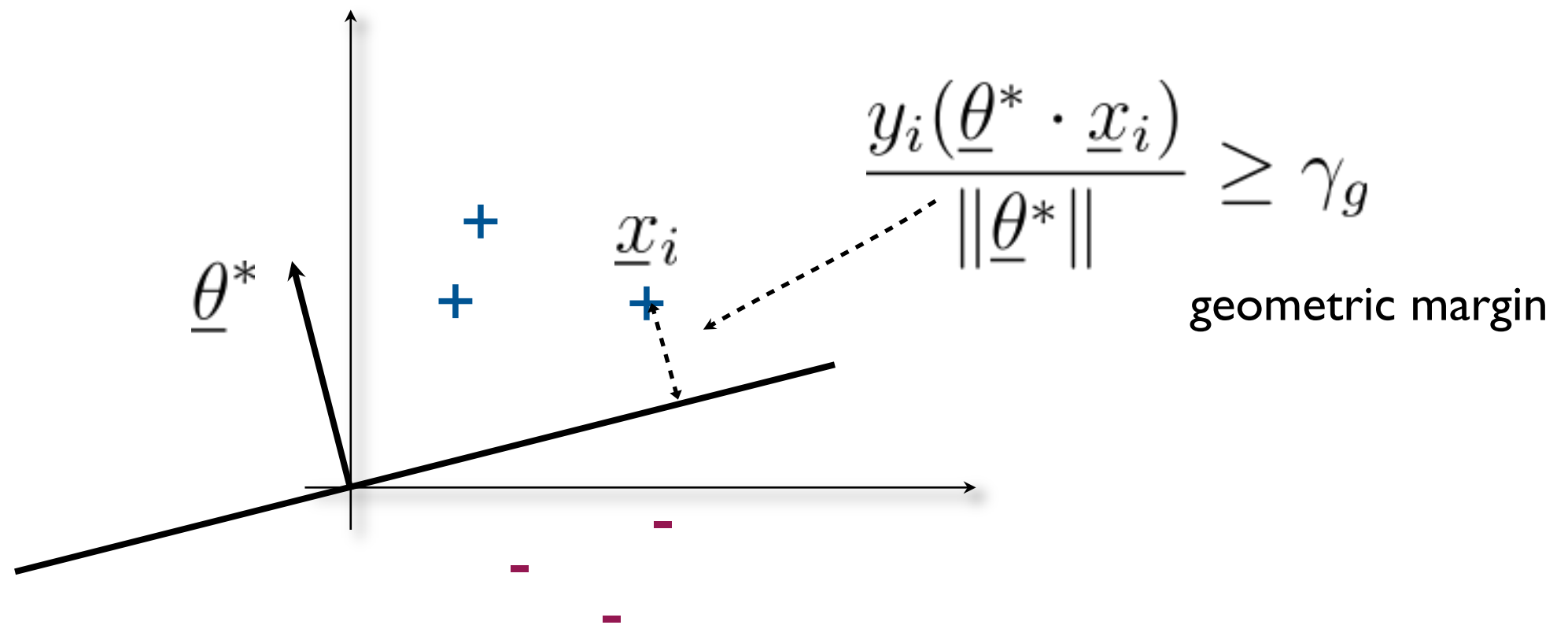
subject to $\begin{cases} z_1 \geq 3 \\ z_2 \leq 2 \end{cases}$ → Constraints

Optimization variable $\underline{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$

- Graphically:

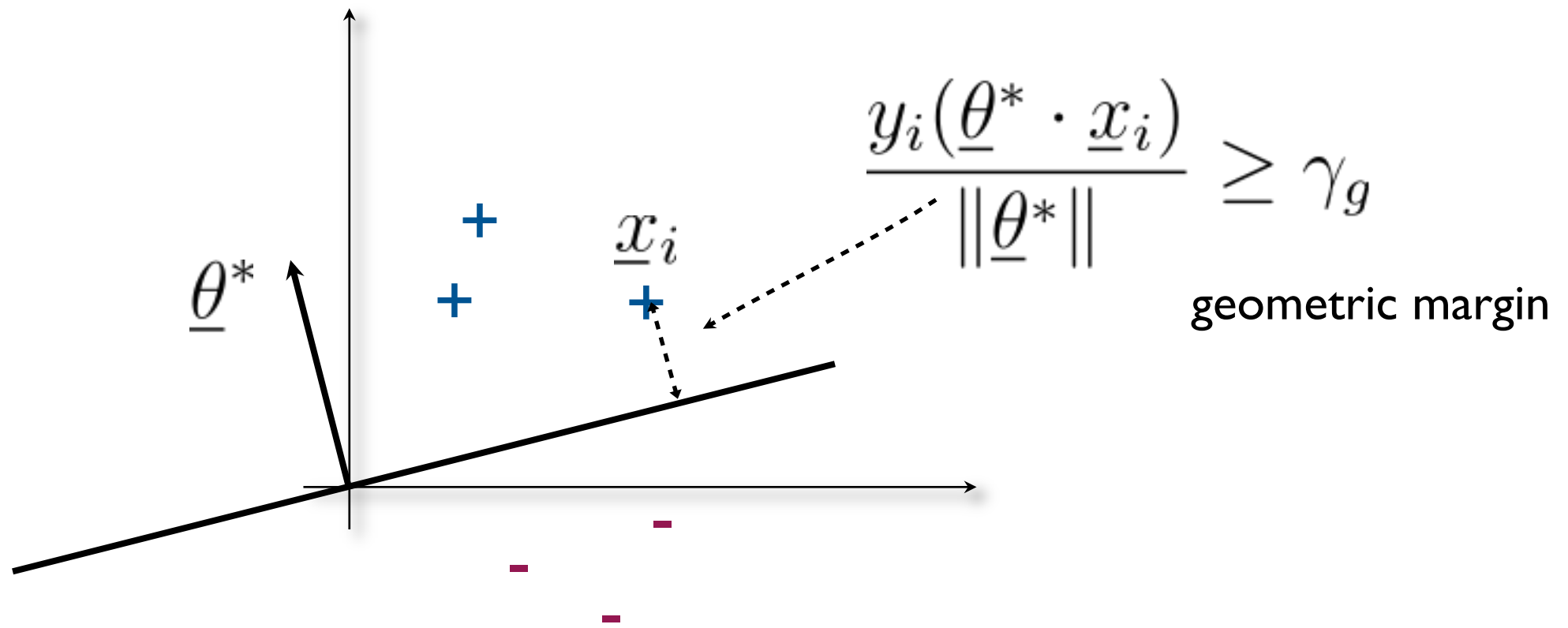


Maximum margin classifier



- Lets maximize the margin γ_g directly

Maximum margin classifier

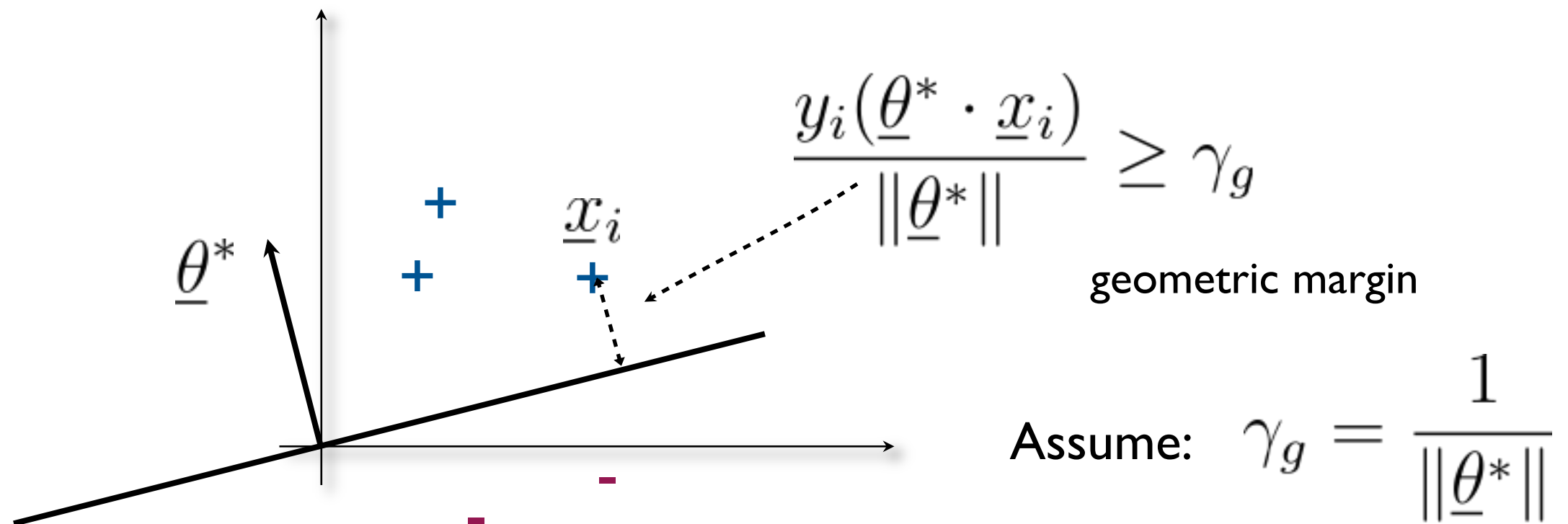


maximize γ_g subject to

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

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

Maximum margin classifier

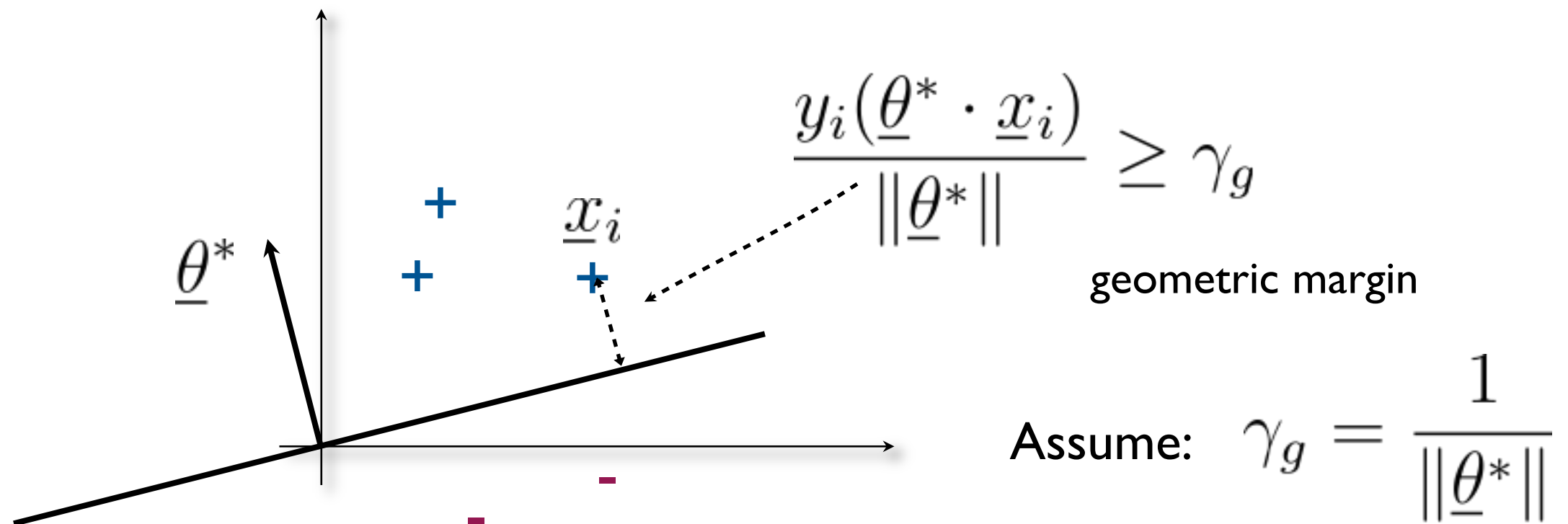


To find $\underline{\theta}^*$:

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

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

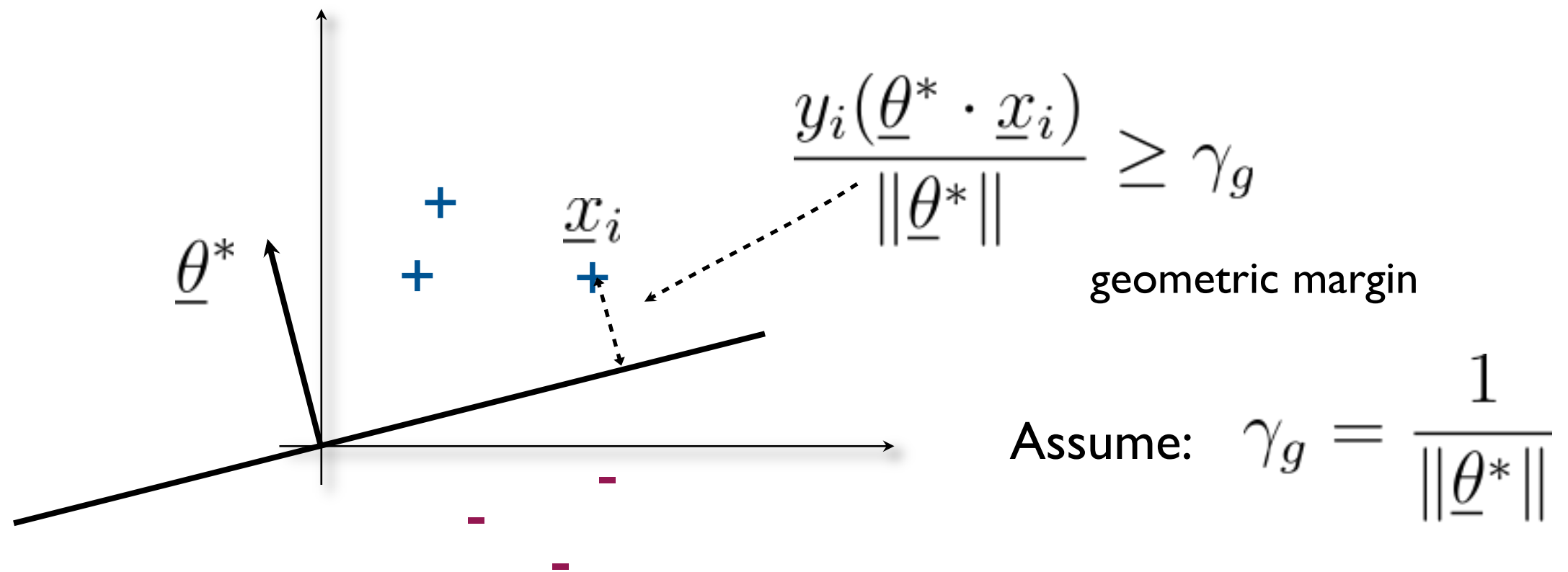
Maximum margin classifier



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

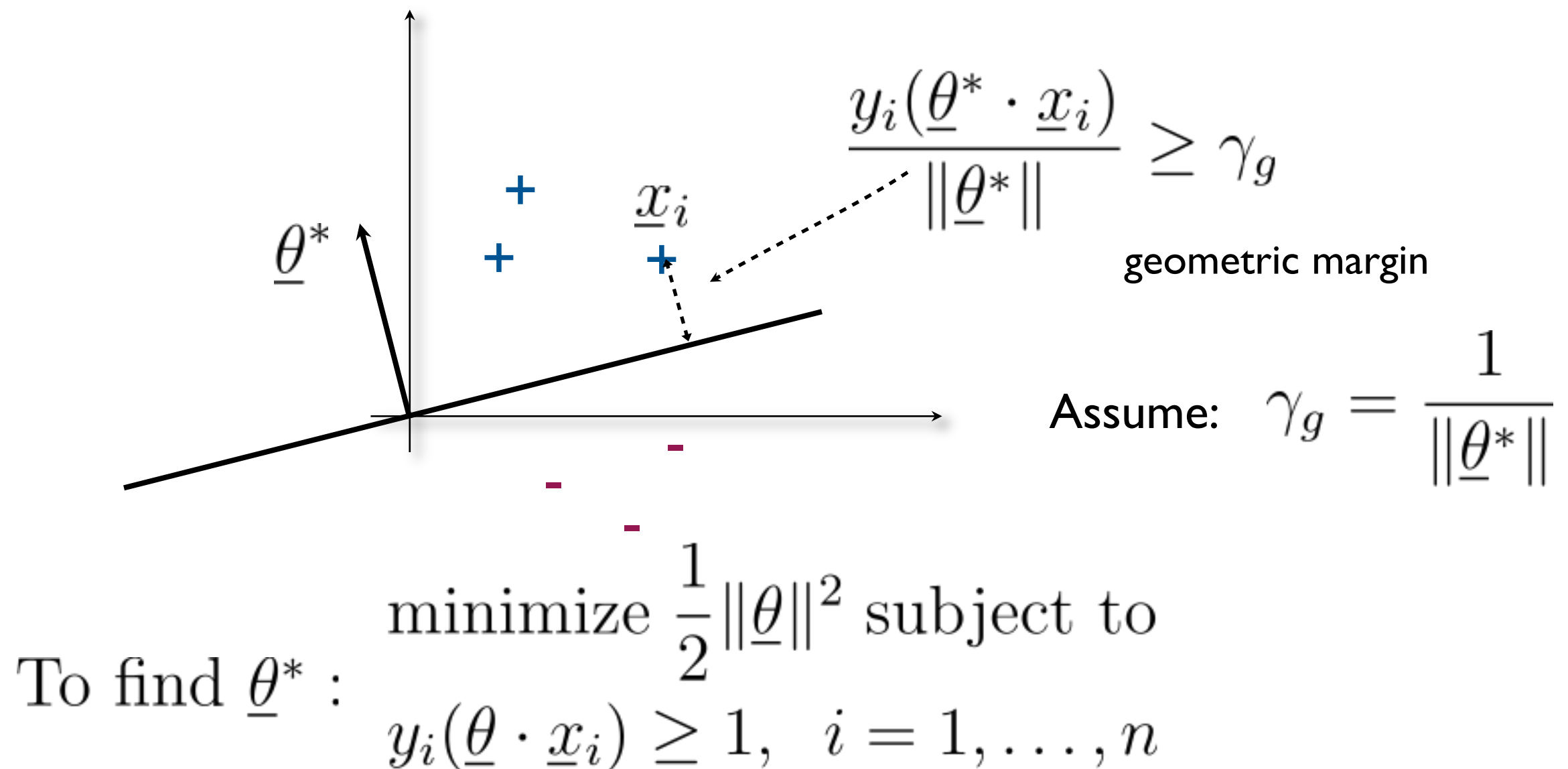
maximize $\frac{1}{\|\underline{\theta}\|}$ subject to
 $y_i(\underline{\theta} \cdot \underline{x}_i) \geq 1, \quad i = 1, \dots, n$

Maximum margin classifier



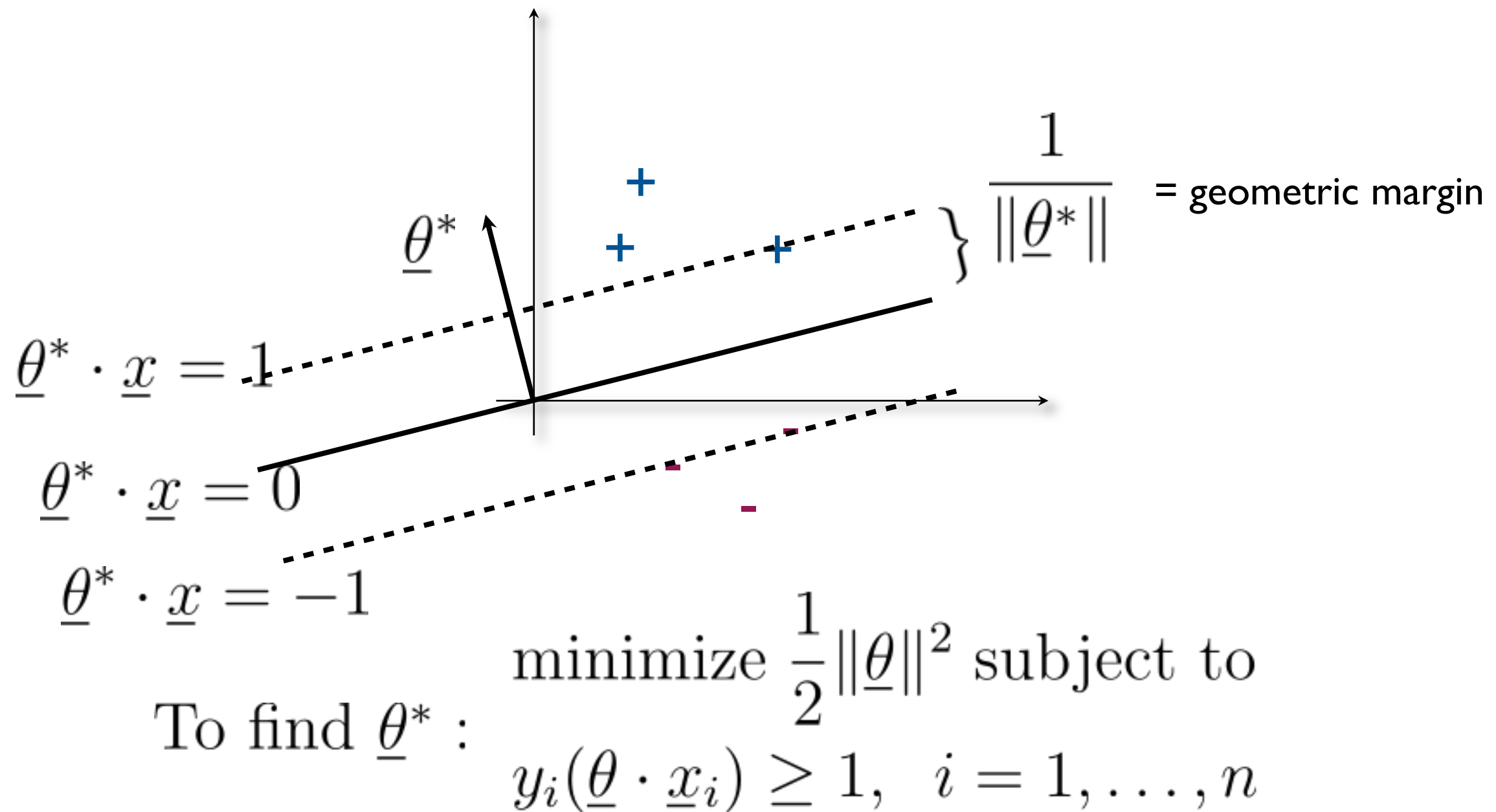
To find $\underline{\theta}^*$: minimize $\|\underline{\theta}\|$ subject to
 $y_i(\underline{\theta} \cdot \underline{x}_i) \geq 1, \quad i = 1, \dots, n$

Support vector machine

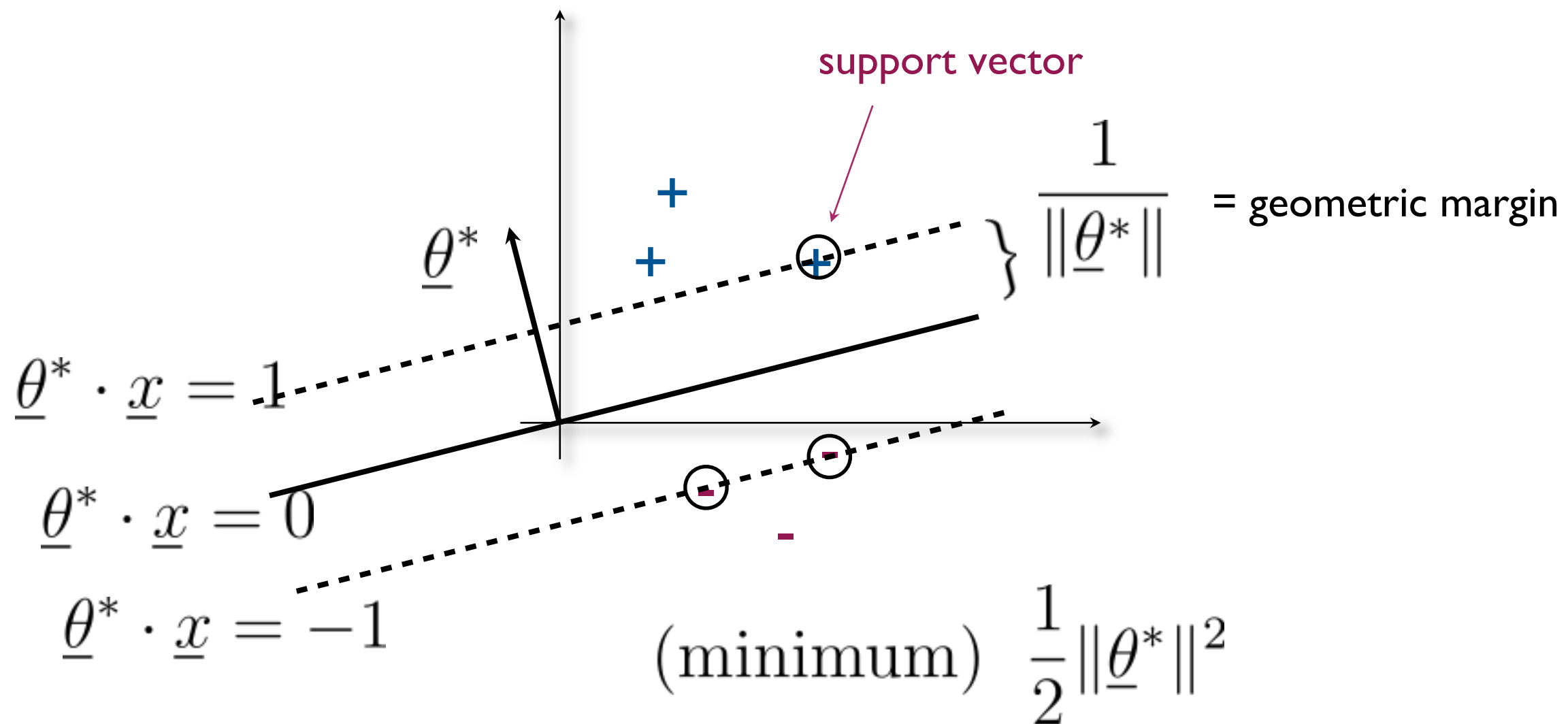


- This is a quadratic programming problem (quadratic objective, linear constraints)

Support vector machine



Support vector machine



The solution is
sparse

$$y_1(\underline{\theta}^* \cdot \underline{x}_1) = 1$$

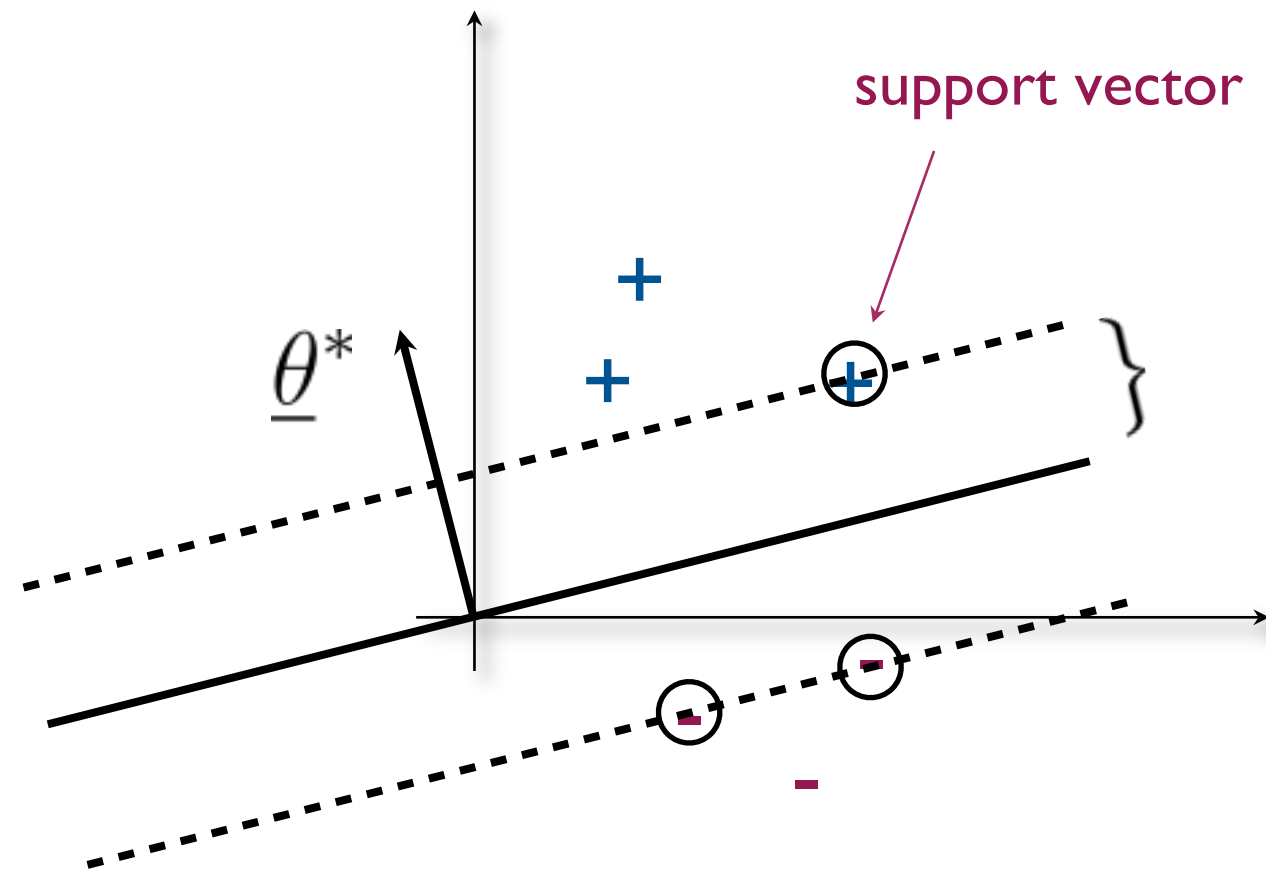
$$y_2(\underline{\theta}^* \cdot \underline{x}_2) > 1$$

$$y_3(\underline{\theta}^* \cdot \underline{x}_3) = 1$$

...

active constraints
= support vectors

Is sparse solution good?



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

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

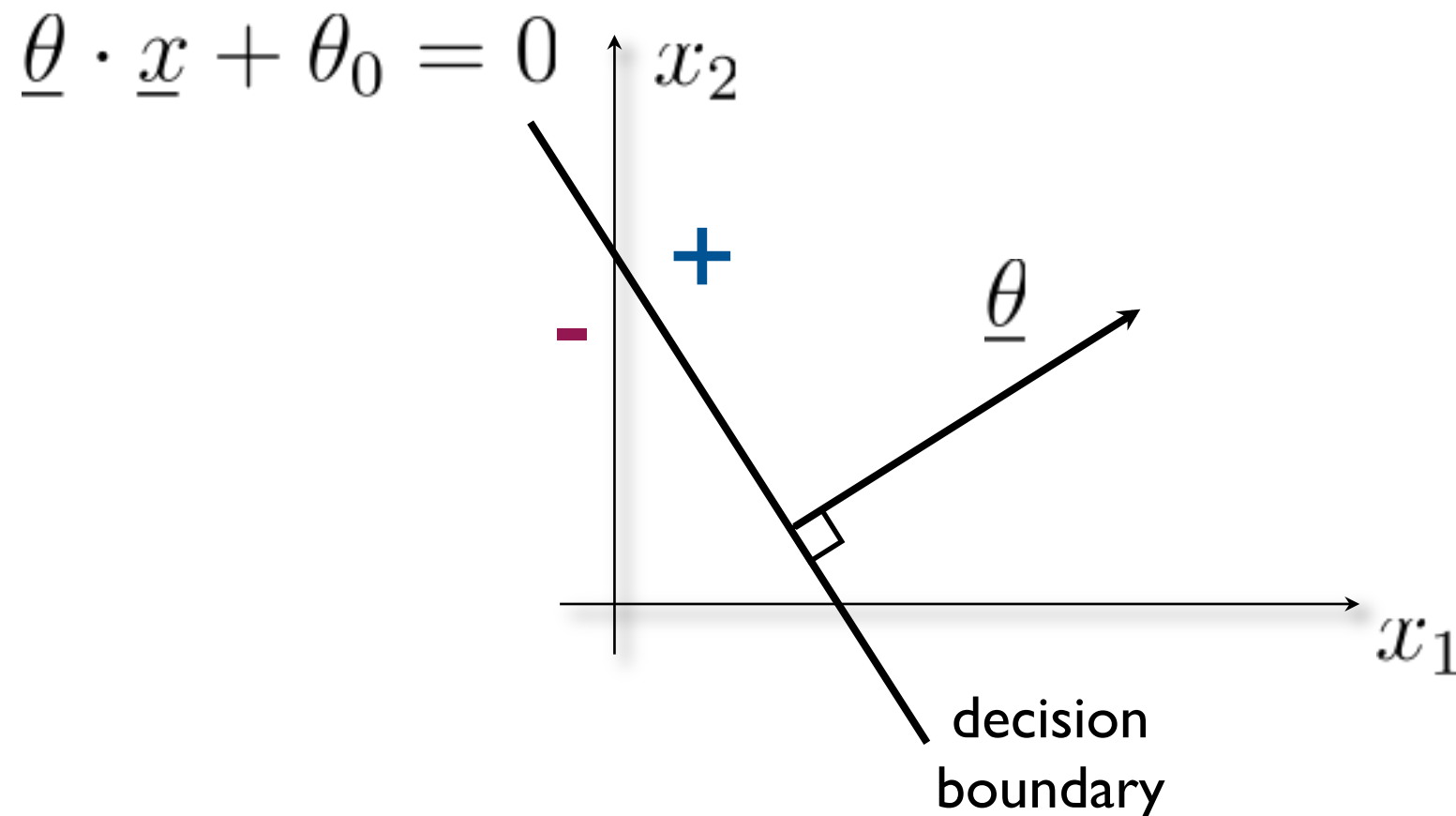
Intuitively:

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

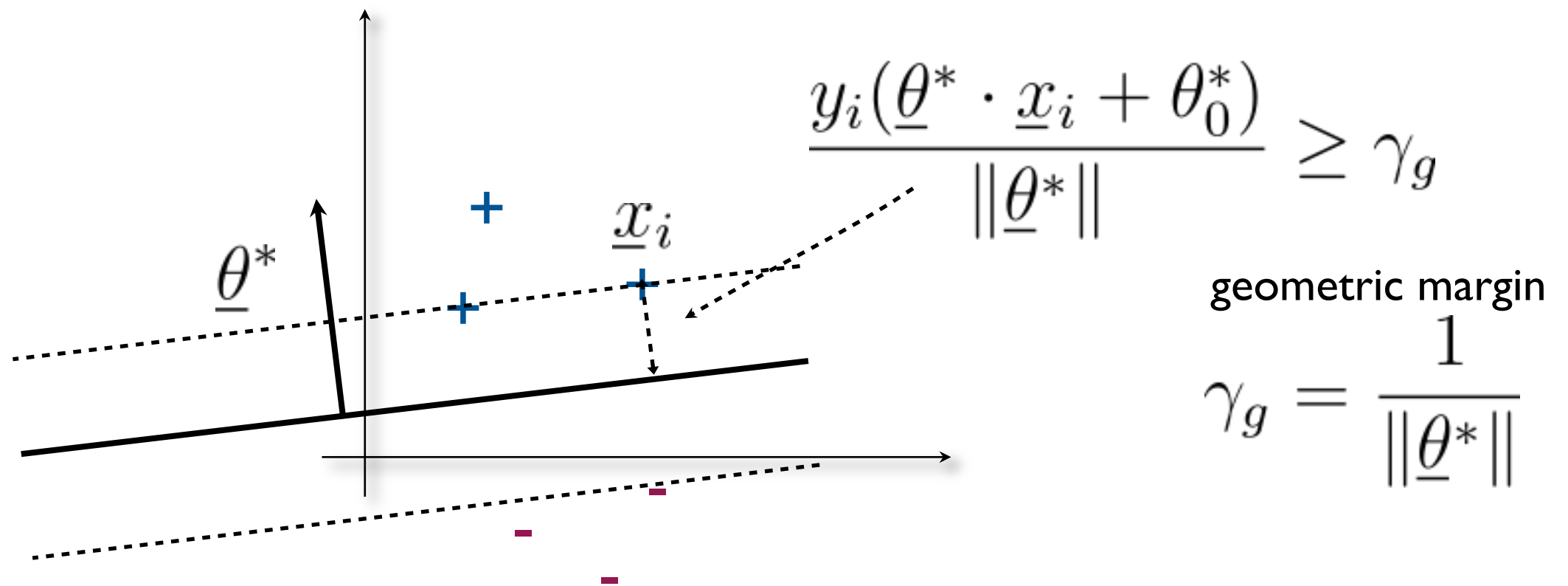
Linear classifiers (with offset)

- 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



To find $\underline{\theta}^*, \theta_0^*$: minimize $\frac{1}{2} \|\underline{\theta}\|^2$ subject to

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

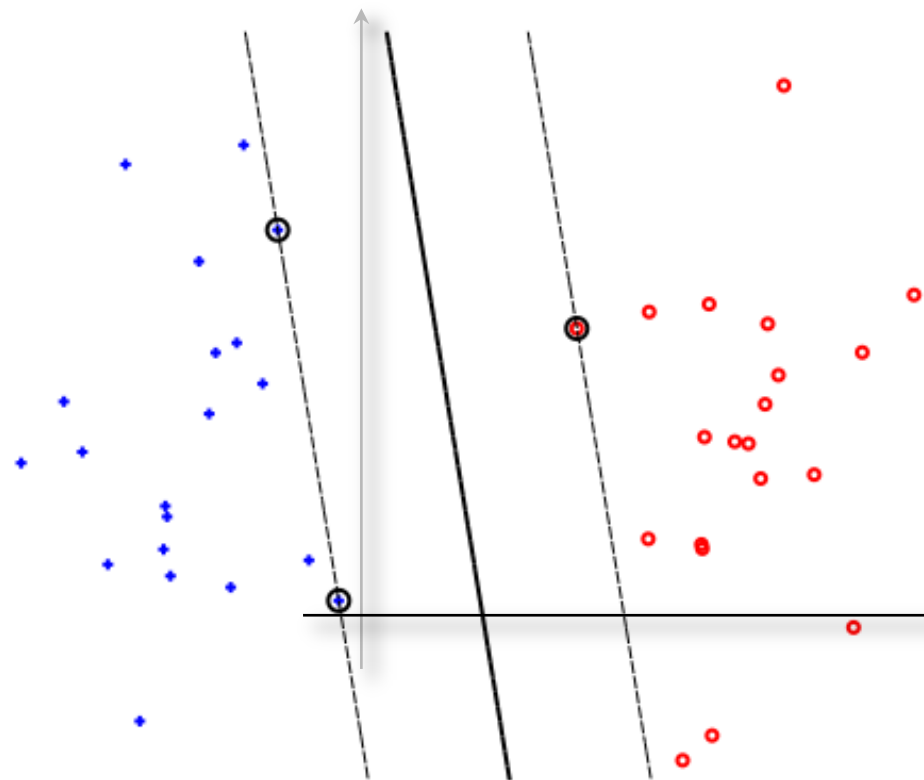
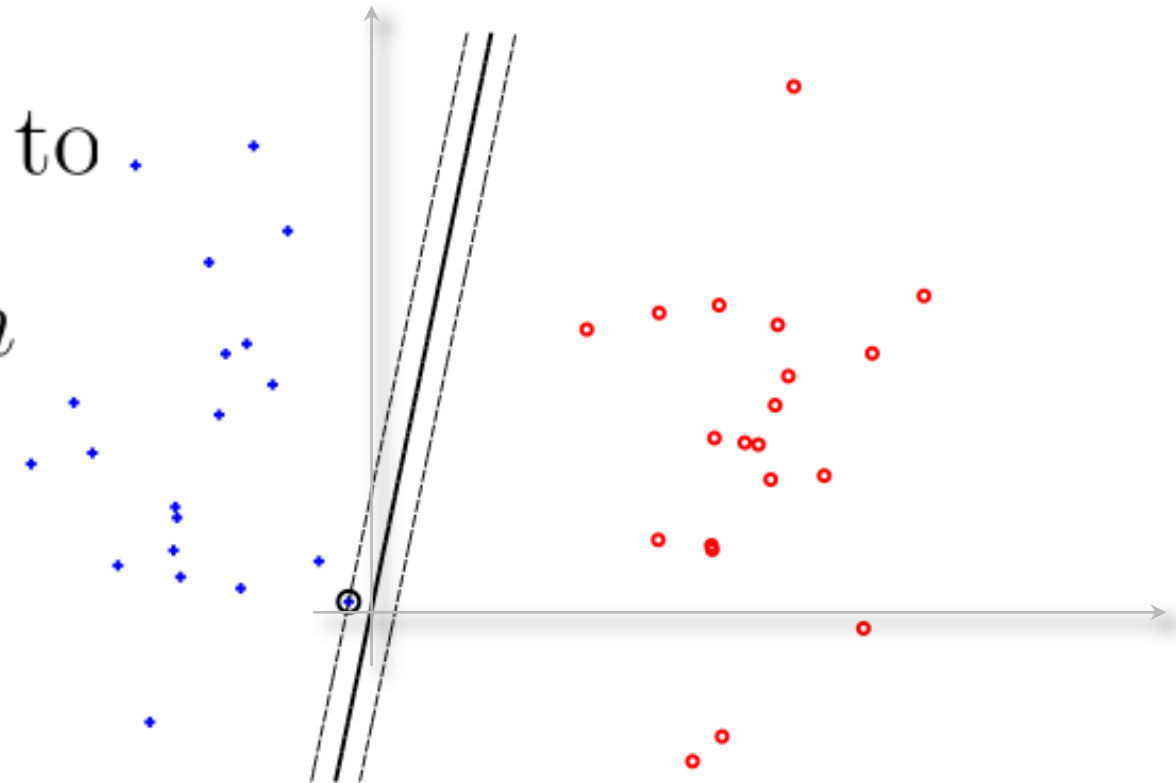
- Still a quadratic programming problem (quadratic objective, linear constraints)

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

- A desirable property
 - maximizes the margin on the training set (\approx good generalization, since it will still predict correctly for points between the decision boundary and the dotted line)
- 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

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

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

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

slack variables
permit us to violate
some of the margin
constraints

Support vector machine

- Relaxed quadratic optimization problem

penalty for constraint violation

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

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

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

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

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

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

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

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 $1/\|\underline{\theta}^*\|$

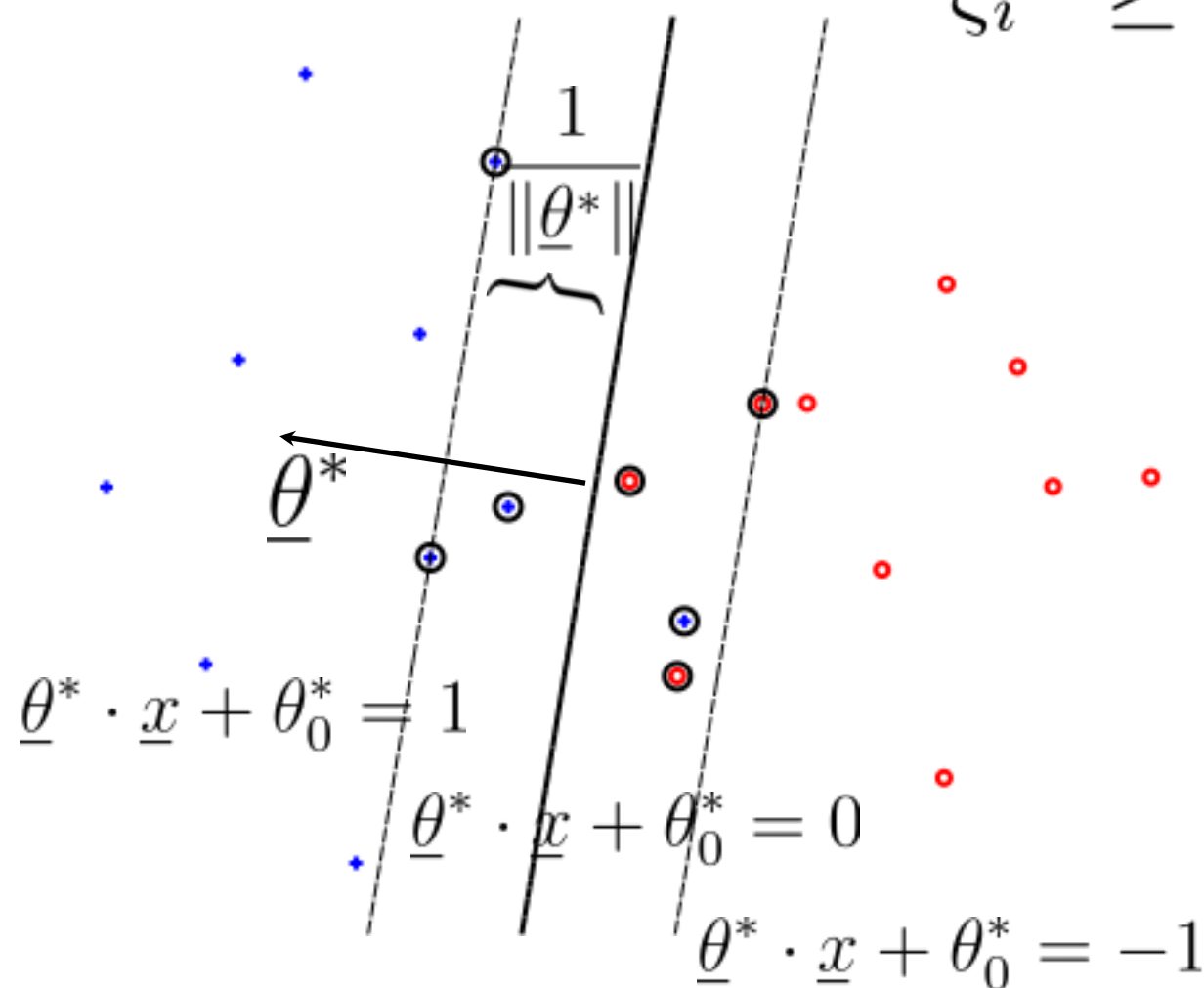
Support vector machine

- Relaxed quadratic optimization problem

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

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

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



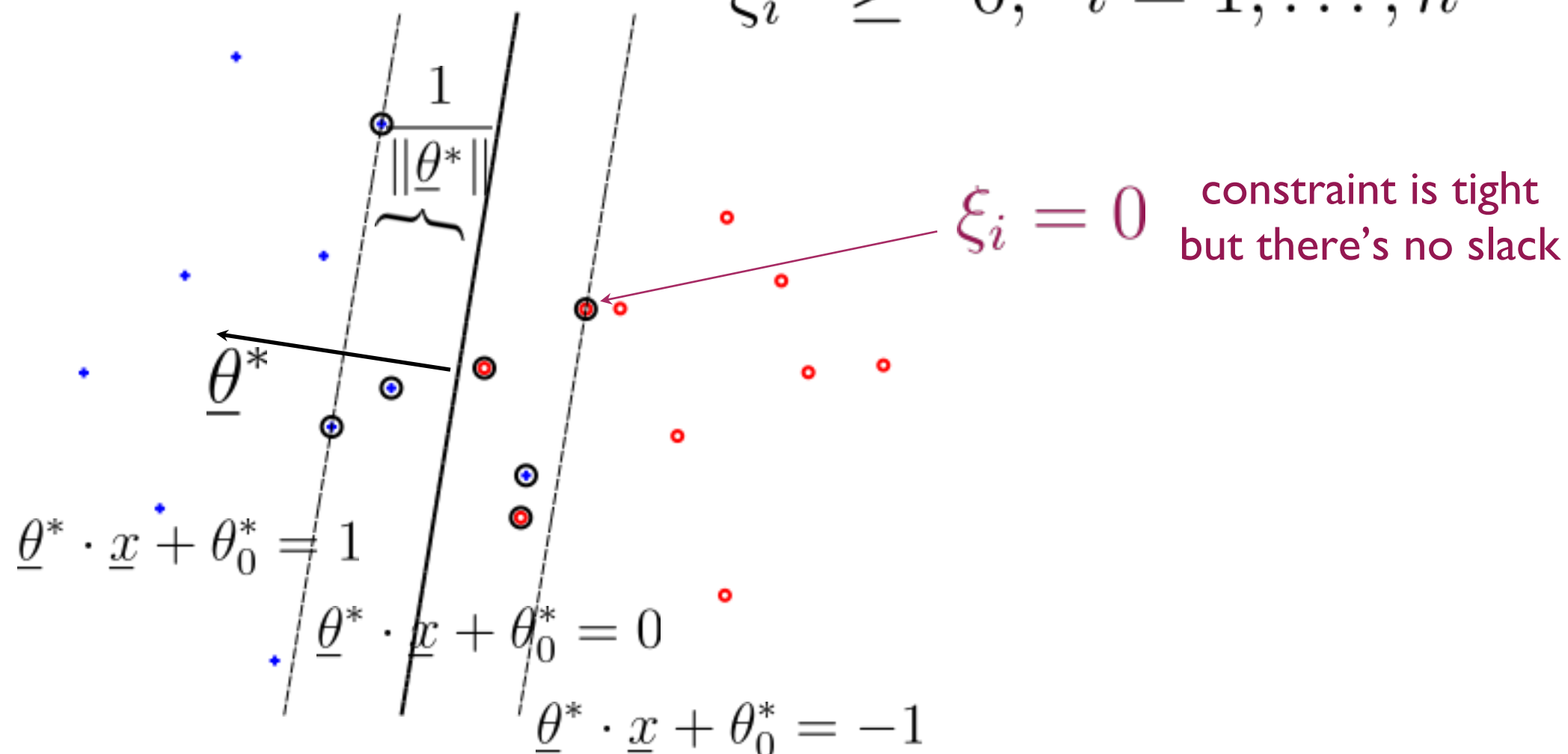
Support vectors and slack

- The solution now has three types of support vectors

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

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

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



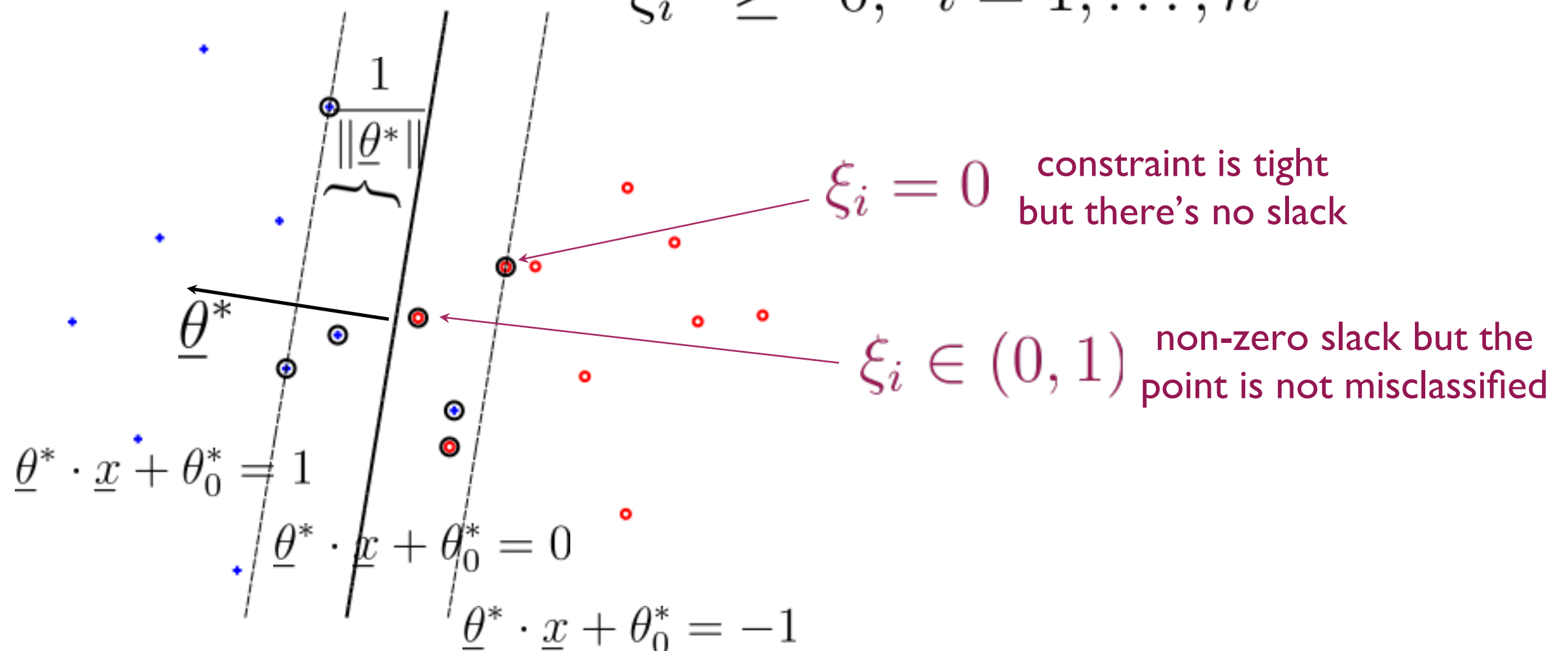
Support vectors and slack

- The solution now has three types of support vectors

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

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

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



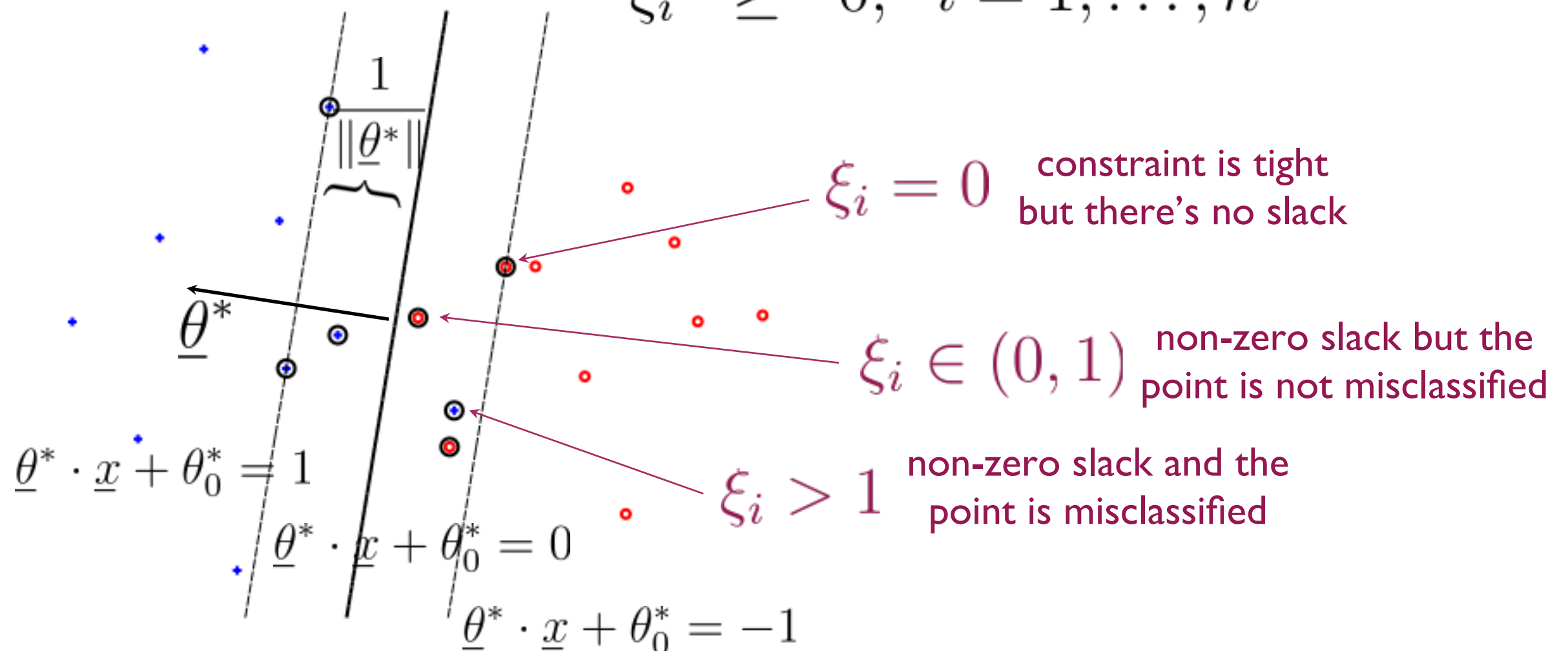
Support vectors and slack

- The solution now has three types of support vectors

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

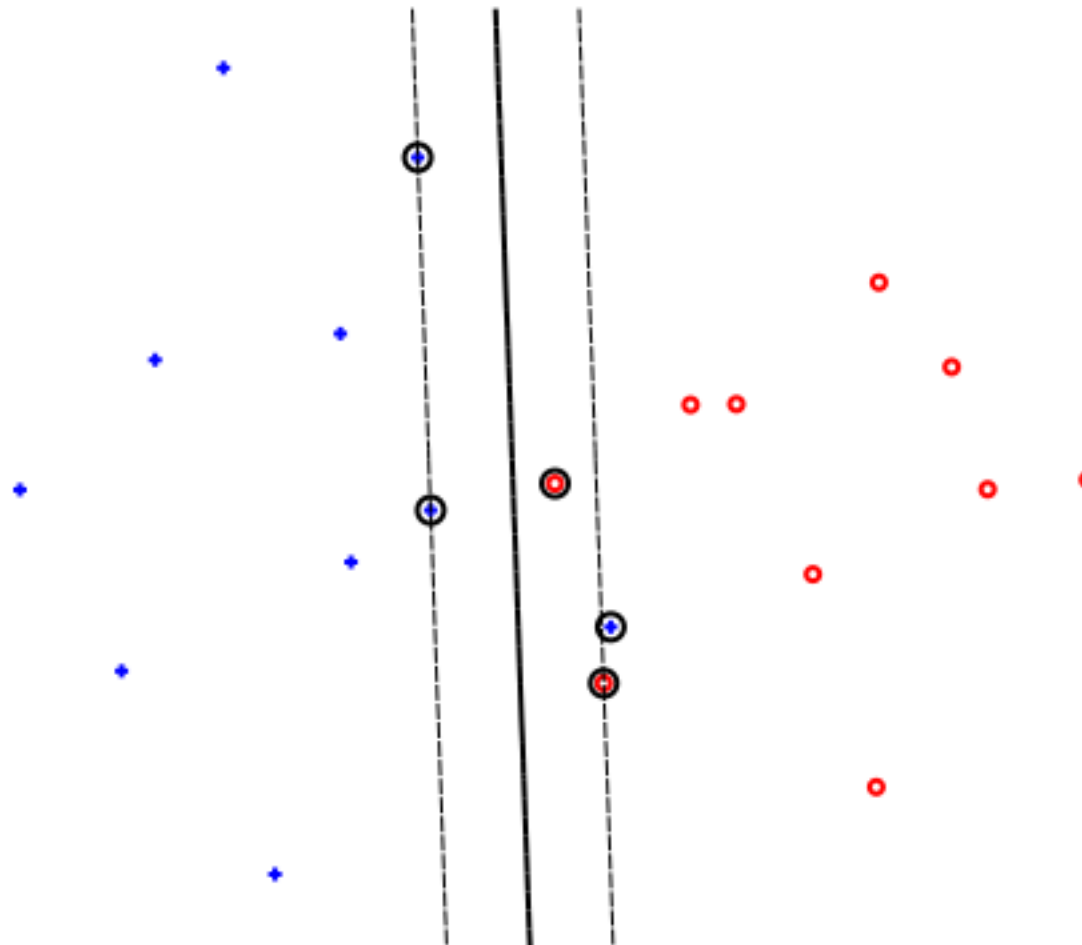
$$y_i(\underline{\theta} \cdot \underline{x}_i + \theta_0) \geq 1 - \xi_i, \quad i = 1, \dots, n$$

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



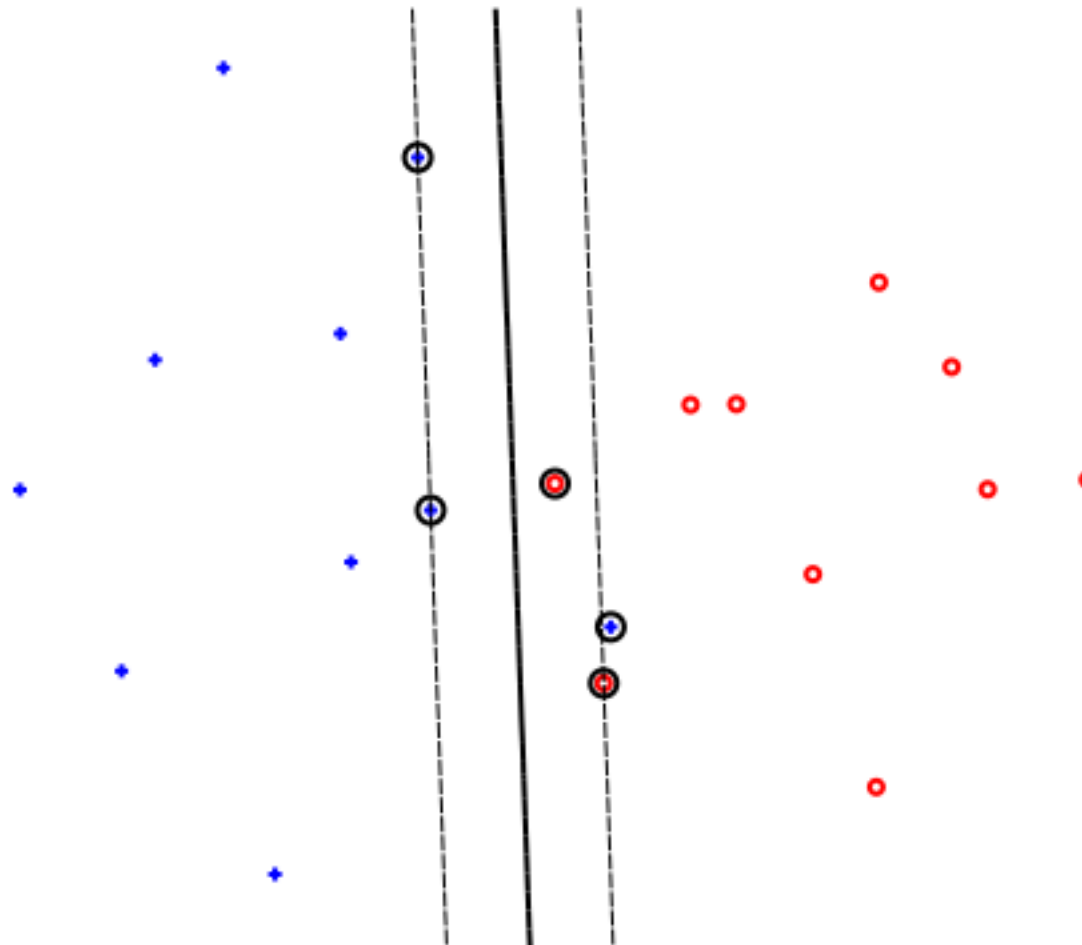
Examples

- $C=100$



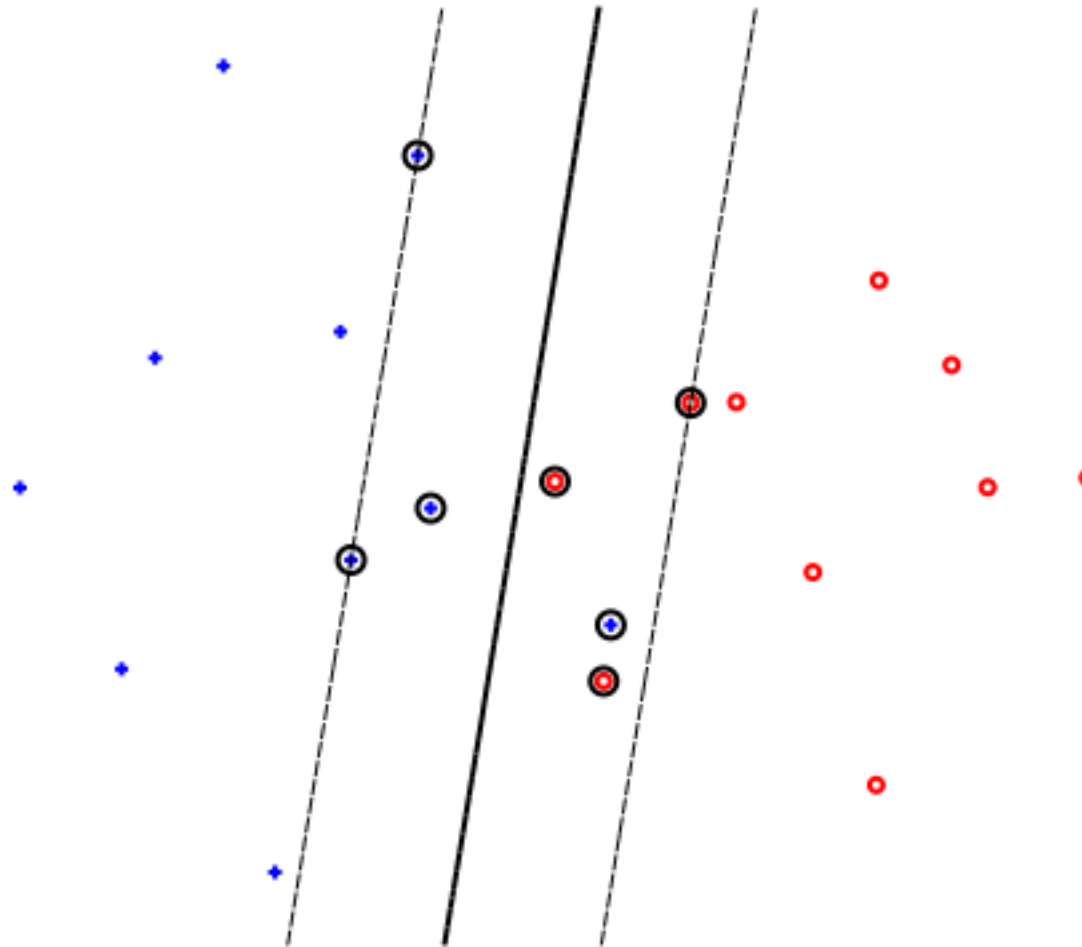
Examples

- $C=10$



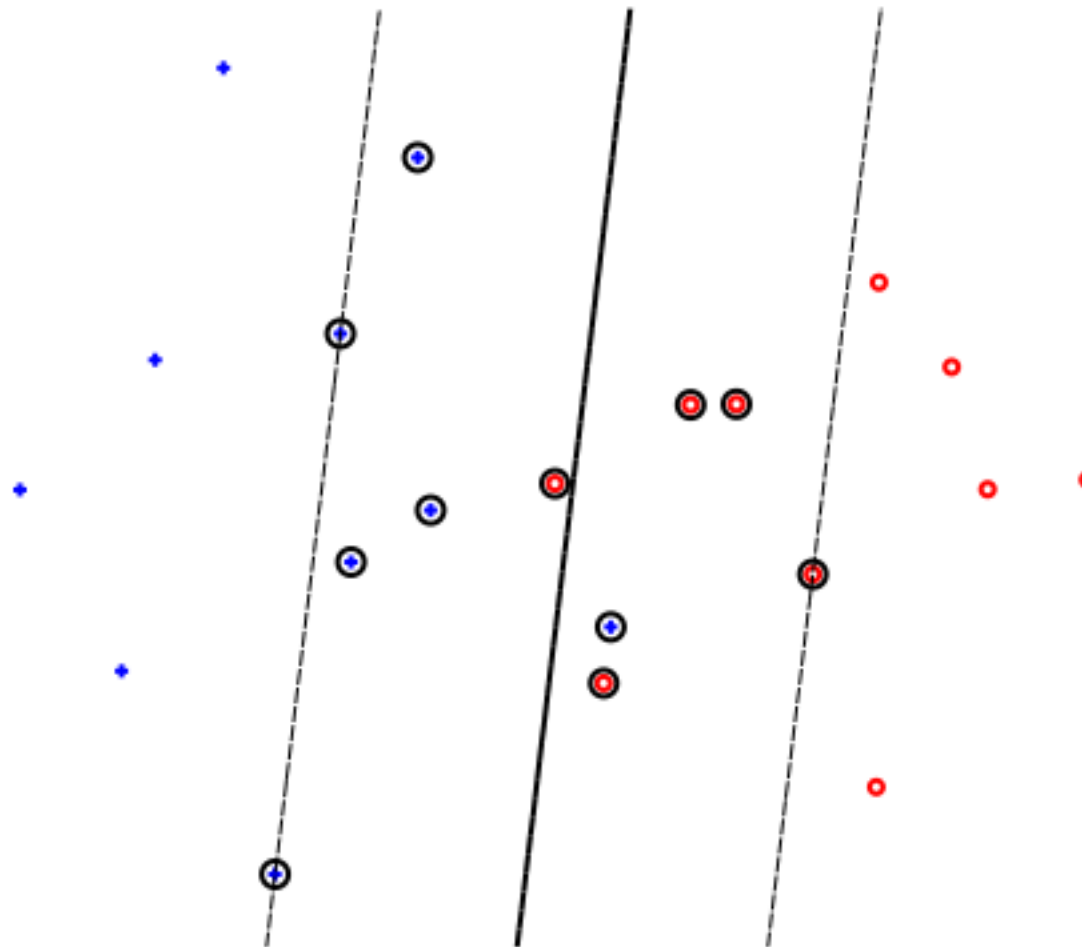
Examples

- $C=1$



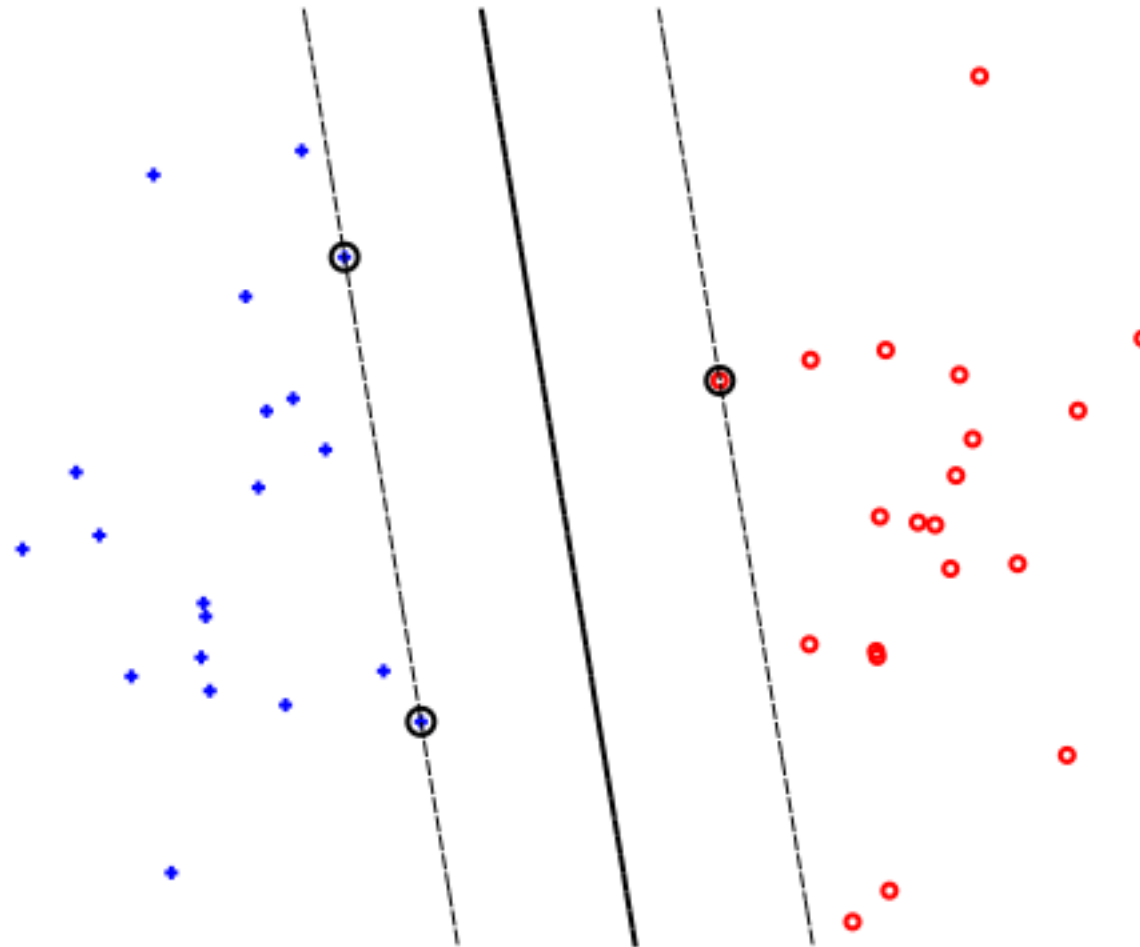
Examples

- $C=0.1$



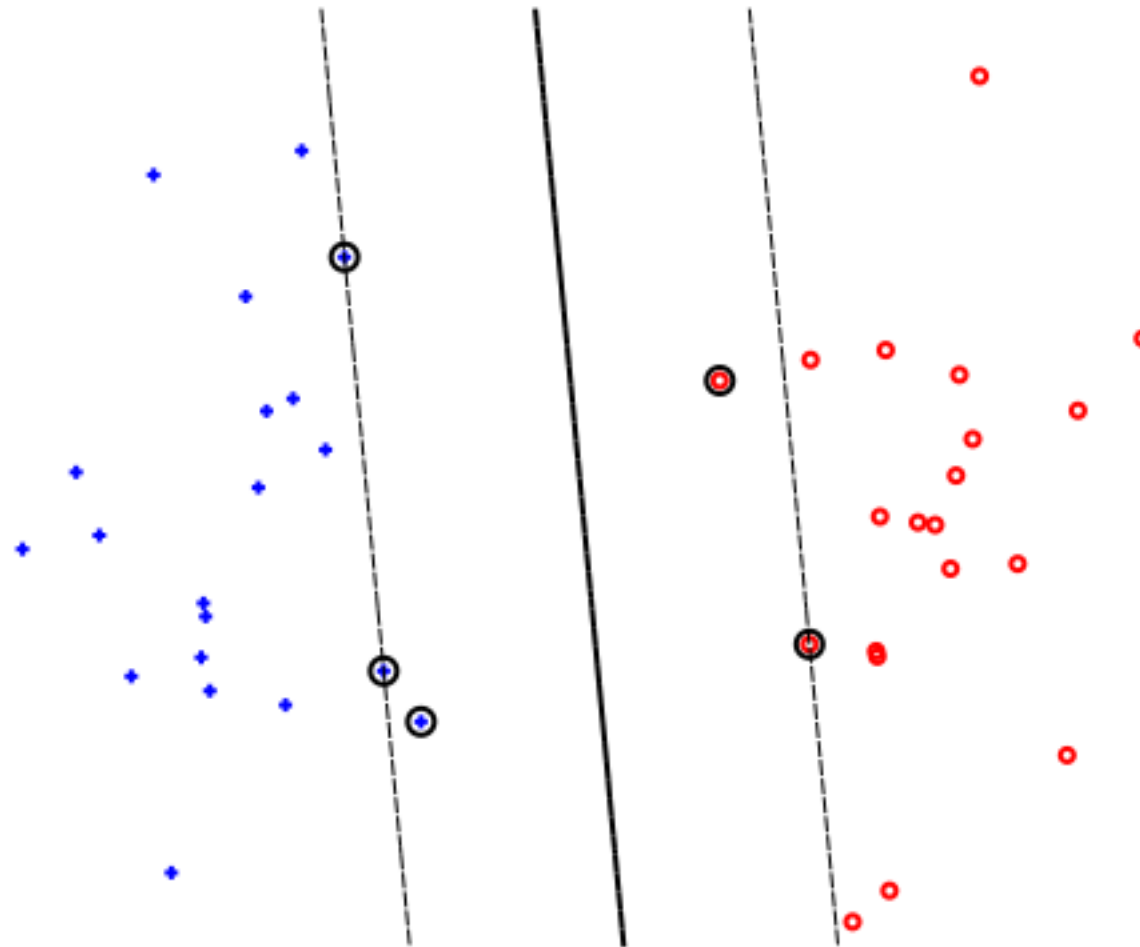
Examples

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



Examples

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



Examples

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

