

CS578 Statistical Machine Learning Lecture I

Jean Honorio
Purdue University

(based on slides by Tommi Jaakkola, MIT CSAIL)

First things first...

- Website: <http://www.cs.purdue.edu/homes/jhonorio/20spring-cs57800.html>
 - TAs
 - Textbooks
 - Assignments
 - Grading
 - Late policy
 - Schedule
- Office hours: Doodle on Friday, we will decide on Monday
- Piazza will be made available soon
- **Answer prerequisite survey sent on January 6**
- **Solve Homework 0, due on January 16**

Course topics

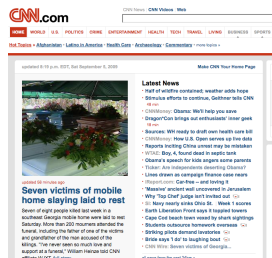
- Supervised learning
 - linear and non-linear classifiers, kernels
 - rating, ranking, collaborative filtering
 - model selection, complexity, generalization
 - conditional Random fields, structured prediction
- Unsupervised learning, modeling
 - mixture models, topic models
 - Hidden Markov Models
 - Bayesian networks
 - Markov Random Fields, factor graphs

Machine learning

- Learning from examples



+ |



- |



+ |



- |

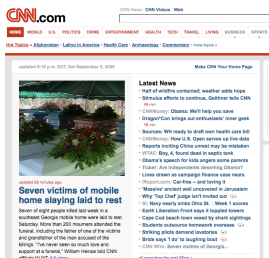
training set

Machine learning

- Learning from examples



+ |



- |



+ |



- |

training set

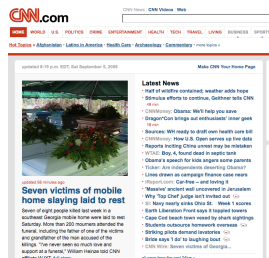
The training set of labeled examples specifies the learning task only implicitly

Machine learning

- Learning from examples



+ |



- |



+ |



- |



?

Our goal is to accurately label new websites that were not part of the training set

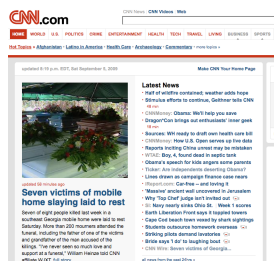
training set

Machine learning

- Learning from examples



+ |



+ |



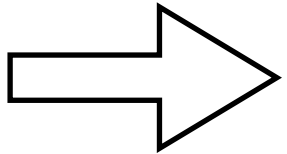
- |



training set

predicted label new website

- |



$$y = f \left(\begin{array}{c} \text{new website} \end{array} \right)$$



classifier = mapping
from websites to labels

“Examples”

- We will have to first represent the examples (websites) in a manner that can be easily mapped to labels

news article

White House
officials consulted
with the Justice
Department in
preparing a list of
U.S. attorneys who
would be removed.

(NYT 03/13/07)

“Examples”

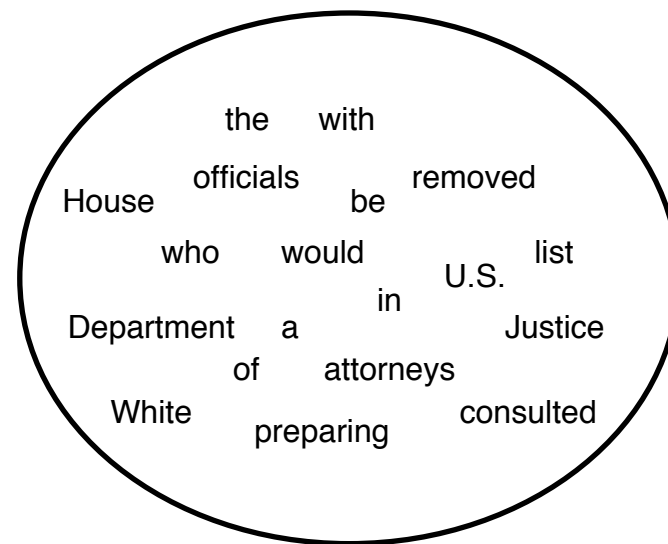
- We will have to first represent the examples (websites) in a manner that can be easily mapped to labels

news article

White House
officials consulted
with the Justice
Department in
preparing a list of
U.S. attorneys who
would be removed.

(NYT 03/13/07)

bag of
words



“Examples”

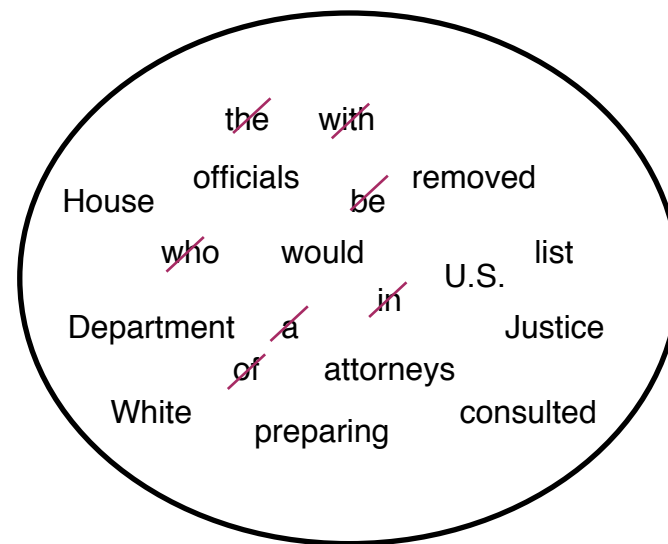
- We will have to first represent the examples (websites) in a manner that can be easily mapped to labels

news article

White House
officials consulted
with the Justice
Department in
preparing a list of
U.S. attorneys who
would be removed.

(NYT 03/13/07)

bag of
words



“Examples”

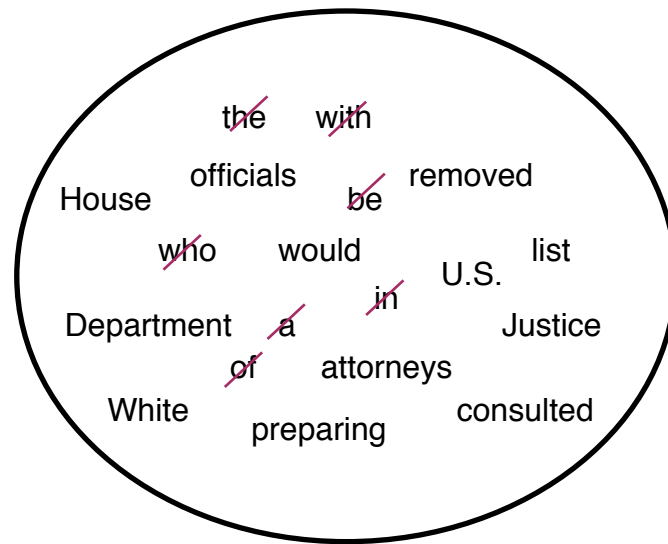
- We will have to first represent the examples (websites) in a manner that can be easily mapped to labels

news article

White House
officials consulted
with the Justice
Department in
preparing a list of
U.S. attorneys who
would be removed.

(NYT 03/13/07)

bag of
words



counts

0	politics
1	Justice
0	government
0	president
1	House
...	...

$$\underline{x} \in \mathcal{R}^d$$

a vector whose coordinates (features)
specify how many times (or whether)
particular words appeared in the article

$$y \in \{-1, 1\}$$

the label of the website

Classifier

- Intuitively:

predicted label new website

$$y = f \left(\begin{array}{c} \text{The Boston Globe} \\ \text{PCB risk feared at older N.E. schools} \\ \text{He could not leave a comrade behind} \\ \text{As Joe Kennedy considers} \end{array} \right)$$

- Programming point of view:

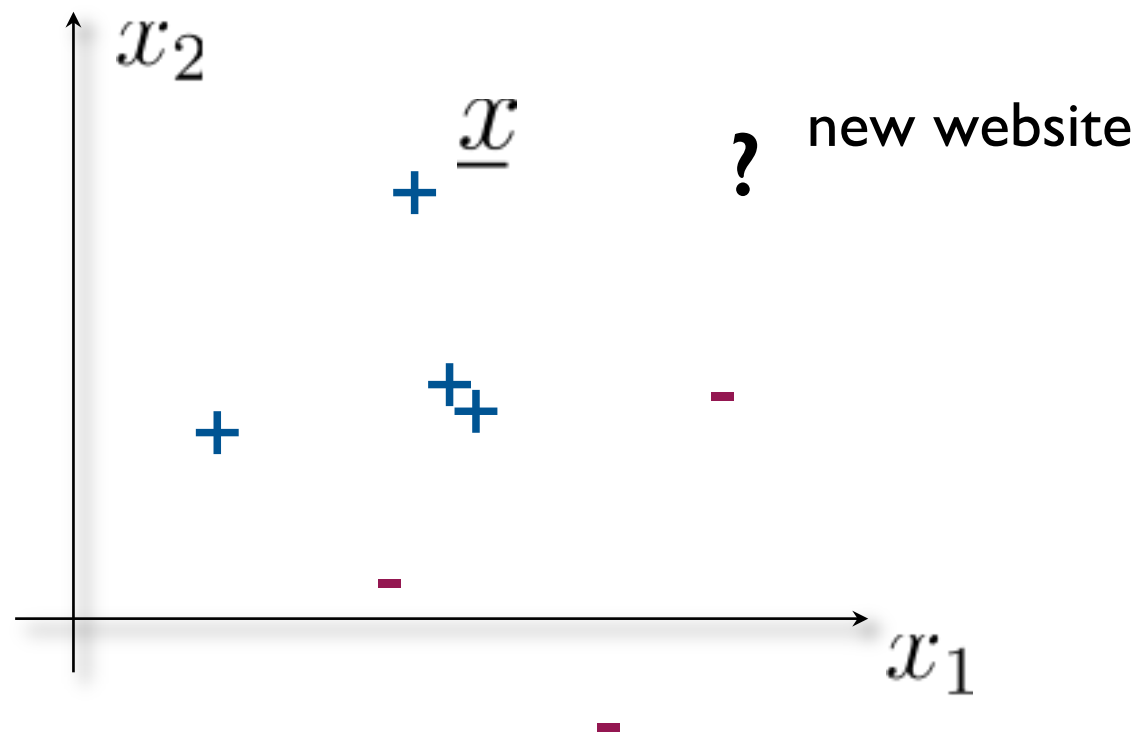
f is a piece of code that receives a vector $\underline{x} \in \mathcal{R}^d$ as input, and outputs a label $y \in \{-1, 1\}$

- Mathematical notation:

$$f : \mathcal{R}^d \rightarrow \{-1, 1\}$$

The learning task

- The training set is now a set of labeled points



- We seek for a “good” classifier $f : \mathcal{R}^d \rightarrow \{-1, 1\}$ based on the training set $D = \{\underline{x}_1, y_1, \underline{x}_2, y_2, \dots, \underline{x}_n, y_n\}$ so that $f(\underline{x})$ correctly labels any new websites \underline{x}

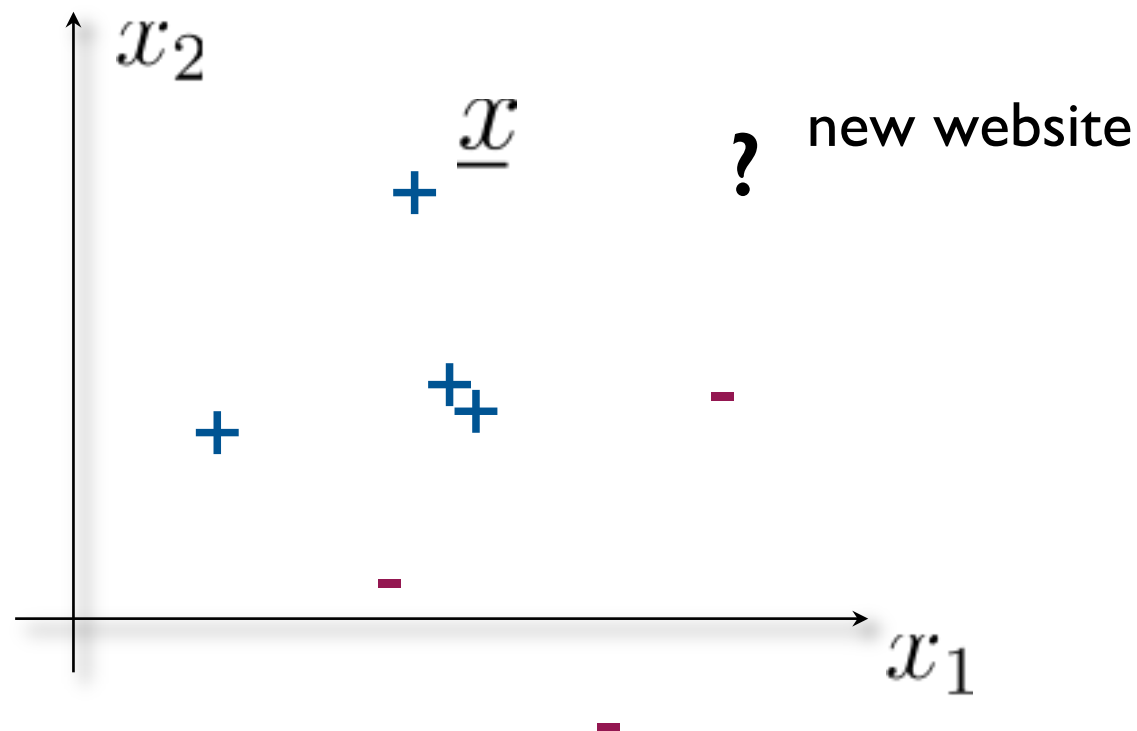
The learning task

- The training set is now a set of labeled points

Part I:

Model selection

what type of classifiers
should we consider?



- We seek for a “good” classifier $f : \mathcal{R}^d \rightarrow \{-1, 1\}$ based on the training set $D = \{\underline{x}_1, y_1, \underline{x}_2, y_2, \dots, \underline{x}_n, y_n\}$ so that $f(\underline{x})$ correctly labels any new websites \underline{x}

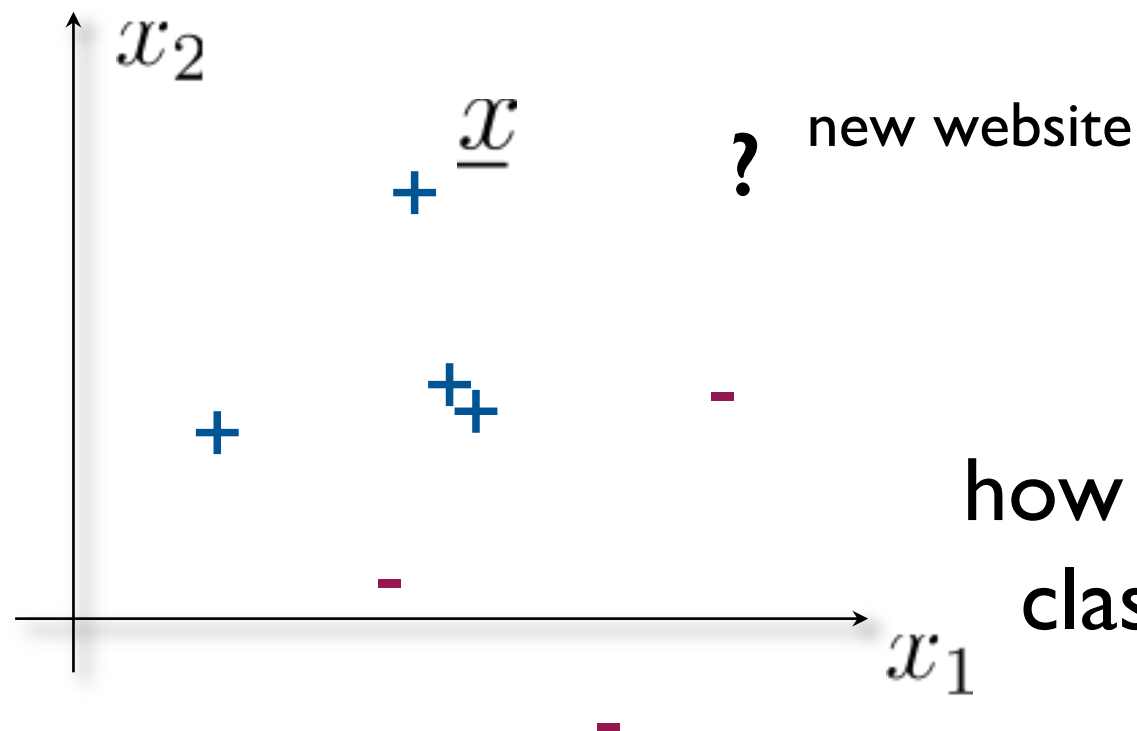
The learning task

- The training set is now a set of labeled points

Part 1:

Model selection

what type of classifiers
should we consider?



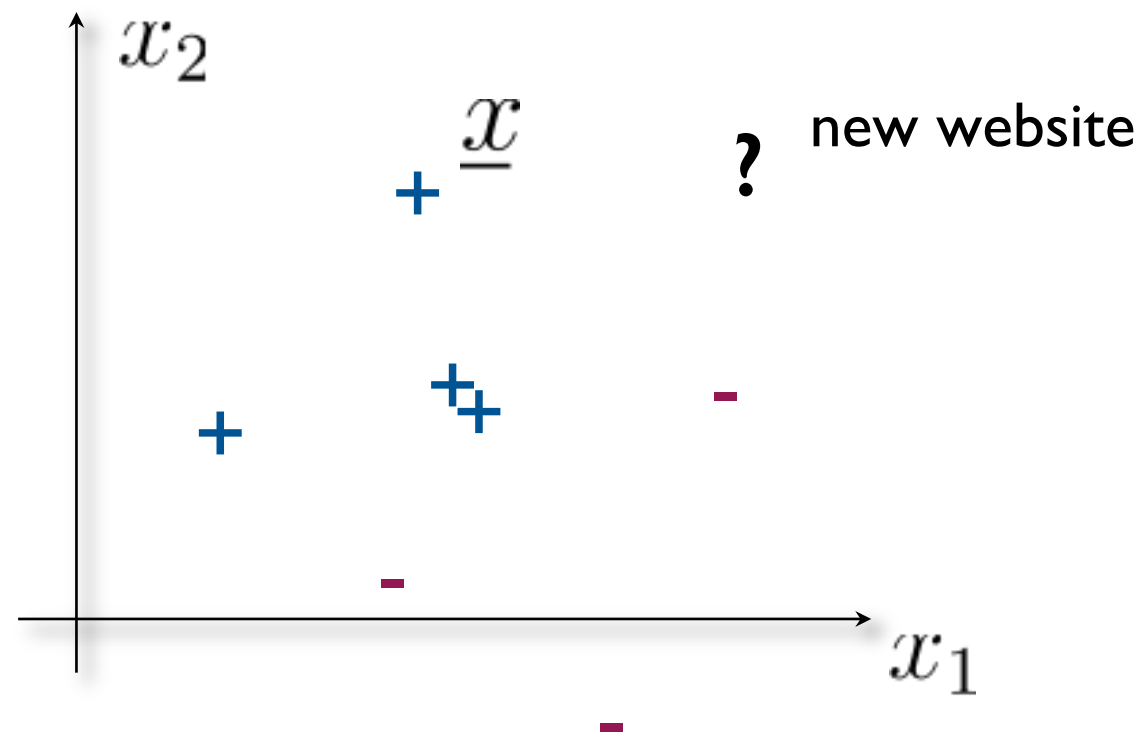
Part 2:

Estimation

how to select the best
classifier in the set?

- We seek for a “good” classifier $f : \mathcal{R}^d \rightarrow \{-1, 1\}$
based on the training set $D = \{\underline{x}_1, y_1, \underline{x}_2, y_2, \dots, \underline{x}_n, y_n\}$
so that $f(\underline{x})$ correctly labels any new websites \underline{x}

Part I: allowing all classifiers?

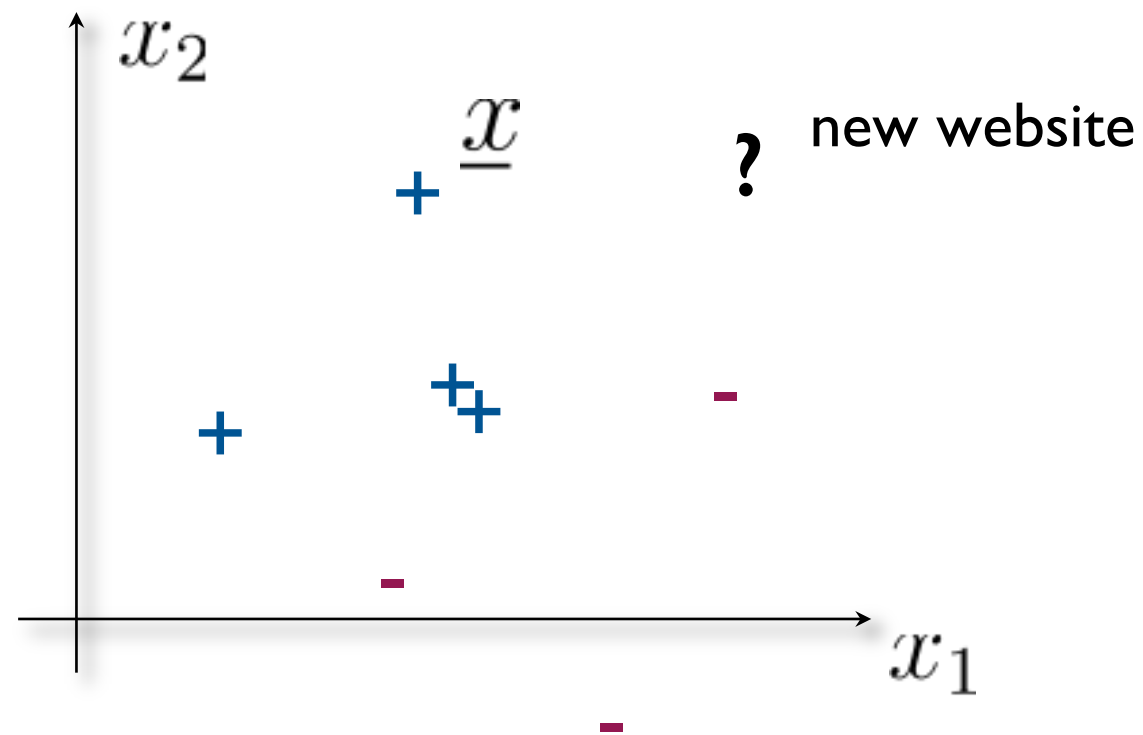


- We can easily construct a “silly classifier” that perfectly classifies any distinct set of training points

$$f(\underline{x}) = \begin{cases} y_i, & \text{if } \underline{x} = \underline{x}_i \text{ for some } i \\ -1, & \text{otherwise} \end{cases}$$

- But it doesn’t “generalize” (it doesn’t classify new points very well)

Part I: allowing few classifiers?

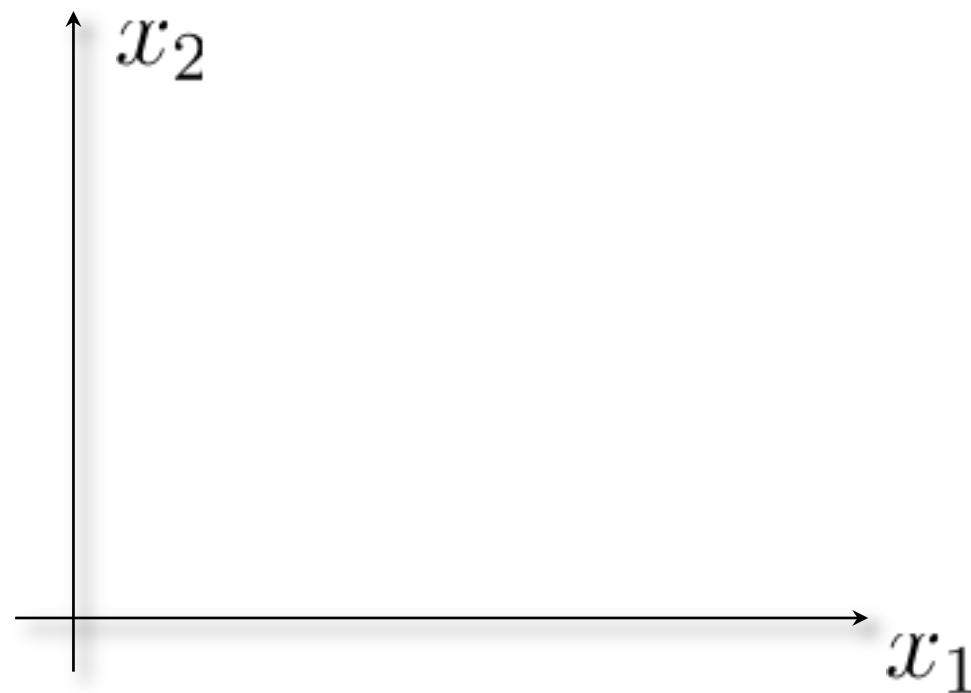


- We could instead consider very few alternatives such as
$$f(\underline{x}) = 1, \quad \text{for all } \underline{x} \in \mathcal{R}^d \text{ or}$$
$$f(\underline{x}) = -1, \quad \text{for all } \underline{x} \in \mathcal{R}^d$$
- But neither one classifies even training points very well

Part I: linear classifiers

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

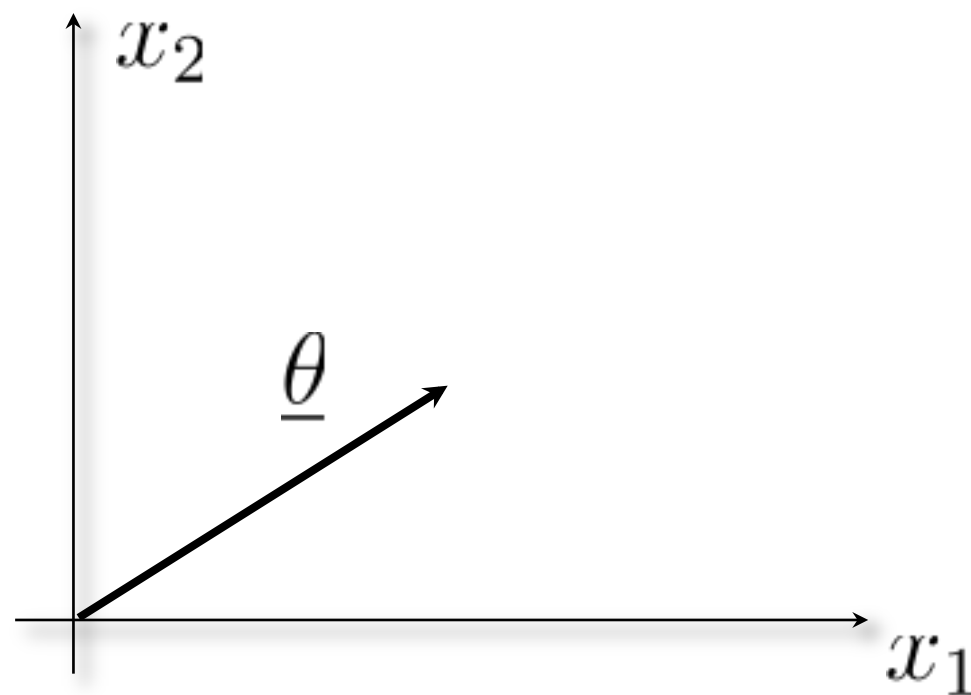
$$\begin{aligned} f(\underline{x}; \underline{\theta}) &= \text{sign}(\underline{\theta} \cdot \underline{x}) = \text{sign}(\theta_1 x_1 + \dots + \theta_d x_d) \\ &= \begin{cases} +1, & \text{if } \underline{\theta} \cdot \underline{x} > 0 \\ -1, & \text{if } \underline{\theta} \cdot \underline{x} \leq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



Part I: linear classifiers

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

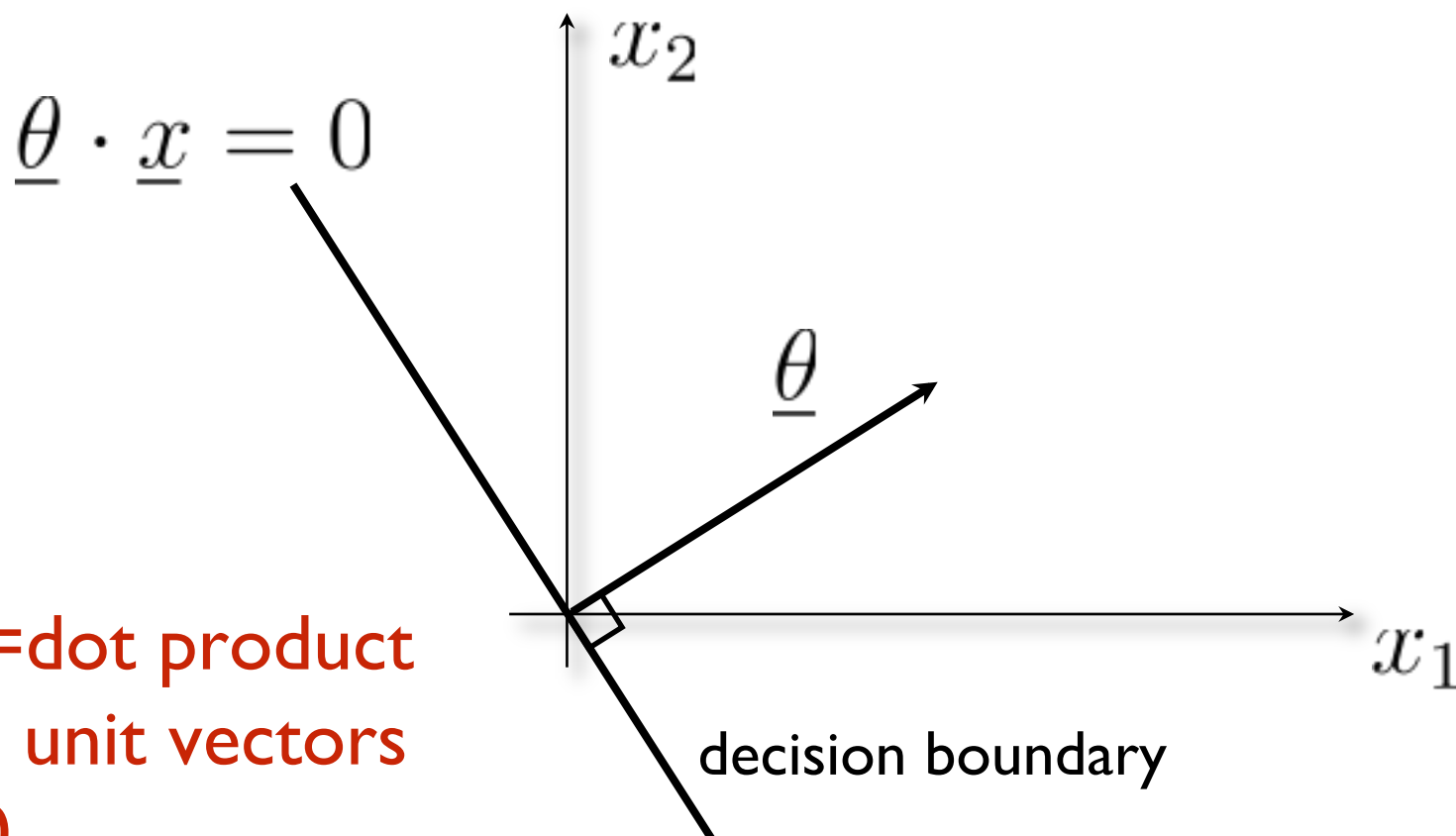
$$\begin{aligned} f(\underline{x}; \underline{\theta}) &= \text{sign}(\underline{\theta} \cdot \underline{x}) = \text{sign}(\theta_1 x_1 + \dots + \theta_d x_d) \\ &= \begin{cases} +1, & \text{if } \underline{\theta} \cdot \underline{x} > 0 \\ -1, & \text{if } \underline{\theta} \cdot \underline{x} \leq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



Part I: linear classifiers

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

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

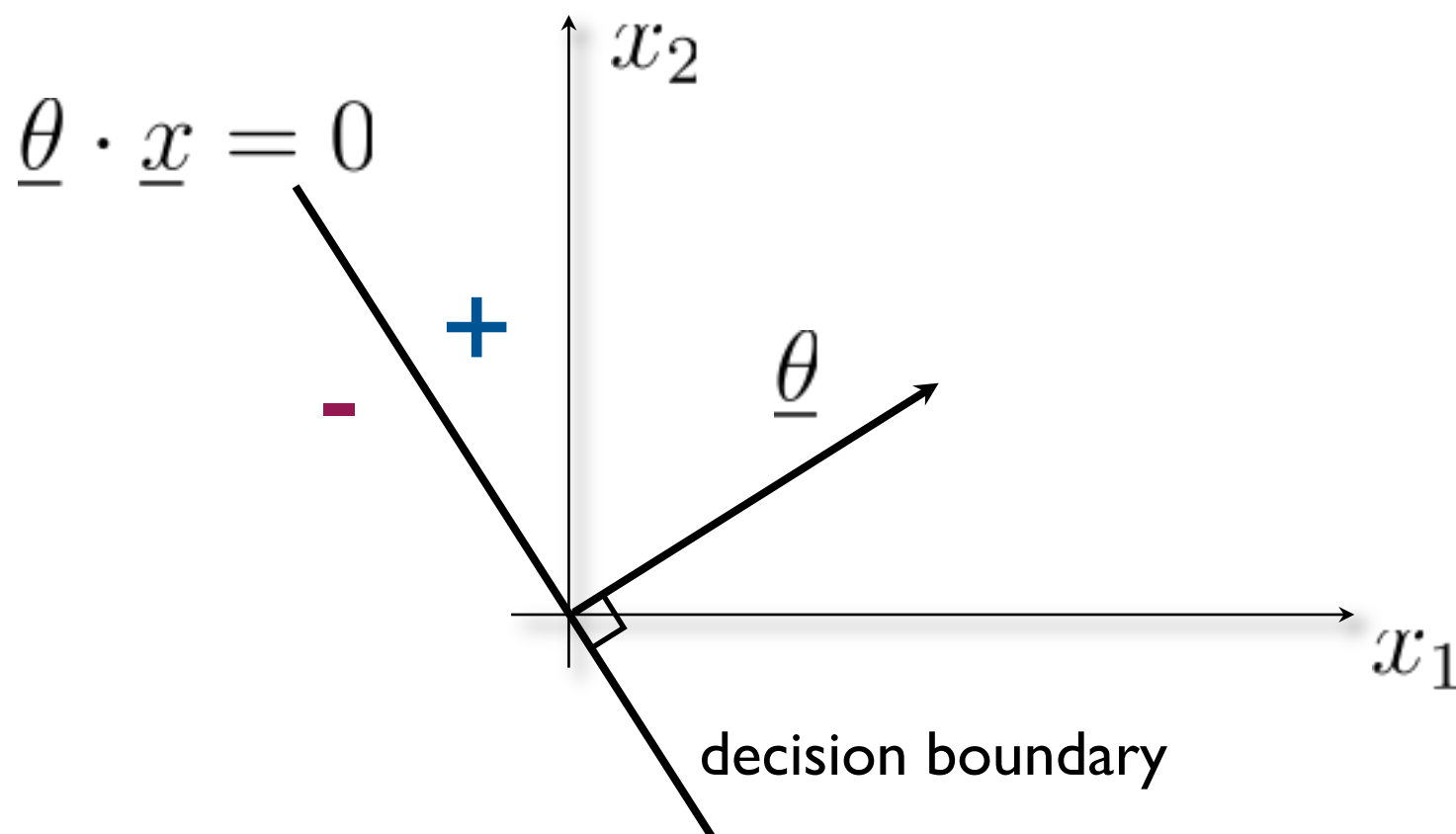


Recall:
 $\cos(\text{angle}) = \text{dot product of two unit vectors}$
 $\cos(90^\circ) = 0$

Part I: linear classifiers

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

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



Part 2: estimation

- We can use the training error as a surrogate criterion for finding the best linear classifier (through origin)

$$\hat{R}_n(\underline{\theta}) = \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, f(\underline{x}_i; \underline{\theta}))$$

where $\text{Loss}(y, y') = \begin{cases} 1, & \text{if } y \neq y' \\ 0, & \text{o.w.} \end{cases}$

- Other choices are possible (and often preferable)

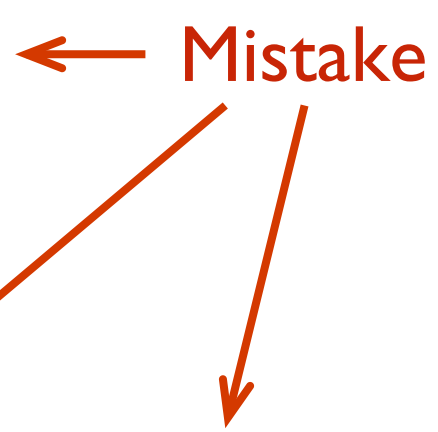
Perceptron algorithm

- We can use the training error as a surrogate criterion for finding the best linear classifier (through origin)

$$\hat{R}_n(\underline{\theta}) = \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, f(\underline{x}_i; \underline{\theta}))$$

$$\text{where } \text{Loss}(y, y') = \begin{cases} 1, & \text{if } y \neq y' \\ 0, & \text{o.w.} \end{cases}$$

- Since $f(\underline{x}; \underline{\theta}) = \text{sign}(\underline{\theta} \cdot \underline{x})$

$$\text{Loss}(y_t, f(\underline{x}_t; \underline{\theta})) = \begin{cases} 1, & y_t \neq \text{sign}(\underline{\theta} \cdot \underline{x}_t) \\ 0, & \text{o.w.} \end{cases}$$


- Checking for the condition $y_t \neq \text{sign}(\underline{\theta} \cdot \underline{x}_t)$ is equivalent to checking for the condition $y_t(\underline{\theta} \cdot \underline{x}_t) \leq 0$

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)

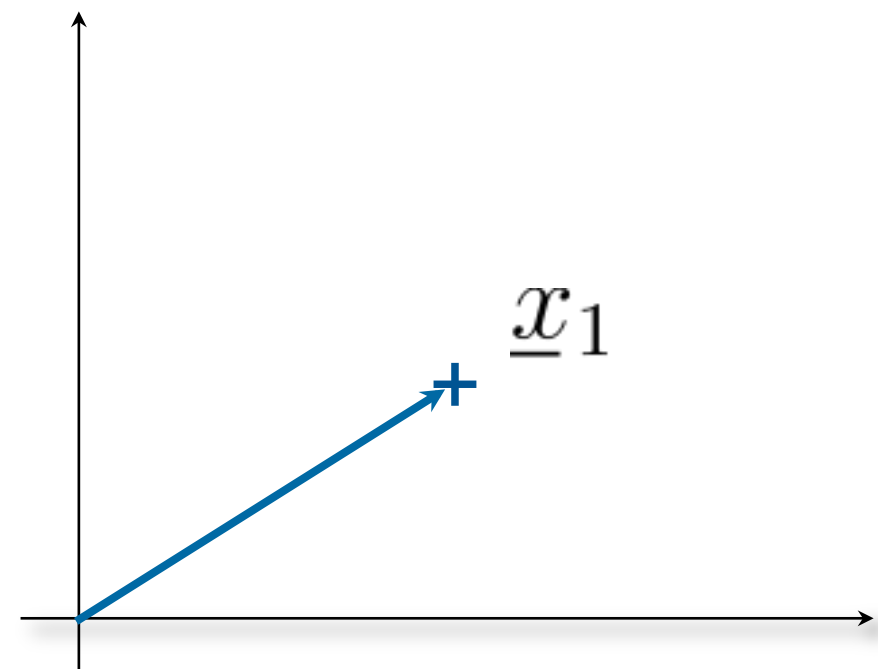
$\underline{\theta} \leftarrow \underline{\theta} + y_t \underline{x}_t$

- The algorithm will converge (no mistakes) if the training points are *linearly separable*, otherwise it won't converge

Perceptron algorithm

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

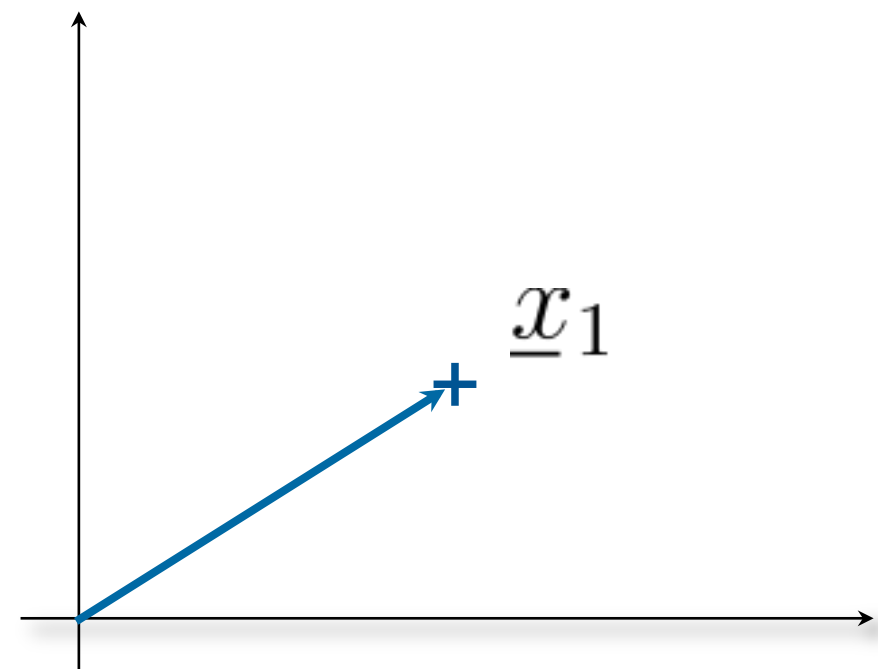


Perceptron algorithm

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$

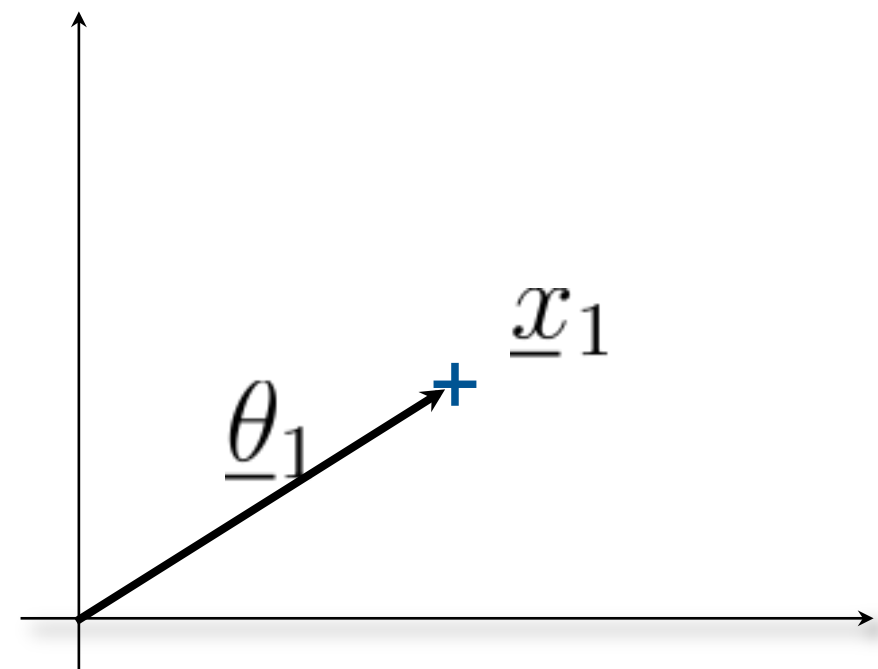


Perceptron algorithm

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$

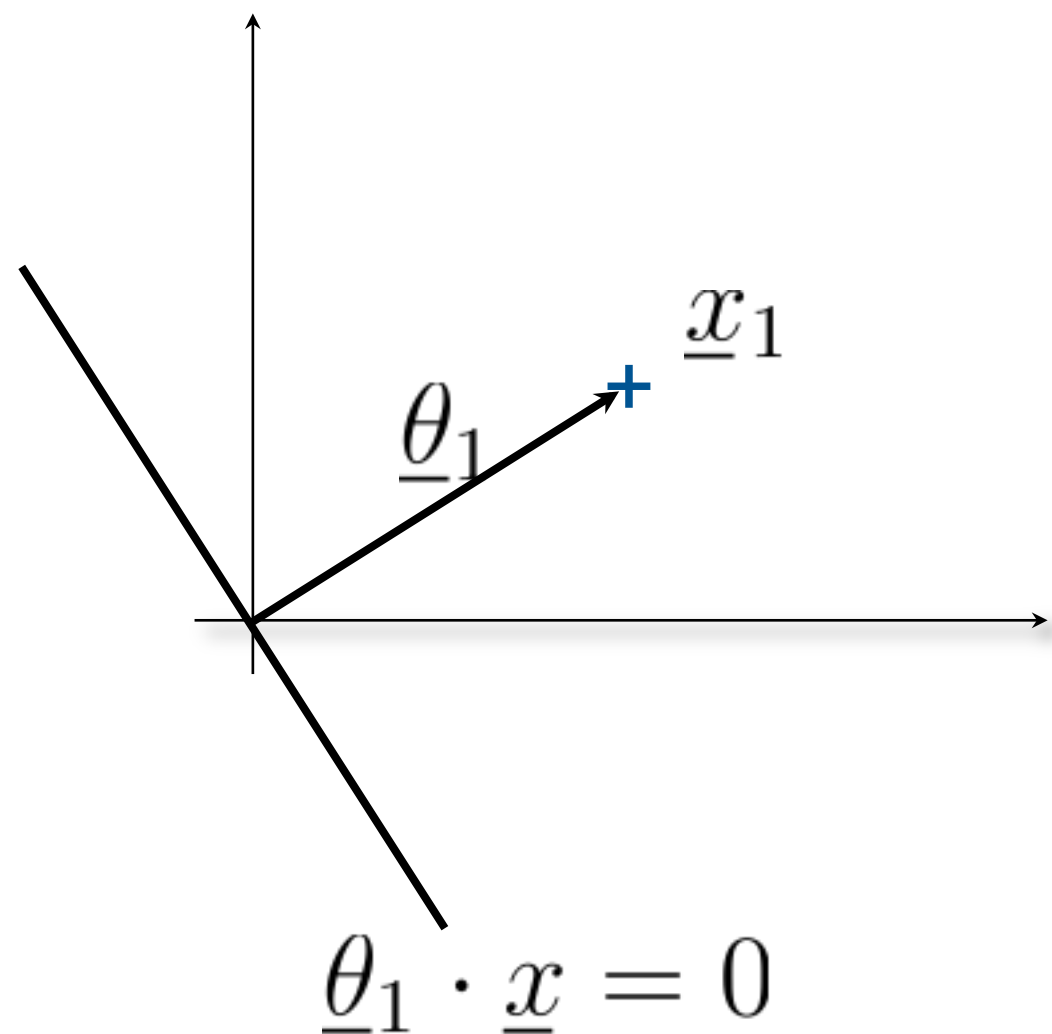


Perceptron algorithm

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$

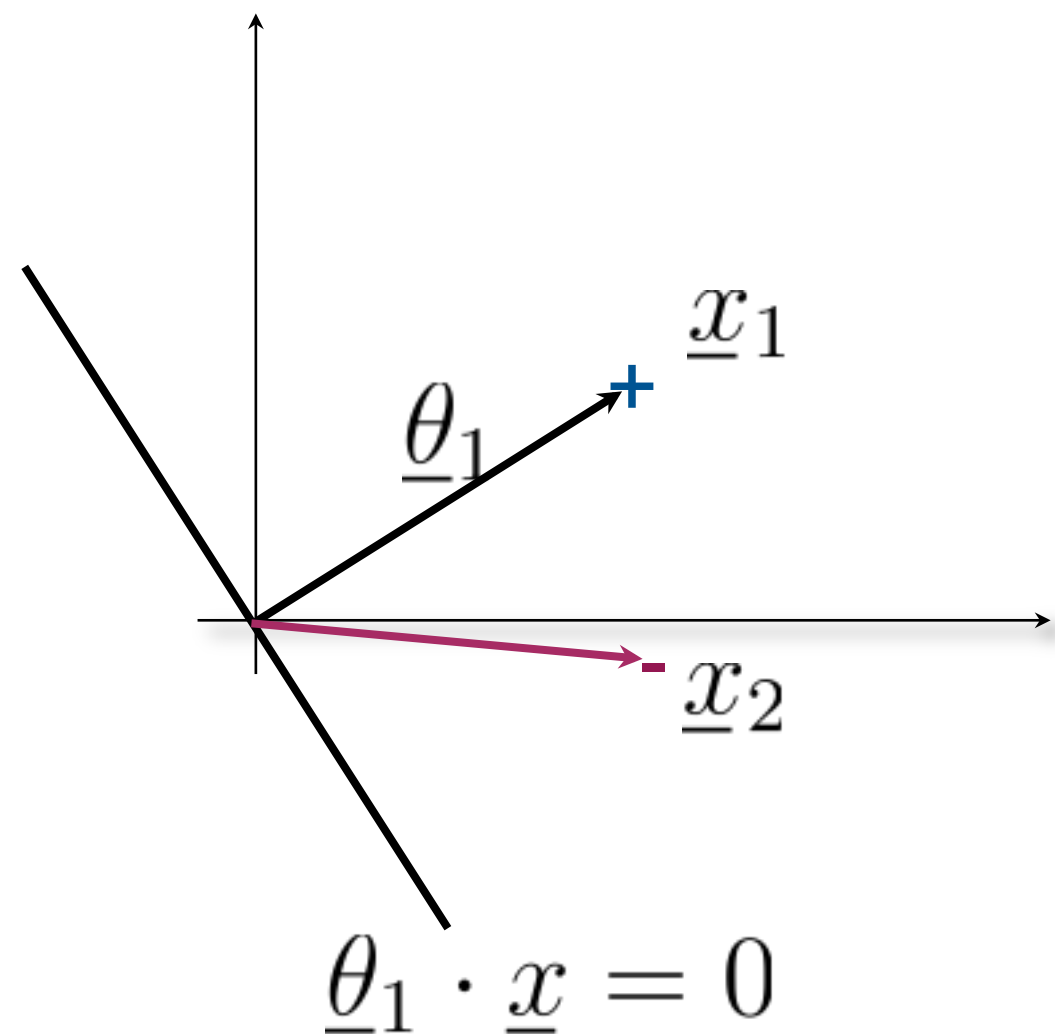


Perceptron algorithm

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$



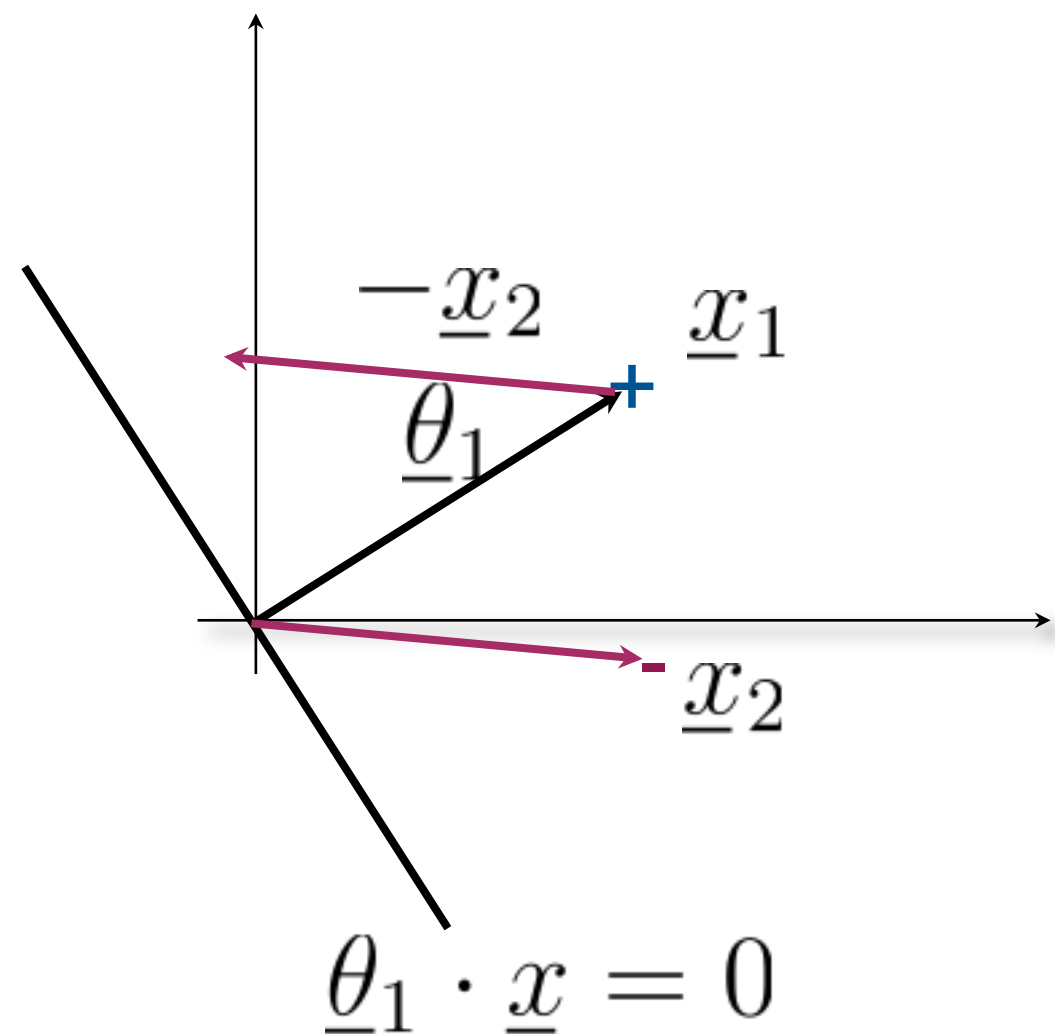
Perceptron algorithm

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$

$$\underline{\theta}_2 = \underline{\theta}_1 + (-1) \underline{x}_2$$



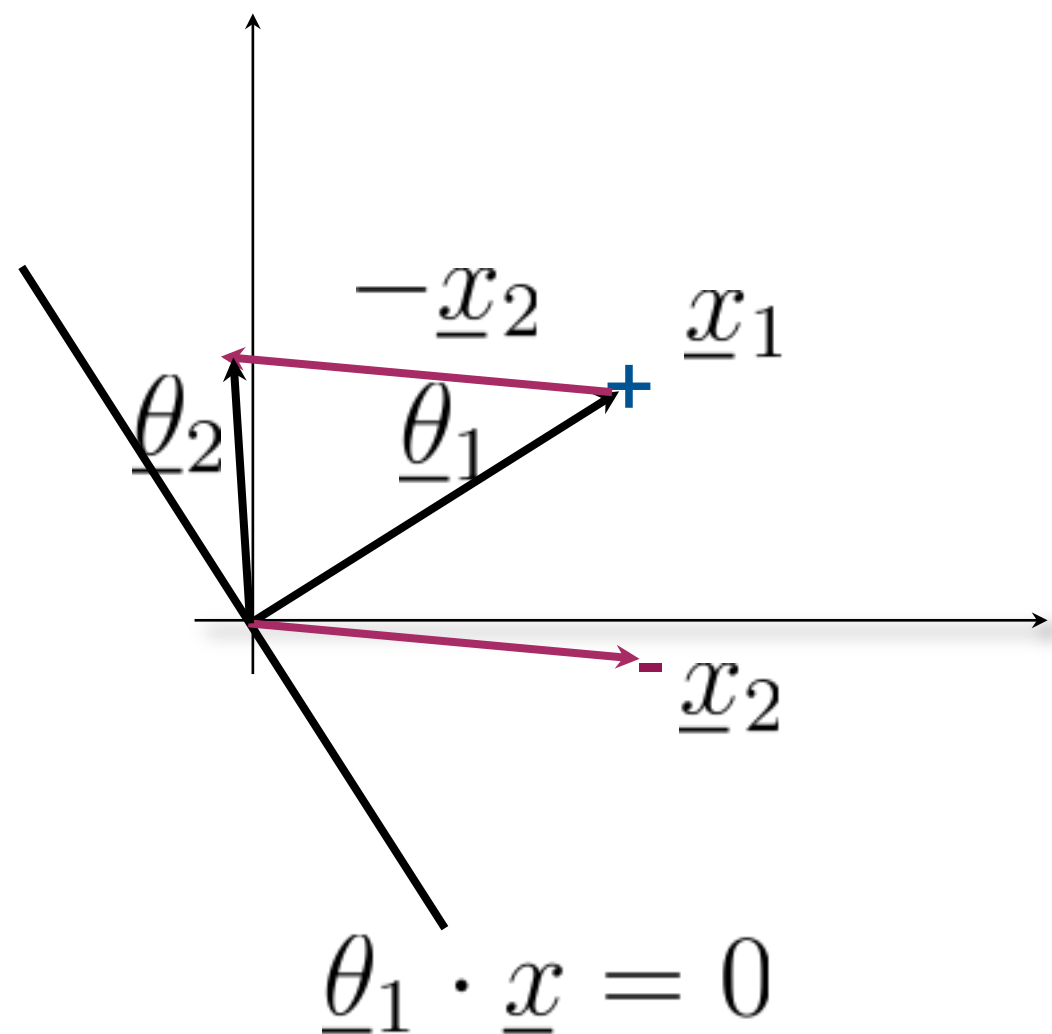
Perceptron algorithm

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$

$$\underline{\theta}_2 = \underline{\theta}_1 + (-1) \underline{x}_2$$



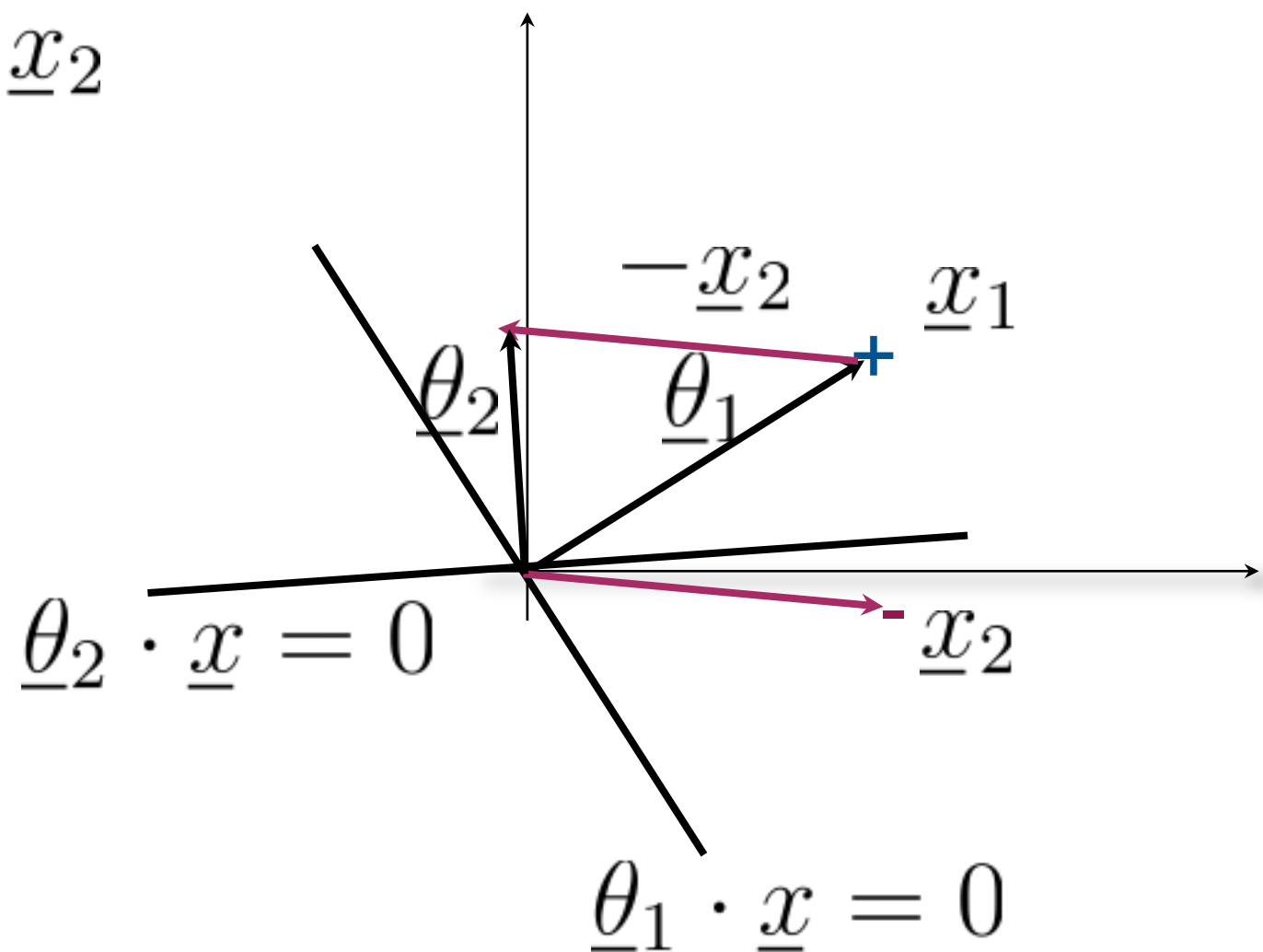
Perceptron algorithm

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$

$$\underline{\theta}_2 = \underline{\theta}_1 + (-1) \underline{x}_2$$



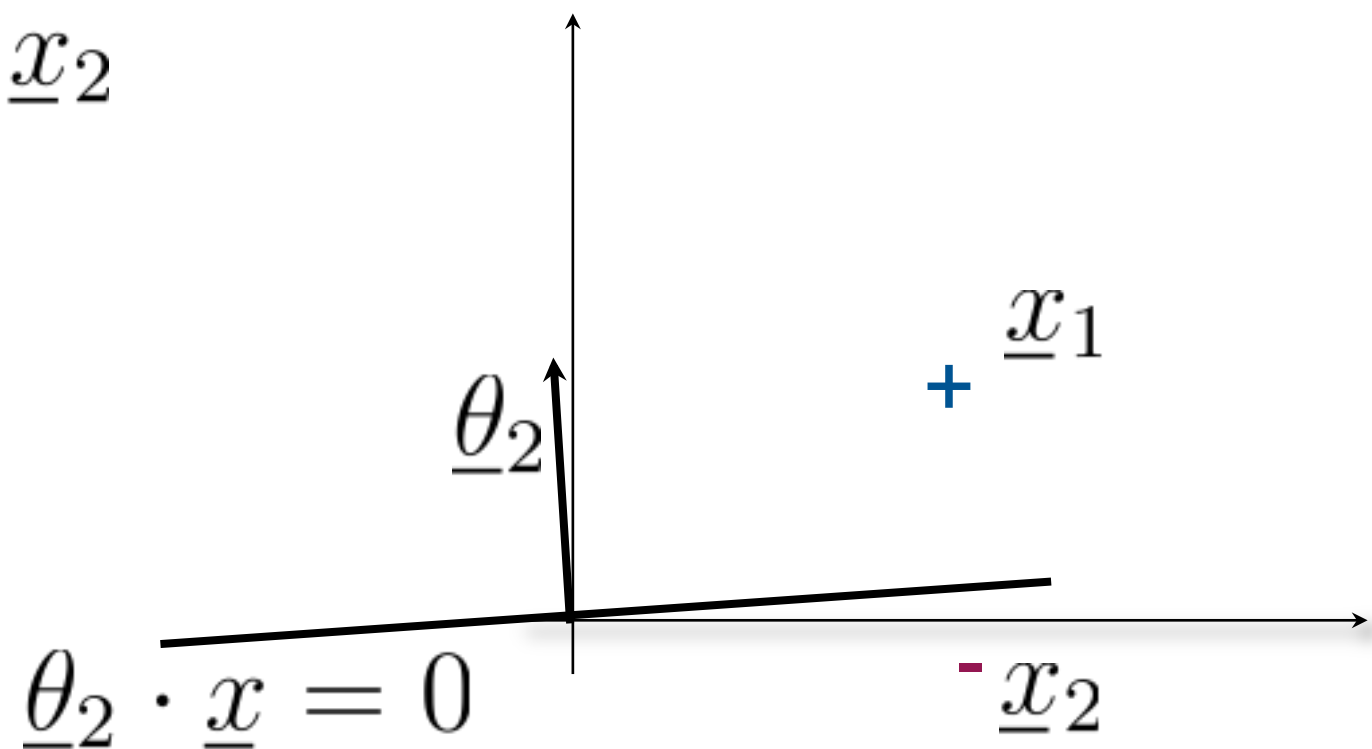
Perceptron algorithm

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$

$$\underline{\theta}_2 = \underline{\theta}_1 + (-1) \underline{x}_2$$

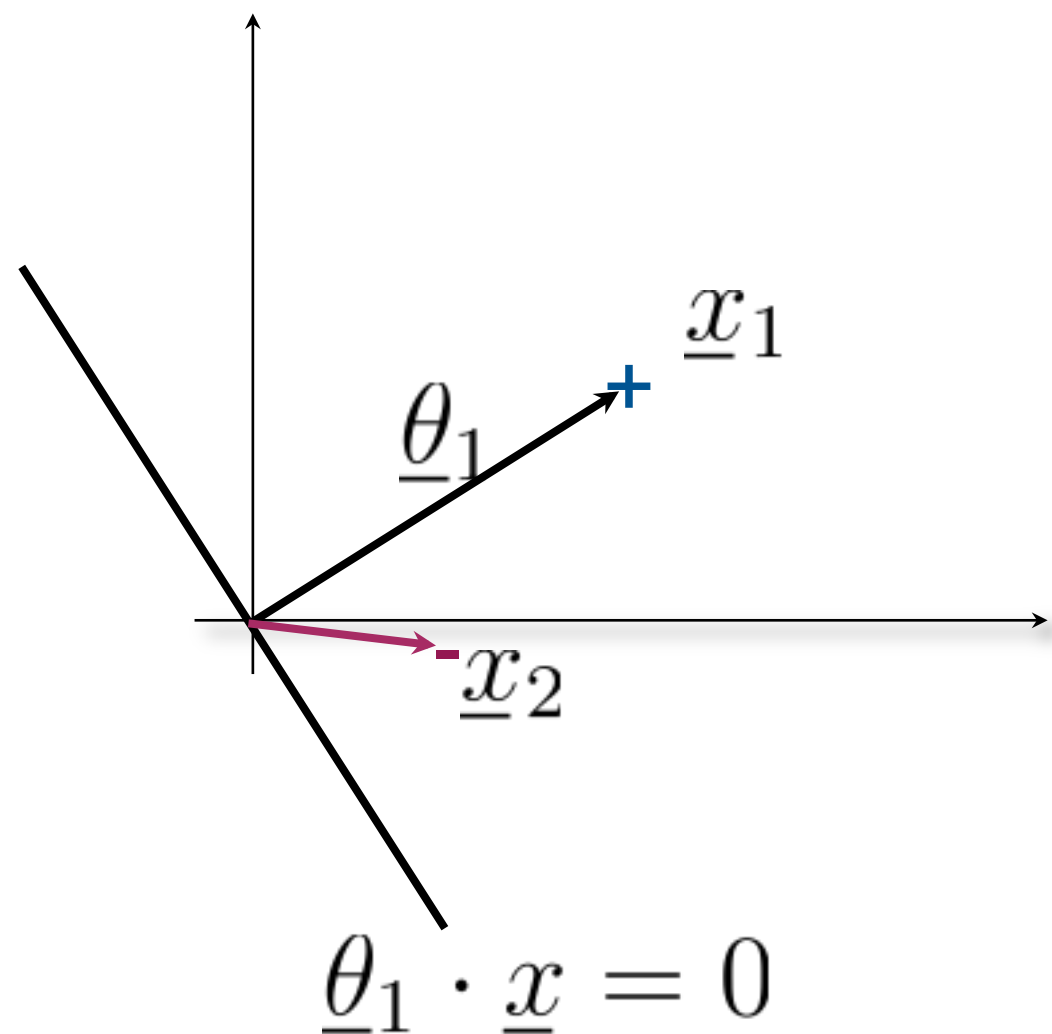


Perceptron algorithm (take 2)

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$



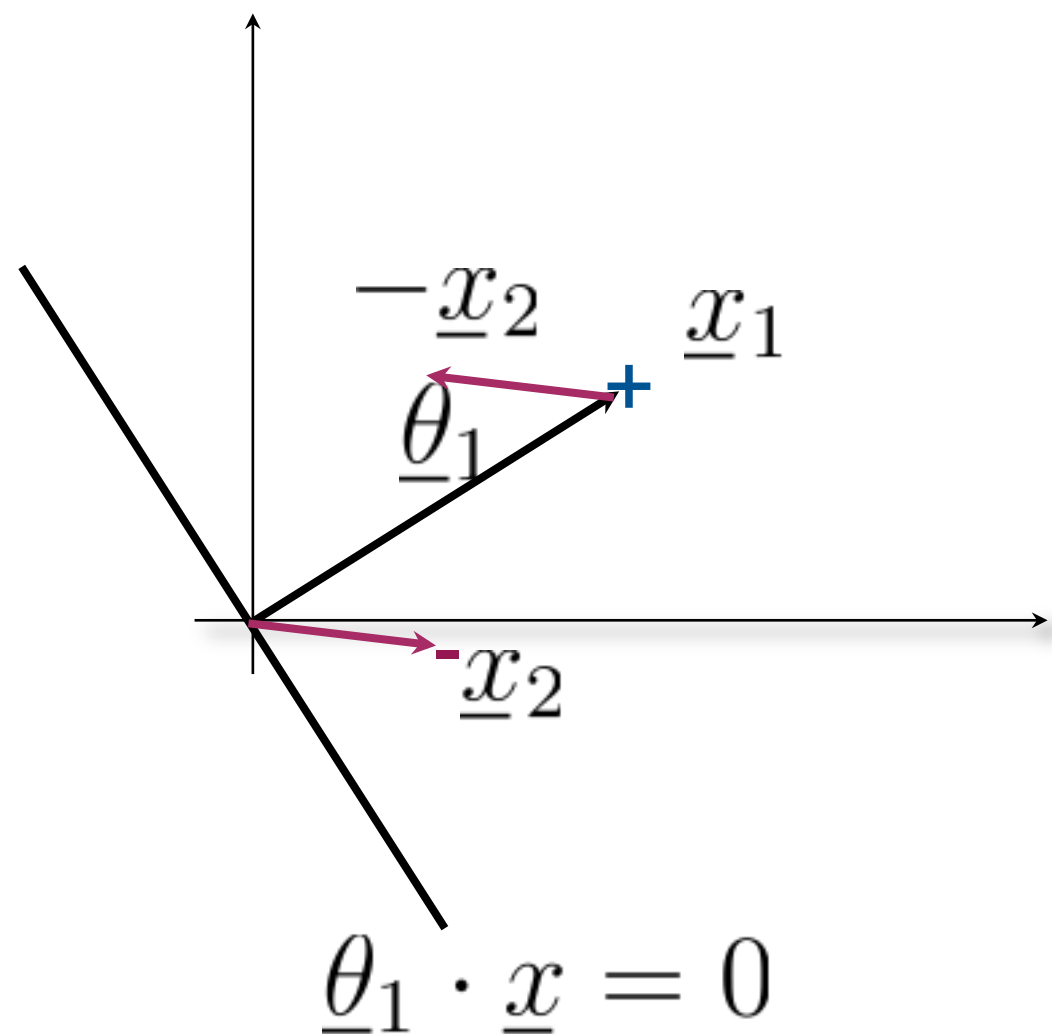
Perceptron algorithm (take 2)

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$

$$\underline{\theta}_2 = \underline{\theta}_1 + (-1) \underline{x}_2$$



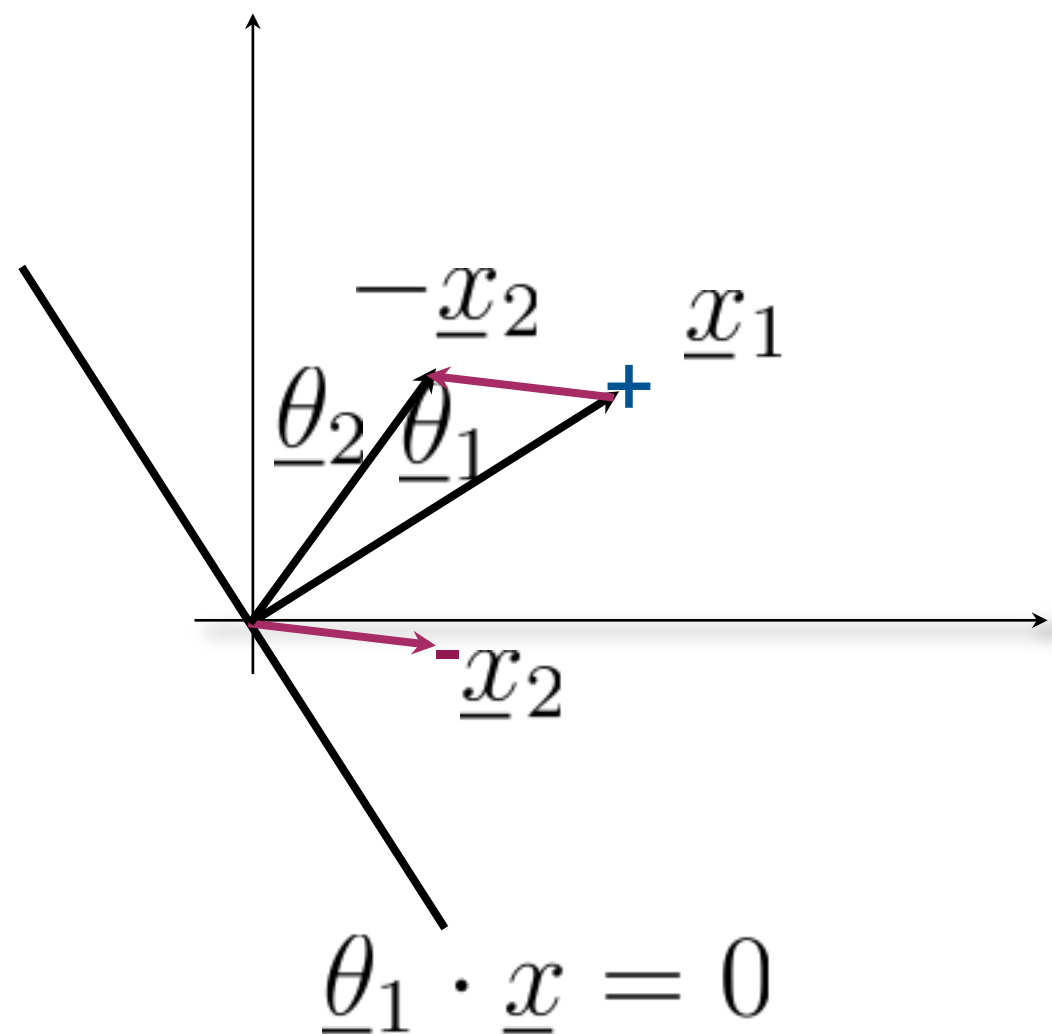
Perceptron algorithm (take 2)

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$

$$\underline{\theta}_2 = \underline{\theta}_1 + (-1) \underline{x}_2$$



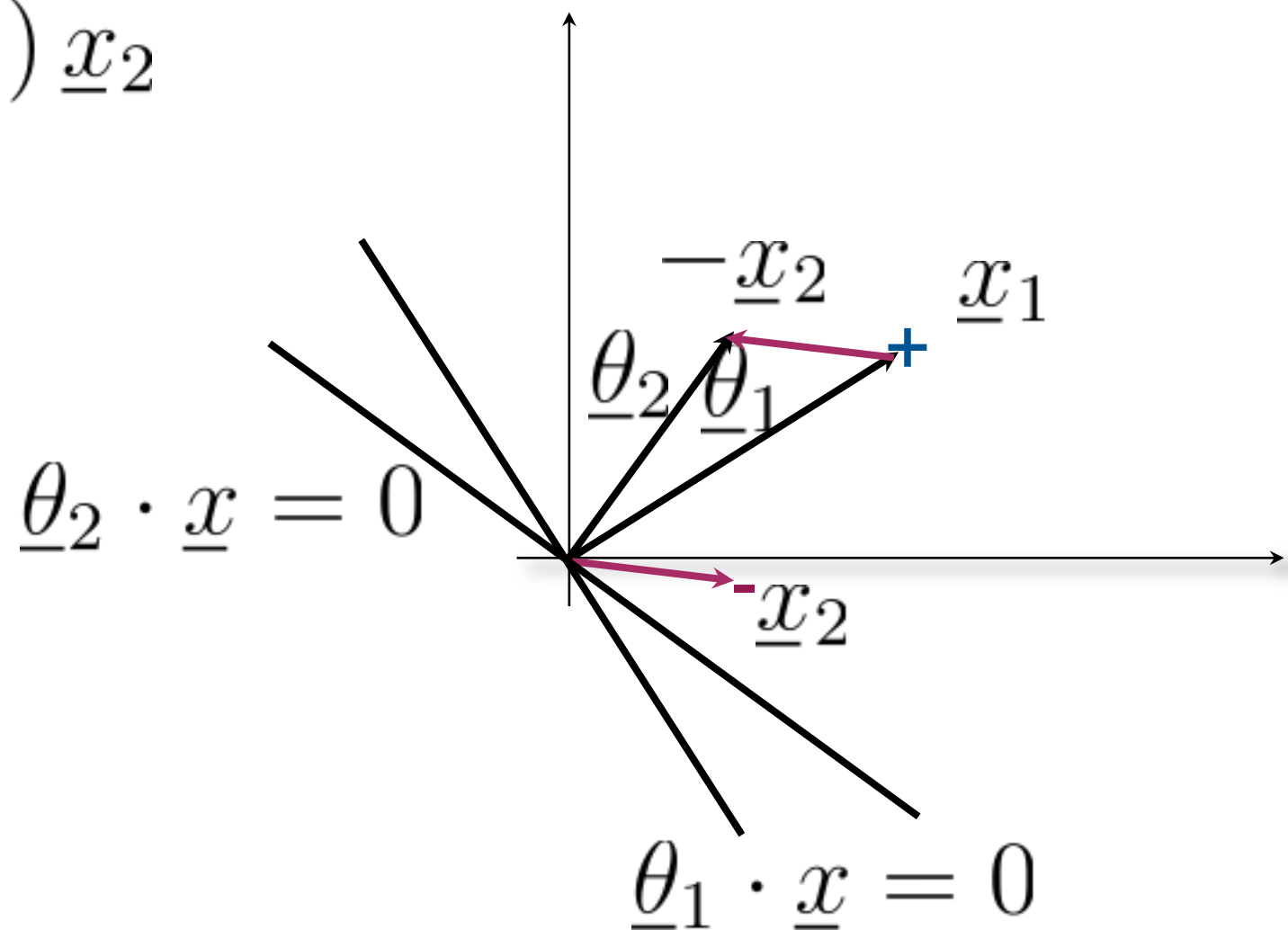
Perceptron algorithm (take 2)

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$

$$\underline{\theta}_2 = \underline{\theta}_1 + (-1) \underline{x}_2$$



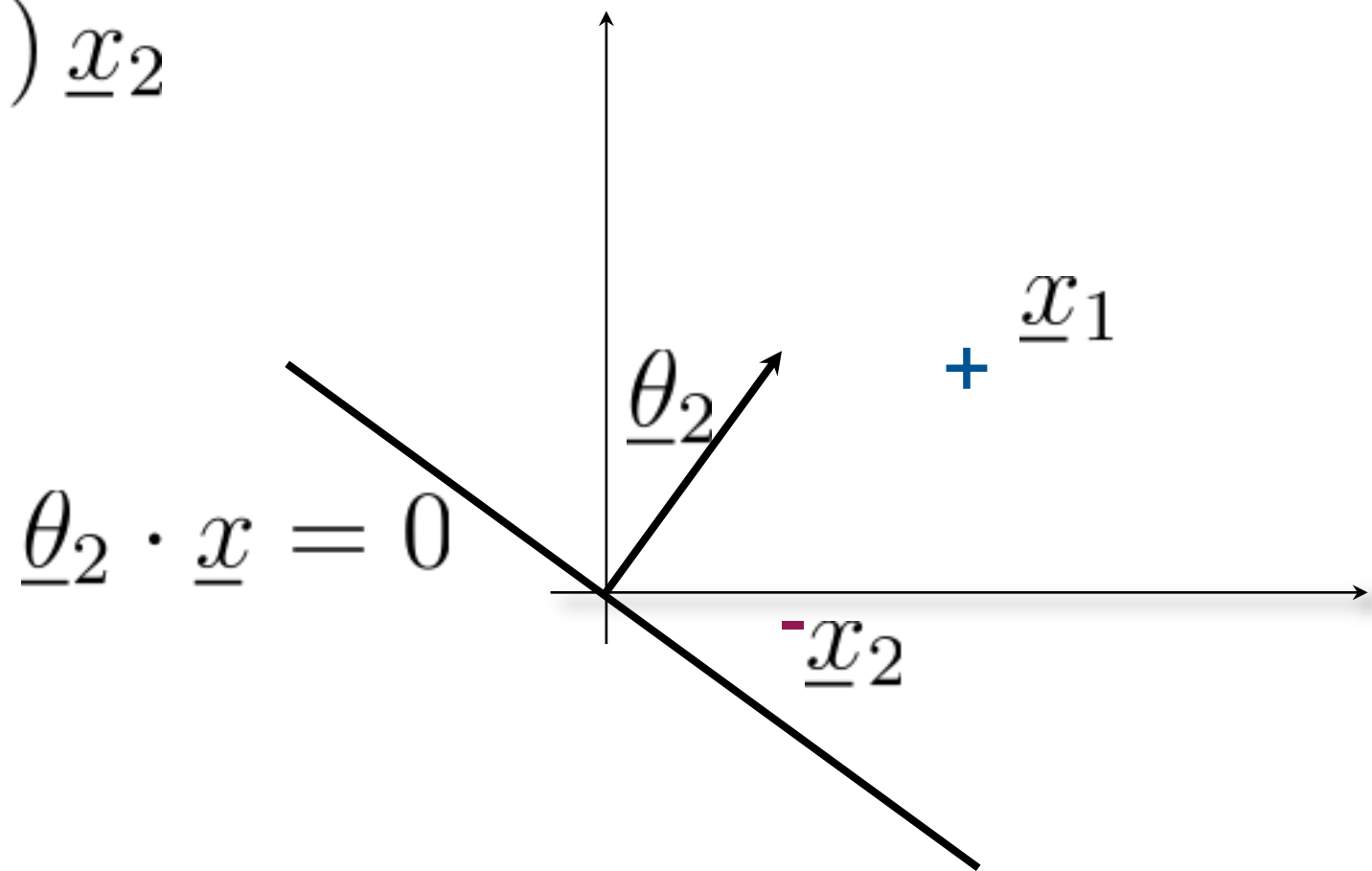
Perceptron algorithm (take 2)

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$

$$\underline{\theta}_2 = \underline{\theta}_1 + (-1) \underline{x}_2$$

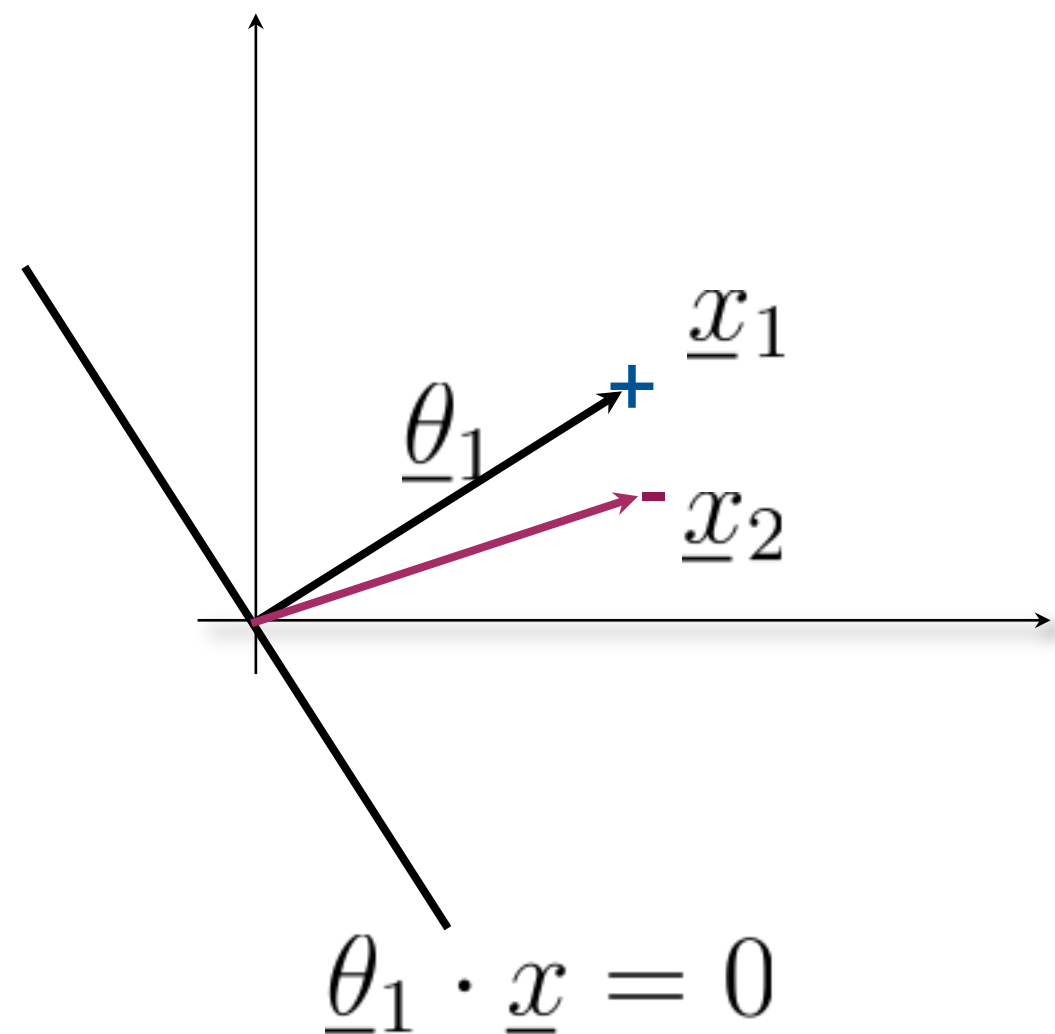


Perceptron algorithm (take 3)

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$



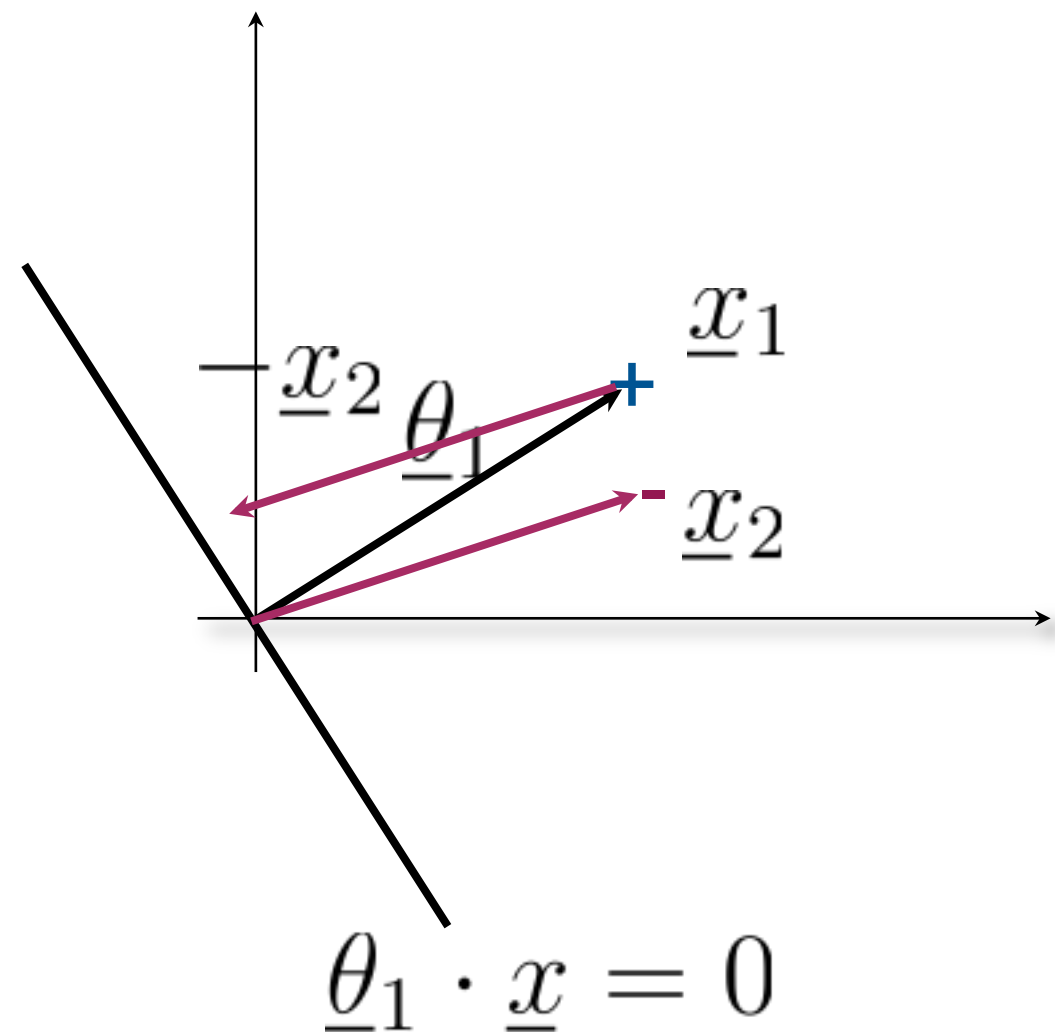
Perceptron algorithm (take 3)

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$

$$\underline{\theta}_2 = \underline{\theta}_1 + (-1) \underline{x}_2$$



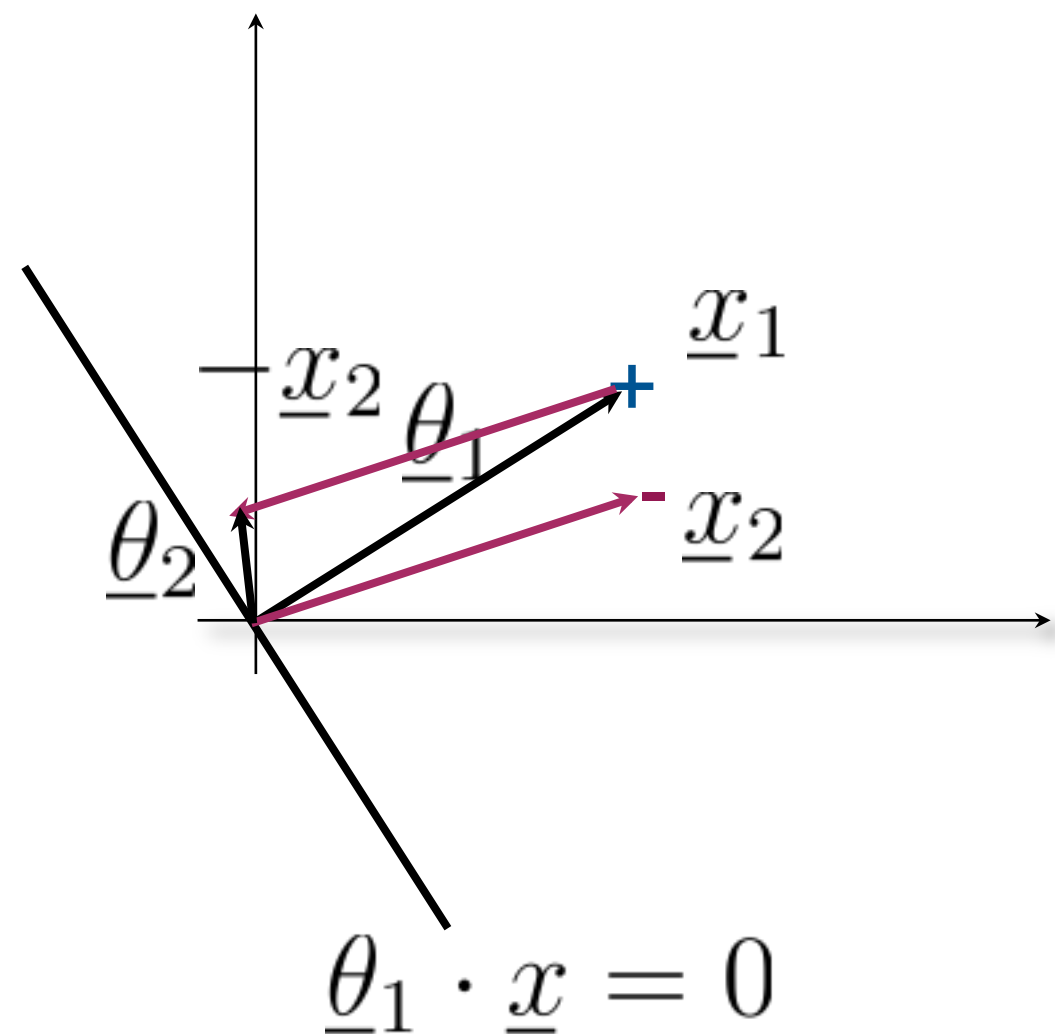
Perceptron algorithm (take 3)

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$

$$\underline{\theta}_2 = \underline{\theta}_1 + (-1) \underline{x}_2$$



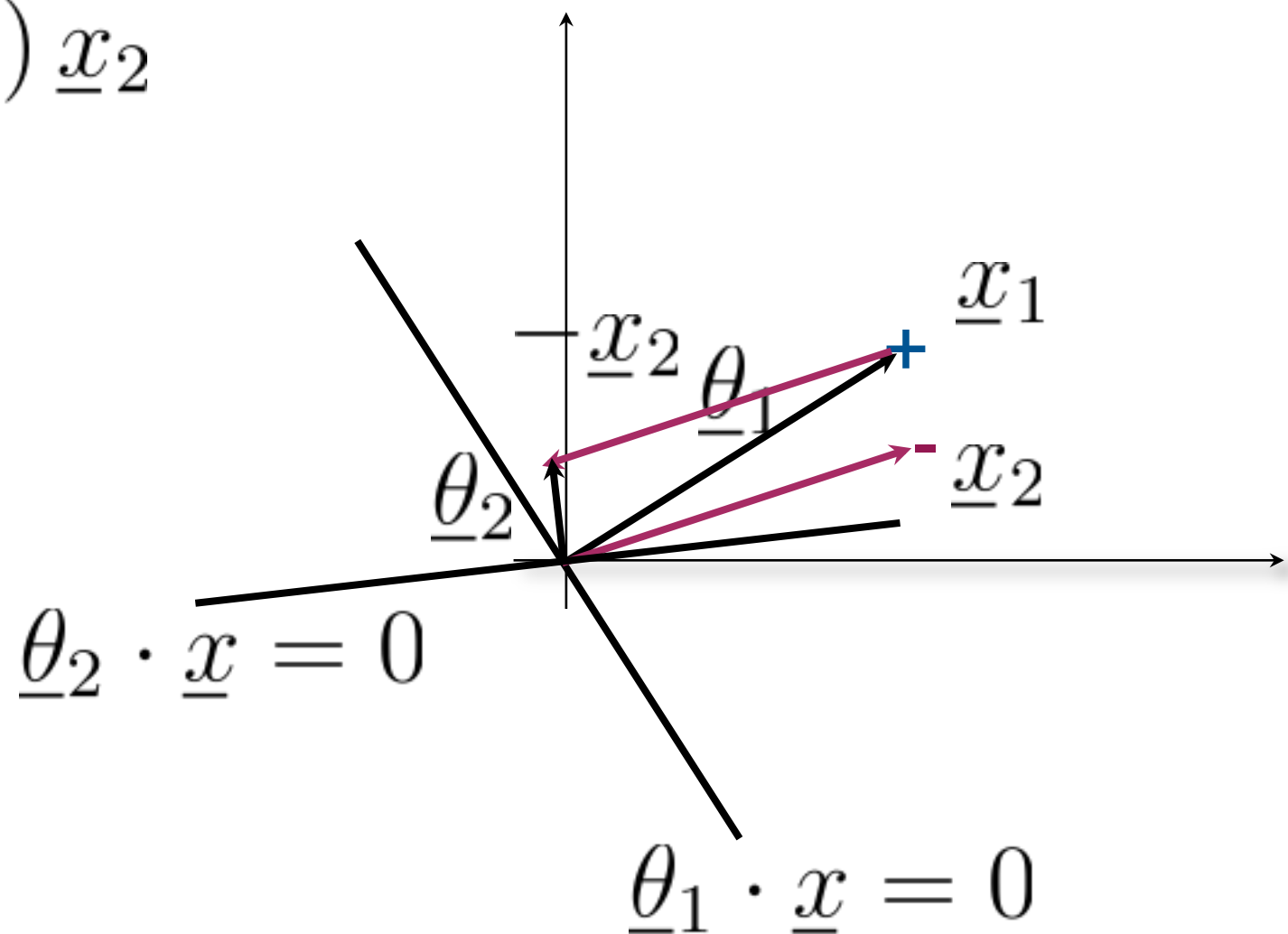
Perceptron algorithm (take 3)

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$

$$\underline{\theta}_2 = \underline{\theta}_1 + (-1) \underline{x}_2$$



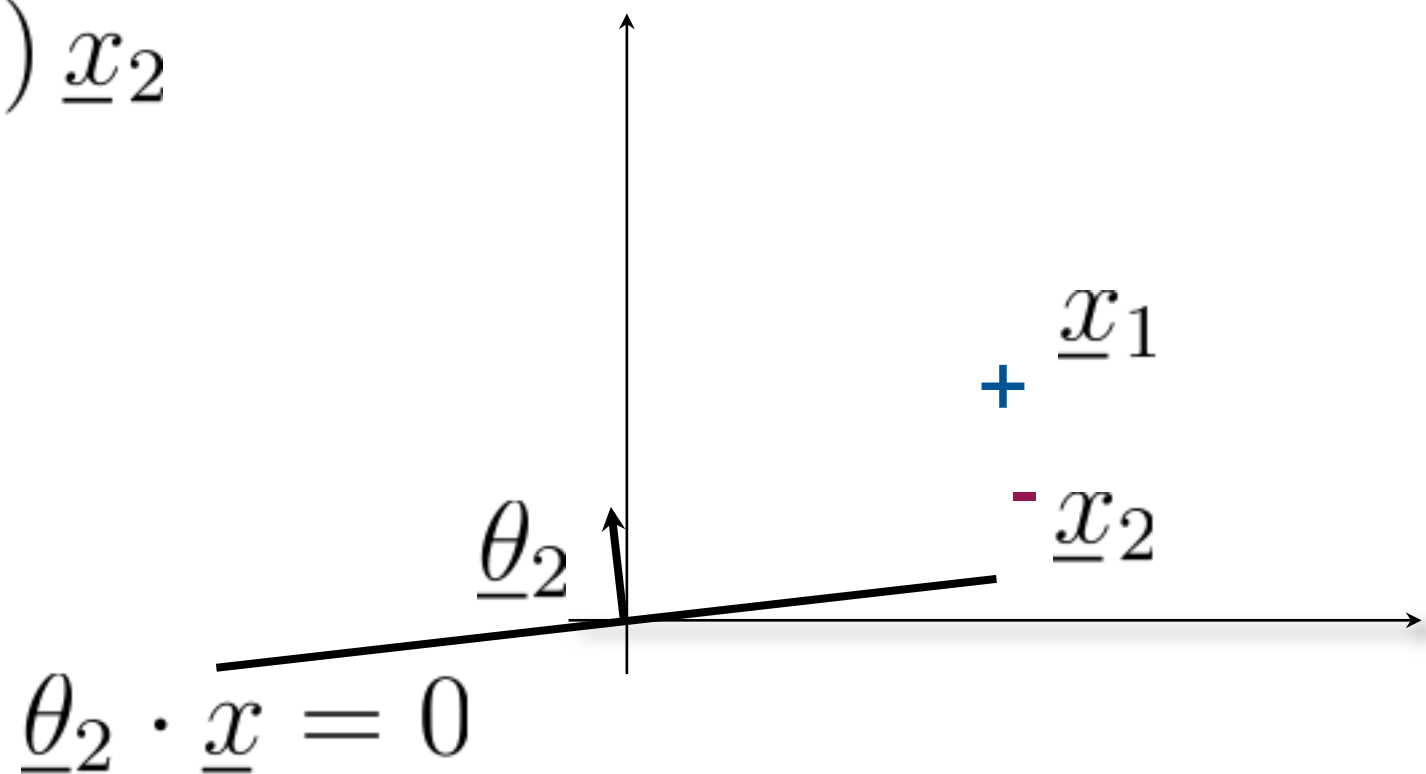
Perceptron algorithm (take 3)

- Iterative updates based on mistakes

$$\underline{\theta}_0 = 0$$

$$\underline{\theta}_1 = \underline{\theta}_0 + 1 \underline{x}_1$$

$$\underline{\theta}_2 = \underline{\theta}_1 + (-1) \underline{x}_2$$



Perceptron algorithm: motivation

- If we make a mistake on the t^{th} training point, then

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

- After the update, we have

$$\underline{\theta}' = \underline{\theta} + y_t \underline{x}_t$$

Perceptron algorithm: motivation

- If we make a mistake on the t^{th} training point, then

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

- After the update, we have

$$\underline{\theta}' = \underline{\theta} + y_t \underline{x}_t$$

$$y_t(\underline{\theta}' \cdot \underline{x}_t) = y_t([\underline{\theta} + y_t \underline{x}_t] \cdot \underline{x}_t)$$

Perceptron algorithm: motivation

- If we make a mistake on the t^{th} training point, then

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

- After the update, we have

$$\underline{\theta}' = \underline{\theta} + y_t \underline{x}_t$$

$$\begin{aligned} y_t(\underline{\theta}' \cdot \underline{x}_t) &= y_t([\underline{\theta} + y_t \underline{x}_t] \cdot \underline{x}_t) \\ &= y_t(\underline{\theta} \cdot \underline{x}_t + y_t \underline{x}_t \cdot \underline{x}_t) \end{aligned}$$

Perceptron algorithm: motivation

- If we make a mistake on the t^{th} training point, then

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

- After the update, we have

$$\underline{\theta}' = \underline{\theta} + y_t \underline{x}_t$$

$$\begin{aligned} y_t(\underline{\theta}' \cdot \underline{x}_t) &= y_t([\underline{\theta} + y_t \underline{x}_t] \cdot \underline{x}_t) \\ &= y_t(\underline{\theta} \cdot \underline{x}_t + y_t \underline{x}_t \cdot \underline{x}_t) \\ &= y_t(\underline{\theta} \cdot \underline{x}_t) + y_t^2 \underline{x}_t \cdot \underline{x}_t \end{aligned}$$

Perceptron algorithm: motivation

- If we make a mistake on the t^{th} training point, then

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

- After the update, we have

$$\underline{\theta}' = \underline{\theta} + y_t \underline{x}_t$$

$$\begin{aligned} y_t(\underline{\theta}' \cdot \underline{x}_t) &= y_t([\underline{\theta} + y_t \underline{x}_t] \cdot \underline{x}_t) \\ &= y_t(\underline{\theta} \cdot \underline{x}_t + y_t \underline{x}_t \cdot \underline{x}_t) \\ &= y_t(\underline{\theta} \cdot \underline{x}_t) + y_t^2 \underline{x}_t \cdot \underline{x}_t \\ &= y_t(\underline{\theta} \cdot \underline{x}_t) + \|\underline{x}_t\|^2 \end{aligned}$$

Perceptron algorithm: motivation

- If we make a mistake on the t^{th} training point, then

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

- After the update, we have

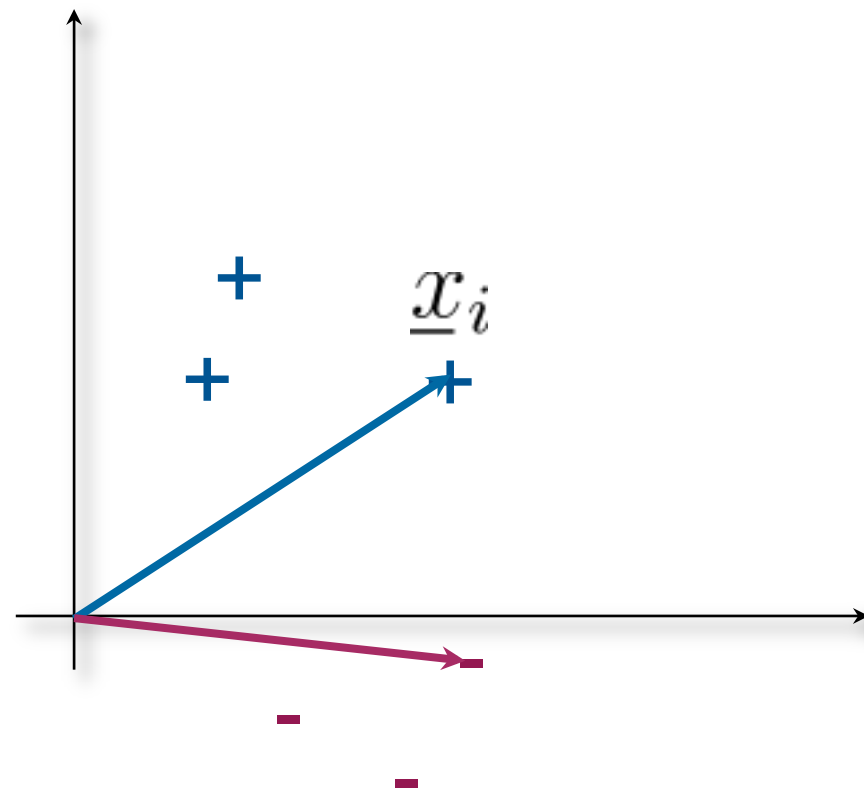
$$\underline{\theta}' = \underline{\theta} + y_t \underline{x}_t$$

$$\begin{aligned} y_t(\underline{\theta}' \cdot \underline{x}_t) &= y_t([\underline{\theta} + y_t \underline{x}_t] \cdot \underline{x}_t) \\ &= y_t(\underline{\theta} \cdot \underline{x}_t + y_t \underline{x}_t \cdot \underline{x}_t) \\ &= y_t(\underline{\theta} \cdot \underline{x}_t) + y_t^2 \underline{x}_t \cdot \underline{x}_t \\ &= y_t(\underline{\theta} \cdot \underline{x}_t) + \|\underline{x}_t\|^2 \end{aligned}$$

- So that $y_t(\underline{\theta}' \cdot \underline{x}_t)$ increases based on the update

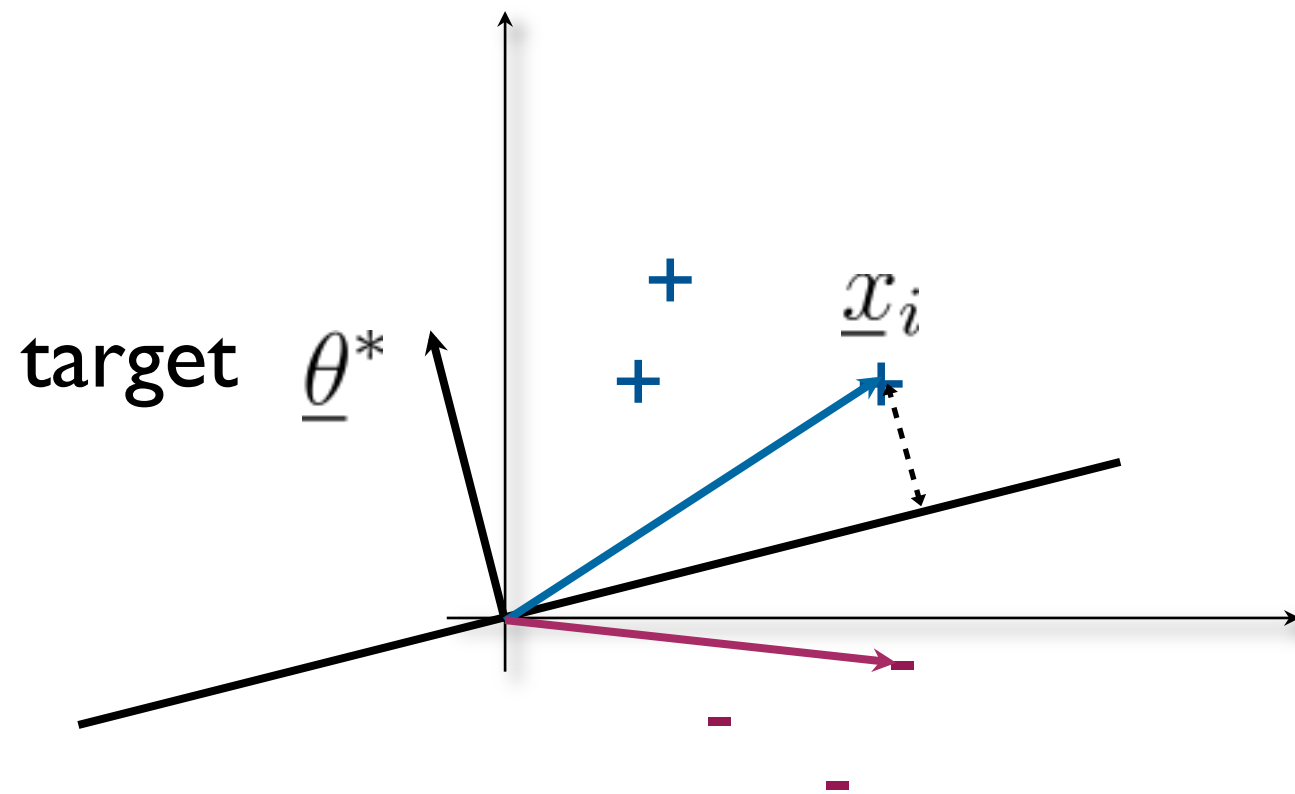
“Margin”

- We can get a handle on convergence by assuming that there exists a target classifier with good properties
- One such property is geometric margin, i.e., how close the separating boundary is to the points



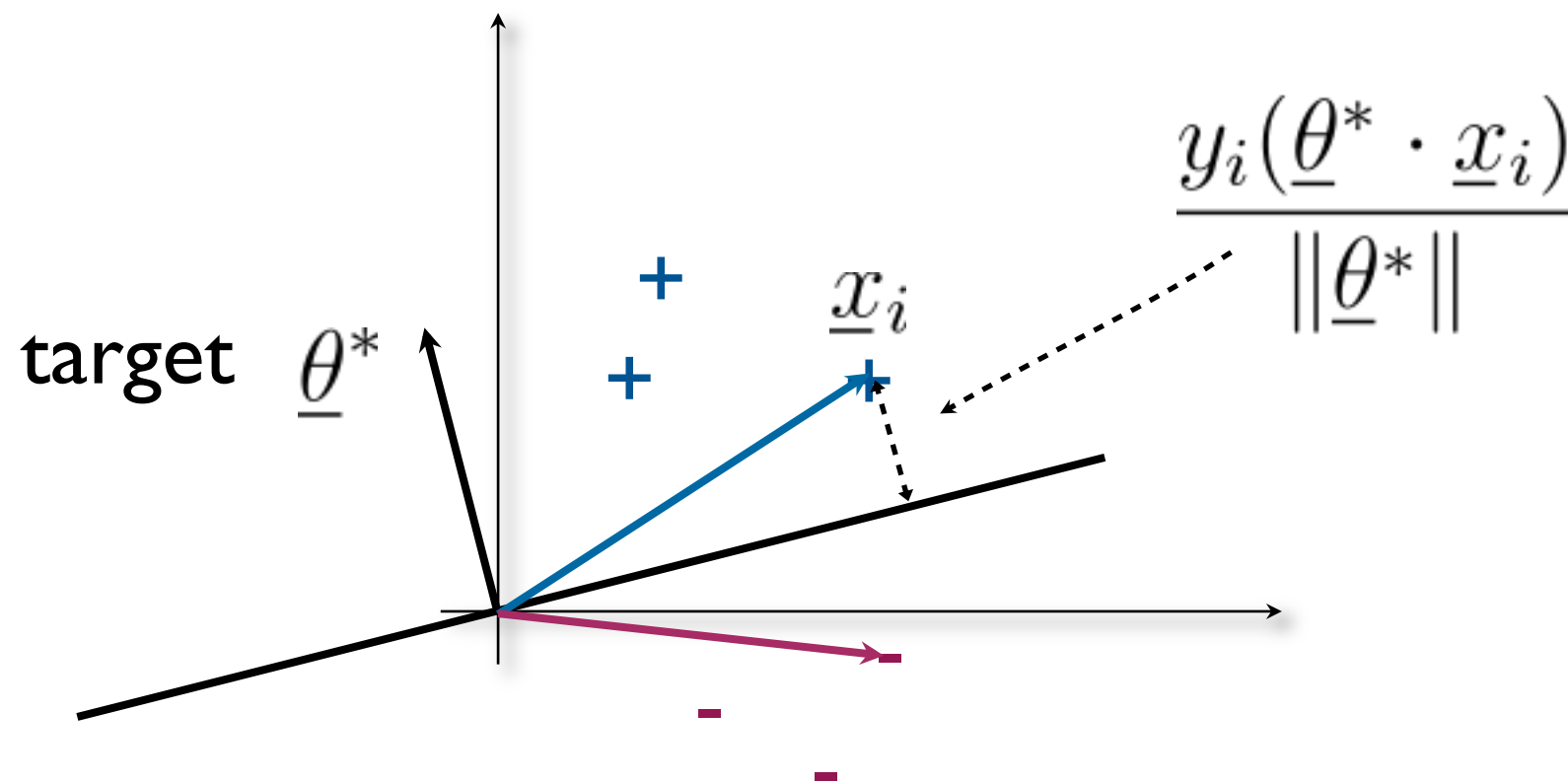
“Margin”

- We can get a handle on convergence by assuming that there exists a target classifier with good properties
- One such property is geometric margin, i.e., how close the separating boundary is to the points



“Margin”

- We can get a handle on convergence by assuming that there exists a target classifier with good properties
- One such property is geometric margin, i.e., how close the separating boundary is to the points

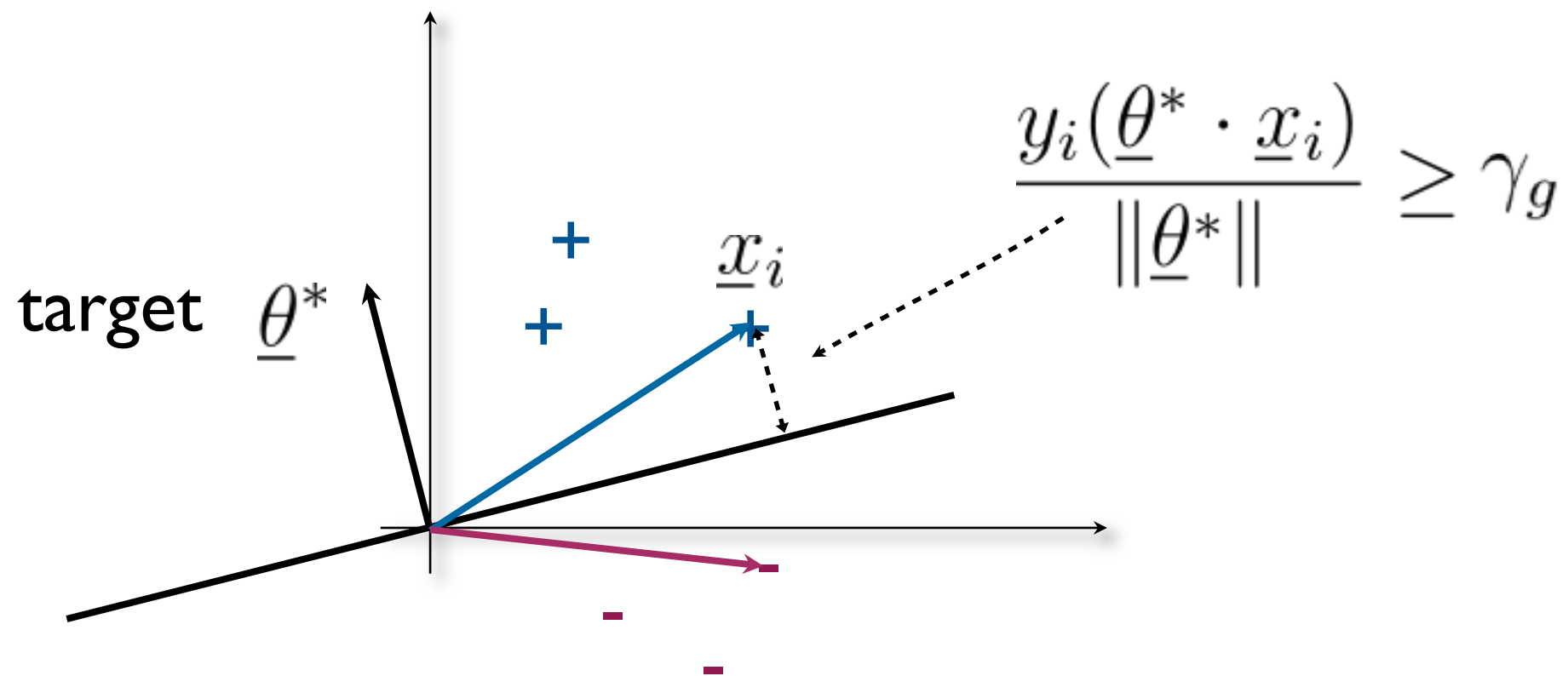


Recall:

$\cos(\text{angle}) = \text{dot product of two unit vectors}$
dotted distance = radius * $\cos(\text{angle})$

“Margin”

- We can get a handle on convergence by assuming that there exists a target classifier with good properties
- One such property is geometric margin, i.e., how close the separating boundary is to the points



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).

- Note that the result does NOT depend on d or n