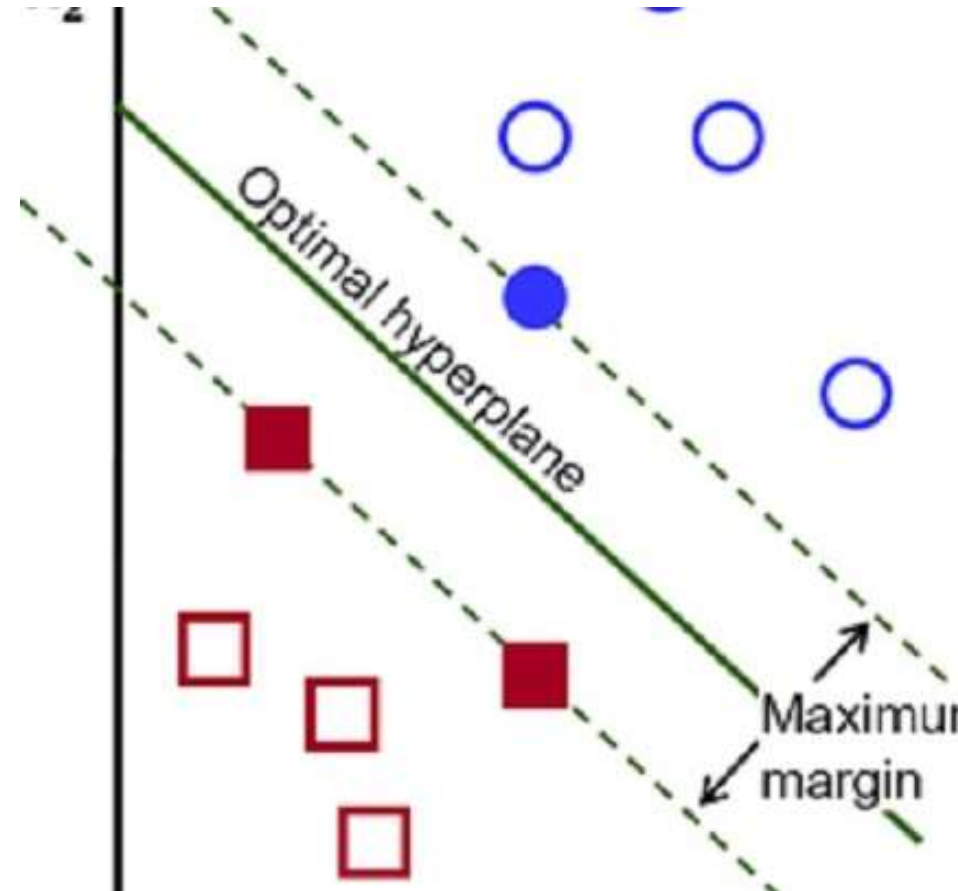


LARGE MARGIN CLASSIFIERS

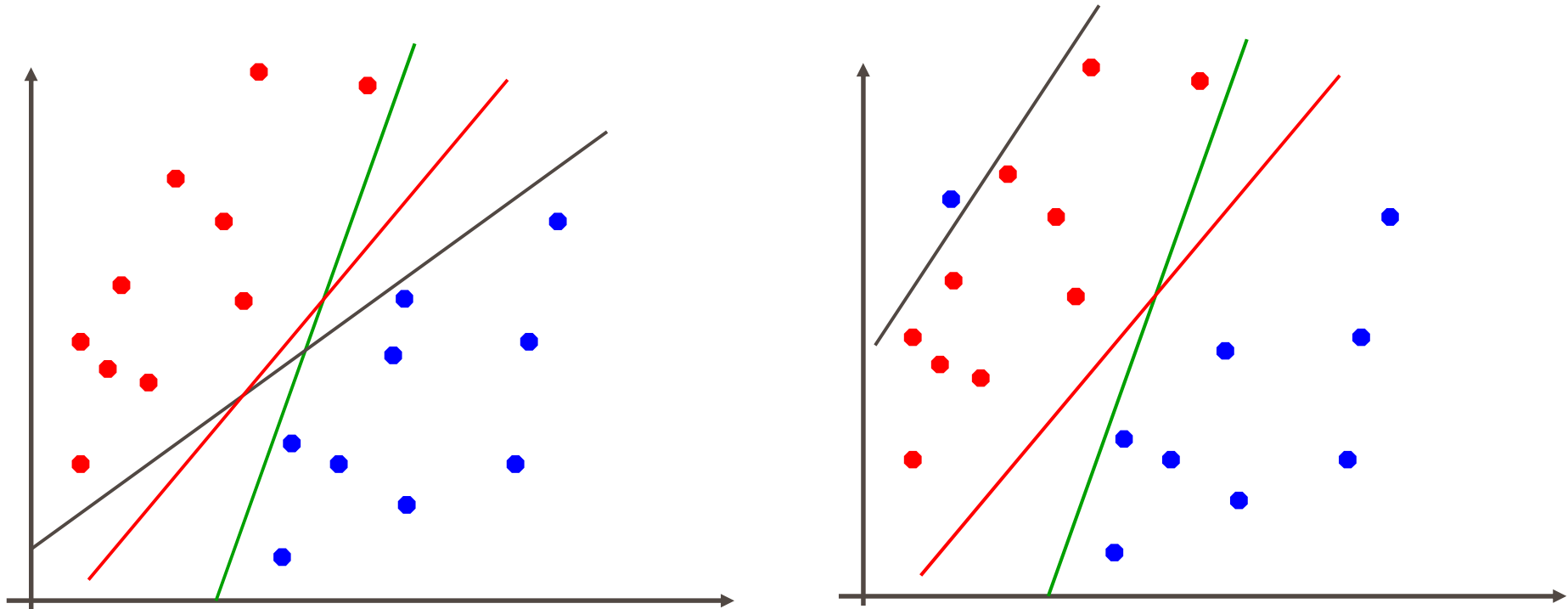
SUPPORT VECTOR MACHINE (SVM)

許志仲 (Chih-Chung Hsu)

Assistant Professor
Institute of Data Science
National Cheng Kung University



Which hyperplane?



Two main variations in linear classifiers:

- which hyperplane they choose when the data is linearly separable
- how they handle data that is not linearly separable

Linear approaches so far

Perceptron:

- separable:
- non-separable:

Gradient descent:

- separable:
- non-separable:

Linear approaches so far

Perceptron:

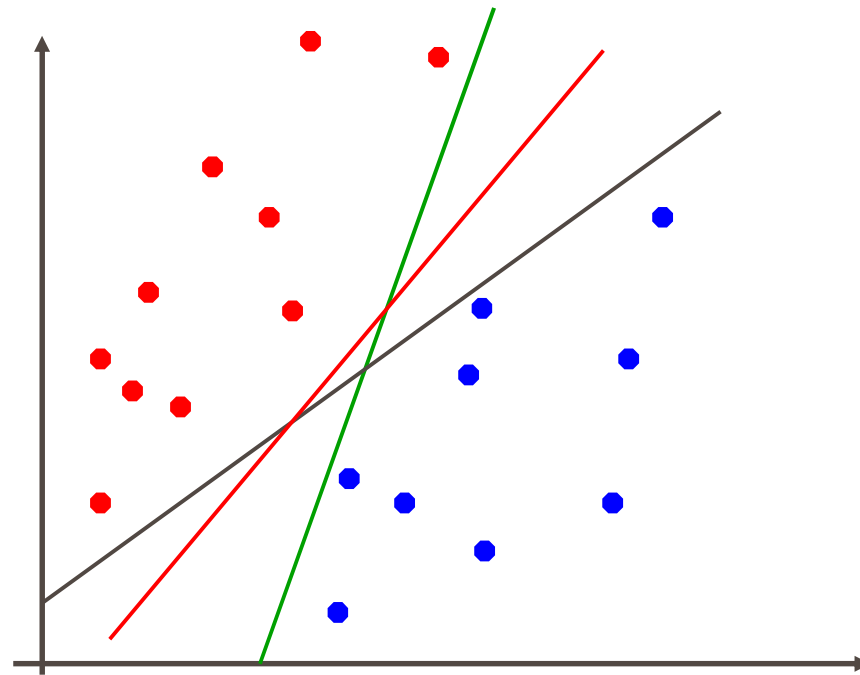
- separable:
 - finds *some* hyperplane that separates the data
- non-separable:
 - will continue to adjust as it iterates through the examples
 - final hyperplane will depend on which examples it saw recently

Gradient descent:

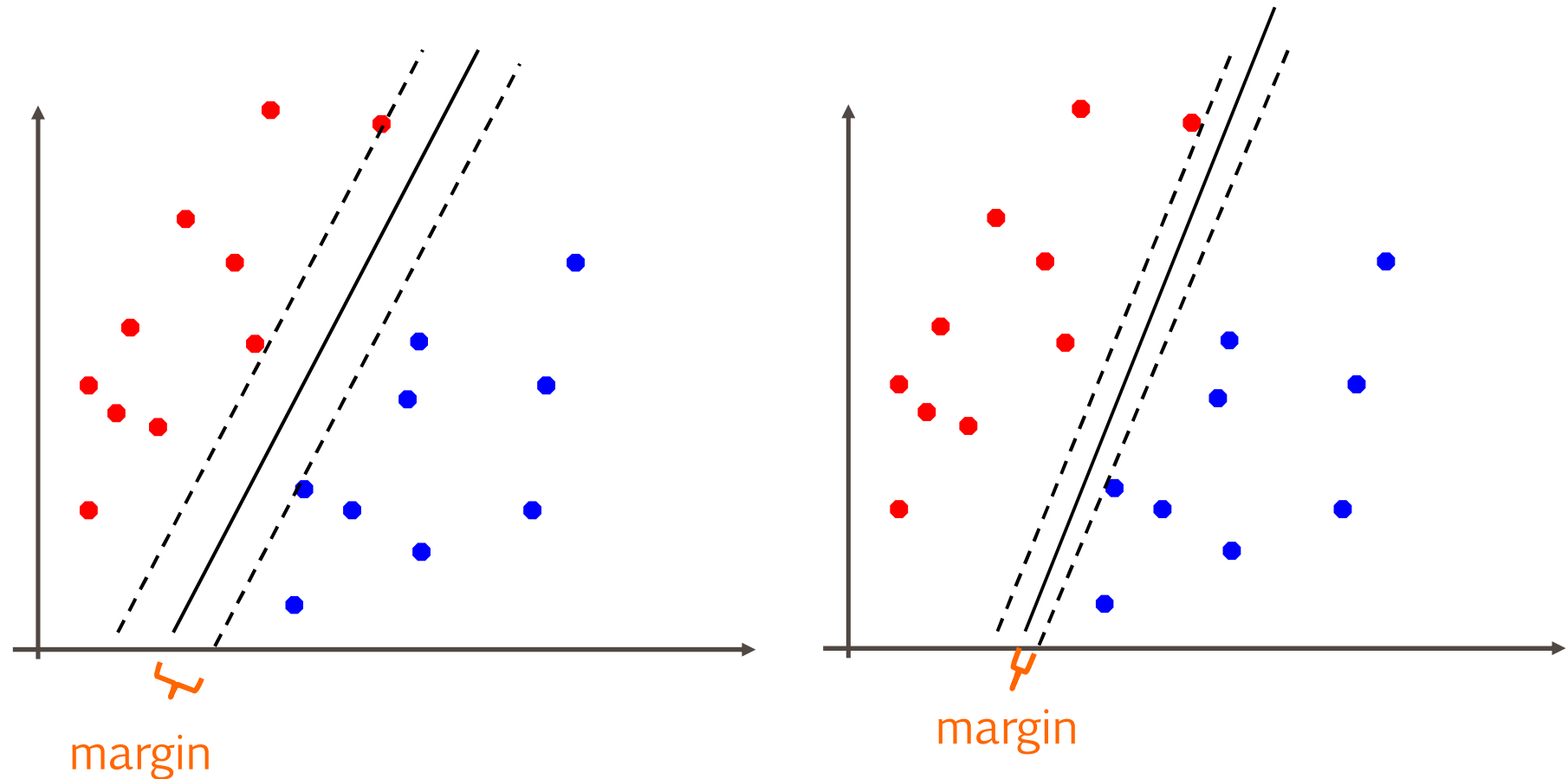
- separable and non-separable
 - finds the hyperplane that minimizes the objective function (loss + regularization)

Which hyperplane is this?

Which hyperplane would you choose?

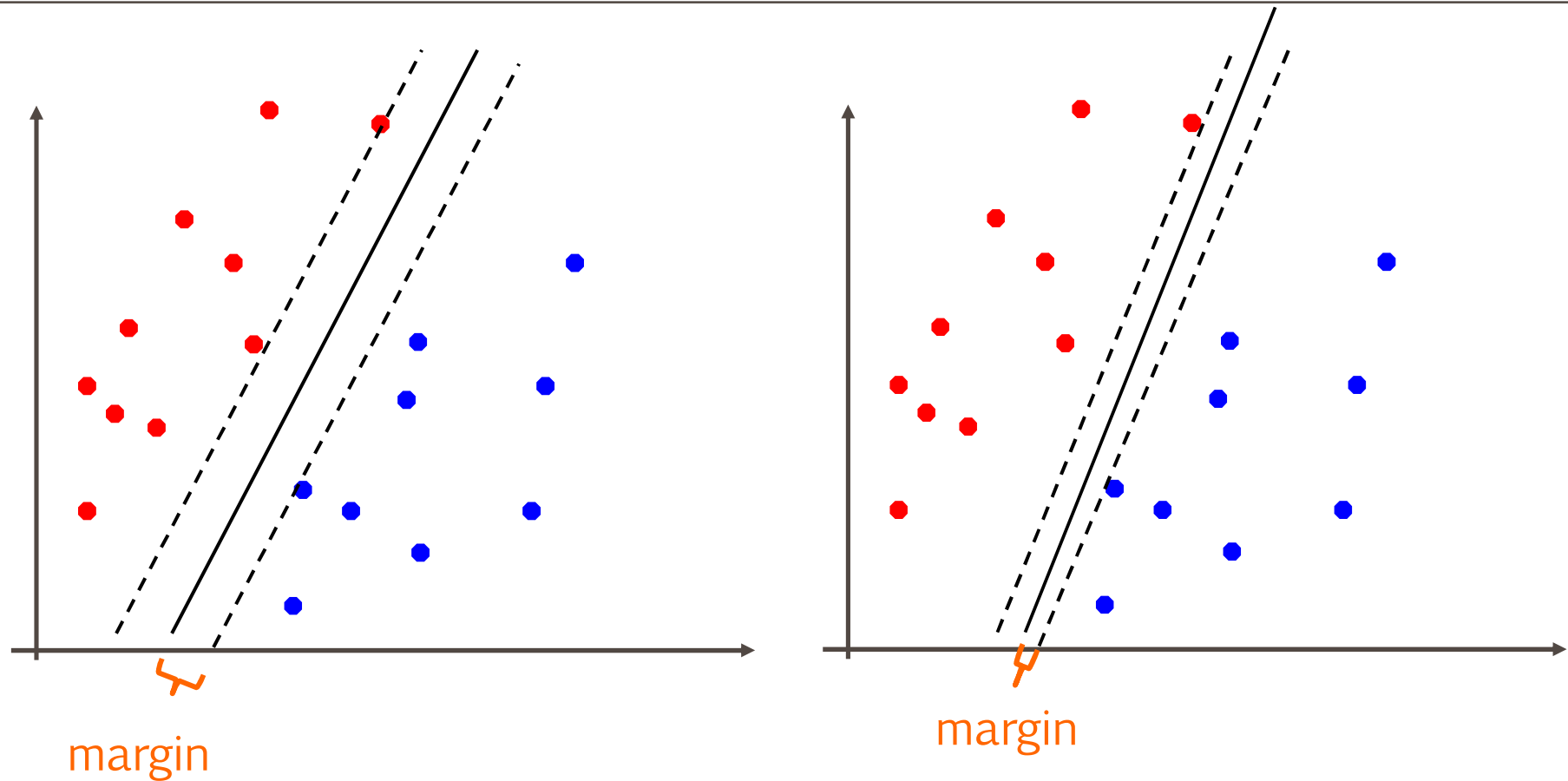


Large margin classifiers



Choose the line where the distance to the nearest point(s) is as large as possible

Large margin classifiers



The **margin** of a classifier is the distance to the closest points of either class
Large **margin** classifiers attempt to maximize this

Large margin classifier setup

Select the hyperplane with the largest margin where the points are classified correctly!

Setup as a **constrained optimization problem**:

$$\begin{aligned} & \max_{w,b} \text{margin}(w,b) \\ & \text{subject to:} \end{aligned}$$

$$y_i(w \cdot x_i + b) > 0 \quad \forall i \quad \text{what does this say?}$$

Large margin classifier setup

$$\max_{w,b} \text{margin}(w,b)$$

subject to:

$$y_i(w \cdot x_i + b) > 0 \quad \forall i$$

$$\max_{w,b} \text{margin}(w,b)$$

subject to:

$$y_i(w \cdot x_i + b) \geq c \quad \forall i$$
$$c > 0$$

Are these equivalent?

Large margin classifier setup

$$\max_{w,b} \text{margin}(w,b)$$

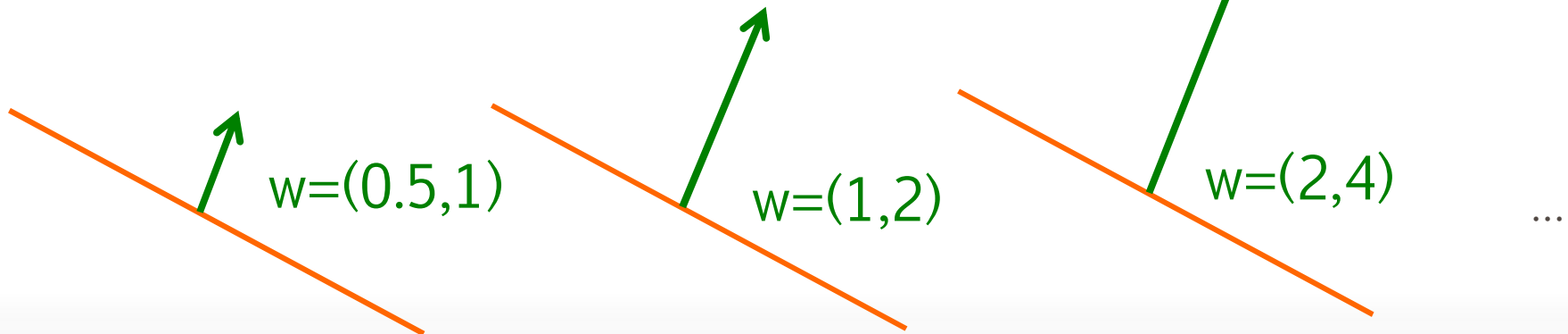
subject to:

$$y_i(w \cdot x_i + b) > 0 \quad \forall i$$

$$\max_{w,b} \text{margin}(w,b)$$

subject to:

$$y_i(w \cdot x_i + b) \geq c \quad \forall i$$
$$c > 0$$



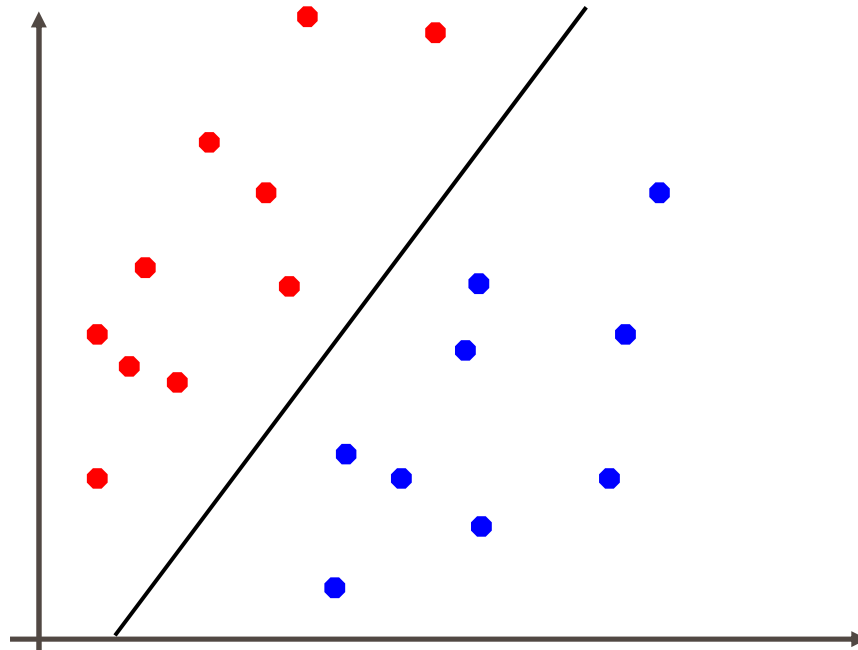
Large margin classifier setup

$$\begin{aligned} & \max_{w,b} \text{margin}(w,b) \\ & \text{subject to:} \\ & \quad y_i(w \cdot x_i + b) \geq 1 \quad \forall i \end{aligned}$$

We'll assume $c = 1$, however, any $c > 0$ works

Measuring the margin

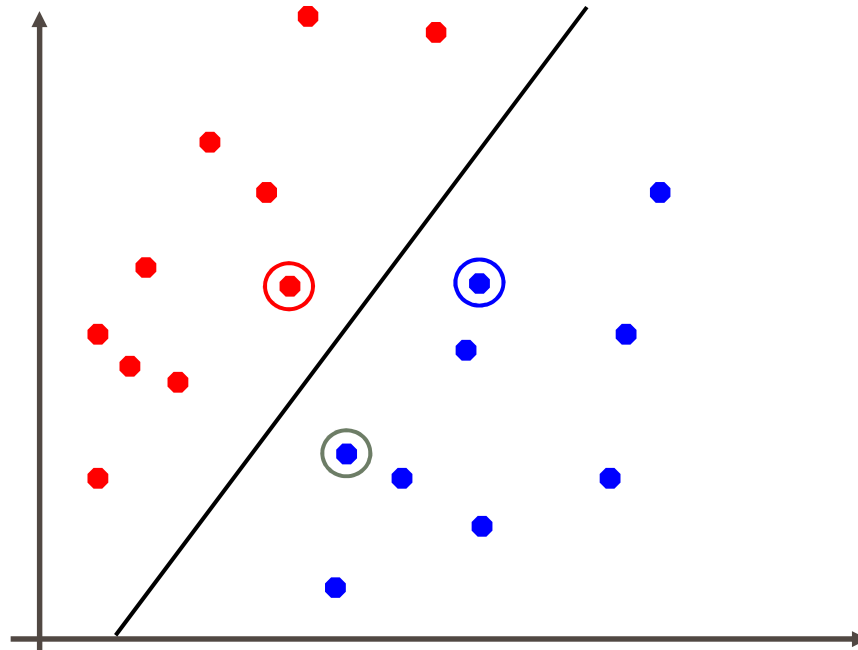
How do we calculate the margin?



Support vectors

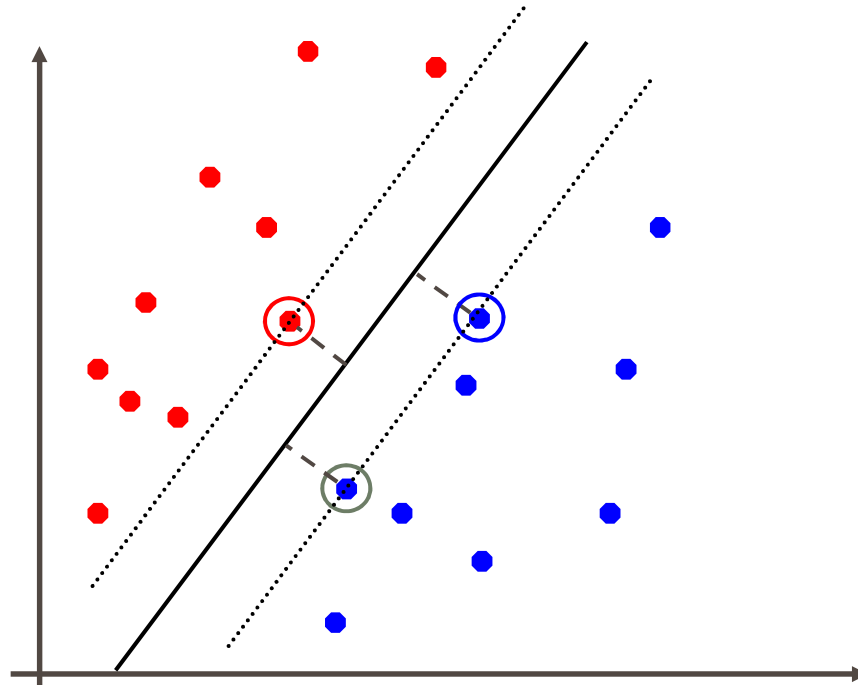
For any separating hyperplane, there exist some set of “closest points”

These are called the support vectors



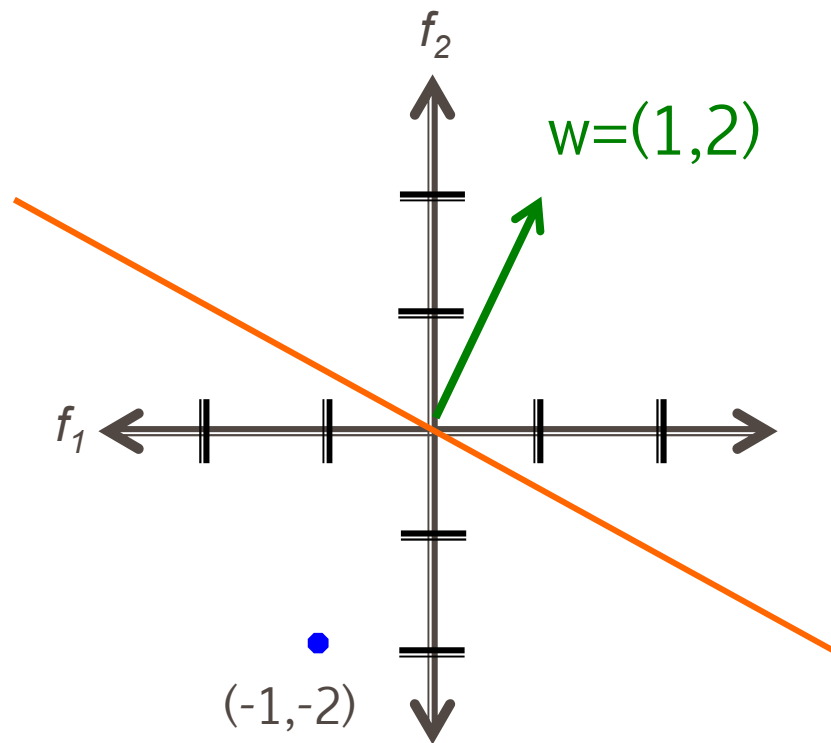
Measuring the margin

The margin is the distance to the support vectors, i.e. the “closest points”, on either side of the hyperplane



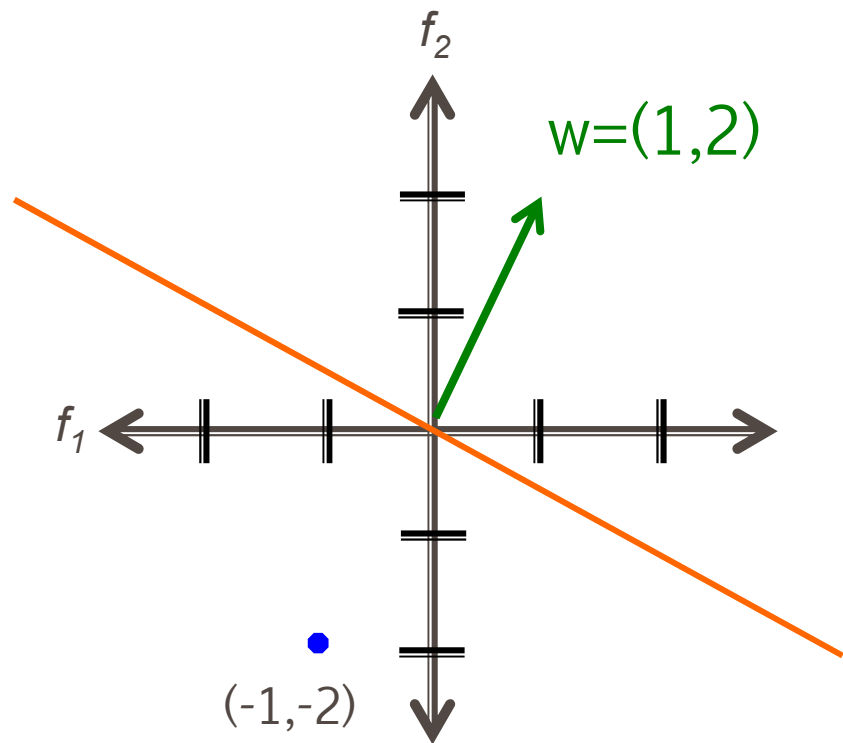
Distance from the hyperplane

How far away is this point from the hyperplane?



Distance from the hyperplane

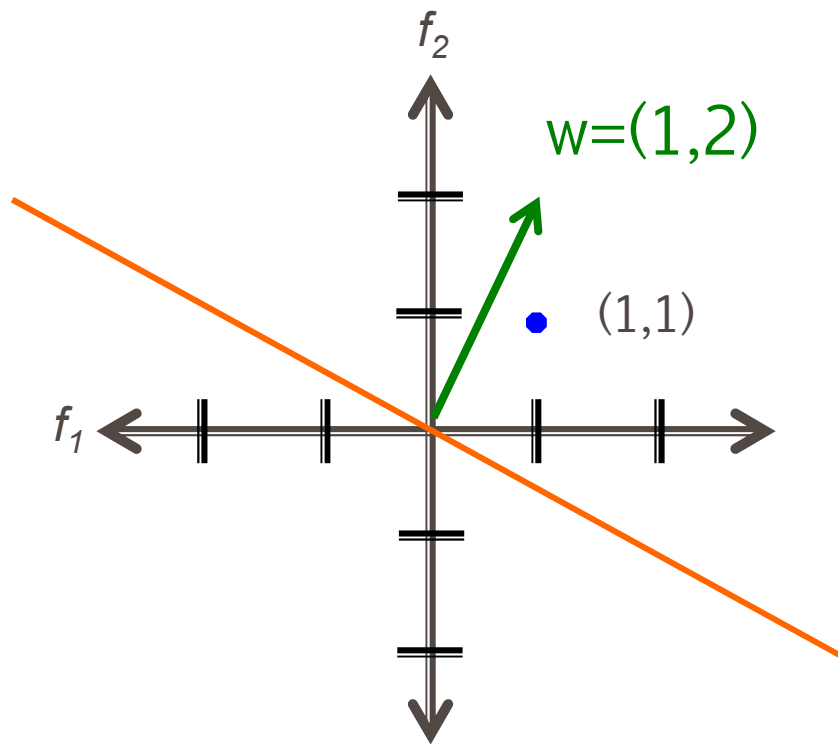
How far away is this point from the hyperplane?



$$d = \sqrt{1^2 + 2^2} = \sqrt{5}$$

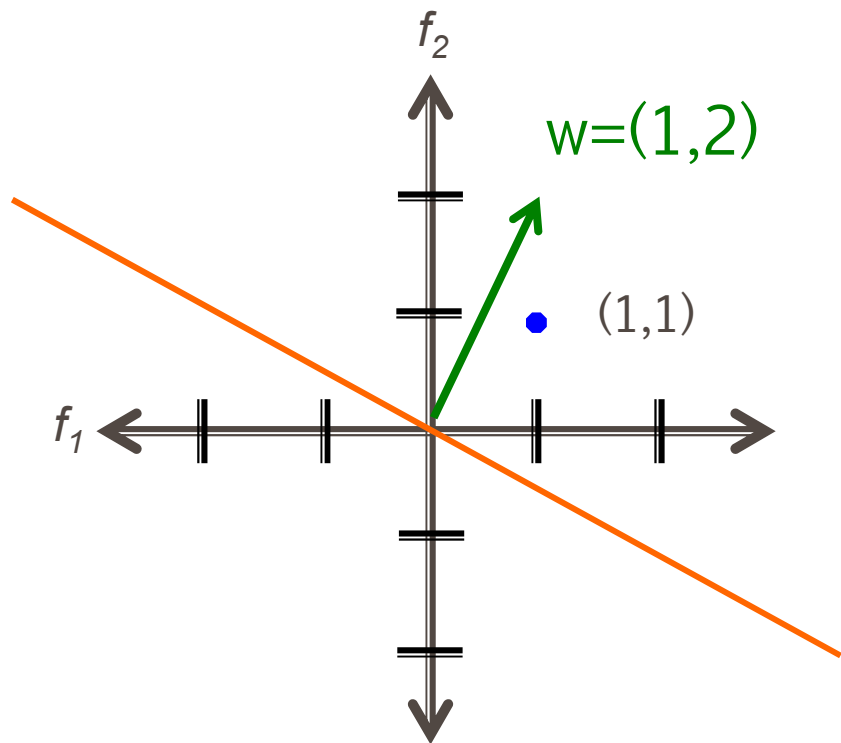
Distance from the hyperplane

How far away is this point from the hyperplane?



Distance from the hyperplane

How far away is this point from the hyperplane?

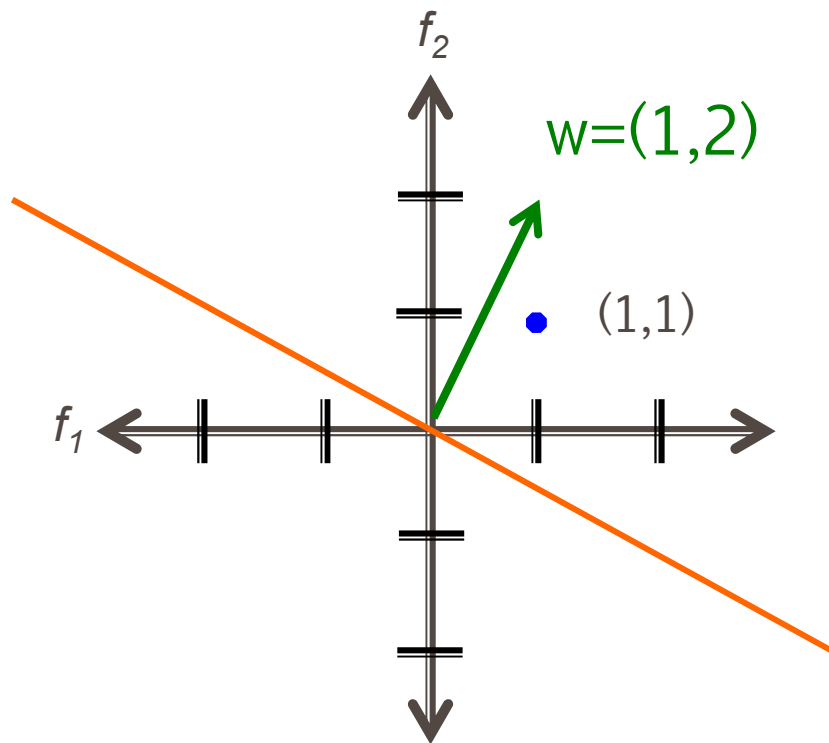


$$d(x) = \frac{w}{\|w\|} \cdot x + b$$

length
normalized weight vectors

Distance from the hyperplane

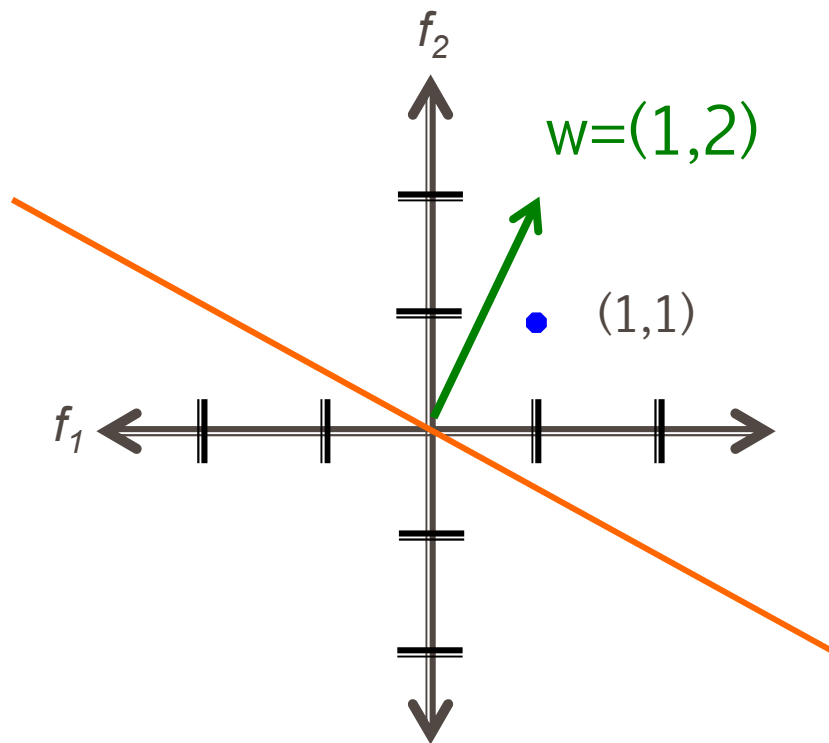
How far away is this point from the hyperplane?



$$\begin{aligned}d(x) &= \frac{w}{\|w\|} \cdot x + b \\&= \frac{1}{\sqrt{5}} (w_1 x_1 + w_2 x_2) + b \\&= \frac{1}{\sqrt{5}} (1 * 1 + 1 * 2) + 0 \\&= 1.34\end{aligned}$$

Distance from the hyperplane

Why length normalized?

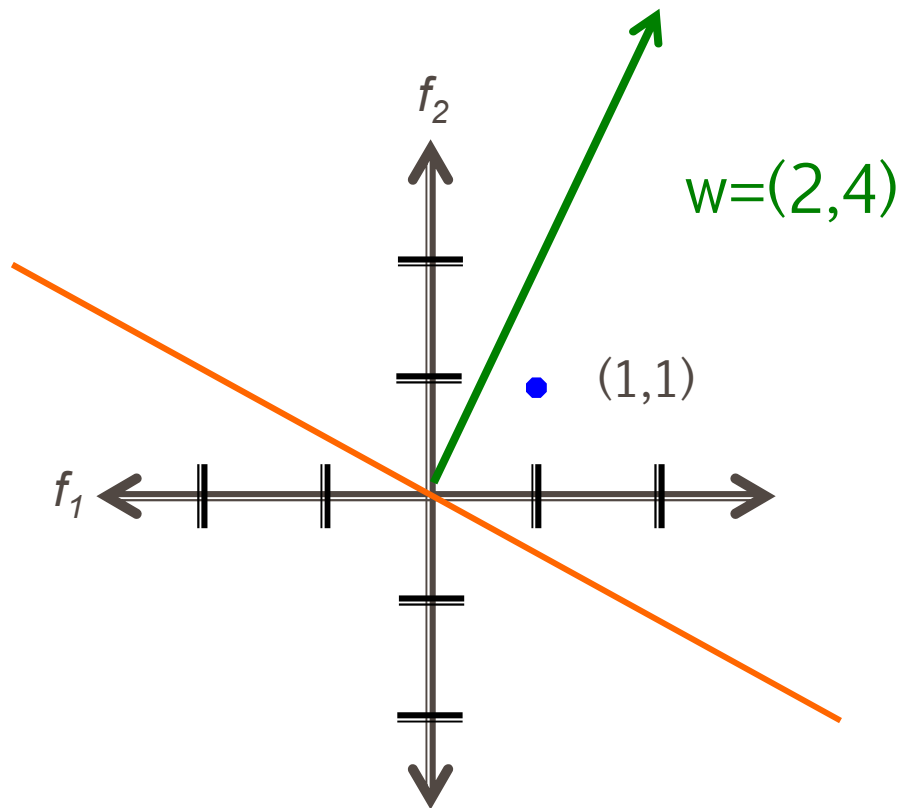


$$d(x) = \frac{w}{\|w\|} \cdot x + b$$

length
normalized
weight vectors

Distance from the hyperplane

Why length normalized?

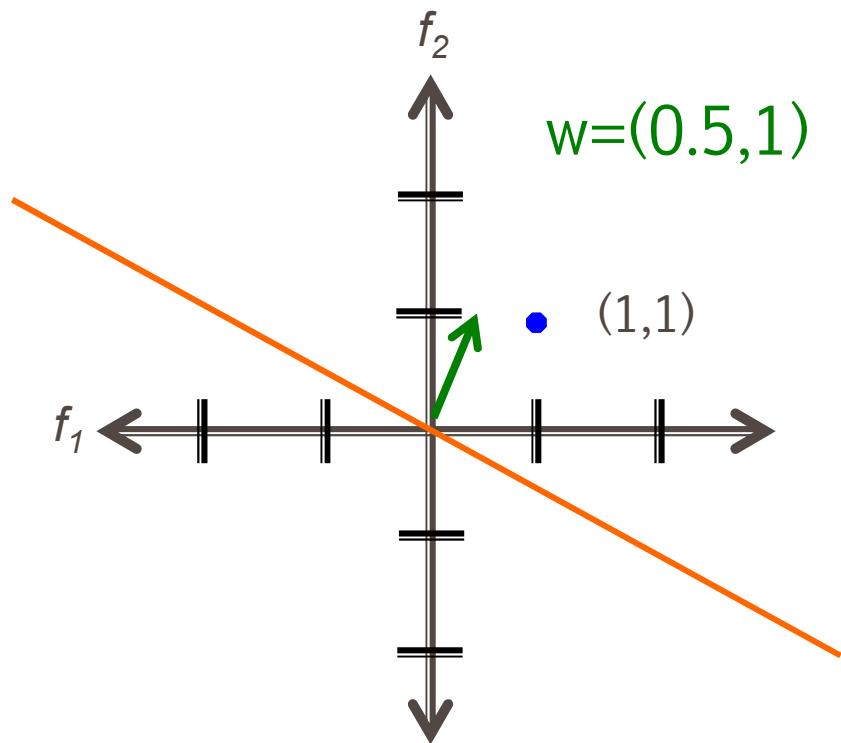


$$d(x) = \frac{w}{\|w\|} \cdot x + b$$

length
normalized
weight vectors

Distance from the hyperplane

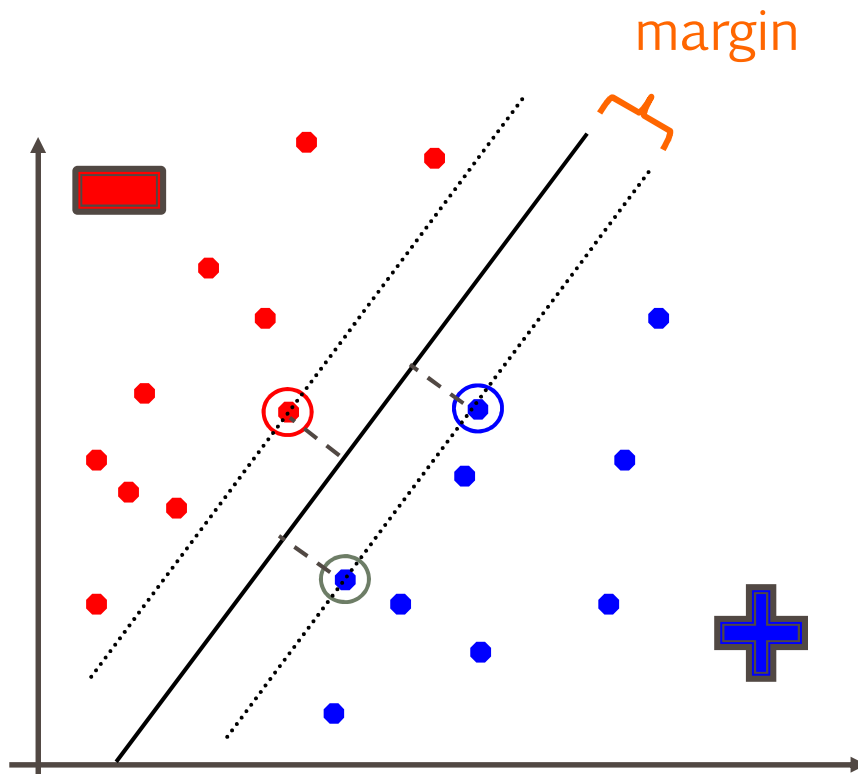
Why length normalized?



$$d(x) = \frac{w}{\|w\|} \cdot x + b$$

length
normalized
weight vectors

Measuring the margin



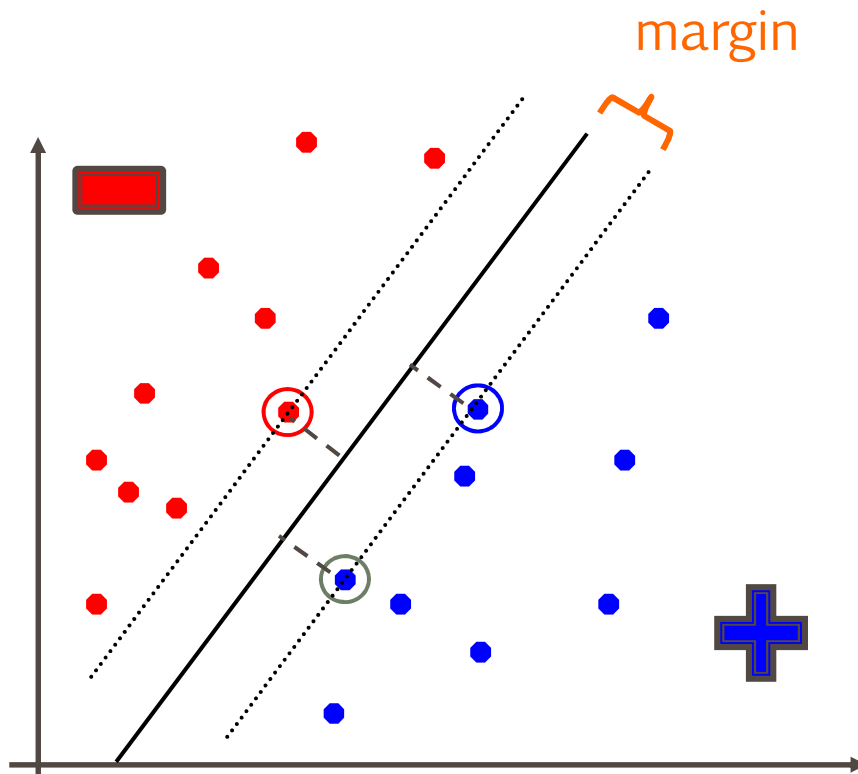
Thought experiment:

Someone gives you the optimal support vectors

Where is the max margin hyperplane?

Measuring the margin

$$d(x) = \frac{w}{\|w\|} \cdot x + b$$



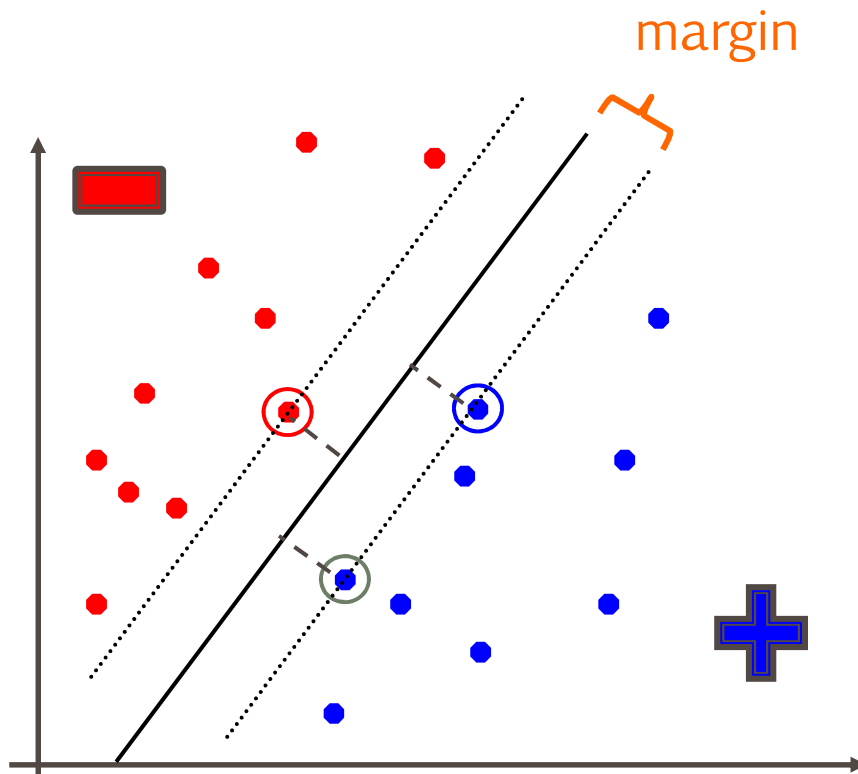
$$\text{Margin} = (d^+ - d^-) / 2$$

Max margin hyperplane is halfway in between the positive support vectors and the negative support vectors

Why?

Measuring the margin

$$d(x) = \frac{w}{\|w\|} \cdot x + b$$



$$\text{Margin} = (d^+ - d^-) / 2$$

Max margin hyperplane is halfway in between the positive support vectors and the negative support vectors

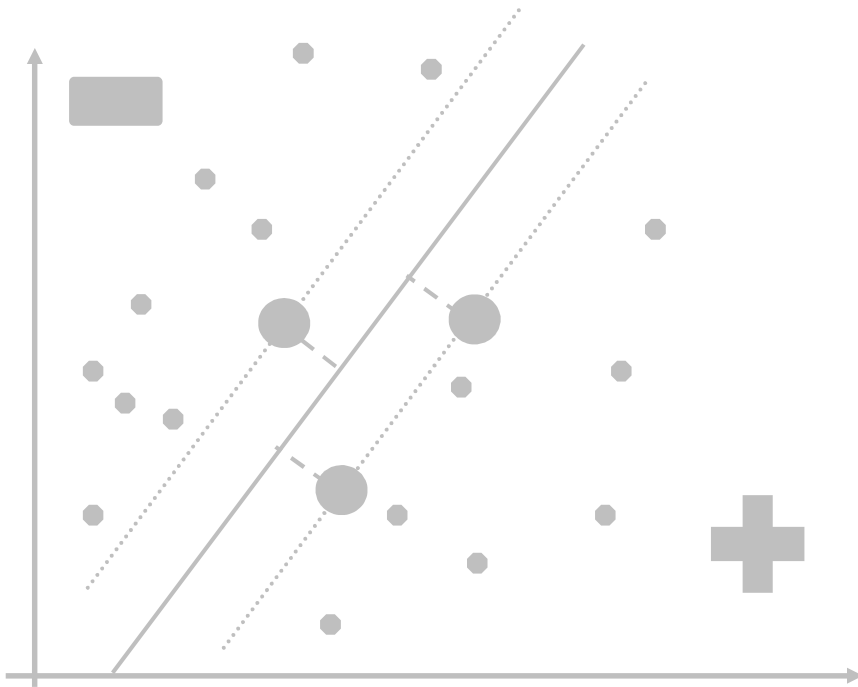
- All support vectors are the same distance
- To maximize, hyperplane should be directly in between

Measuring the margin

$$d(x) = \frac{w}{\|w\|} \cdot x + b$$

$$\text{Margin} = (d^+ - d^-) / 2$$

$$\text{margin} = \frac{1}{2} \left(\frac{w}{\|w\|} \cdot x^+ + b - \left(\frac{w}{\|w\|} \cdot x^- + b \right) \right)$$



What is $w \cdot x + b$ for support vectors?

Hint:

$$\max_{w,b} \text{margin}(w,b)$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

Measuring the margin

$$\max_{w,b} \text{margin}(w,b)$$

subject to:

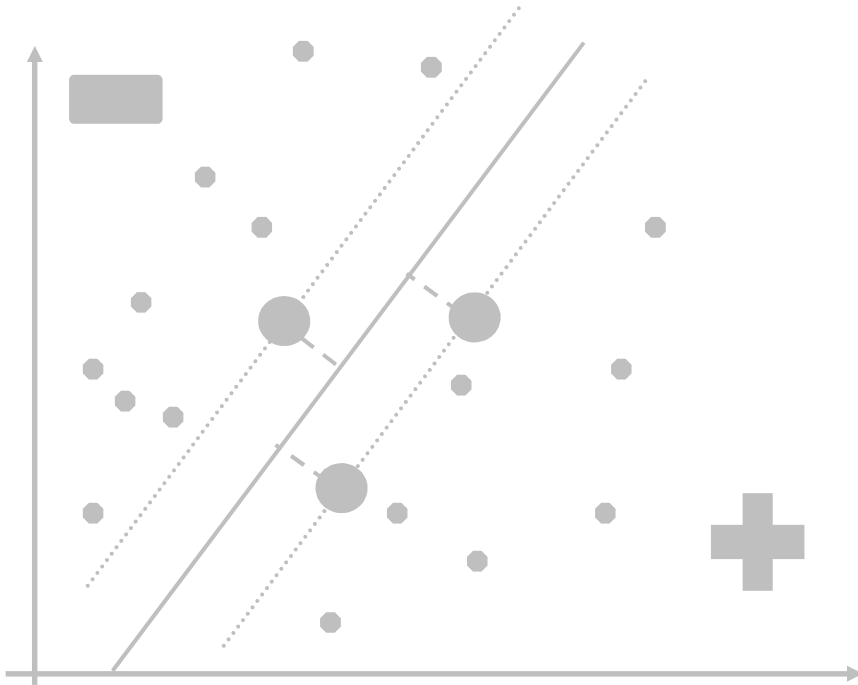
$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

The support vectors have $y_i(w \cdot x_i + b) = 1$

Otherwise, we could make the margin larger!

Measuring the margin

$$d(x) = \frac{w}{\|w\|} \cdot x + b$$



$$\text{Margin} = (d^+ - d^-) / 2$$

$$\text{margin} = \frac{1}{2} \left(\frac{w}{\|w\|} \cdot x^+ + b - \left(\frac{w}{\|w\|} \cdot x^- + b \right) \right)$$

$$= \frac{1}{2} \left(\frac{1}{\|w\|} - \frac{-1}{\|w\|} \right) \quad \leftarrow \text{negative example}$$

$$= \frac{1}{\|w\|}$$

Maximizing the margin

$$\begin{aligned} & \max_{w,b} \frac{1}{\|w\|} \\ & \text{subject to:} \\ & y_i(w \cdot x_i + b) \geq 1 \quad \forall i \end{aligned}$$

Maximizing the margin is equivalent to minimizing $\|w\|$!
(subject to the separating constraints)

Maximizing the margin

$$\min_{w,b} \|w\|$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

Maximizing the margin is equivalent to minimizing $\|w\|$
(subject to the separating constraints)

Maximizing the margin

The minimization criterion wants w to be as small as possible

$$\min_{w,b} \|w\|$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

The constraints:

1. make sure the data is separable
2. encourages w to be larger (once the data is separable)

Maximizing the margin: the real problem

$$\begin{aligned} & \min_{w,b} \quad \|w\|^2 \\ & \text{subject to:} \\ & \quad y_i(w \cdot x_i + b) \geq 1 \quad \forall i \end{aligned}$$

Why the squared?

Maximizing the margin: the real problem

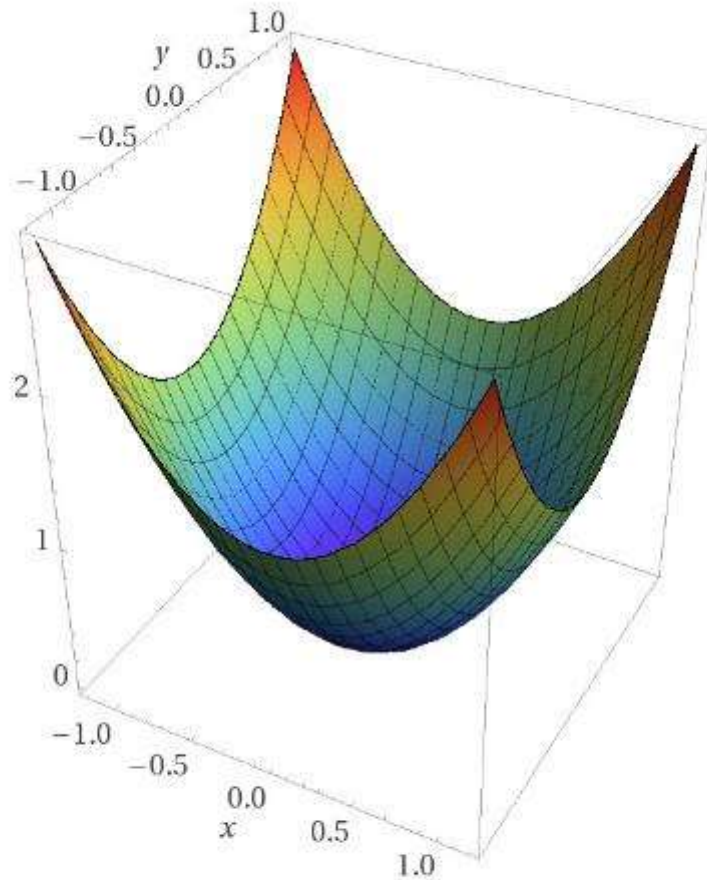
$$\begin{array}{ll} \min_{w,b} & \|w\| = \sqrt{\sum_i w_i^2} \\ \text{subject to:} & \\ & y_i(w \cdot x_i + b) \geq 1 \quad \forall i \end{array}$$

$$\begin{array}{ll} \min_{w,b} & \|w\|^2 = \sum_i w_i^2 \\ \text{subject to:} & \\ & y_i(w \cdot x_i + b) \geq 1 \quad \forall i \end{array}$$

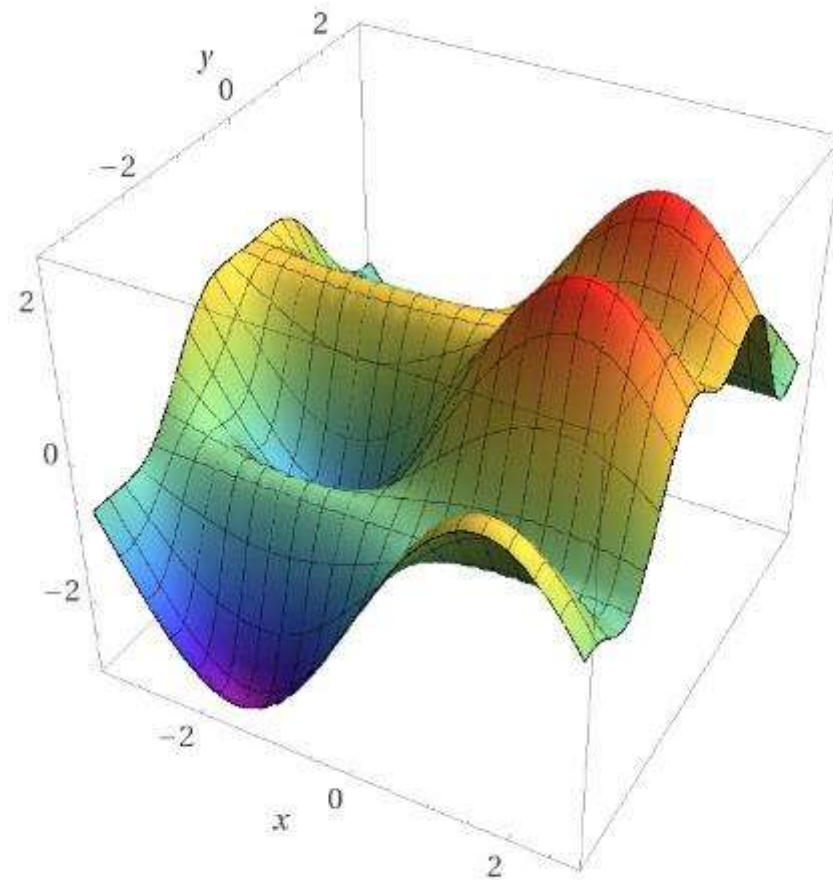
Minimizing $\|w\|$ is equivalent to minimizing $\|w\|^2$

The sum of the squared weights is a convex function!

Convex and Nonconvex



Computed by Wolfram|Alpha



Computed by Wolfram|Alpha

Source: <https://www.oreilly.com/ideas/the-hard-thing-about-deep-learning>

Support vector machine problem

$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 \\ \text{subject to:} \quad & y_i(w \cdot x_i + b) \geq 1 \quad \forall i \end{aligned}$$

This is a version of a **quadratic optimization problem**

Maximize/minimize a quadratic function

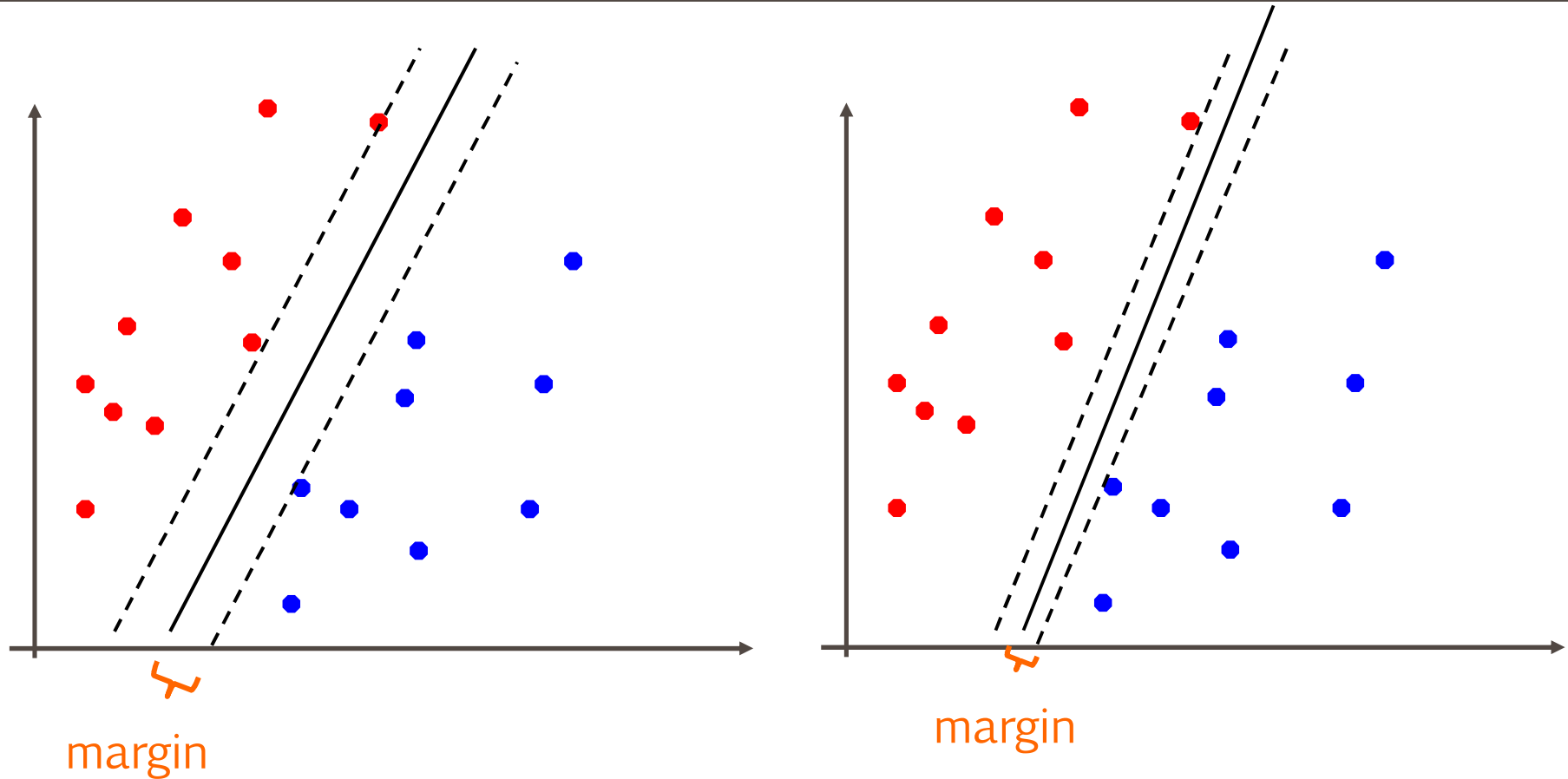
Subject to a set of linear constraints

Many, many variants of solving this problem (we'll see one in a bit)



SOFT LARGE MARGIN CLASSIFIERS

Large margin classifiers



The **margin** of a classifier is the distance to the closest points of either class

Large margin classifiers attempt to maximize this

Support vector machine problem

$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 \\ \text{subject to:} \quad & y_i(w \cdot x_i + b) \geq 1 \quad \forall i \end{aligned}$$

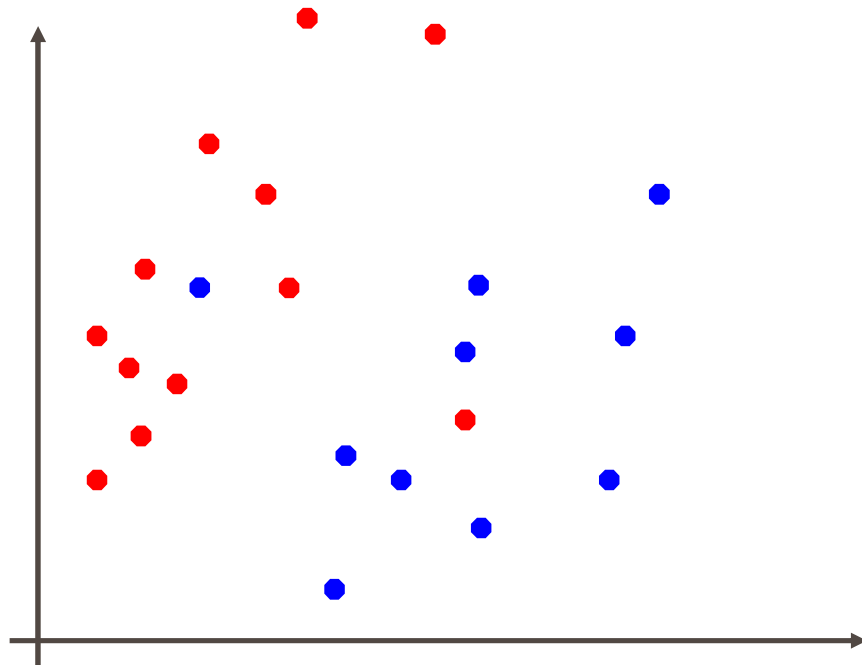
This is a a **quadratic optimization problem**

Maximize/minimize a quadratic function

Subject to a set of linear constraints

Many, many variants of solving this problem (we'll see one in a bit)

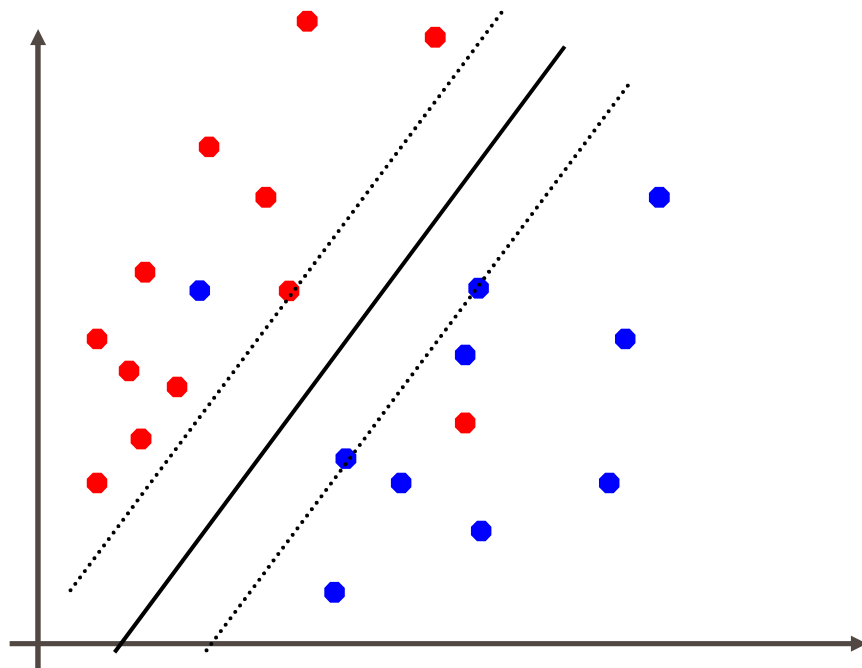
Soft Margin Classification



$$\begin{aligned} & \min_{w,b} \|w\|^2 \\ & \text{subject to:} \\ & y_i(w \cdot x_i + b) \geq 1 \quad \forall i \end{aligned}$$

What about this problem?

Soft Margin Classification

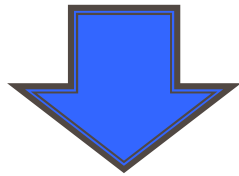


$$\begin{aligned} & \min_{w,b} \quad \|w\|^2 \\ & \text{subject to:} \\ & \quad y_i(w \cdot x_i + b) \geq 1 \quad \forall i \end{aligned}$$

We'd like to learn something like this,
but our constraints won't allow it.

Slack variables

$$\begin{aligned} & \min_{w,b} \quad \|w\|^2 \\ & \text{subject to:} \\ & \quad y_i(w \cdot x_i + b) \geq 1 \quad \forall i \end{aligned}$$

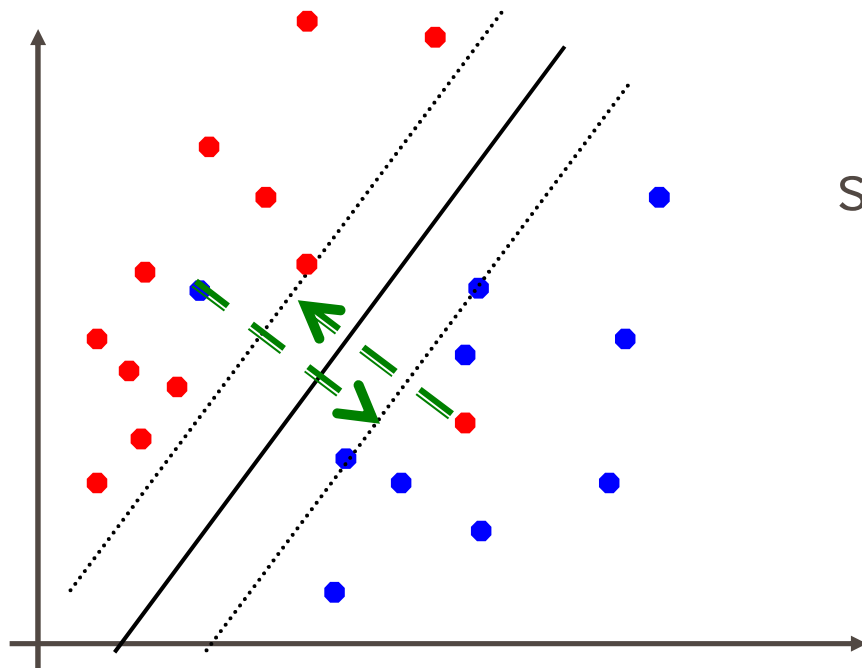


$$\begin{aligned} & \min_{w,b} \quad \|w\|^2 + C \sum_i \zeta_i \\ & \text{subject to:} \\ & \quad y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i \\ & \quad \zeta_i \geq 0 \end{aligned}$$

slack variables
(one for each example)

What effect does this have?

Slack variables



slack penalties

$$\min_{w,b} \quad \|w\|^2 + C \sum_i \zeta_i$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i$$
$$\zeta_i \geq 0$$

Slack variables

margin

trade-off between margin maximization and penalization

penalized by how far from “correct”

subject to:

$$\min_{w,b} \|w\|^2 + C \sum_i \zeta_i$$
$$y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i$$
$$\zeta_i \geq 0$$

allowed to make a mistake

The diagram illustrates the components of the SVM optimization problem. A blue arrow points from the word 'margin' to the $\|w\|^2$ term in the objective function. Another blue arrow points from the text 'trade-off between margin maximization and penalization' to the coefficient C in the same term. A third blue arrow points from the text 'penalized by how far from “correct”' to the slack variable ζ_i in the summation. Below the objective function, the text 'subject to:' is followed by two constraints. A blue arrow points from the text 'allowed to make a mistake' to the slack variable ζ_i in the first constraint, indicating that ζ_i represents the margin violation for each data point.

Soft margin SVM

$$\min_{w,b} \quad \|w\|^2 + C \sum_i \zeta_i$$

subject to:

$$y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i$$
$$\zeta_i \geq 0$$

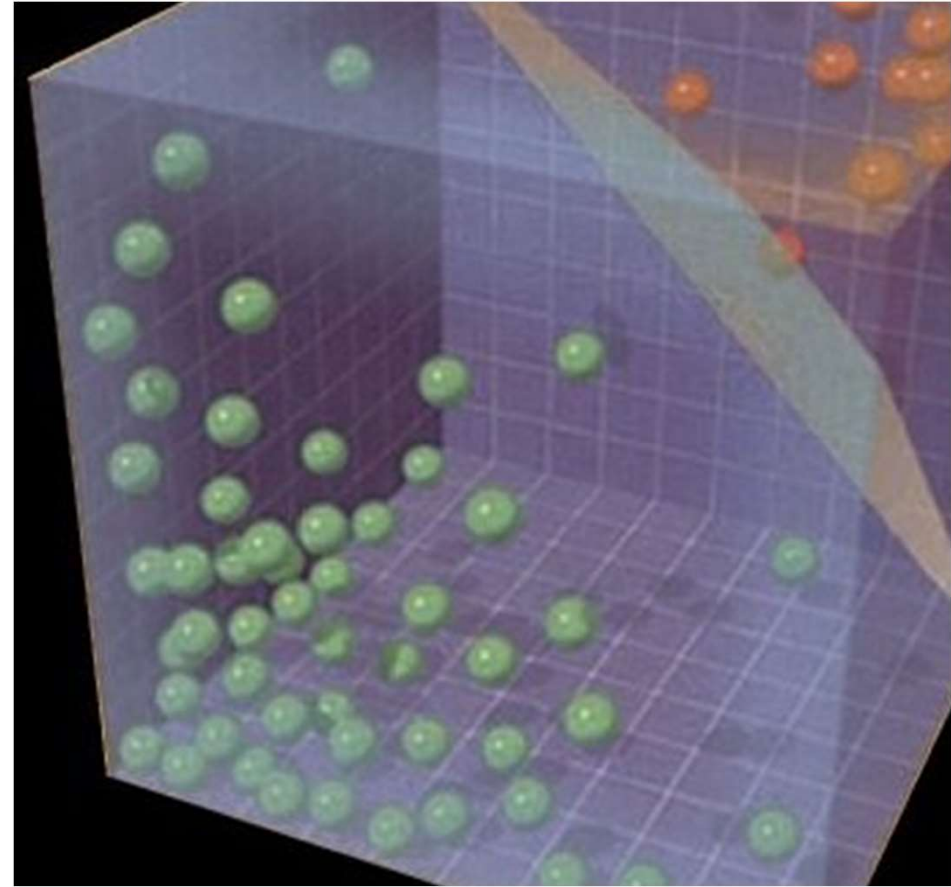
Still a **quadratic optimization problem!**

Demo

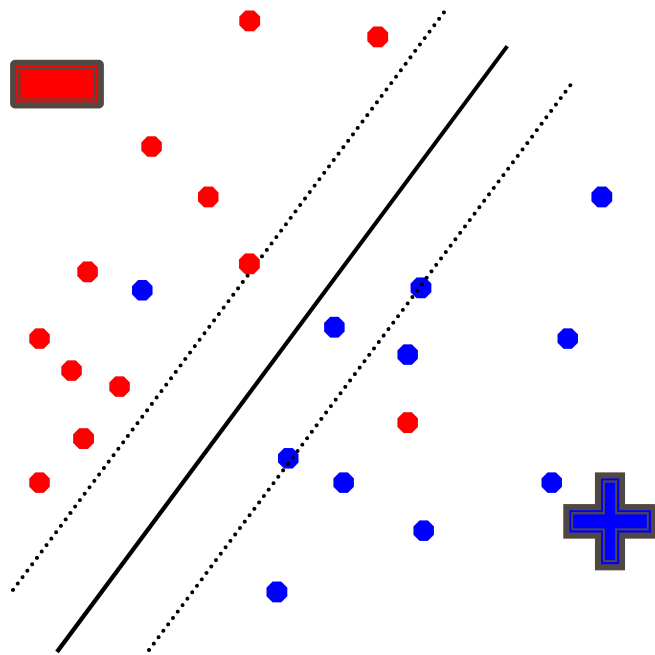
Stanford CS:

[http://cs.stanford.edu/people/karpathy/svmjs/
demo/](http://cs.stanford.edu/people/karpathy/svmjs/demo/)

SOLVING THE SVM PROBLEM



Understanding the Soft Margin SVM

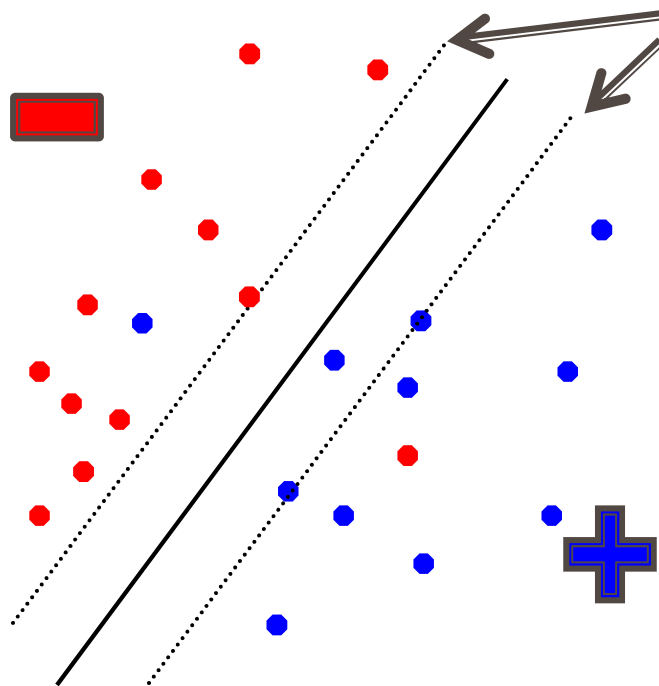


$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 + C \sum_i \zeta_i \\ \text{subject to:} \quad & y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i \\ & \zeta_i \geq 0 \end{aligned}$$

Given the optimal solution, w , b :

Can we figure out what the slack penalties are for each point?

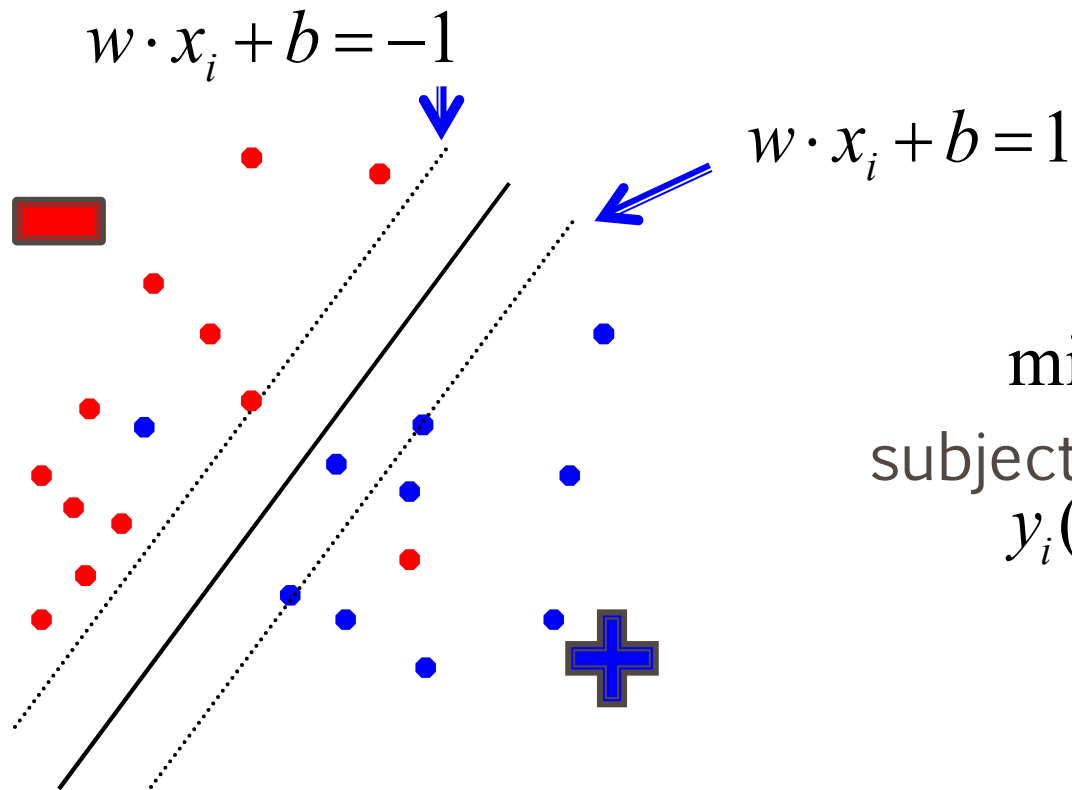
Understanding the Soft Margin SVM



What do the margin lines
represent wrt w, b ?

$$\begin{aligned} & \min_{w,b} \quad \|w\|^2 + C \sum_i \zeta_i \\ & \text{subject to:} \\ & \quad y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i \\ & \quad \zeta_i \geq 0 \end{aligned}$$

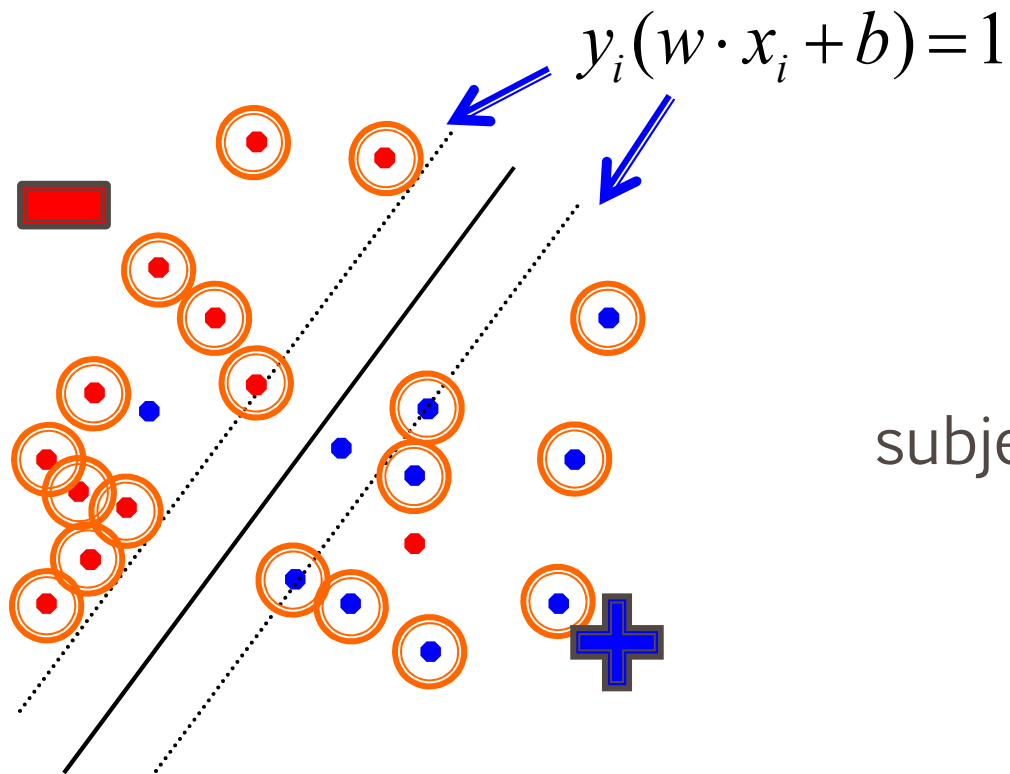
Understanding the Soft Margin SVM



$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 + C \sum_i \zeta_i \\ \text{subject to:} \quad & y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i \\ & \zeta_i \geq 0 \end{aligned}$$

Or: $y_i(w \cdot x_i + b) = 1$

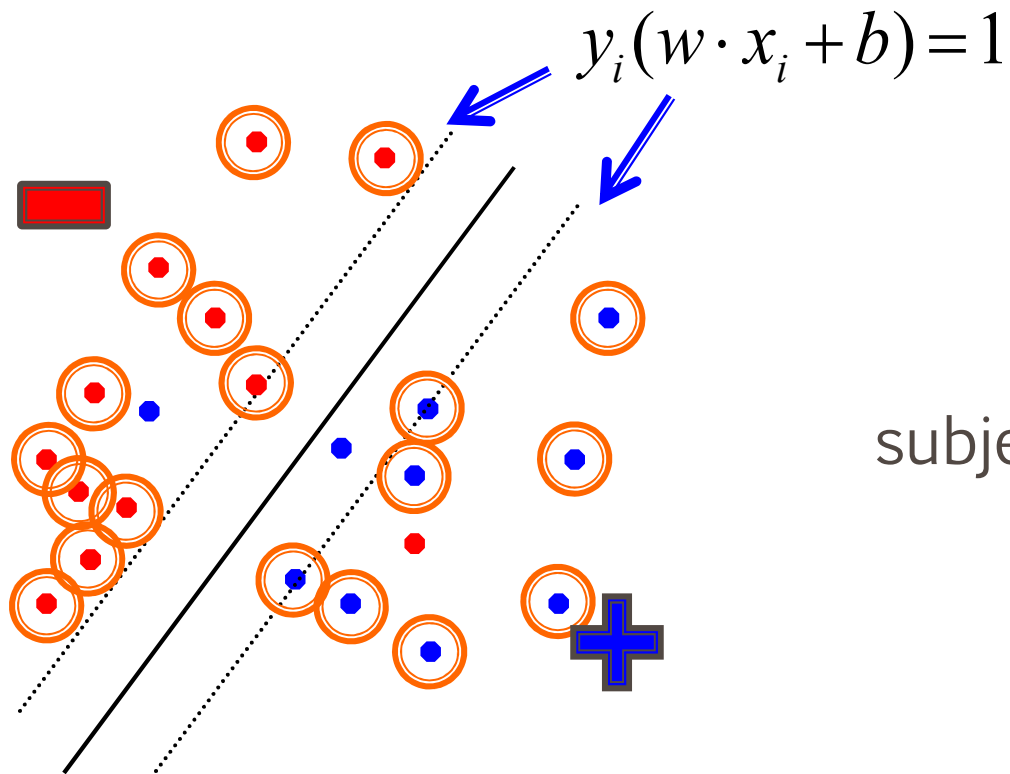
Understanding the Soft Margin SVM



$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 + C \sum_i \zeta_i \\ \text{subject to:} \quad & y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i \\ & \zeta_i \geq 0 \end{aligned}$$

What are the slack values for points outside (or on) the margin AND correctly classified?

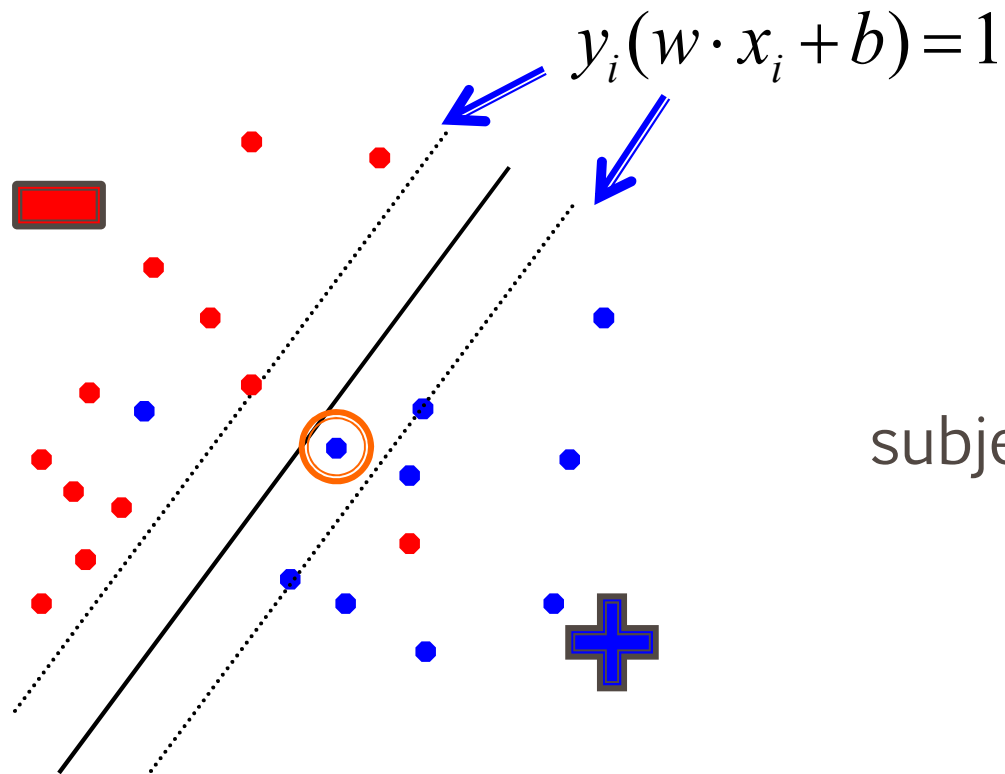
Understanding the Soft Margin SVM



$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 + C \sum_i \zeta_i \\ \text{subject to:} \quad & y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i \\ & \zeta_i \geq 0 \end{aligned}$$

0! The slack variables have to be greater than or equal to zero and if they're on or beyond the margin then $y_i(w \cdot x_i + b) \geq 1$ already

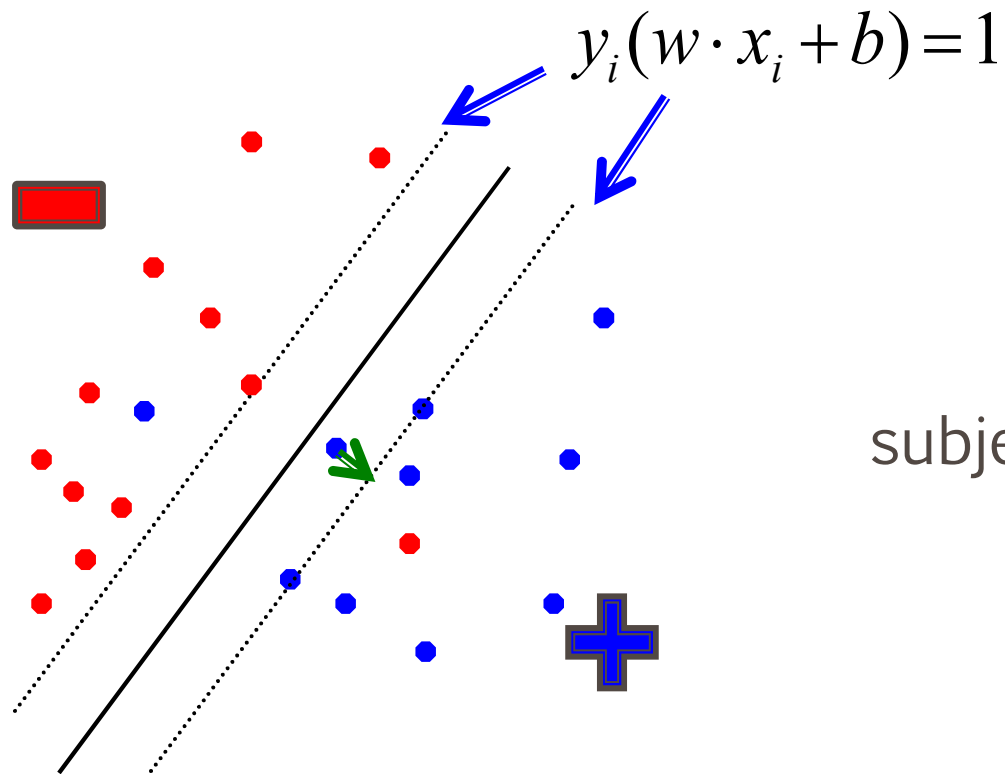
Understanding the Soft Margin SVM



$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 + C \sum_i \zeta_i \\ \text{subject to:} \quad & y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i \\ & \zeta_i \geq 0 \end{aligned}$$

What are the slack values for points inside the margin AND classified correctly?

Understanding the Soft Margin SVM

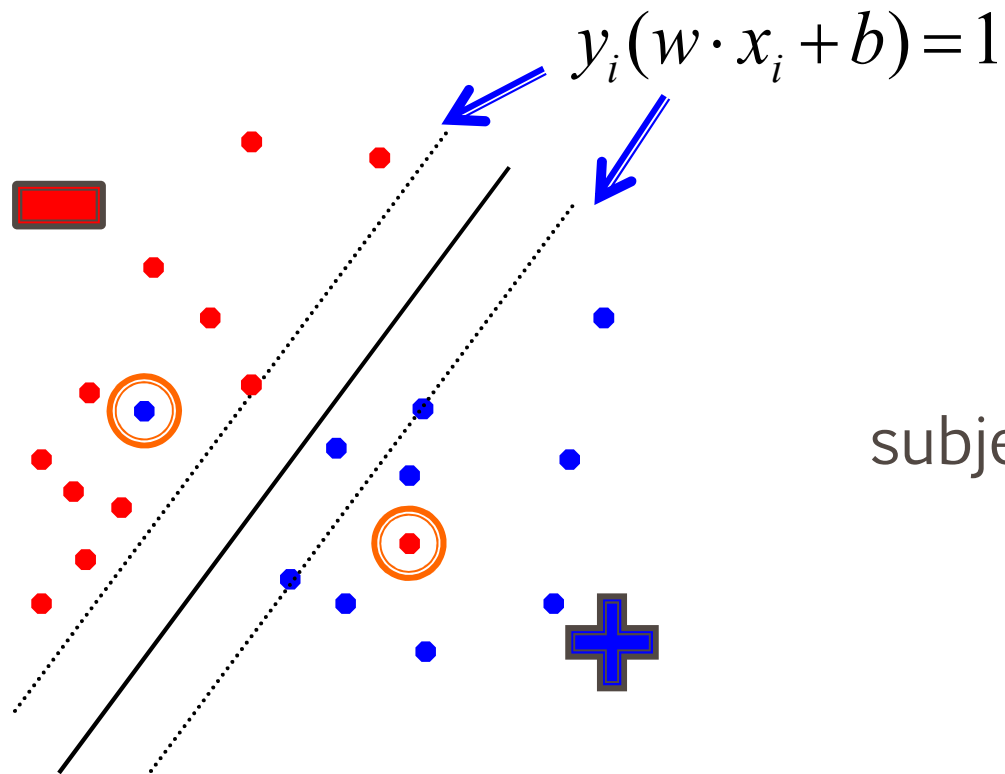


$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 + C \sum_i \zeta_i \\ \text{subject to:} \quad & y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i \\ & \zeta_i \geq 0 \end{aligned}$$

Difference from point to the margin. Which is?

$$\zeta_i = 1 - y_i(w \cdot x_i + b)$$

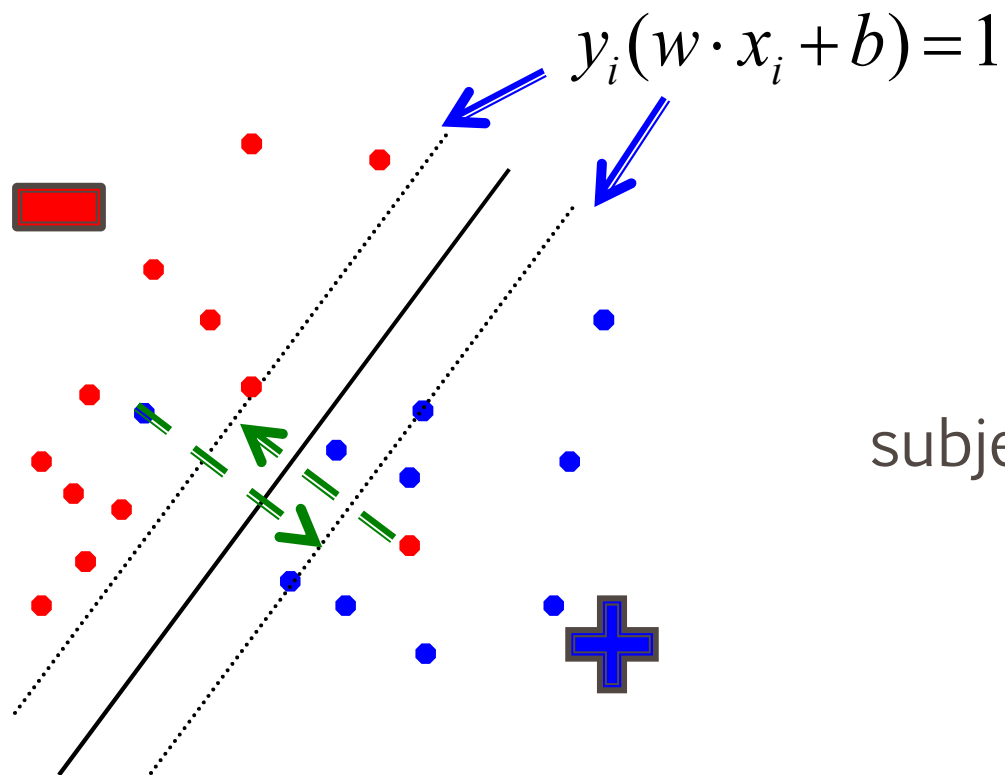
Understanding the Soft Margin SVM



$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 + C \sum_i \zeta_i \\ \text{subject to:} \quad & y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i \\ & \zeta_i \geq 0 \end{aligned}$$

What are the slack values for points that are incorrectly classified?

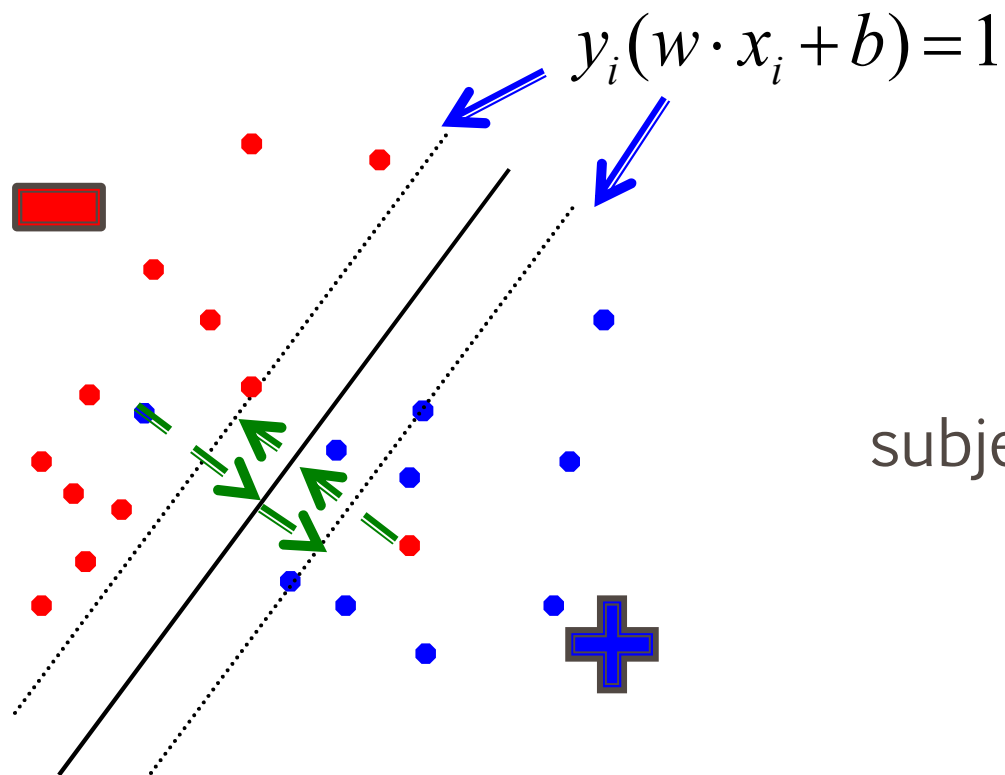
Understanding the Soft Margin SVM



$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 + C \sum_i \zeta_i \\ \text{subject to:} \quad & y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i \\ & \zeta_i \geq 0 \end{aligned}$$

Which is?

Understanding the Soft Margin SVM

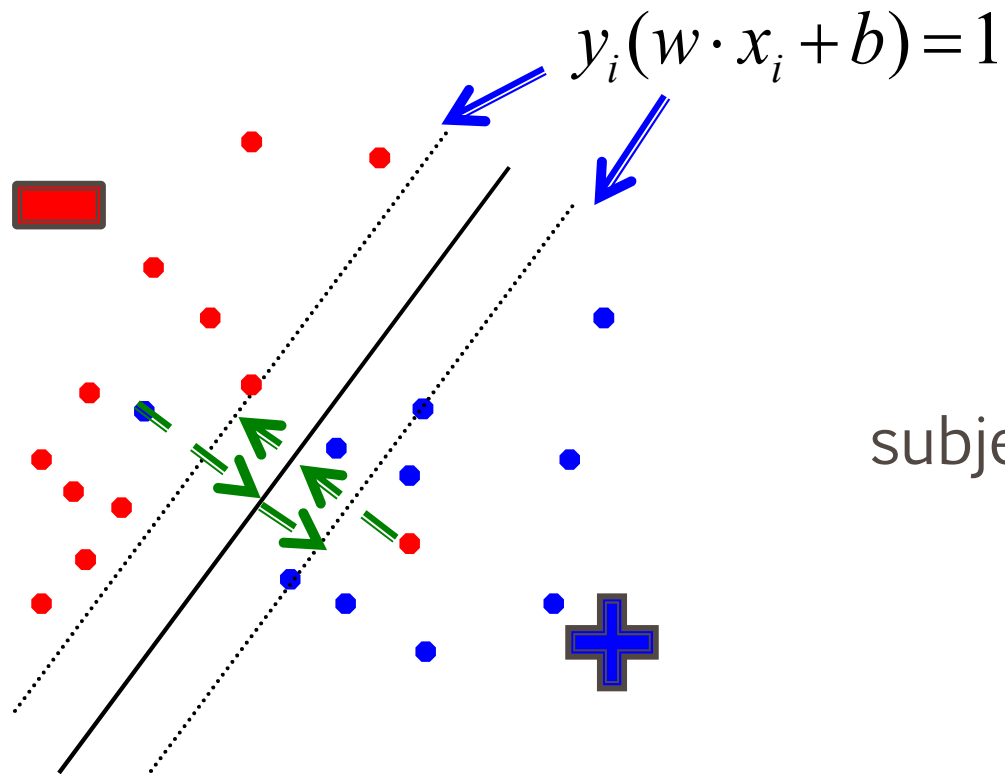


$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 + C \sum_i \zeta_i \\ \text{subject to:} \quad & y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i \\ & \zeta_i \geq 0 \end{aligned}$$

“distance” to the hyperplane *plus* the “distance” to the margin

?

Understanding the Soft Margin SVM

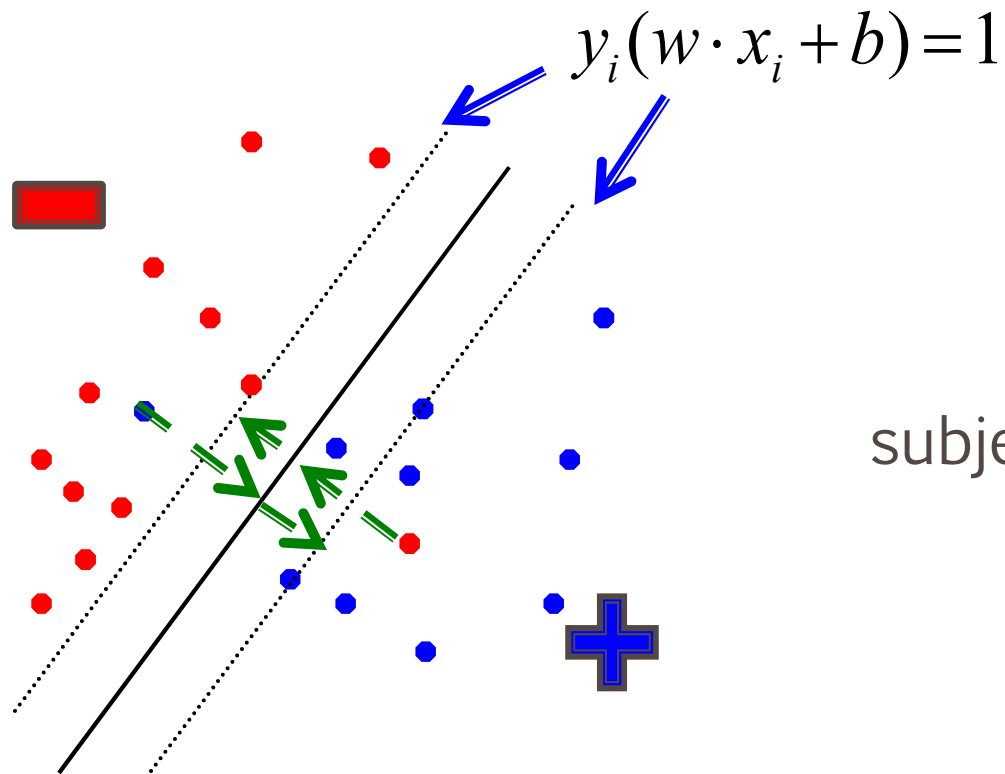


$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 + C \sum_i \zeta_i \\ \text{subject to:} \quad & y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i \\ & \zeta_i \geq 0 \end{aligned}$$

“distance” to the hyperplane *plus* the “distance” to the margin

$$-y_i(w \cdot x_i + b) \quad \text{Why -?}$$

Understanding the Soft Margin SVM



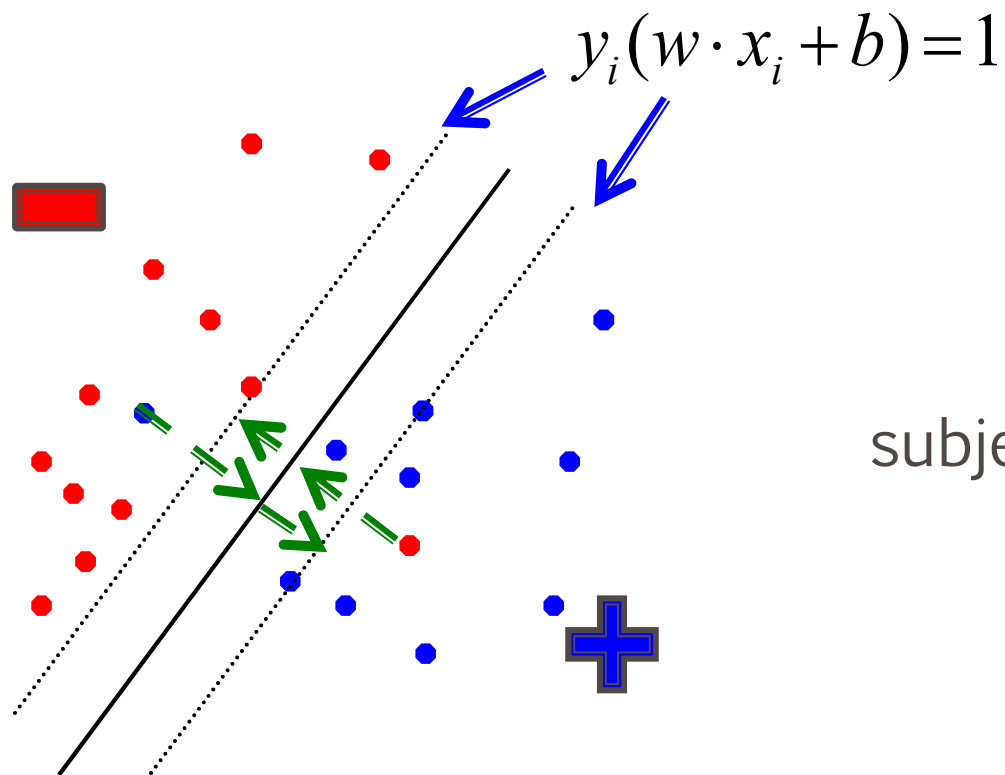
$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 + C \sum_i \xi_i \\ \text{subject to:} \quad & y_i(w \cdot x_i + b) \geq 1 - \xi_i \quad \forall i \\ & \xi_i \geq 0 \end{aligned}$$

“distance” to the hyperplane *plus* the “distance” to the margin

$$-y_i(w \cdot x_i + b)$$

?

Understanding the Soft Margin SVM



$$\min_{w,b} \quad \|w\|^2 + C \sum_i \xi_i$$

subject to:

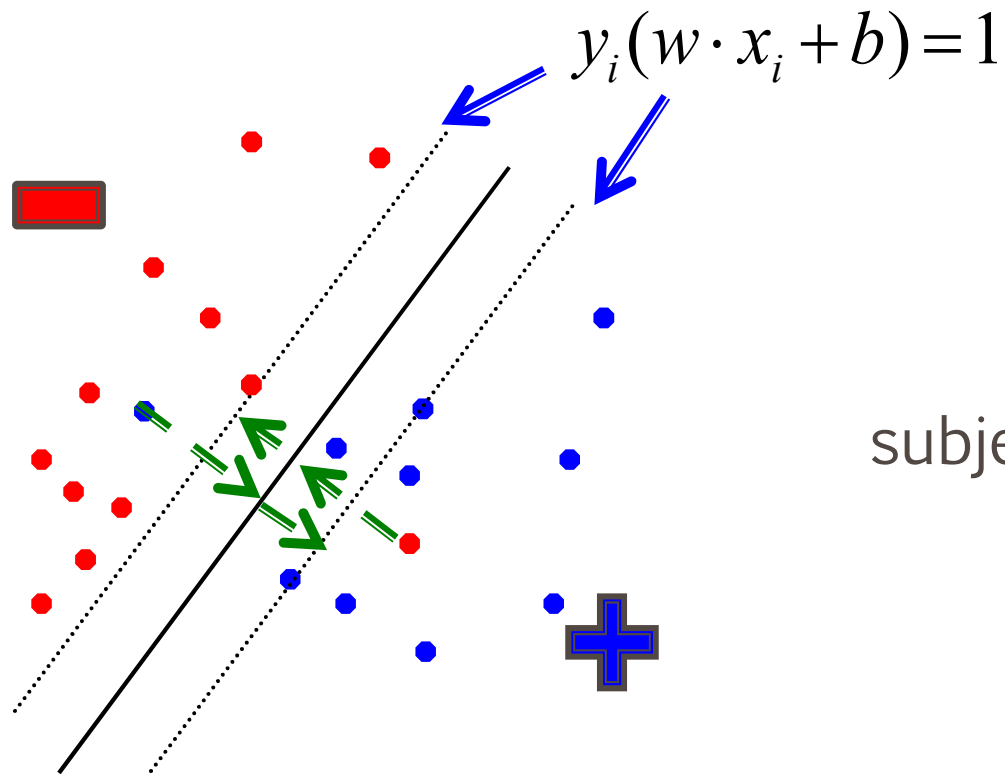
$$y_i(w \cdot x_i + b) \geq 1 - \xi_i \quad \forall i$$

$$\xi_i \geq 0$$

“distance” to the hyperplane *plus* the “distance” to the margin

$$-y_i(w \cdot x_i + b)$$

Understanding the Soft Margin SVM



$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 + C \sum_i \zeta_i \\ \text{subject to:} \quad & y_i(w \cdot x_i + b) \geq 1 - \zeta_i \quad \forall i \\ & \zeta_i \geq 0 \end{aligned}$$

“distance” to the hyperplane *plus* the “distance” to the margin

$$\zeta_i = 1 - y_i(w \cdot x_i + b)$$

Understanding the Soft Margin SVM

$$\begin{aligned} \min_{w,b} \quad & \|w\|^2 + C \sum_i \varsigma_i \\ \text{subject to:} \quad & y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i \\ & \varsigma_i \geq 0 \end{aligned}$$

$$\varsigma_i = \begin{cases} 0 & \text{if } y_i(w \cdot x_i + b) \geq 1 \\ 1 - y_i(w \cdot x_i + b) & \text{otherwise} \end{cases}$$

Understanding the Soft Margin SVM

$$\zeta_i = \begin{cases} 0 & \text{if } y_i(w \cdot x_i + b) \geq 1 \\ 1 - y_i(w \cdot x_i + b) & \text{otherwise} \end{cases}$$



$$\begin{aligned} \zeta_i &= \max(0, 1 - y_i(w \cdot x_i + b)) \\ &= \max(0, 1 - yy') \end{aligned}$$

Hinge loss!

0/1 loss: $l(y, y') = 1[y y' \leq 0]$

Hinge: $l(y, y') = \max(0, 1 - y y')$

Exponential: $l(y, y') = \exp(-y y')$

Squared loss: $l(y, y') = (y - y')^2$

Understanding the Soft Margin SVM

$$\begin{aligned} & \min_{w,b} \quad \|w\|^2 + C \sum_i \varsigma_i \\ & \text{subject to:} \\ & \quad y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i \\ & \quad \varsigma_i \geq 0 \end{aligned}$$

$$\varsigma_i = \max(0, 1 - y_i(w \cdot x_i + b))$$

Do we need the constraints still?

Understanding the Soft Margin SVM

$$\begin{aligned} & \min_{w,b} \quad \|w\|^2 + C \sum_i \varsigma_i \\ & \text{subject to:} \\ & \quad y_i(w \cdot x_i + b) \geq 1 - \varsigma_i \quad \forall i \\ & \quad \varsigma_i \geq 0 \end{aligned}$$

$$\varsigma_i = \max(0, 1 - y_i(w \cdot x_i + b))$$



$$\min_{w,b} \quad \|w\|^2 + C \sum_i \max(0, 1 - y_i(w \cdot x_i + b))$$

Unconstrained problem!

Understanding the Soft Margin SVM

$$\min_{w,b} \quad \|w\|^2 + C \sum_i \text{loss}_{\text{hinge}}(y_i, y_i')$$

Does this look like something we've seen before?

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \text{loss}(y_i y_i') + \lambda \text{ regularizer}(w, b)$$

Gradient descent problem!

Soft margin SVM as gradient descent

$$\min_{w,b} \|w\|^2 + C \sum_i \text{loss}_{\text{hinge}}(y_i, y_i')$$

multiply through by 1/C
and rearrange

$$\min_{w,b} \sum_i \text{loss}_{\text{hinge}}(y_i, y_i') + \frac{1}{C} \|w\|^2$$

let $\lambda = 1/C$

$$\min_{w,b} \sum_i \text{loss}_{\text{hinge}}(y_i, y_i') + \lambda \|w\|^2$$

What type of gradient descent problem?

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \text{loss}(y_i, y_i') + \lambda \text{ regularizer}(w, b)$$

Soft margin SVM as gradient descent

One way to solve the soft margin SVM problem is using gradient descent

$$\min_{w,b} \sum_i \text{loss}_{\text{hinge}}(y_i, y_i') + \lambda \|w\|^2$$

Diagram illustrating the components of the soft margin SVM objective function:

- The term $\sum_i \text{loss}_{\text{hinge}}(y_i, y_i')$ is labeled "hinge loss" (indicated by a blue arrow).
- The term $\lambda \|w\|^2$ is labeled "L2 regularization" (indicated by a blue arrow).

Gradient descent SVM solver

- pick a starting point (w)
- repeat until loss doesn't decrease in all dimensions:
 - pick a dimension
 - move a small amount in that dimension towards decreasing loss (using the derivative)

$$w_i = w_i - \eta \frac{d}{dw_i} (\text{loss}(w) + \text{regularizer}(w, b))$$

$$w_j = w_j + \eta \sum_{i=1}^n y_i x_i 1[y_i(w \cdot x + b) < 1] - \eta \lambda w_j$$

hinge loss

L2 regularization

Finds the largest margin hyperplane while allowing for a soft margin

Support vector machines

One of the most successful (if not the most successful) classification approach:

decision tree

About 2,160,000 results (0.05 sec)

Support vector machine

About 1,960,000 results (0.04 sec)

k nearest neighbor

About 746,000 results (0.04 sec)

perceptron algorithm

About 84,300 results (0.04 sec)

