

UVA CS 4501:

Machine Learning

Lecture 16: Support Vector Machine

Dr. Yanjun Qi

University of Virginia
Department of Computer Science

Where are we ? →

Five major sections of this course

- Regression (supervised)
- Classification (supervised)
- Unsupervised models
- Learning theory
- Graphical models

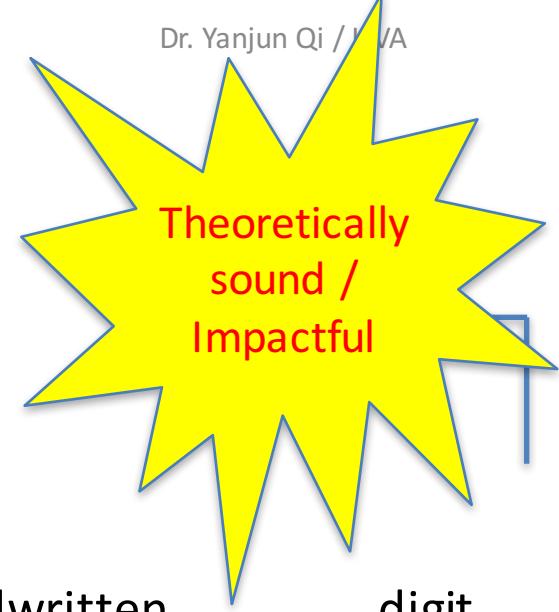
Today

❑ Support Vector Machine (SVM)

- ✓ History of SVM
- ✓ Large Margin Linear Classifier
- ✓ Define Margin (M) in terms of model parameter
- ✓ Optimization to learn model parameters (w , b)
- ✓ Linearly Non-separable case
- ✓ Optimization with dual form
- ✓ Nonlinear decision boundary
- ✓ Multiclass SVM

History of SVM

- SVM is inspired from statistical learning theory [3]
- SVM was first introduced in 1992 [1]
- SVM becomes popular because of its success in handwritten recognition (1994)
 - 1.1% test error rate for SVM.
 - The same as the error rates of a carefully constructed neural network, LeNet 4.
 - Section 5.11 in [2] or the discussion in [3] for details
- Regarded as an important example of “kernel methods”, **arguably the hottest area in machine learning 20 years ago**



- [1] B.E. Boser *et al.* A Training Algorithm for Optimal Margin Classifiers. Proceedings of the Fifth Annual Workshop on Computational Learning Theory 5 144-152, Pittsburgh, 1992.
- [2] L. Bottou *et al.* Comparison of classifier methods: a case study in handwritten digit recognition. Proceedings of the 12th IAPR International Conference on Pattern Recognition, vol. 2, pp. 77-82, 1994.
- [3] V. Vapnik. The Nature of Statistical Learning Theory. 2nd edition, Springer, 1999.

Applications of SVMs

- Computer Vision
- Text Categorization
- Ranking (e.g., Google searches)
- Handwritten Character Recognition
- Time series analysis
- Bioinformatics
-

→ Lots of very successful applications!!!

Handwritten digit recognition

→ MNIST

(SVM)

1994



3-nearest-neighbor = 2.4% error

400–300–10 unit MLP = 1.6% error

LeNet: 768–192–30–10 unit MLP = 0.9% error

best (kernel machines, vision algorithms) \approx 0.6% error

X ₁	X ₂	X ₃	Y

A Dataset for **binary/** classification

$$f : \boxed{X} \longrightarrow \boxed{Y}$$

Output as
Binary Class:
only two
possibilities

- **Data/points/instances/examples/samples/records:** [rows]
- **Features/attributes/dimensions/independent variables/covariates/predictors/regressors:** [columns, except the last]
- **Target/outcome/response/label/dependent variable:** special column to be predicted [last column]

Affine Hyperplanes

- <https://en.wikipedia.org/wiki/Hyperplane>
- Any hyperplane can be given in coordinates as the solution of a single linear (algebraic) equation of degree 1.

$$[a_1x_1 + a_2x_2 + a_3x_3 + \dots + a_px_p = b], \text{ at least one } a_i \neq 0$$

⇒ e.g. classification Boundary $w^T x + b = 0$

$$\begin{cases} x \in R^P \\ b \in R \end{cases}$$

Review :

Vector Product, Orthogonal, and Norm

For two vectors x and y ,

$$x^T y$$

is called the *(inner) vector product*.

x and y are called *orthogonal* if

$$x^T y = 0$$

The square root of the product of a vector with itself,

$$\sqrt{x^T x}$$

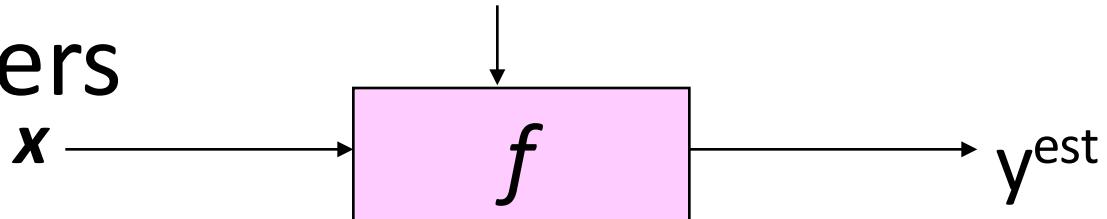
is called the *2-norm* ($\|x\|_2$), can also write as $|x|$

Today

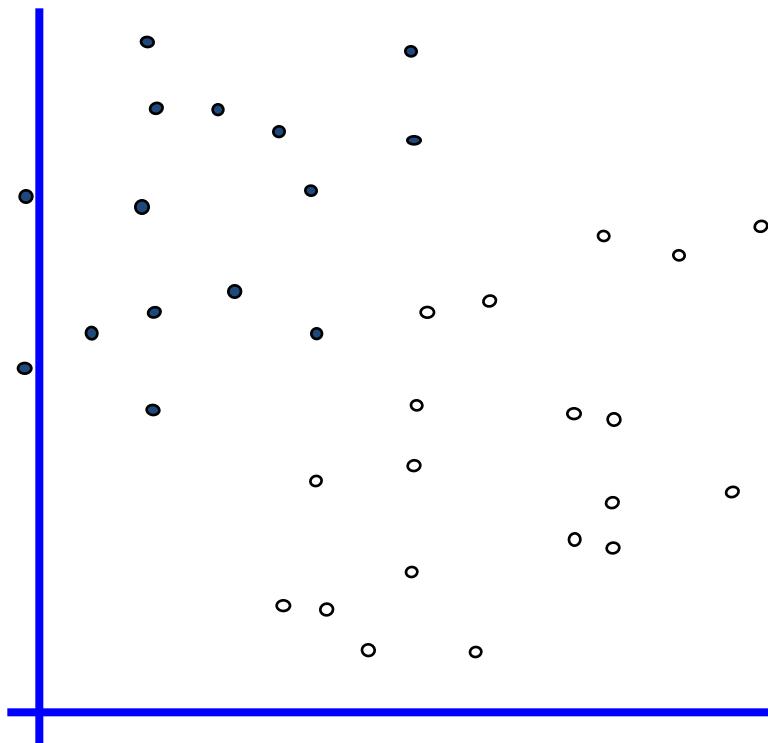
❑ Support Vector Machine (SVM)

- ✓ History of SVM
- ✓ Large Margin Linear Classifier
- ✓ Define Margin (M) in terms of model parameter
- ✓ Optimization to learn model parameters (w, b)
- ✓ Linearly Non-separable case
- ✓ Optimization with dual form
- ✓ Nonlinear decision boundary
- ✓ Multiclass SVM

Linear Classifiers

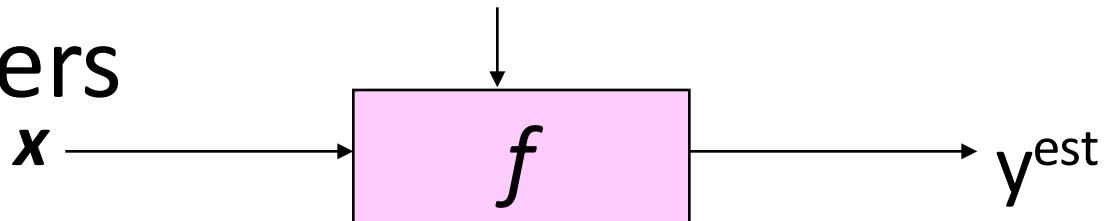


- denotes +1
- denotes -1

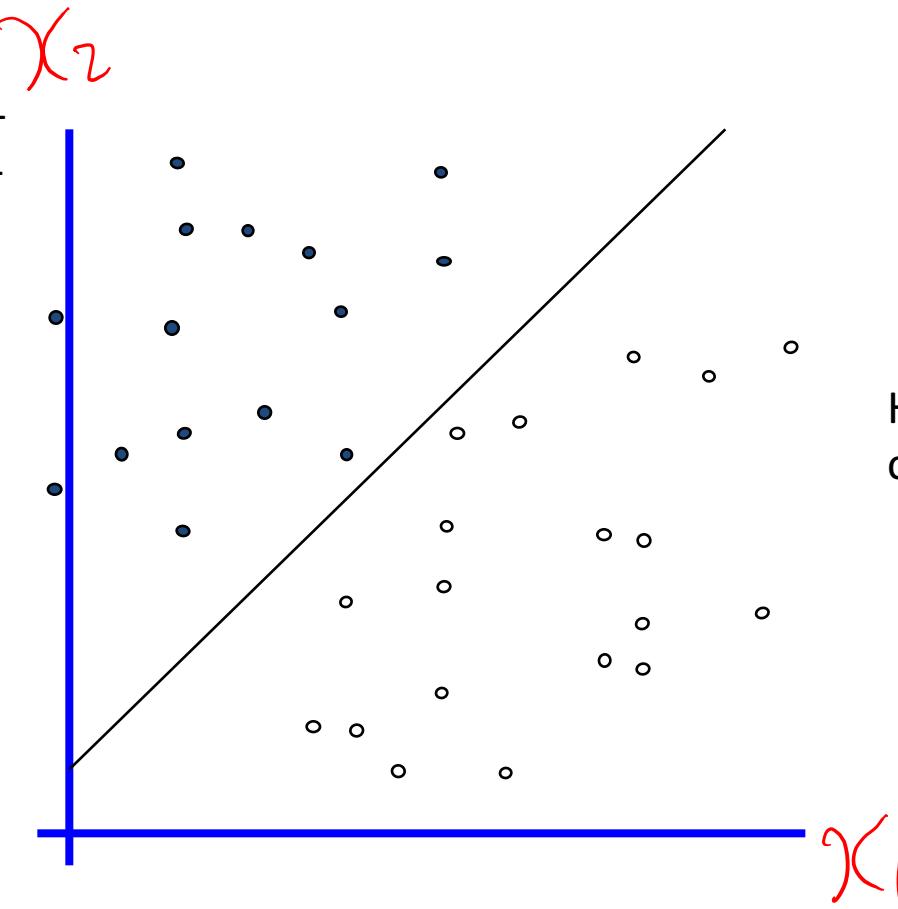


How would you
classify this data?

Linear Classifiers

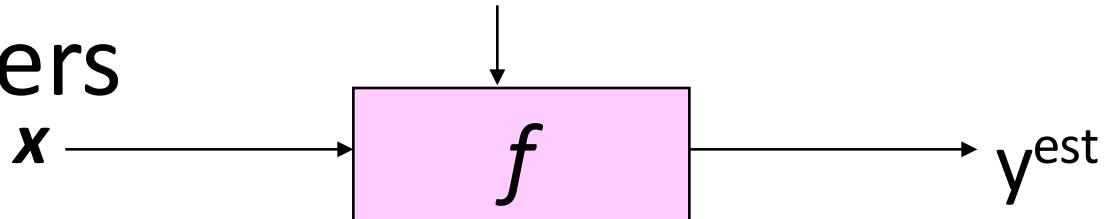


- denotes +1
- denotes -1

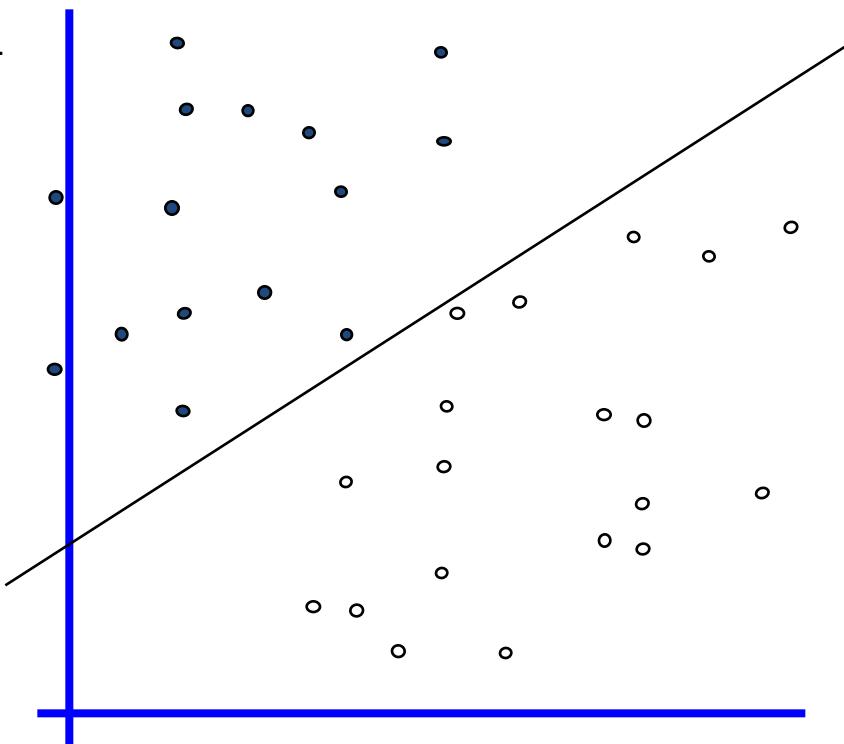


How would you
classify this data?

Linear Classifiers

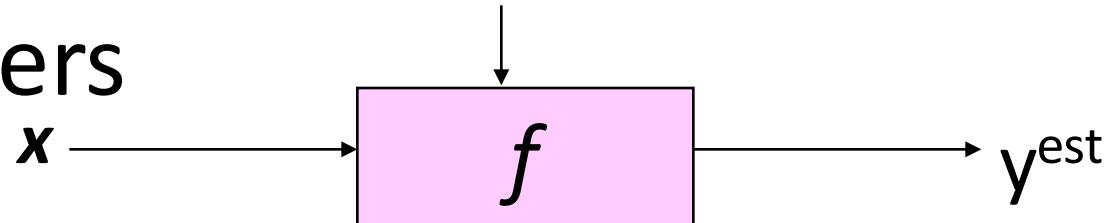


- denotes +1
- denotes -1

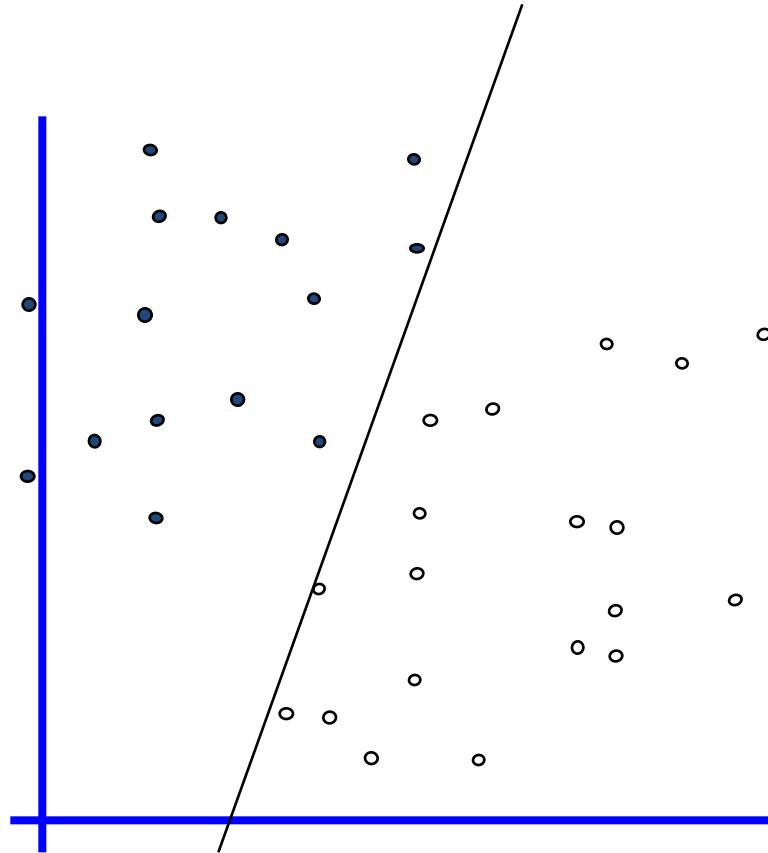


How would you
classify this data?

Linear Classifiers

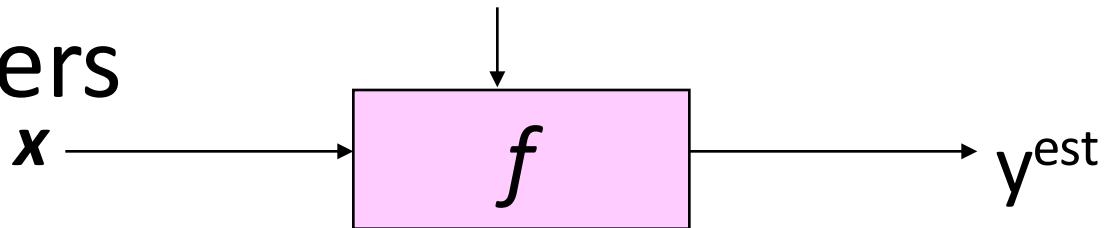


- denotes +1
- denotes -1

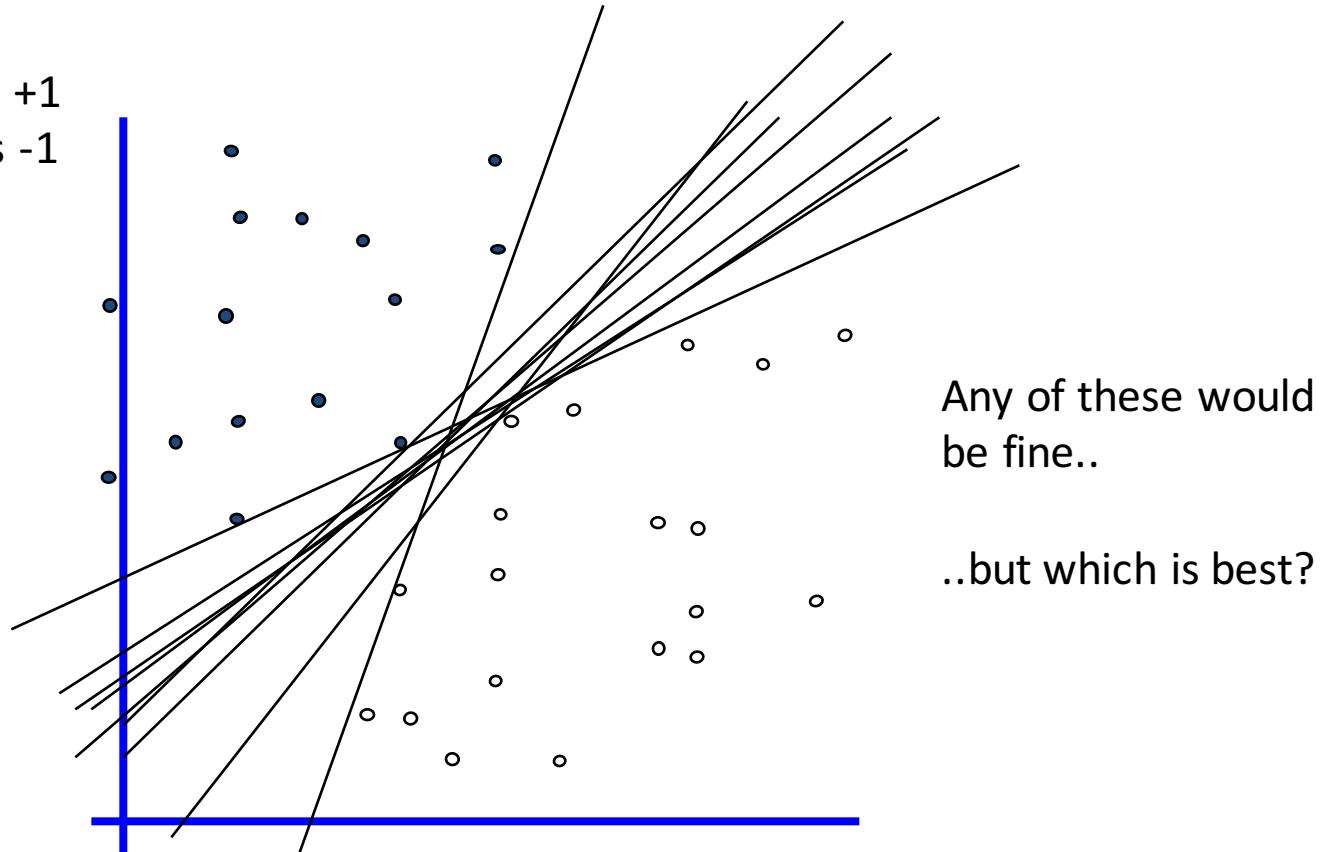


How would you
classify this data?

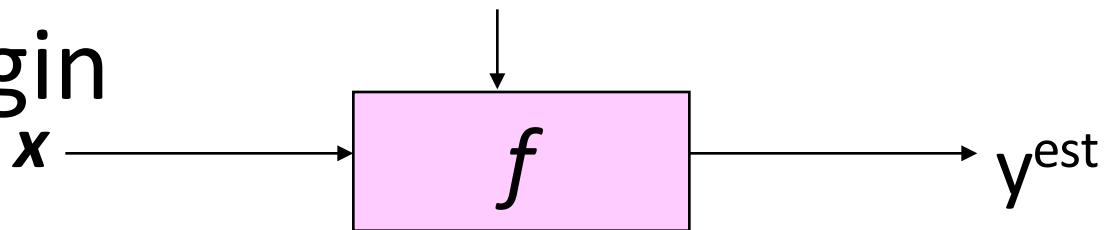
Linear Classifiers



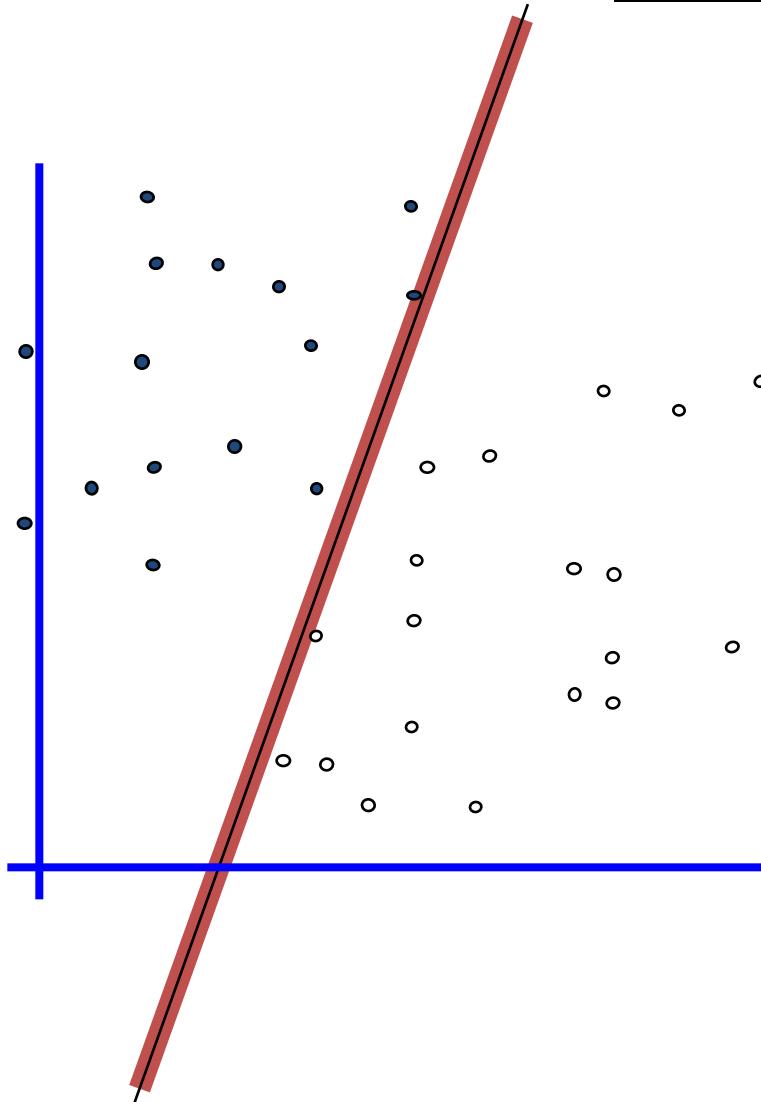
- denotes +1
- denotes -1



Classifier Margin

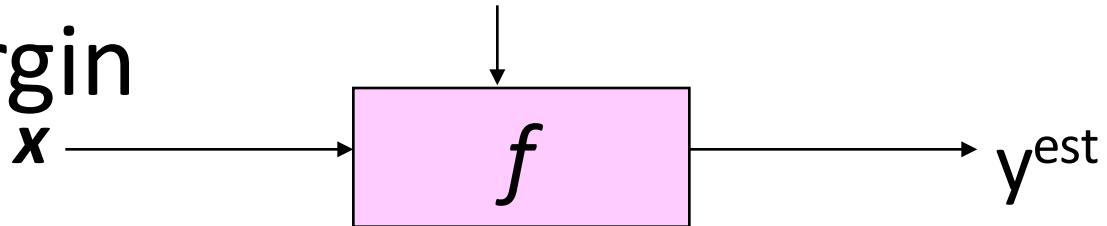


- denotes +1
- denotes -1

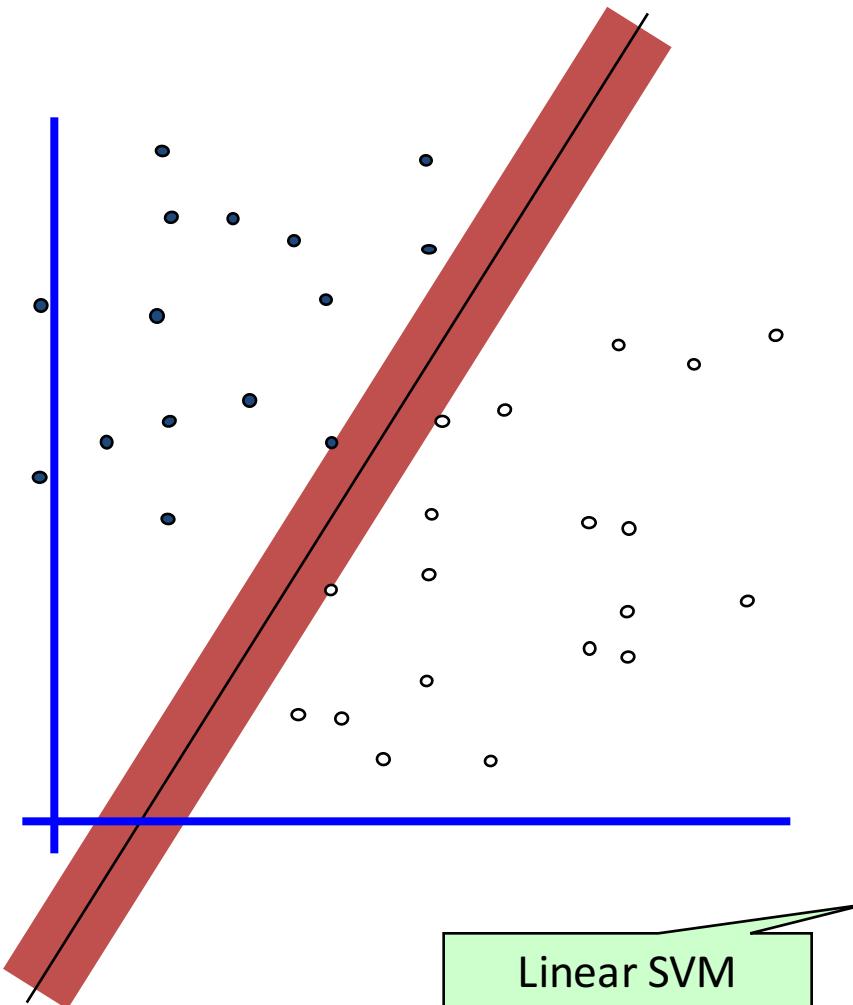


Define the **margin** of a linear classifier as the width that the **boundary could be increased by** before hitting a datapoint.

Maximum Margin

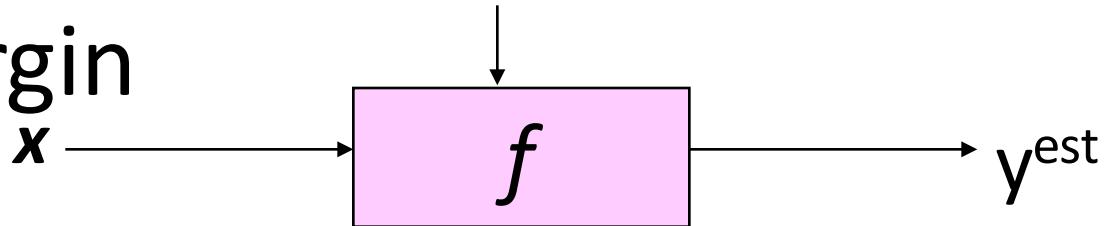


- denotes +1
- denotes -1



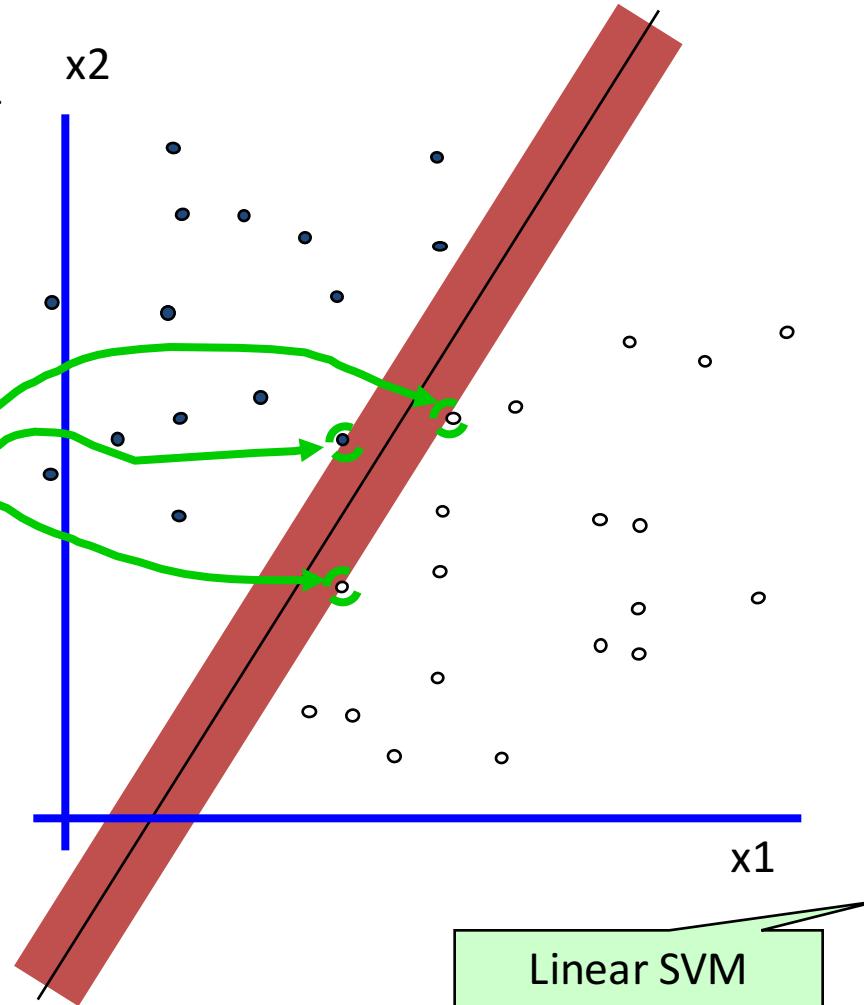
The maximum margin linear classifier is the linear classifier with the maximum margin. This is the simplest kind of SVM (Called an LSVM)

Maximum Margin



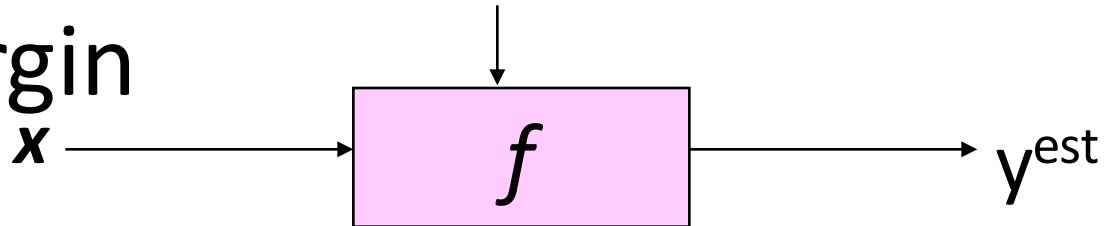
- denotes +1
- denotes -1

Support Vectors are those datapoints that the margin pushes up against



The **maximum margin linear classifier** is the linear classifier with the maximum margin.
This is the simplest kind of SVM (Called an LSVM)

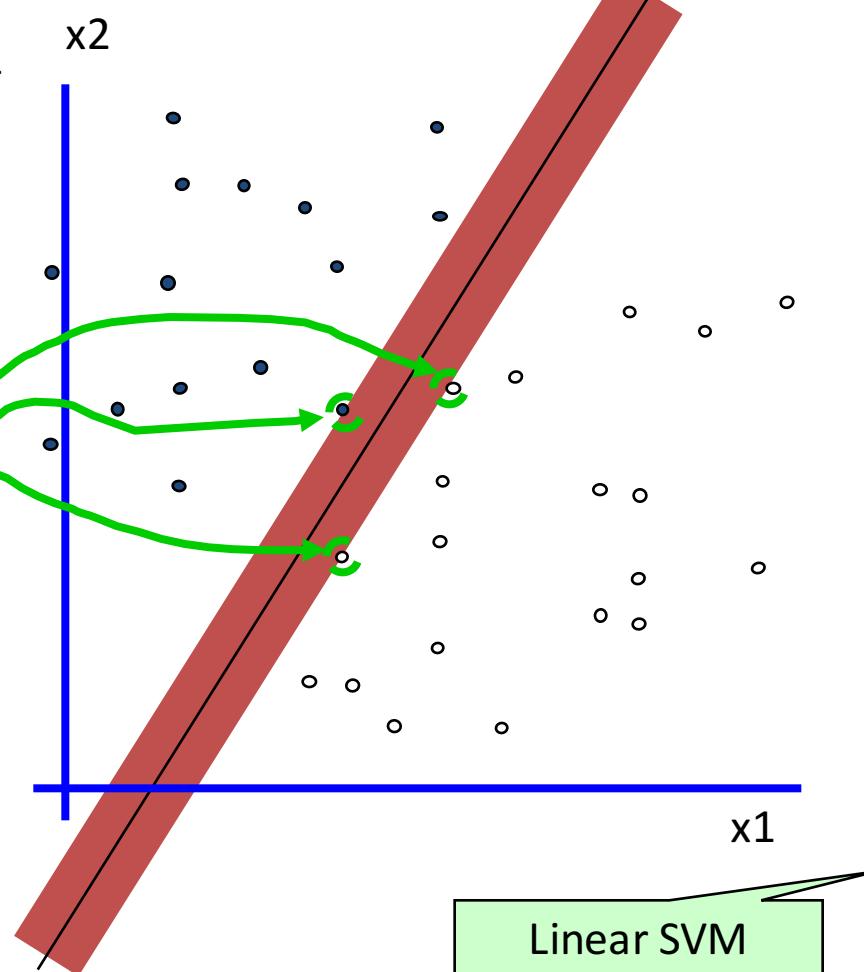
Maximum Margin



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

- denotes +1
- denotes -1

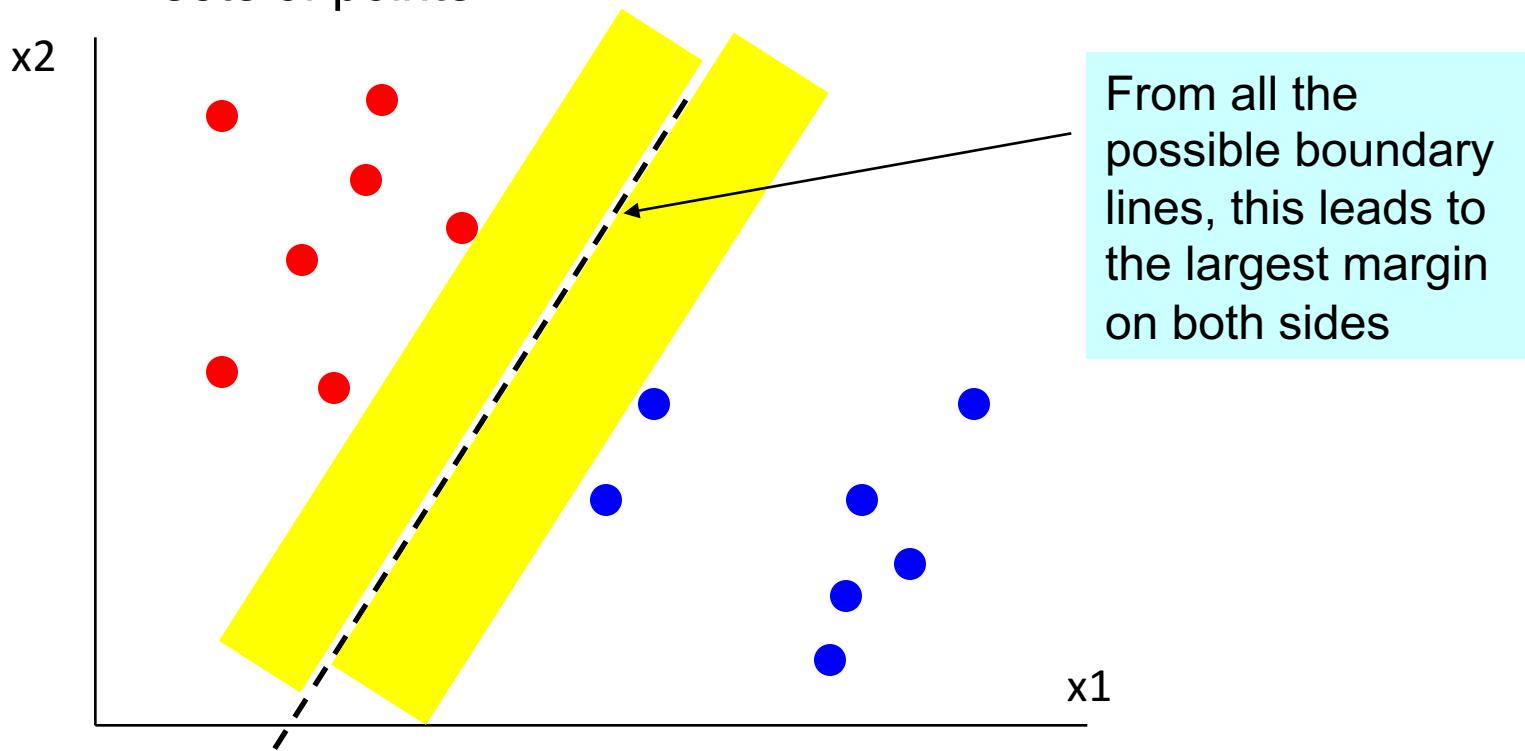
Support Vectors are those datapoints that the margin pushes up against



The maximum margin linear classifier is the linear classifier with the maximum margin. This is the simplest kind of SVM (Called an LSVM)

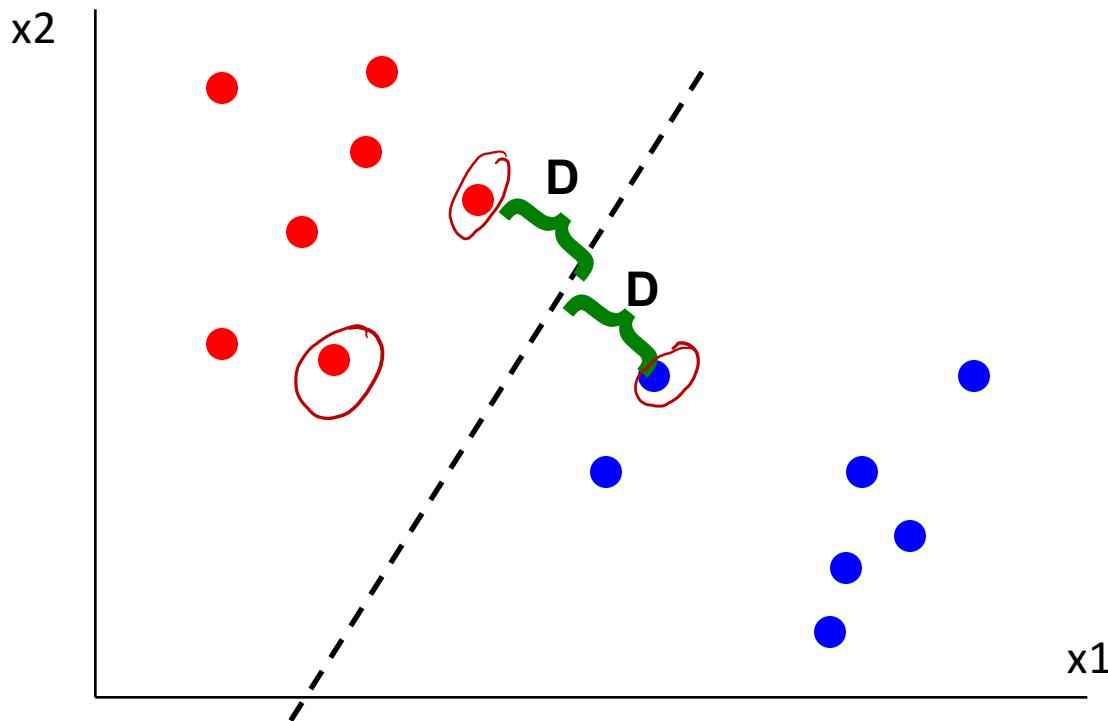
Max margin classifiers

- Instead of fitting all points, focus on boundary points
- Learn a boundary that leads to the largest margin from both sets of points



Max margin classifiers

- Instead of fitting all points, focus on boundary points
- Learn a boundary that leads to the largest margin from points on both sides

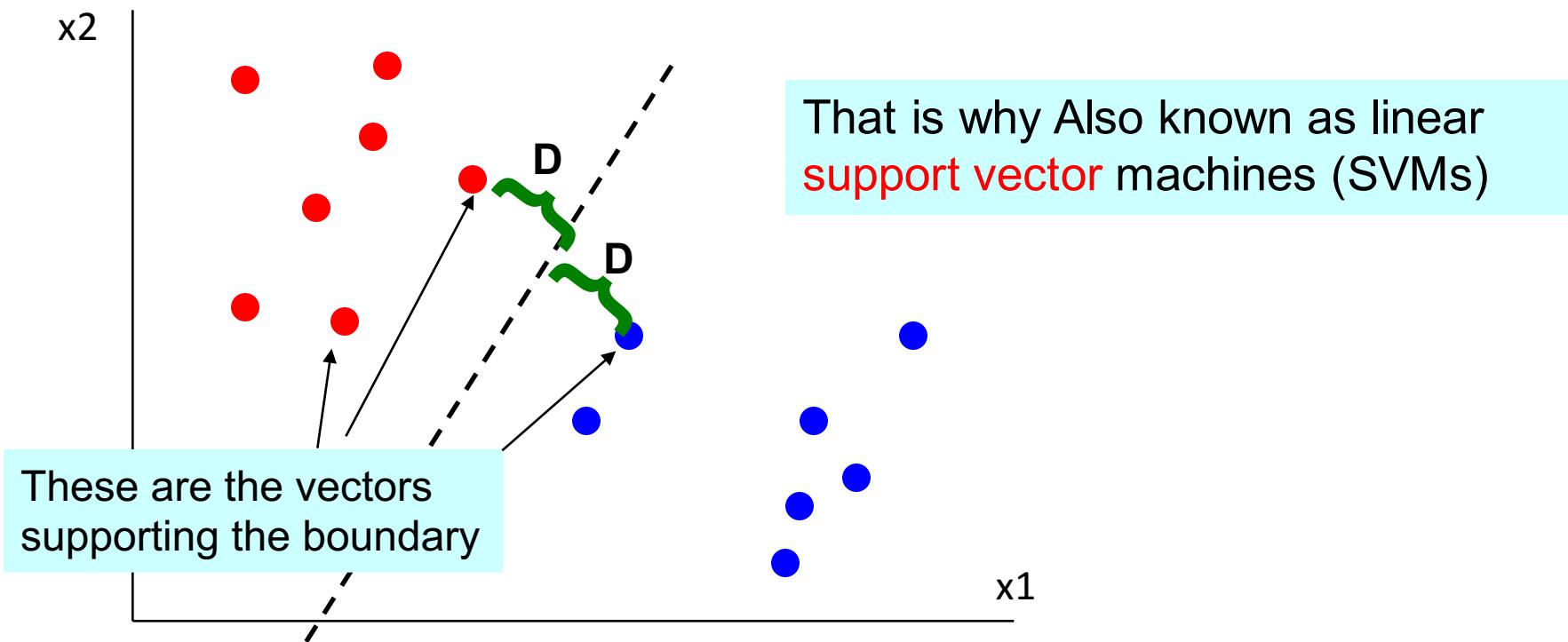


Why MAX margin?

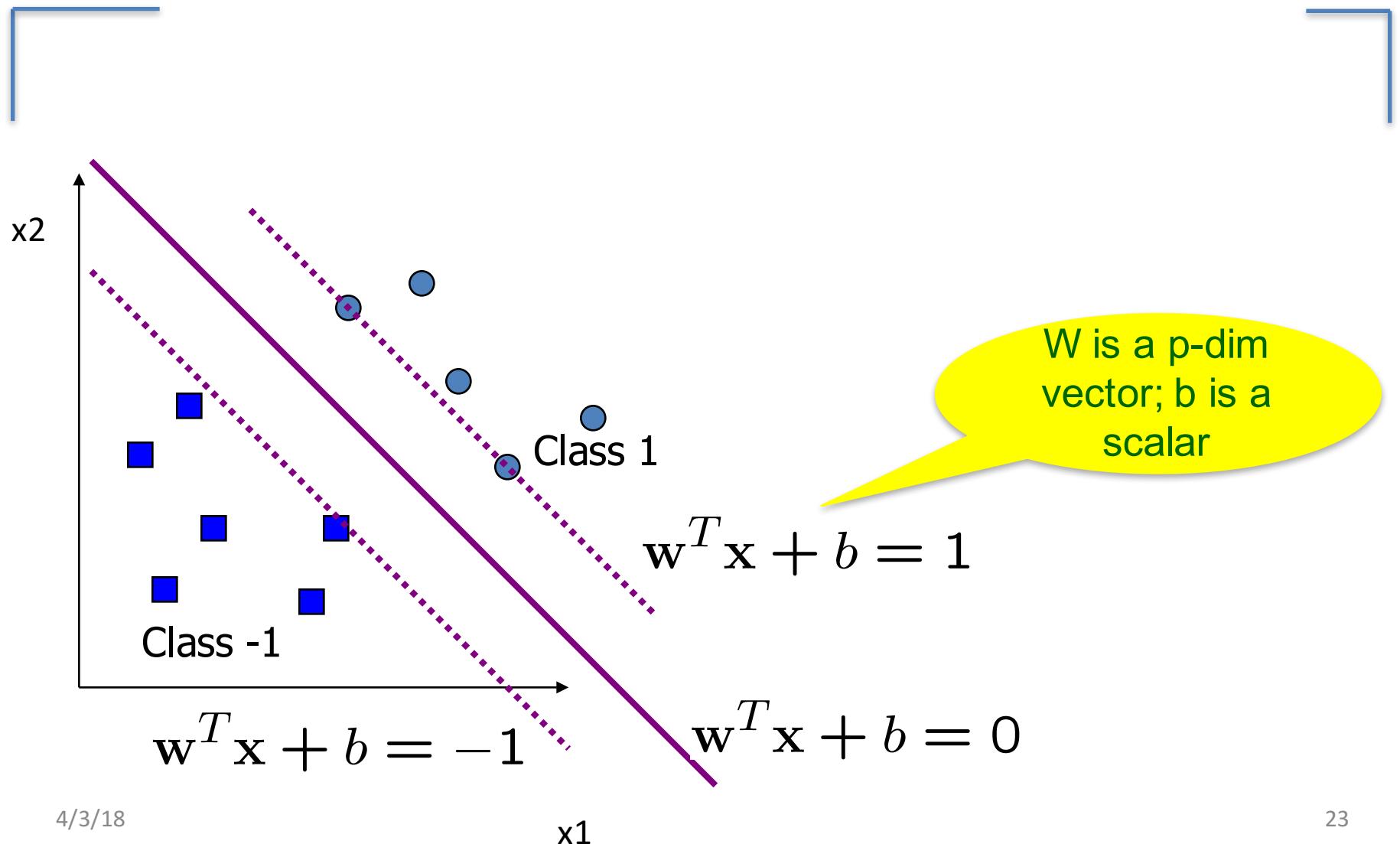
- Intuitive, ‘makes sense’
- Some theoretical support
- Works well in practice

Max margin classifiers

- Instead of fitting all points, focus on boundary points
- Learn a boundary that leads to the largest margin from points on both sides

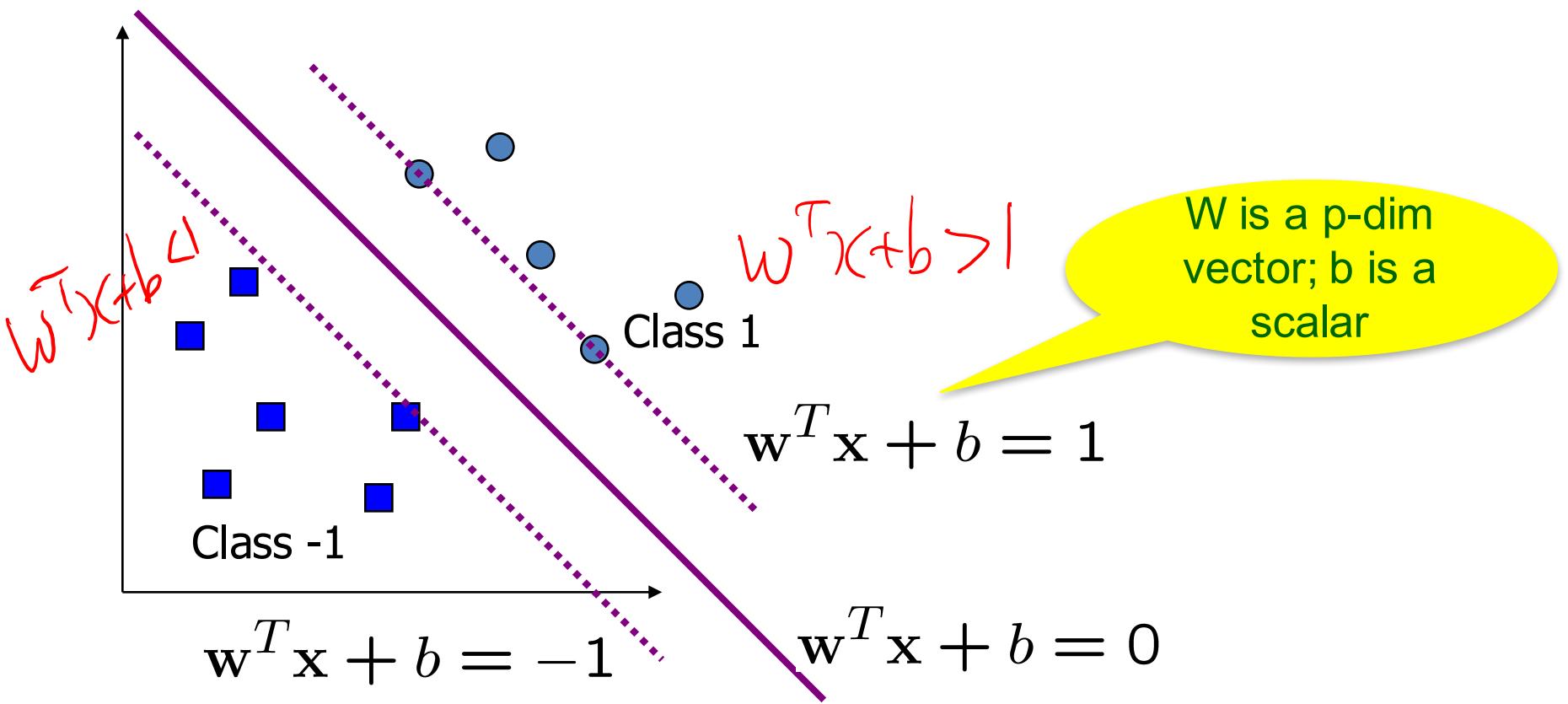


Max-margin & Decision Boundary

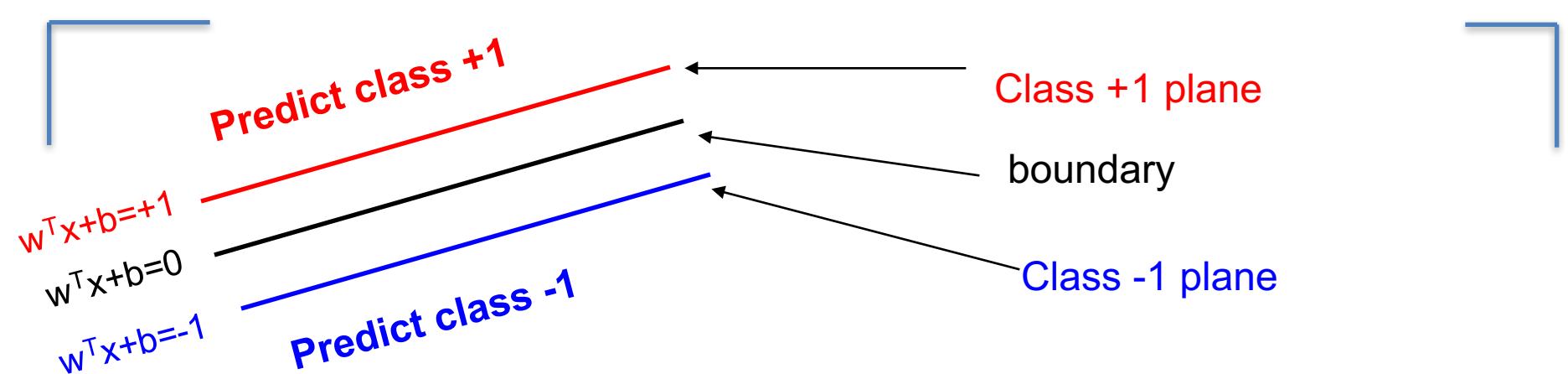


Max-margin & Decision Boundary

- The decision boundary should be as far away from the data of both classes as possible



Specifying a max margin classifier

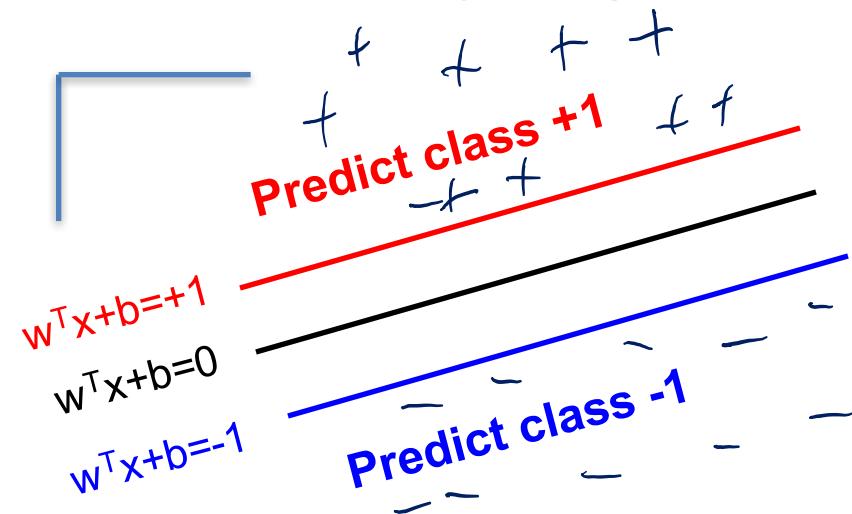


Classify as +1 if $w^T x + b \geq 1$

Classify as -1 if $w^T x + b \leq -1$

Undefined if $-1 < w^T x + b < 1$

Specifying a max margin classifier



Classify as +1 if

Classify as -1 if

Undefined if

Is the linear separation assumption realistic?

We will deal with this shortly, but let's assume it for now

$$w^T x + b \geq 1$$

$$w^T x + b \leq -1$$

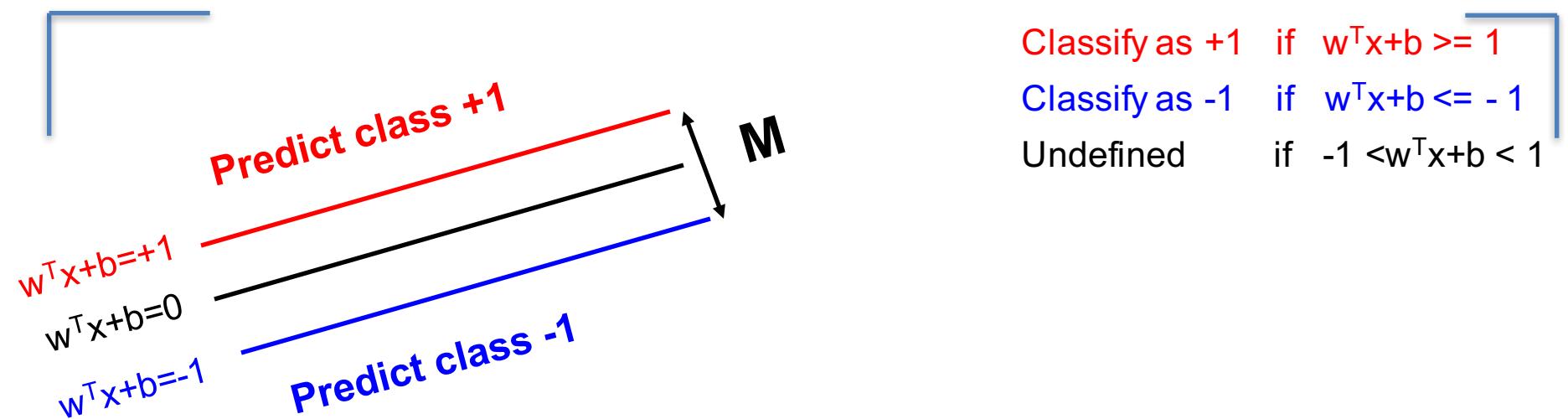
$$-1 < w^T x + b < 1$$

Now assuming such lines exist in our train

Today

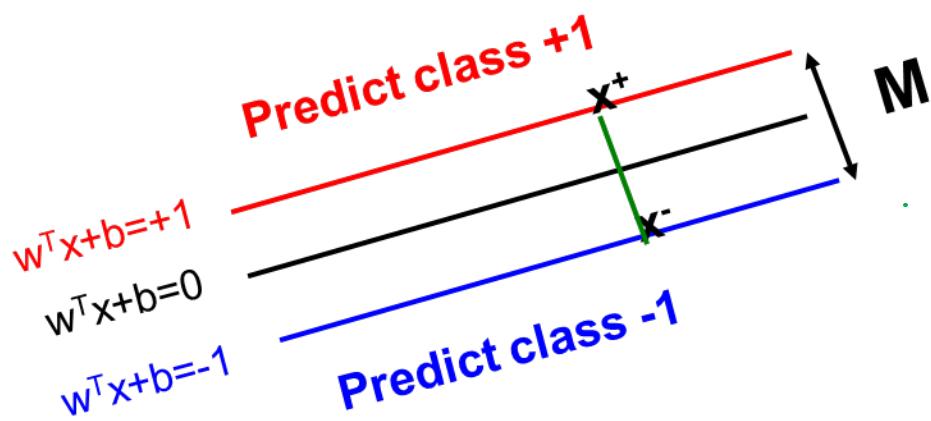
- ❑ Supervised Classification
- ❑ Support Vector Machine (SVM)
 - ✓ History of SVM
 - ✓ Large Margin Linear Classifier
 - ✓ Define Margin (M) in terms of model parameter
 - ✓ Optimization to learn model parameters (w, b)
 - ✓ Linearly Non-separable case
 - ✓ Optimization with dual form
 - ✓ Nonlinear decision boundary
 - ✓ Multiclass SVM

Maximizing the margin



- Lets define the width of the margin by M
- How can we encode our goal of maximizing M in terms of our parameters (w and b)?
- Lets start with a few obsevations

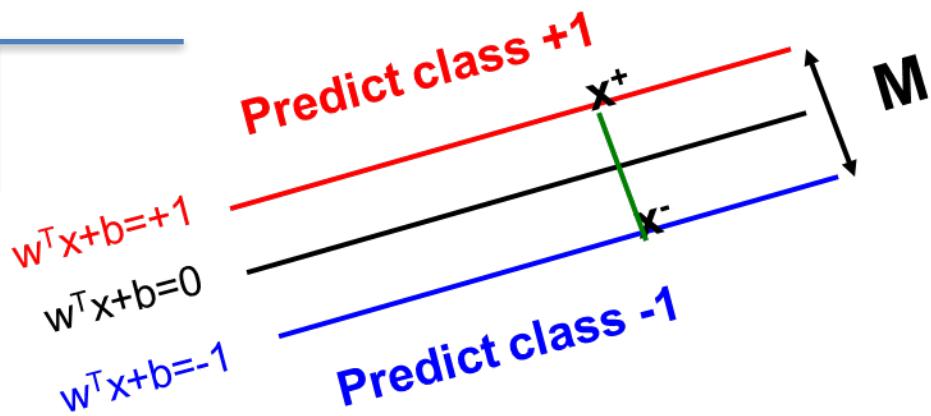
Margin M



Classify as +1	if	$w^T x + b \geq 1$
Classify as -1	if	$w^T x + b \leq -1$
Undefined	if	$-1 < w^T x + b < 1$

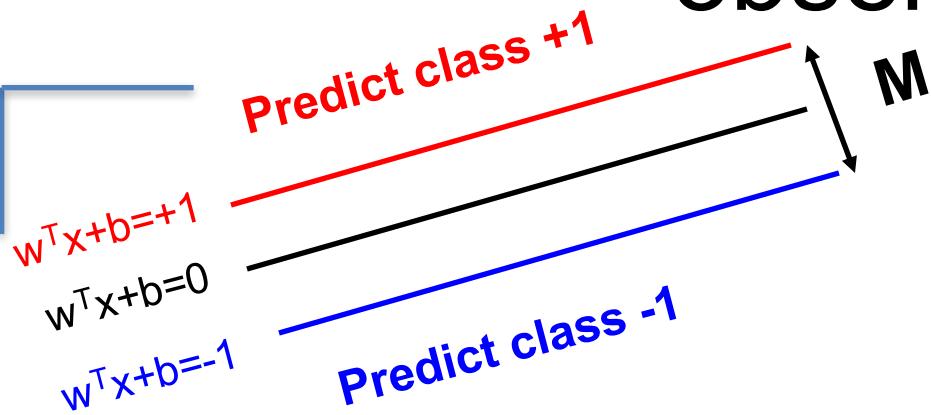
$$M = |x^+ - x^-| \quad \begin{matrix} \text{length of} \\ \text{Vector } (x^+ - x^-) \end{matrix}$$

⇒ How to represent $(x^+ - x^-)$???



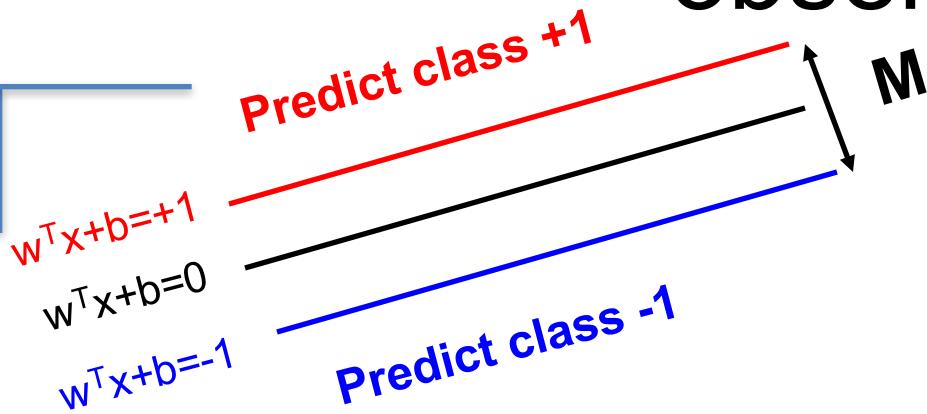
- $w^T x^+ + b = +1$
- $w^T x^- + b = -1$
- $M = |x^+ - x^-| = ?$

Maximizing the margin: observation-1



- Observation 1: the vector w is orthogonal to the +1 plane
- Why?

Maximizing the margin: observation-1

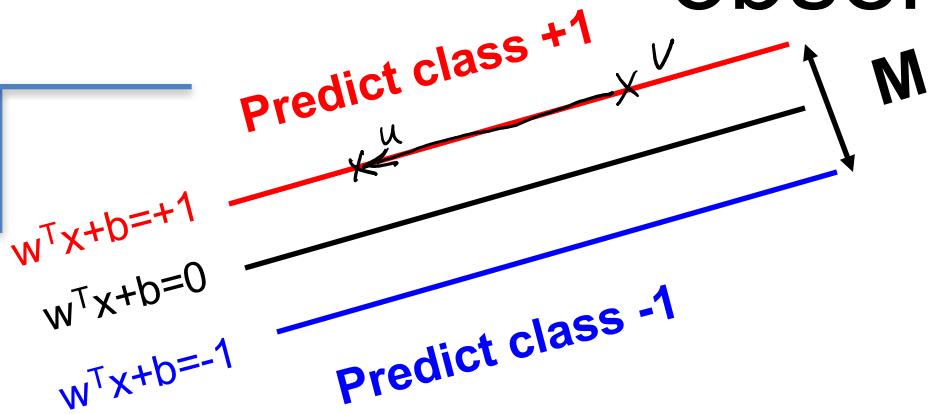


Classify as +1	if	$w^T x + b \geq 1$
Classify as -1	if	$w^T x + b \leq -1$
Undefined	if	$-1 < w^T x + b < 1$

- Observation 1: the vector w is orthogonal to the +1 plane
- Why?

Let u and v be two points on the +1 plane,
then for the vector defined by u and v we have
 $w^T(u-v) = 0$

Maximizing the margin: observation-1



Classify as +1	if	$w^T x + b \geq 1$
Classify as -1	if	$w^T x + b \leq -1$
Undefined	if	$-1 < w^T x + b < 1$

- Observation 1: the vector w is orthogonal to the $+1$ plane

- Why? \rightarrow Vector $(u-v)$ shown above

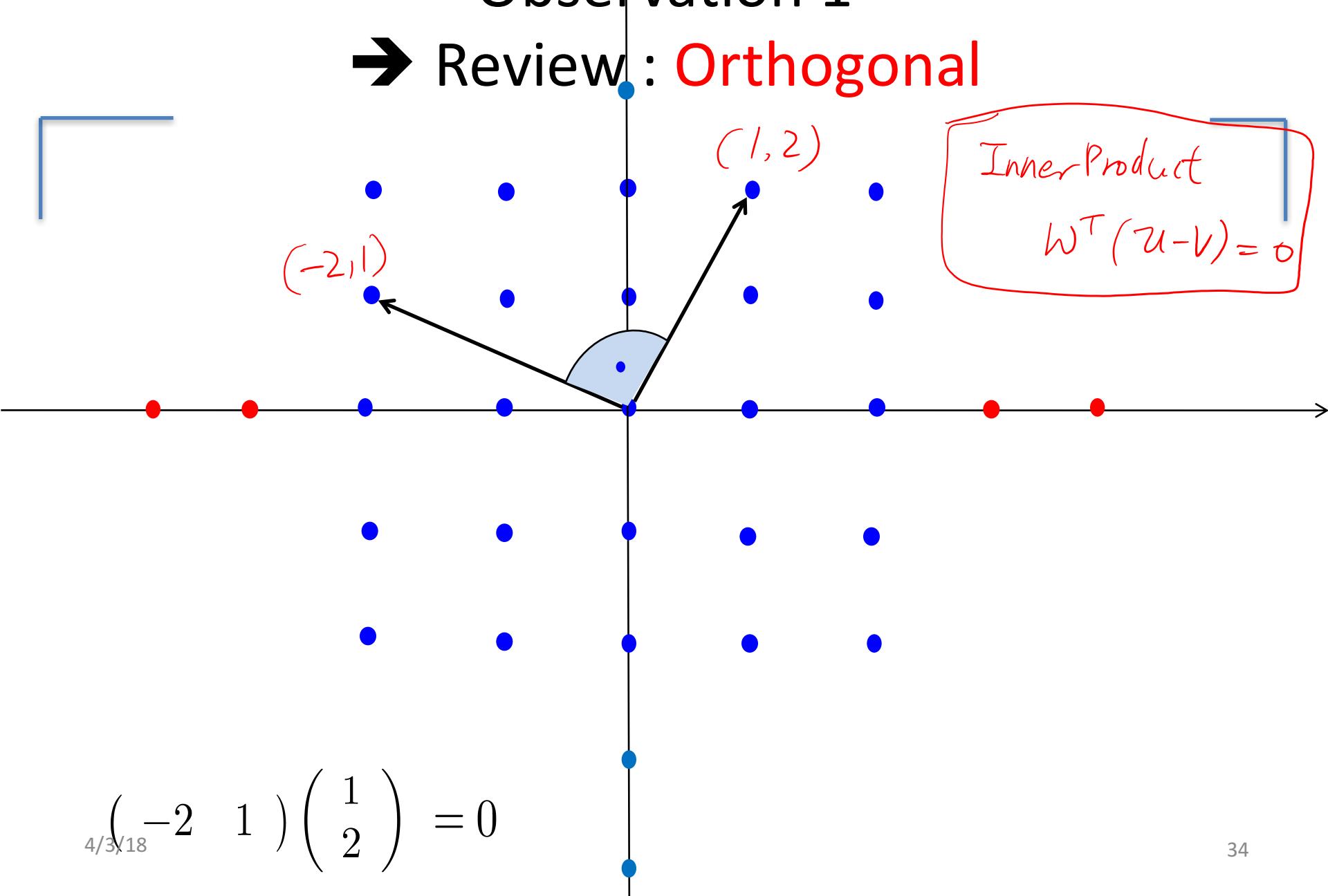
$$\rightarrow w^T(u-v) = w^T u - w^T v = (-b) - (1-b) = 0$$

$\Rightarrow w$ orthogonal
to $(u-v)$

Let u and v be two points on the $+1$ plane,
then for the vector defined by u and v we have
 $w^T(u-v) = 0$

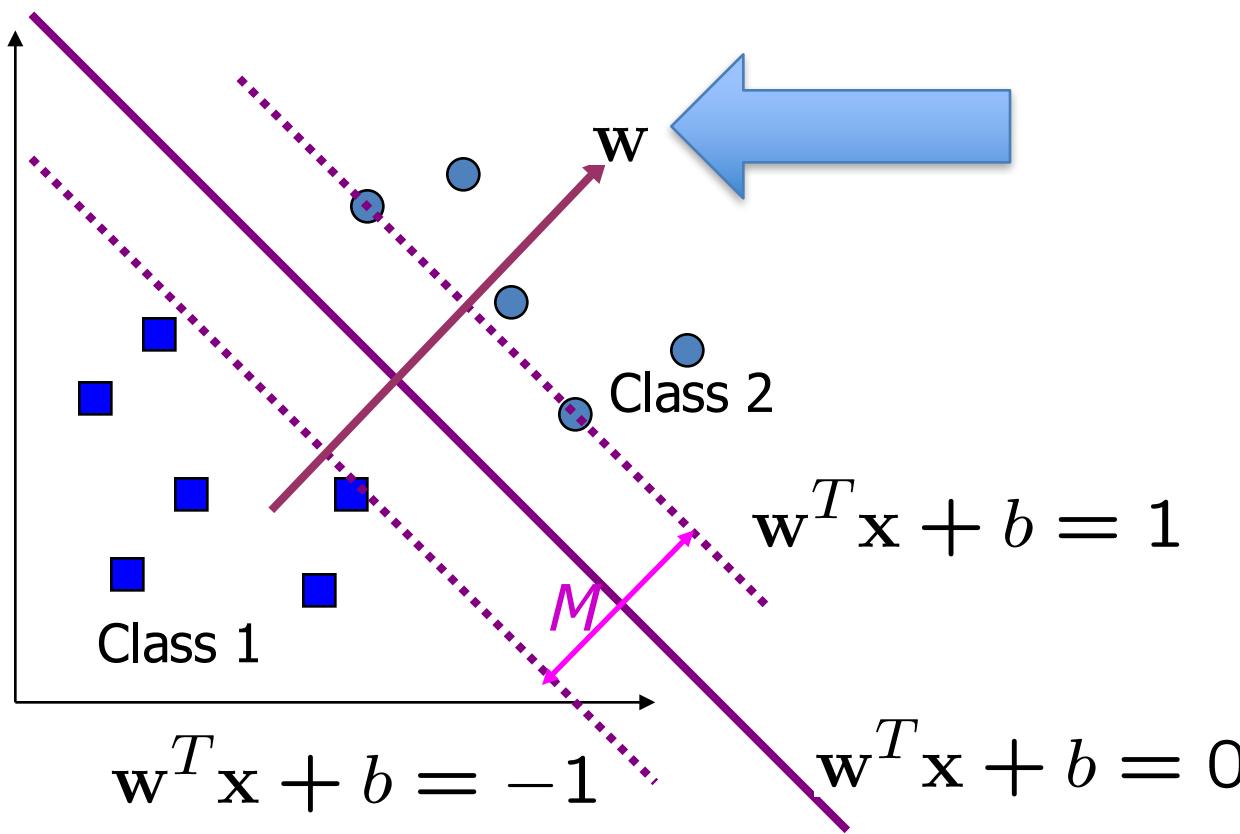
Observation 1

→ Review: Orthogonal



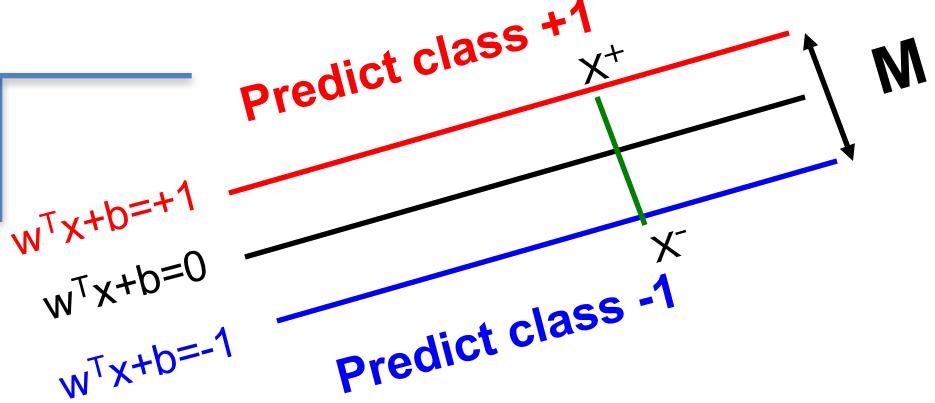
Maximizing the margin: observation-1

- Observation 1: the vector w is orthogonal to the +1 plane



+1 plane
-1 plane
0-plane

Maximizing the margin: observation-2



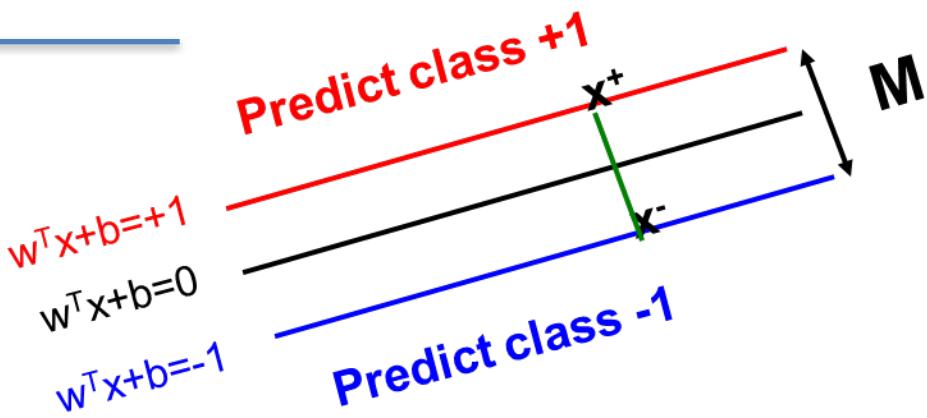
Classify as +1	if	$w^T x + b \geq 1$
Classify as -1	if	$w^T x + b \leq -1$
Undefined	if	$-1 < w^T x + b < 1$

- Observation 1: the vector w is orthogonal to the +1 and -1 planes
- Observation 2: if x^+ is a point on the +1 plane and x^- is the closest point to x^+ on the -1 plane then

$$x^+ = \lambda w + x^-$$

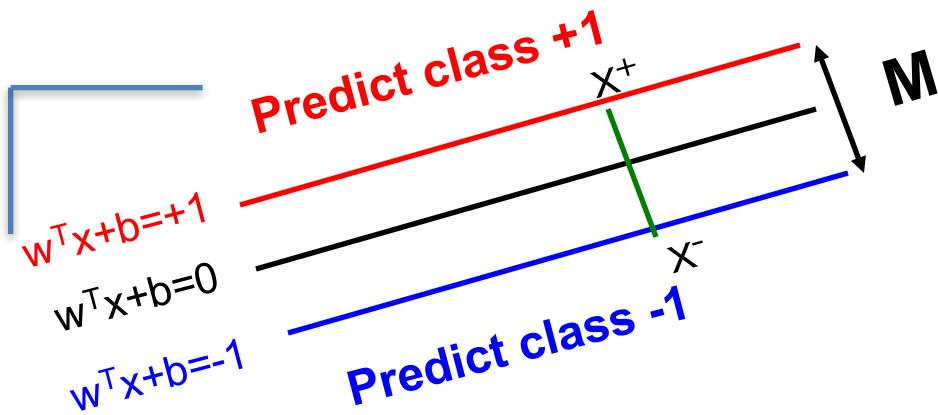
Since w is orthogonal to both planes we need to ‘travel’ some distance along w to get from x^+ to x^-

Putting it together



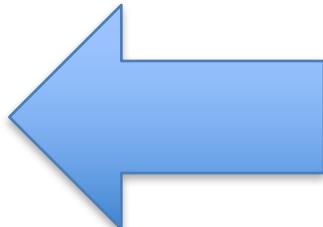
- $w^T x^+ + b = +1$
- $w^T x^- + b = -1$
- $M = |x^+ - x^-| = ?$
- $x^+ = \lambda w + x^-$

Putting it together



- $w^T x^+ + b = +1$
- $w^T x^- + b = -1$
- $x^+ = \lambda w + x^-$
- $|x^+ - x^-| = M$

We can now define M in terms of w and b



$$M = |x^+ - x^-|$$

$$= |\lambda w|$$

$$= \lambda |w|$$

$$= \lambda \sqrt{w^T w}$$

$$= \frac{2}{\sqrt{w^T w}} \sqrt{w^T w}$$

$$= \frac{2}{\sqrt{w^T w}}$$

$$w^T x^+ + b = 1$$

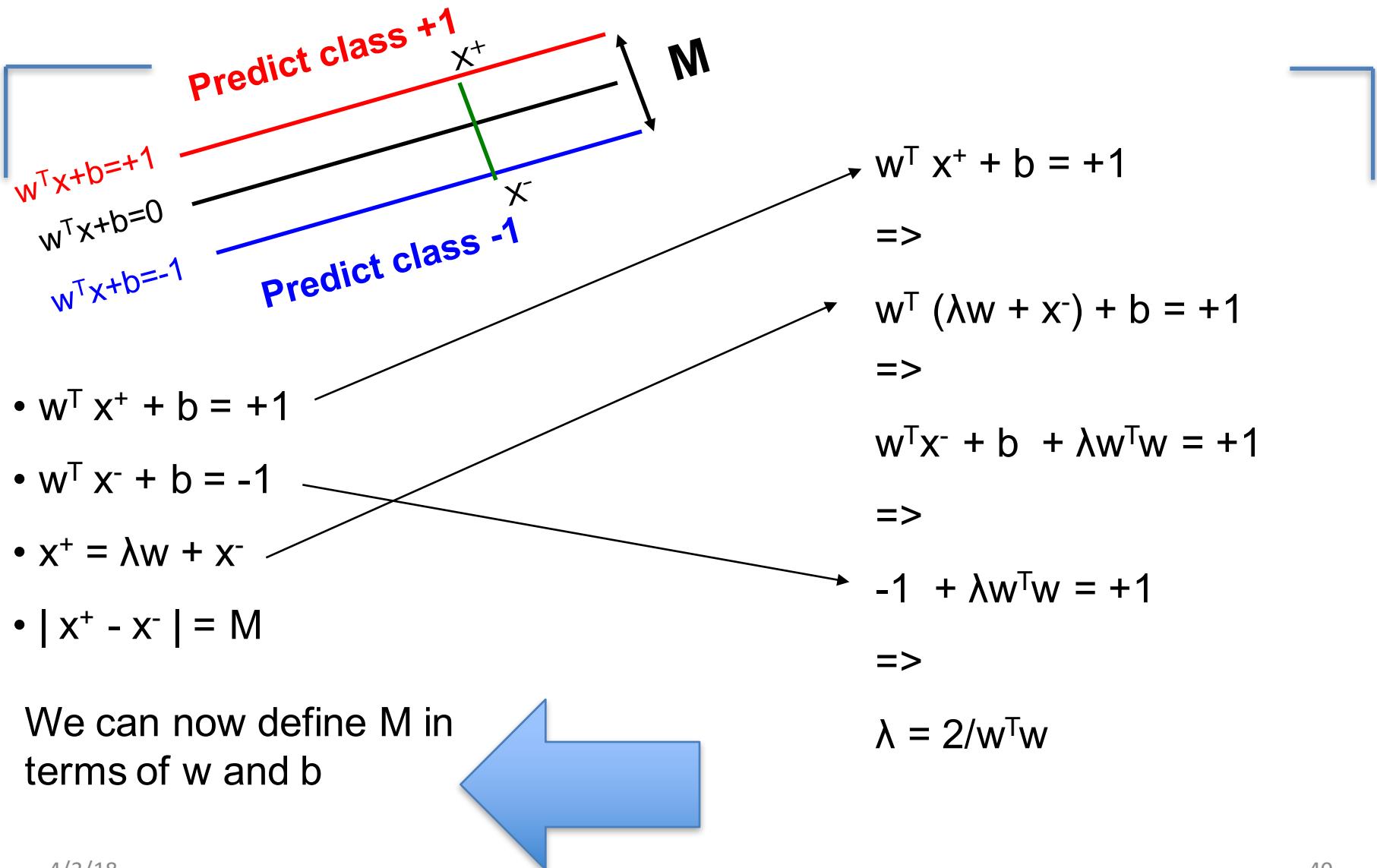
$$w^T(\lambda w + x^-) + b = +1$$

$$\lambda w^T w + \underbrace{w^T x^-}_{\Rightarrow -1} + b = 1$$

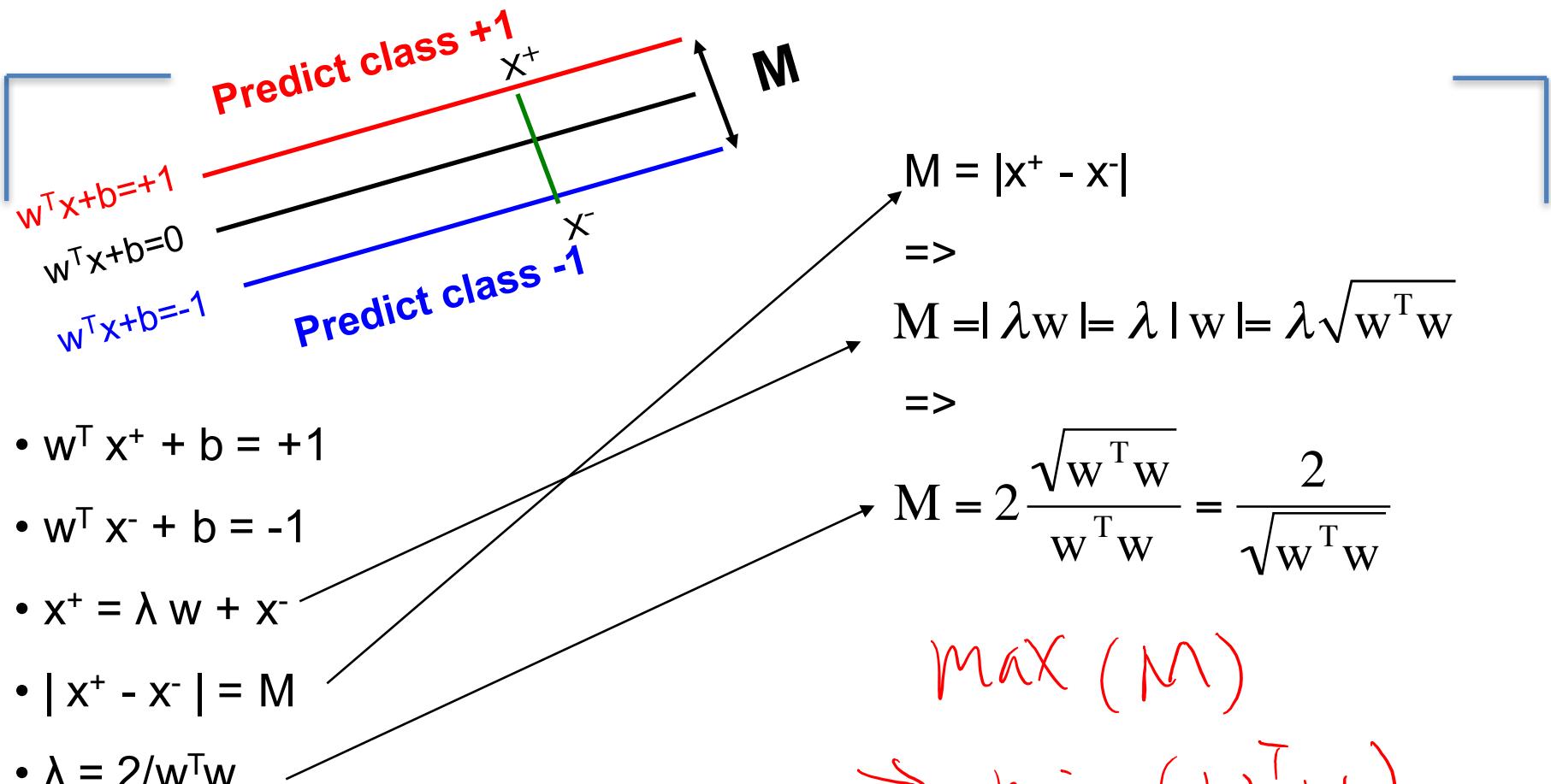
$$\lambda w^T w = 2$$

$$\Rightarrow \lambda = \frac{2}{w^T w}$$

Putting it together

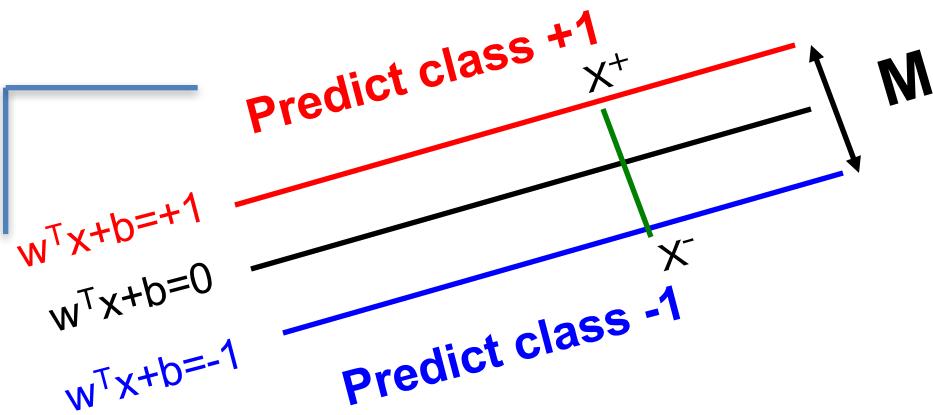


Putting it together

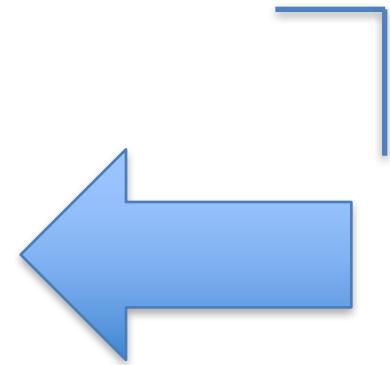


We can now define M in terms of w and b

Finding the optimal parameters



$$\begin{aligned} M &= \frac{2}{\sqrt{w^T w}} \\ &= \frac{2}{\|w\|} \end{aligned}$$



We can now search for the optimal parameters by finding a solution that:

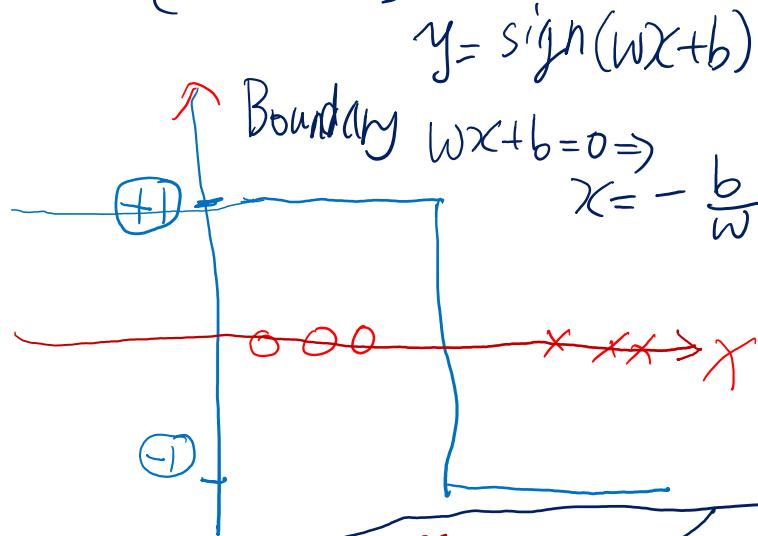
1. Correctly classifies all points
2. Maximizes the margin (or equivalently minimizes $w^T w$)

Several optimization methods can be used:
Gradient descent, OR SMO (see extra slides)

Binary Classification

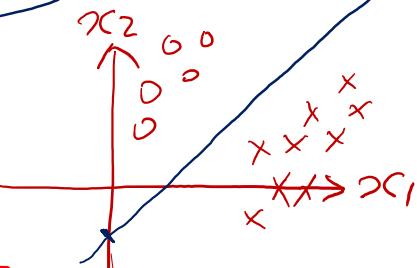
$$y \in \{-1, 1\}$$

\Rightarrow 1D $[x \in \mathbb{R}]$



\Rightarrow 2D $[x \in \mathbb{R}^2]$

$$\text{Boundary } w^T x + b = 0$$



\Rightarrow P dim $[x \in \mathbb{R}^P]$

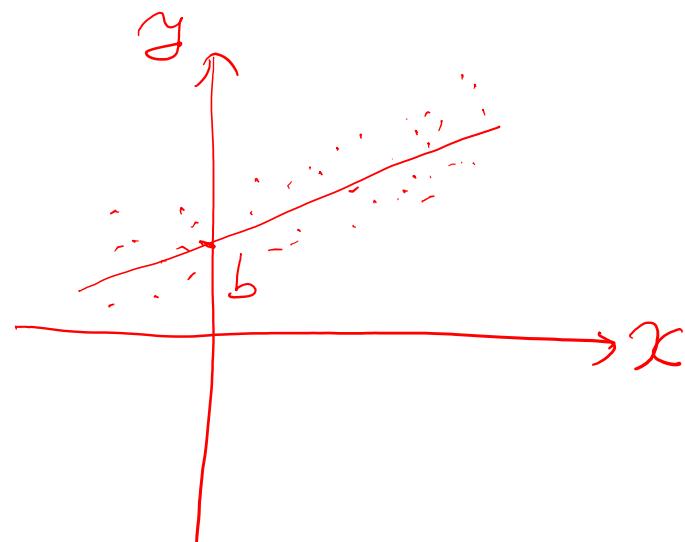
Boundary: hyperplane

4/3/18

Regression

$$y \in \mathbb{R}$$

\Rightarrow 1D $x \in \mathbb{R}$

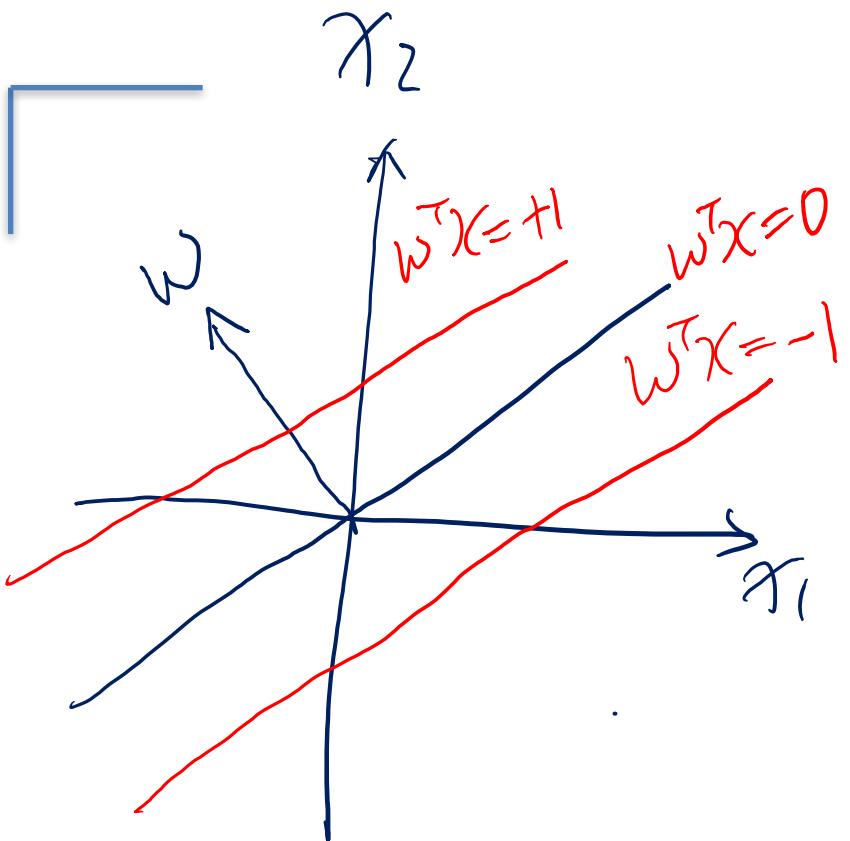


$$y = wx + b$$

$$y = w^T x + b \quad \text{if } x \in \mathbb{R}^P$$

43

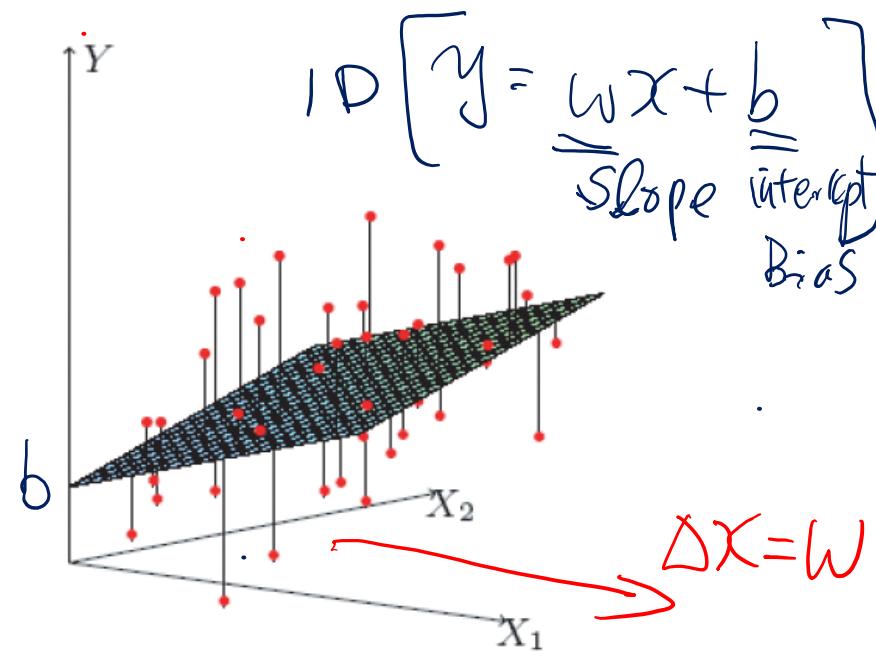
[Classification]



[Regression]

$$y = w^T x + b$$

$$\Delta x = \frac{\partial y}{\partial x} = w$$



The gradient points in the **direction** of the greatest rate of increase of the function and its **magnitude** is the slope of the graph in that direction

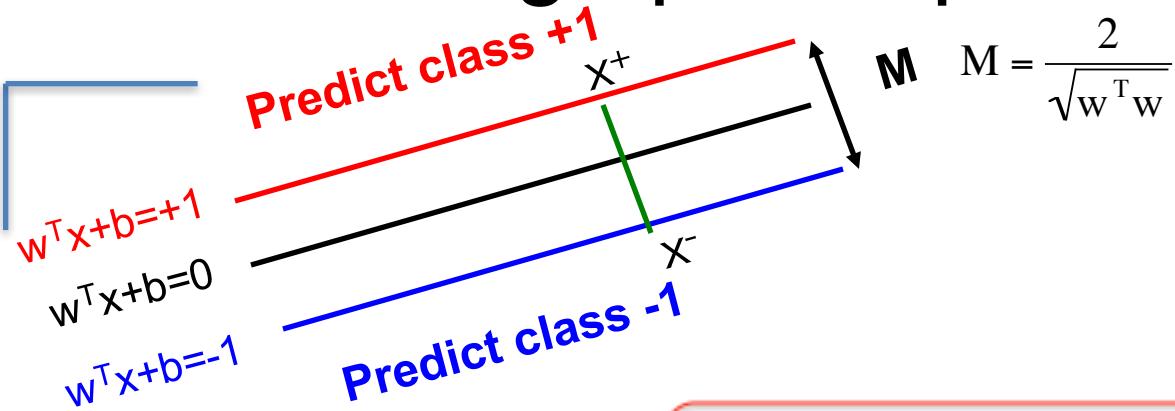
Today

❑ Support Vector Machine (SVM)

- ✓ History of SVM
- ✓ Large Margin Linear Classifier
- ✓ Define Margin (M) in terms of model parameter
-  ✓ Optimization to learn model parameters (w, b)
- ✓ Linearly Non-separable case
- ✓ Optimization with dual form
- ✓ Nonlinear decision boundary
- ✓ Practical Guide

Optimization Step

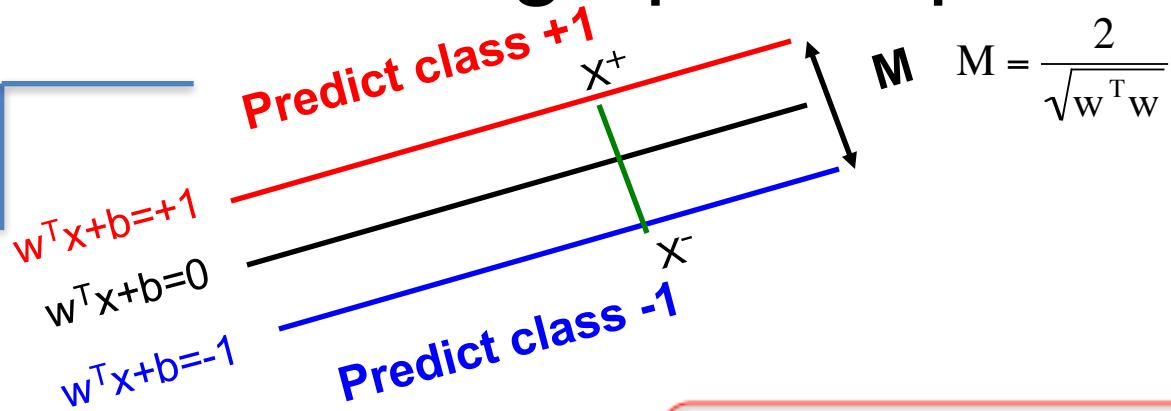
i.e. learning optimal parameter for SVM



1. Correctly classifies all points
2. Maximizes the margin (or equivalently minimizes $w^T w$)

Optimization Step

i.e. learning optimal parameter for SVM



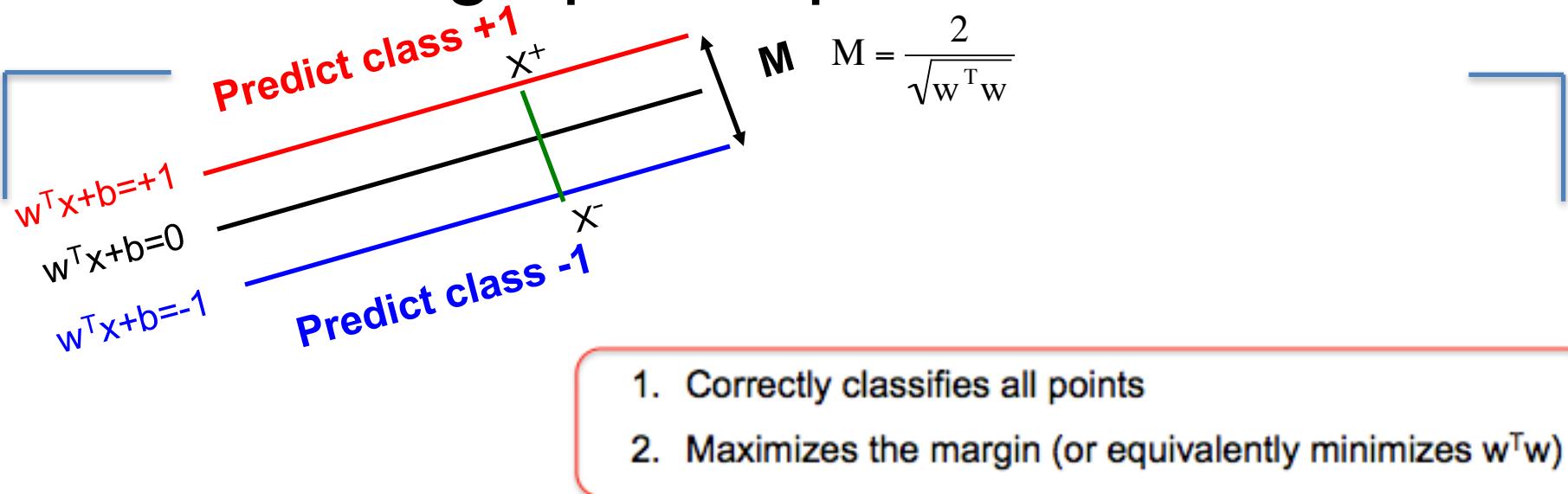
1. Correctly classifies all points
2. Maximizes the margin (or equivalently minimizes $w^T w$)

$$\text{Min } (w^T w)/2$$

subject to the following constraints:

Optimization Step

i.e. learning optimal parameter for SVM



$$\text{Min } (w^T w)/2$$

subject to the following constraints:

For all x in class + 1

$$w^T x + b \geq 1$$

For all x in class - 1

$$w^T x + b \leq -1$$

}

A total of n constraints if we have n training samples

Optimization Reformulation

1. Correctly classifies all points
2. Maximizes the margin (or equivalently minimizes $w^T w$)

$$\text{Min } (w^T w)/2$$

subject to the following constraints:

For all x in class +1

$$w^T x + b \geq 1 \quad y_i = 1$$

For all x in class -1

$$w^T x + b \leq -1 \quad y_i = -1$$

A total of n
constraints if
we have n
input
samples

$$\rightarrow \text{Pos } y_i = 1, w^T x_i + b \geq 1$$

$$y_i (w^T x_i + b) \geq 1$$

$$\rightarrow \text{Neg } y_i = -1, w^T x_i + b \leq -1$$

$$y_i (w^T x_i + b) \geq 1$$

Optimization Reformulation

- 1. Correctly classifies all points
- 2. Maximizes the margin (or equivalently minimizes $\mathbf{w}^T \mathbf{w}$)

$$\text{Min } (\mathbf{w}^T \mathbf{w})/2$$

subject to the following constraints:

For all x in class + 1

$$\mathbf{w}^T \mathbf{x} + b \geq 1$$

For all x in class - 1

$$\mathbf{w}^T \mathbf{x} + b \leq -1$$



A total of n
constraints if
we have n
input samples



$$\underset{\mathbf{w}, b}{\operatorname{argmin}} \sum_{i=1}^p w_i^2$$

$$\text{subject to } \forall \mathbf{x}_i \in D_{\text{train}} : y_i(\mathbf{x}_i \cdot \mathbf{w} + b) \geq 1$$

Optimization Reformulation

- 1. Correctly classifies all points
- 2. Maximizes the margin (or equivalently minimizes $\mathbf{w}^T \mathbf{w}$)

$$\text{Min } (\mathbf{w}^T \mathbf{w})/2$$

subject to the following constraints:

For all x in class + 1

$$\mathbf{w}^T \mathbf{x} + b \geq 1$$

For all x in class - 1

$$\mathbf{w}^T \mathbf{x} + b \leq -1$$

A total of n
constraints if
we have n
inputsamples



Quadratic Objective

$$\underset{\mathbf{w}, b}{\operatorname{argmin}} \sum_{i=1}^p w_i^2$$

$$\text{subject to } \forall \mathbf{x}_i \in D_{\text{train}} : y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

Quadratic programming
i.e.,

- Quadratic objective
- Linear constraints

Today

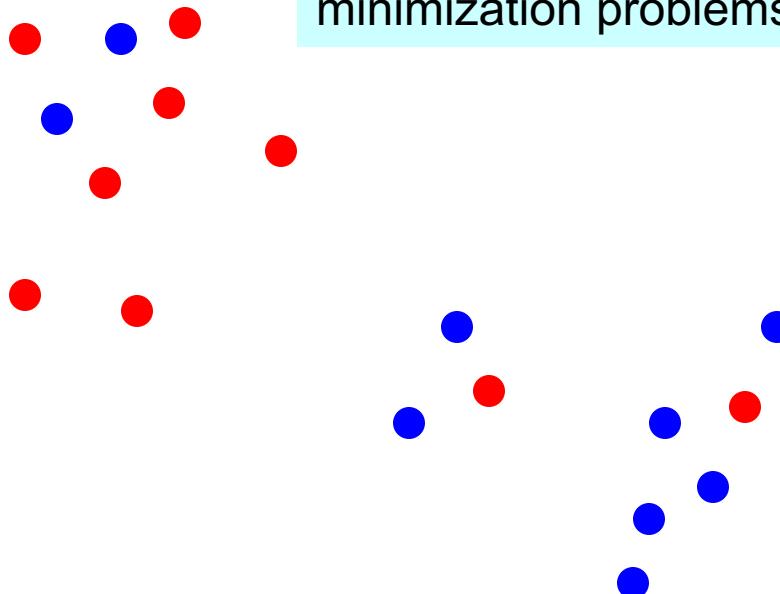
❑ Support Vector Machine (SVM)

- ✓ History of SVM
- ✓ Large Margin Linear Classifier
- ✓ Define Margin (M) in terms of model parameter
- ✓ Optimization to learn model parameters (w , b)
- ✓ Linearly Non-separable case (soft SVM)
- ✓ Optimization with dual form
- ✓ Nonlinear decision boundary
- ✓ Practical Guide

Linearly Non separable case

- So far we assumed that a linear hyperplane can perfectly separate the points
- But this is not usually the case
 - noise, outliers

How can we convert this to a QP problem?

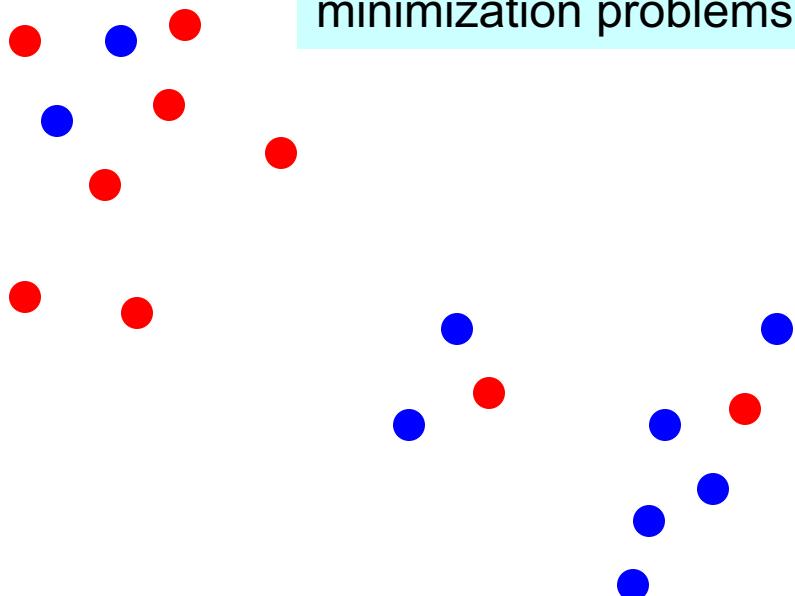


- Minimize training errors?

$$\left. \begin{array}{l} \min w^T w \\ \min \# \text{errors} \end{array} \right\}$$

Linearly Non separable case

- So far we assumed that a linear plane can perfectly separate the points
- But this is not usually the case
 - noise, outliers



How can we convert this to a QP problem?

- Minimize training errors?

$$\min w^T w$$

$$\min \# \text{errors}$$

- Penalize training errors:

$$\min w^T w + C * (\# \text{errors})$$

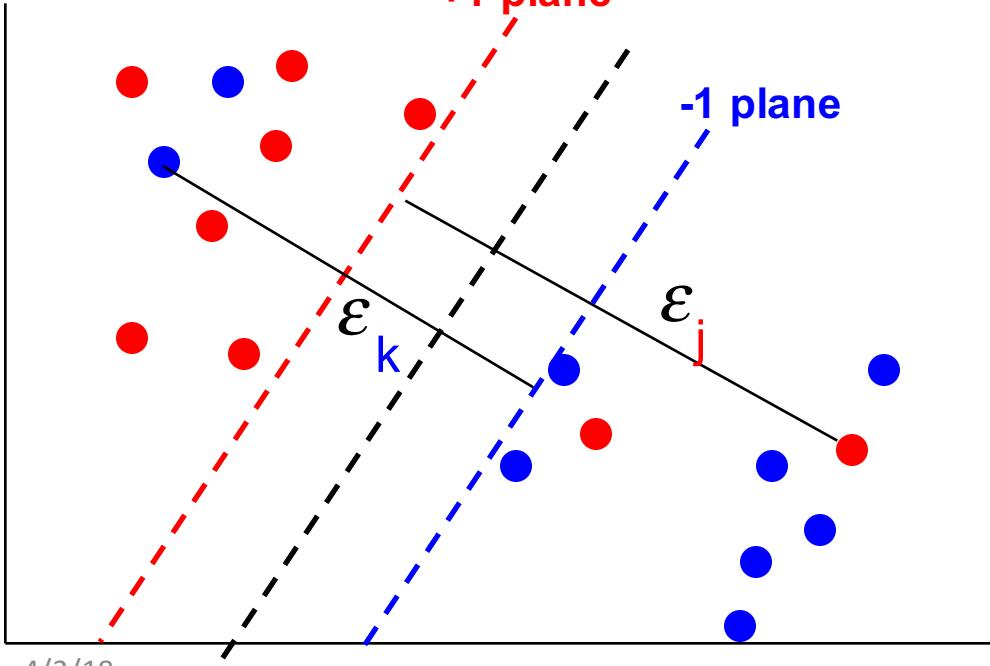
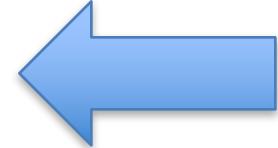
Hard to encode in a QP problem

Linearly Non separable case

- Instead of minimizing the number of misclassified points we can minimize the **distance** between these points and their correct plane

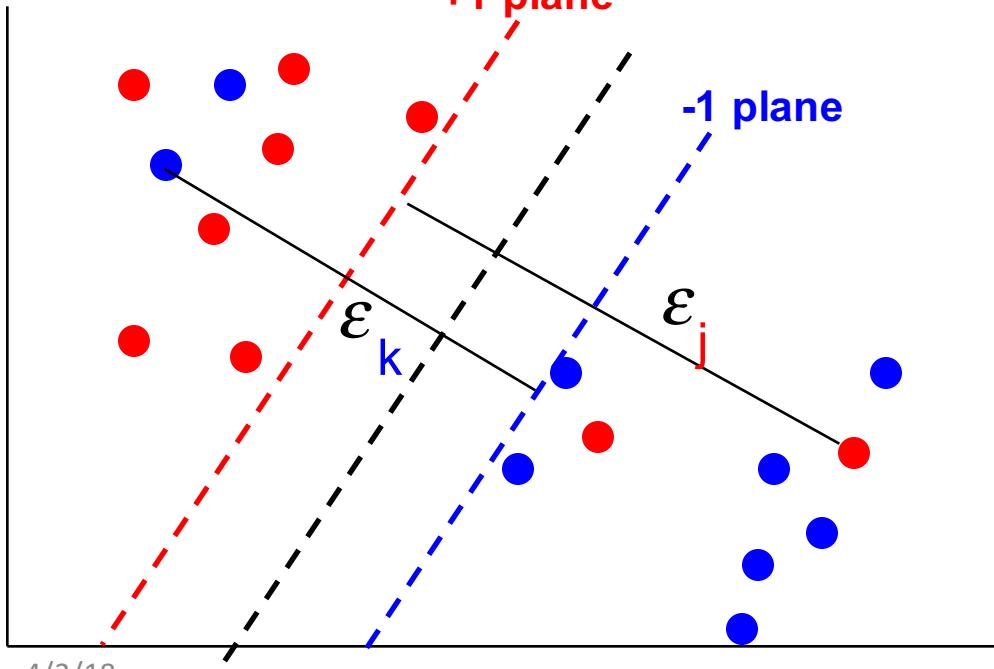
The new optimization problem is:

$$\min_w \frac{w^T w}{2} + \sum_{i=1}^n C \varepsilon_i$$



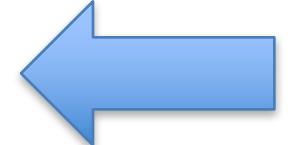
Linearly Non separable case

- Instead of minimizing the number of misclassified points we can minimize the ***distance*** between these points and their correct plane



The new optimization problem is:

$$\min_w \frac{w^T w}{2} + C \sum_{i=1}^n \varepsilon_i$$



subject to the following inequality constraints:

For all x_i in class +1

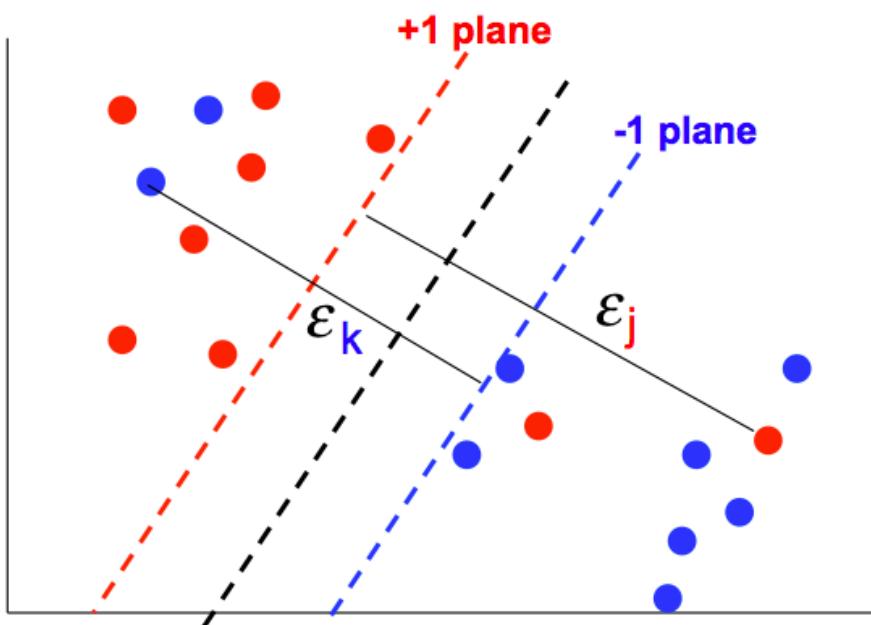
$$w^T x_i + b \geq 1 - \varepsilon_i$$

For all x_i in class -1

$$w^T x_i + b \leq -1 + \varepsilon_i$$

Wait. Are we missing something?

Final optimization for linearly non-separable case



The new optimization problem is:

$$\min_w \frac{w^T w}{2} + C \sum_{i=1}^n \varepsilon_i$$

subject to the following inequality constraints:

For all x_i in class +1

$$w^T x_i + b \geq 1 - \varepsilon_i$$

For all x_i in class -1

$$w^T x_i + b \leq -1 + \varepsilon_i$$

For all i

$$\varepsilon_i \geq 0$$

A total of n constraints

Another n constraints

Two optimization problems:

For the separable and non separable cases

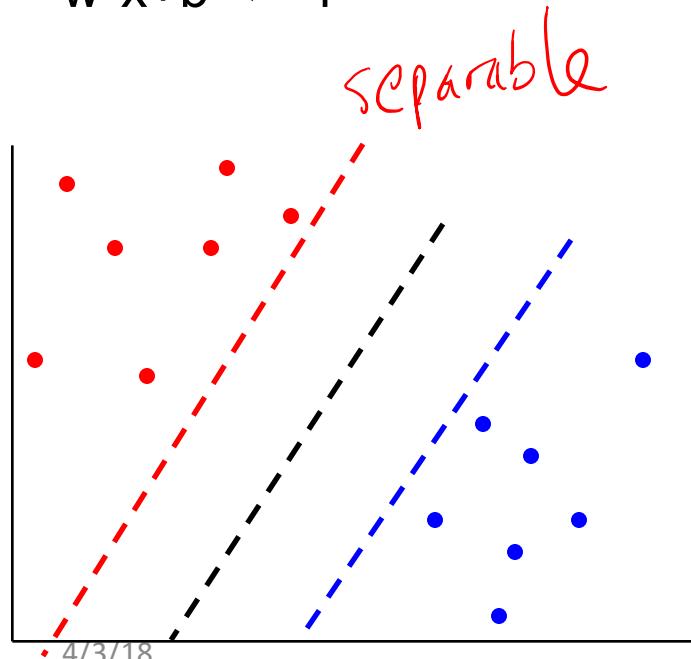
$$\min_w \frac{w^T w}{2}$$

For all x in class + 1

$$w^T x + b \geq 1$$

For all x in class - 1

$$w^T x + b \leq -1$$



$$\min_w \frac{w^T w}{2} + C \sum_{i=1}^n \varepsilon_i$$

For all x_i in class + 1

$$w^T x_i + b \geq 1 - \varepsilon_i$$

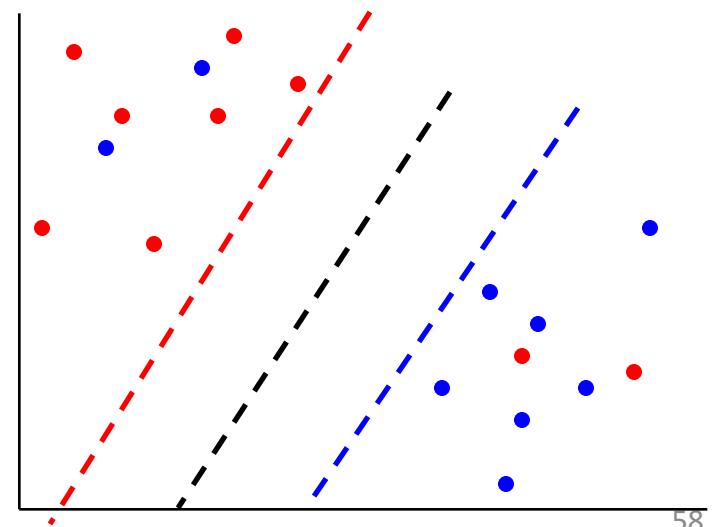
For all x_i in class - 1

$$w^T x_i + b \leq -1 + \varepsilon_i$$

For all i

$$\varepsilon_i \geq 0$$

Non Separable



Hinge Loss for Soft SVM

$$\min_w \frac{w^T w}{2} + C \sum_{i=1}^n \varepsilon_i$$

For all x_i in class + 1

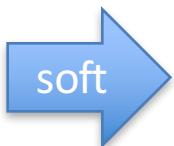
$$w^T x_i + b \geq 1 - \varepsilon_i$$

For all x_i in class - 1

$$w^T x_i + b \leq -1 + \varepsilon_i$$

For all i

$$\varepsilon_i \geq 0$$



$$\operatorname{argmin}_{w,b} \frac{w^T w}{2} + C \sum_{i=1}^n \max(0, 1 - y_i(w^T x_i + b))$$

subject to:

$$y_i(w^T x_i + b) \geq 1 - \varepsilon_i$$

$$\varepsilon_i \geq 0$$

vs. Hard
SVM

$$\operatorname{argmin}_{w,b} \sum_{i=1}^p w_i^2$$

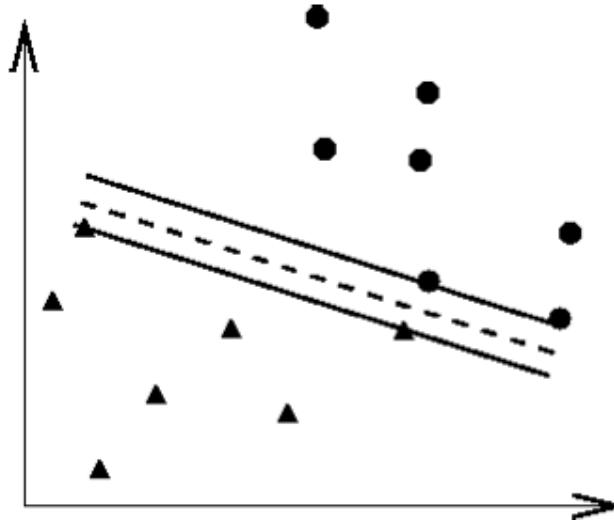
$$\text{subject to } \forall x_i \in D_{train}: y_i(w^T x_i + b) \geq 1$$

Model Selection, find right C

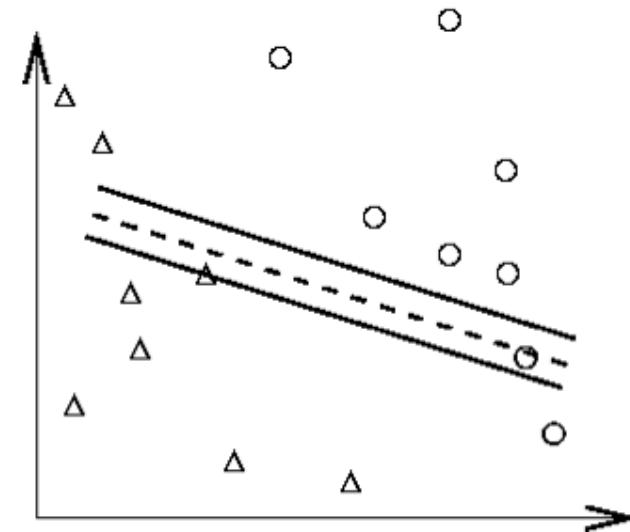
large C

Select the
right
penalty
parameter
C

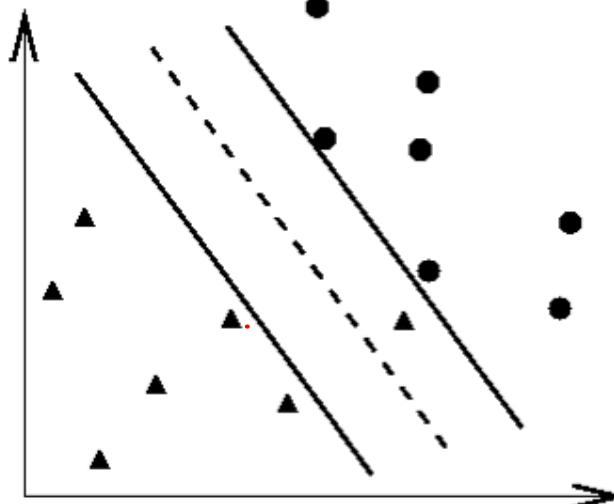
small C



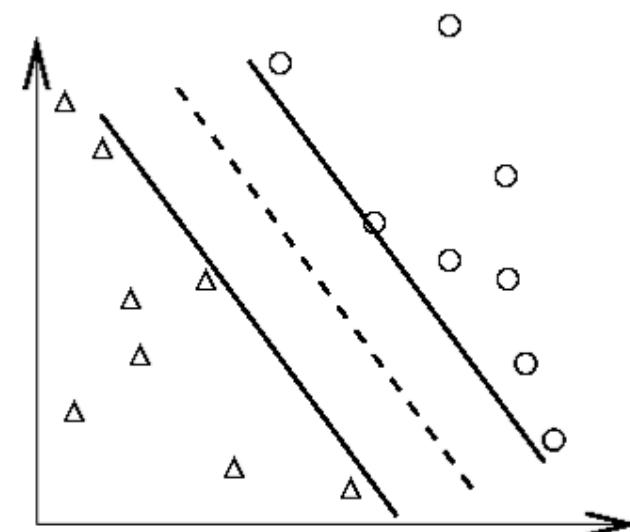
(a) Training data and an overfitting classifier



(b) Applying an overfitting classifier on testing data



(c) Training data and a better classifier



(d) Applying a better classifier on testing data

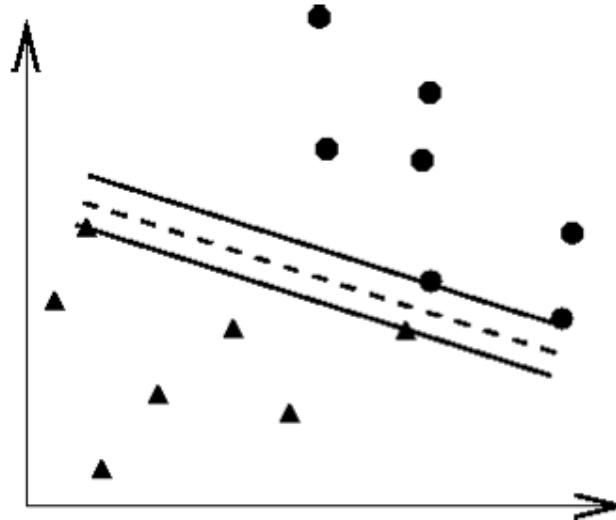
Model Selection, find right C

large C

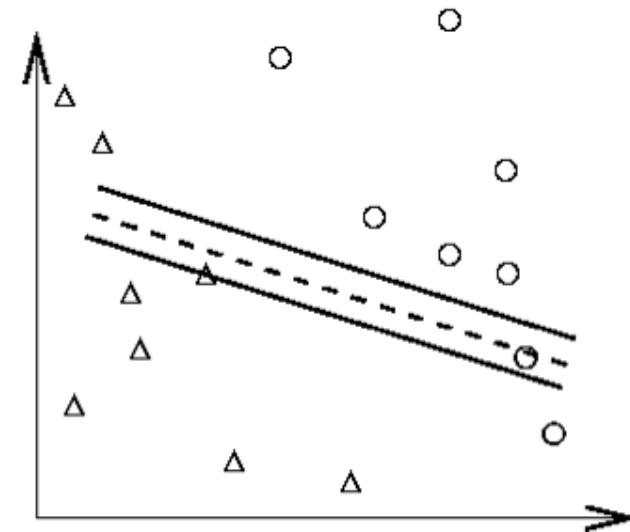
A large value of C means that misclassifications are bad - resulting in smaller margins and less training error (but more expected true error).

A small C results in more training error, hopefully better true error.

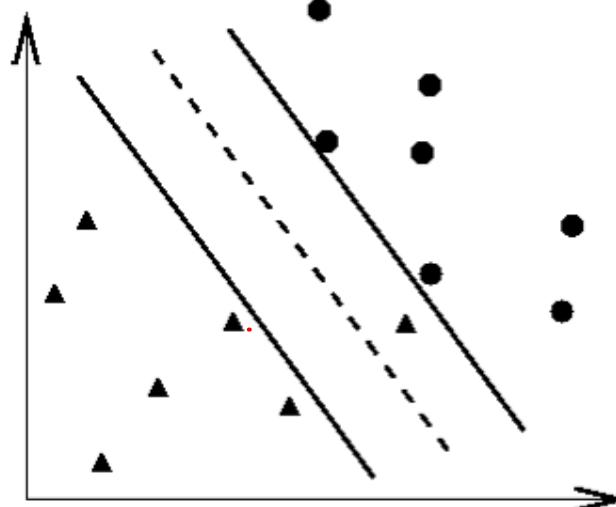
small C



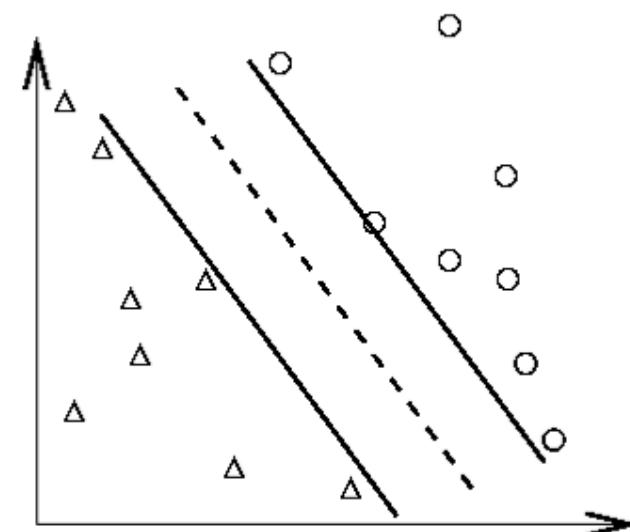
(a) Training data and an overfitting classifier



(b) Applying an overfitting classifier on testing data



(c) Training data and a better classifier



(d) Applying a better classifier on testing data

Today

❑ Support Vector Machine (SVM)

- ✓ History of SVM
- ✓ Large Margin Linear Classifier
- ✓ Define Margin (M) in terms of model parameter
- ✓ Optimization to learn model parameters (w , b)
- ✓ Linearly non-separable case
- ✓ Optimization with dual form
- ✓ Nonlinear decision boundary
- ✓ Practical Guide

Two optimization problems: For the **separable** and **non separable** cases

$$\text{Min } (\mathbf{w}^T \mathbf{w})/2$$

For all x in class + 1

$$\mathbf{w}^T \mathbf{x} + b \geq 1$$

For all x in class - 1

$$\mathbf{w}^T \mathbf{x} + b \leq -1$$

$$\min_w \frac{\mathbf{w}^T \mathbf{w}}{2} + C \sum_{i=1}^n \varepsilon_i$$

For all x_i in class + 1

$$\mathbf{w}^T \mathbf{x}_i + b \geq 1 - \varepsilon_i$$

For all x_i in class - 1

$$\mathbf{w}^T \mathbf{x}_i + b \leq -1 + \varepsilon_i$$

For all i

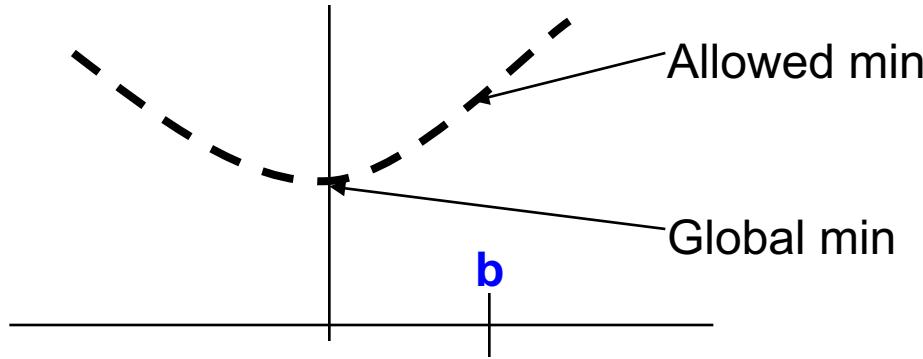
$$\varepsilon_i \geq 0$$

- Instead of solving these QPs directly we will solve a dual formulation of the SVM optimization problem
- The main reason for switching to this type of representation is that it would allow us to use a neat trick that will make our lives easier (and the run time faster)

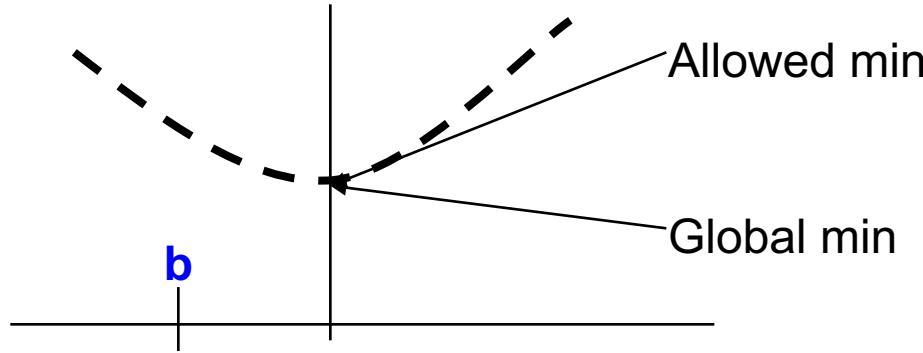
Optimization Review: Constrained Optimization

$$\begin{aligned} & \min_u u^2 \\ & \text{s.t. } u \geq b \end{aligned}$$

Case 1:



Case 2:



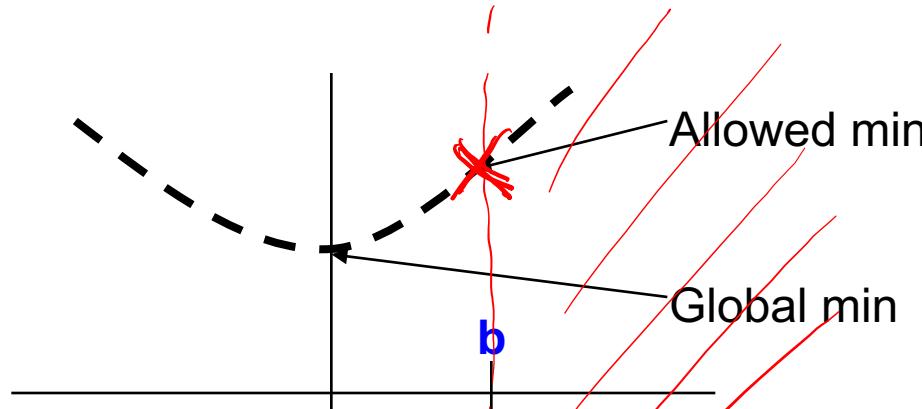
Optimization Review:

Constrained Optimization

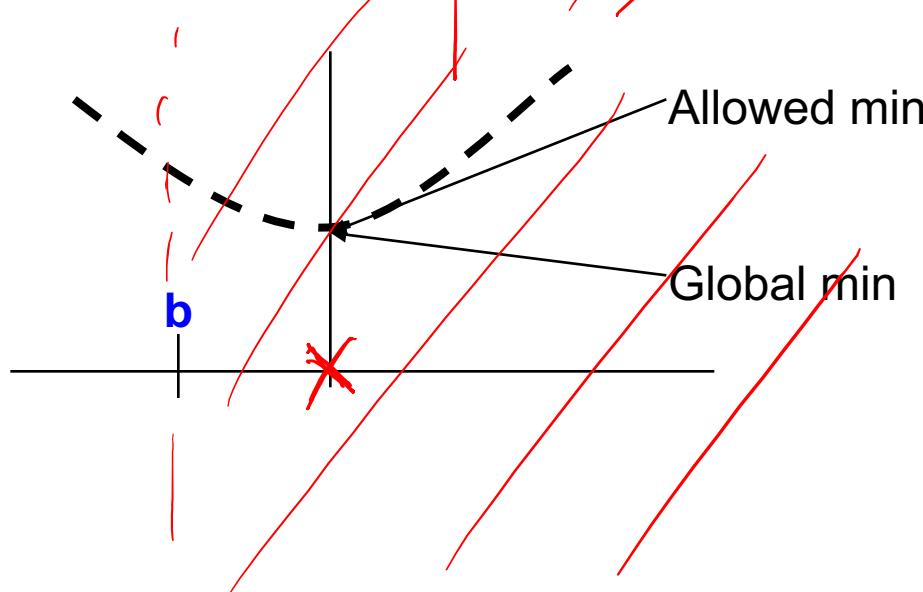
$$f(u) = u^2$$

$$\begin{aligned} \min_u & u^2 \\ \text{s.t. } & u \geq b \end{aligned}$$

Case 1:



Case 2:



Optimization Review:

Constrained Optimization

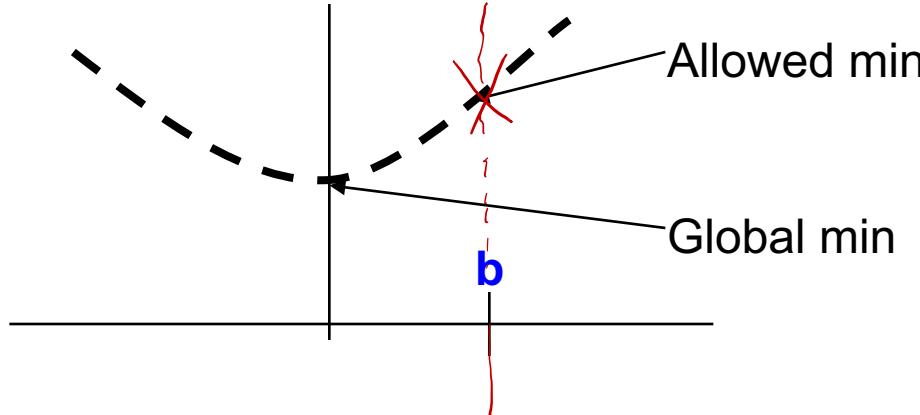
$$f(u)$$

$$\min_u u^2$$

$$\text{s.t. } u \geq b$$

Subject to

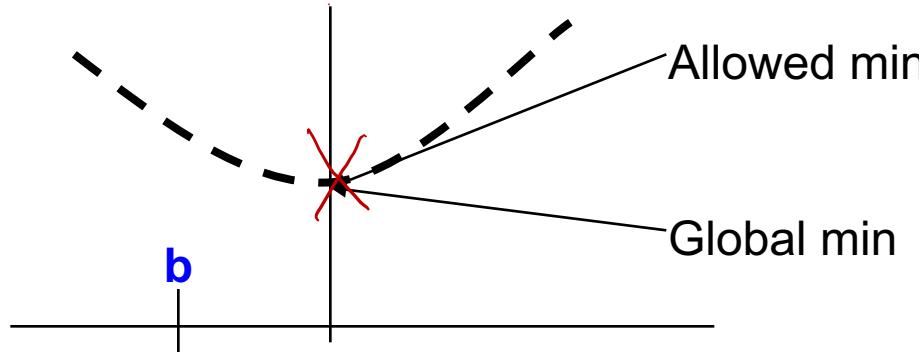
Case 1:



$$b > 0$$

$$f(u) = b^2$$

Case 2:



$$b < 0$$

$$f(u) = 0$$

Optimization Review:

Ingredients

- Objective function
- Variables
- Constraints

**Find values of the variables
that minimize or maximize the objective function
while satisfying the constraints**

Optimization Review:

Lagrangian Duality (Extra)

- The Primal Problem

Primal:

$$\begin{aligned} \min_w \quad & f_0(w) \\ \text{s.t.} \quad & f_i(w) \leq 0, \quad i = 1, \dots, k \end{aligned}$$

The generalized Lagrangian:

“Method of Lagrange multipliers”
convert to a higher-dimensional problem

$$L(w, \alpha) = f_0(w) + \sum_{i=1}^k \alpha_i f_i(w)$$

the α 's ($\alpha_i \geq 0$) are called the Lagrangian multipliers

Lemma:

$$\max_{\alpha, \alpha_i \geq 0} L(w, \alpha) = \begin{cases} f_0(w) & \text{if } w \text{ satisfies primal constraints} \\ \infty & \text{o/w} \end{cases}$$

A re-written Primal:

$$\min_w \max_{\alpha, \alpha_i \geq 0} L(w, \alpha)$$

Optimization Review: Lagrangian Duality, cont. (Extra)

- Recall the Primal Problem:

$$\min_w \max_{\alpha, \alpha_i \geq 0} L(w, \alpha)$$

- The Dual Problem:

$$\max_{\alpha, \alpha_i \geq 0} \min_w L(w, \alpha)$$

- Theorem (weak duality):**

$$d^* = \max_{\alpha, \alpha_i \geq 0} \min_w L(w, \alpha) \leq \min_w \max_{\alpha, \alpha_i \geq 0} L(w, \alpha) = p^*$$

- Theorem (strong duality):**

Iff there exist a saddle point of $L(w, \alpha)$

we have

$$d^* = p^*$$

An alternative representation of the SVM QP

- We will start with the linearly separable case
- Instead of encoding the correct classification rule and constraint we will use Lagrange multiplies to encode it as part of the our minimization problem

$$\begin{aligned} & \text{Min } (\mathbf{w}^T \mathbf{w})/2 \\ & \text{s.t.} \\ & (\mathbf{w}^T \mathbf{x}_i + b) y_i \geq 1 \end{aligned}$$

Recall that Lagrange multipliers can be applied to turn the following problem:

$$L_{\text{primal}}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i (y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1)$$

$\nearrow \mathbf{w}^T \mathbf{w}$

The Dual Problem (Extra)

$$\max_{\alpha_i \geq 0} \underbrace{\min_{w,b} L(w,b,\alpha)}_{\text{Dual formulation}}$$

Dual formulation

- We minimize L with respect to w and b first:

$$\nabla_w L(w,b,\alpha) = w - \sum_{i=1}^{\text{train}} \alpha_i y_i x_i = 0, \quad (*)$$

$$\nabla_b L(w,b,\alpha) = \sum_{i=1}^{\text{train}} \alpha_i y_i = 0, \quad (**)$$

Note that $(*)$ implies: $w = \sum_{i=1}^{\text{train}} \alpha_i y_i x_i$ (***)

- Plus $(***)$ back to L , and using $(**)$, we have:

$$L(w,b,\alpha) = \sum_{i=1} \alpha_i - \frac{1}{2} \sum_{i,j=1} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j)$$

Summary: Dual for SVM (Extra)

Solving for \mathbf{w} that gives maximum margin:

1. Combine objective function and constraints into new objective function, using **Lagrange multipliers** α_i ,

$$L_{primal} = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i (y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1)$$

2. To minimize this **Lagrangian**, we take derivatives of \mathbf{w} and b and set them to 0:

Summary: Dual for SVM (Extra)

3. Substituting and rearranging gives the **dual** of the Lagrangian:

$$L_{dual} = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

which we try to maximize (not minimize).

4. Once we have the α_i , we can substitute into previous equations to get \mathbf{w} and b .
5. This defines \mathbf{w} and b as **linear combinations of the training data**.

$$\mathbf{w} = \sum_{i=1}^{train} \alpha_i y_i \mathbf{x}_i$$

Summary: Dual SVM for linearly separable case

Dual formulation

$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0 \quad \forall i$$

$n \times i$



$$\text{Min } (\mathbf{w}^T \mathbf{w})/2$$

subject to the following inequality constraints:

For all x in class +1

$$\mathbf{w}^T \mathbf{x} + b \geq 1$$

For all x in class -1

$$\mathbf{w}^T \mathbf{x} + b \leq -1$$

}

A total of n constraints if we have n input samples

Easier than original QP, more efficient algorithms exist to find α_i

Dual SVM for linearly separable case – Training

Our dual target function: $\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$

$$\sum_i \alpha_i y_i = 0$$

Dot product for all training samples

$$\alpha_i \geq 0 \quad \forall i$$

$$\left(\begin{array}{c|ccccc|c}
1 & & & & & & n \\
2 & & & & & & \\
3 & & & & & & \\
\vdots & & & & & & \\
i & & & & \xrightarrow{\mathbf{x}_i^T \mathbf{x}_j} & & \\
\vdots & & & & & & \\
n & & & & & &
\end{array} \right) \quad n \times n$$

matrix

Support vectors: non-zero α_i

- only a few α_i 's can be nonzero!!

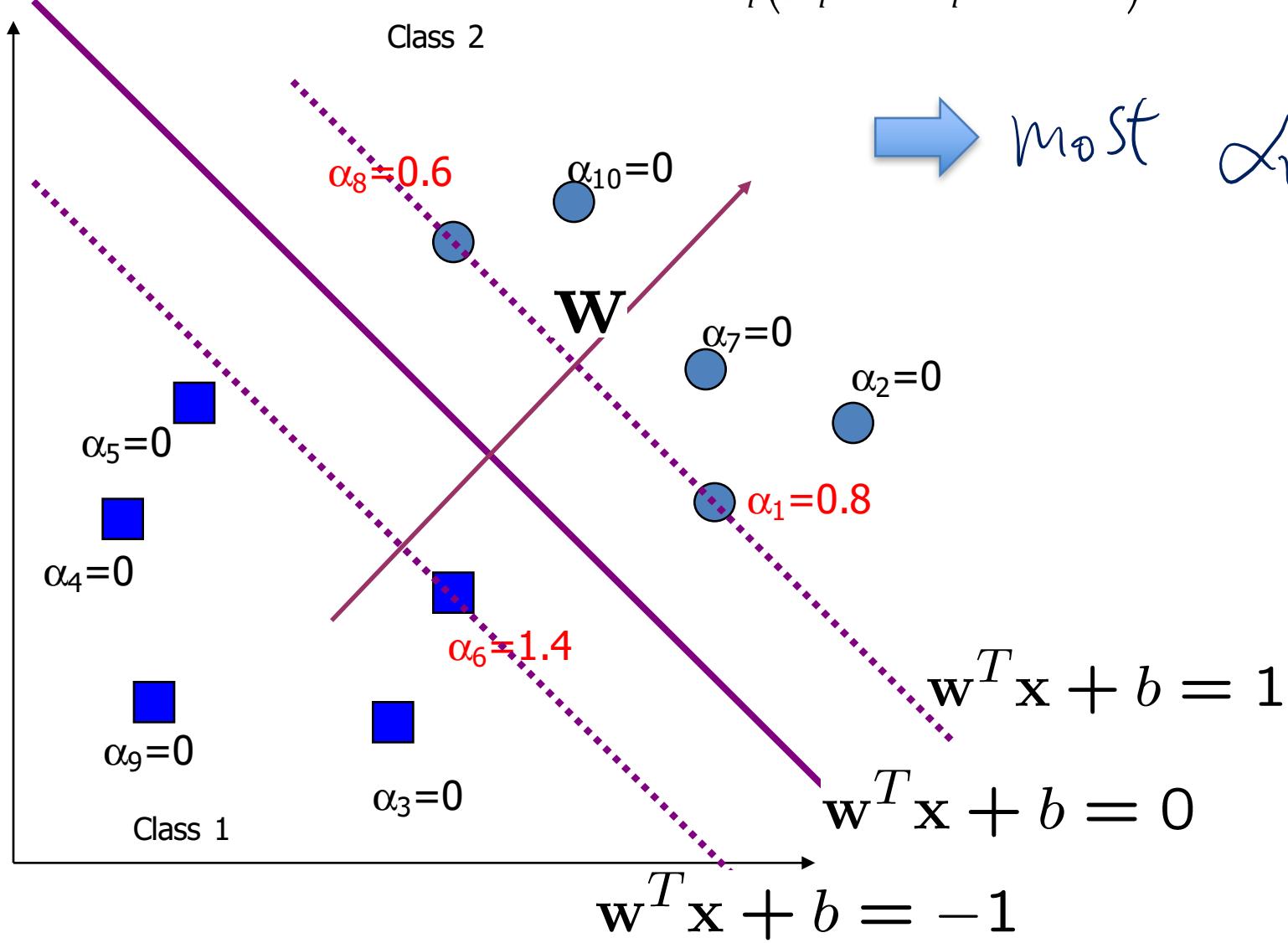
$$\alpha_i(y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1) = 0, \quad i = 1, \dots, n$$

for most $\alpha_i \Rightarrow \alpha_i = 0$

We call the training data points whose α_i 's are nonzero the **support vectors** (SV)

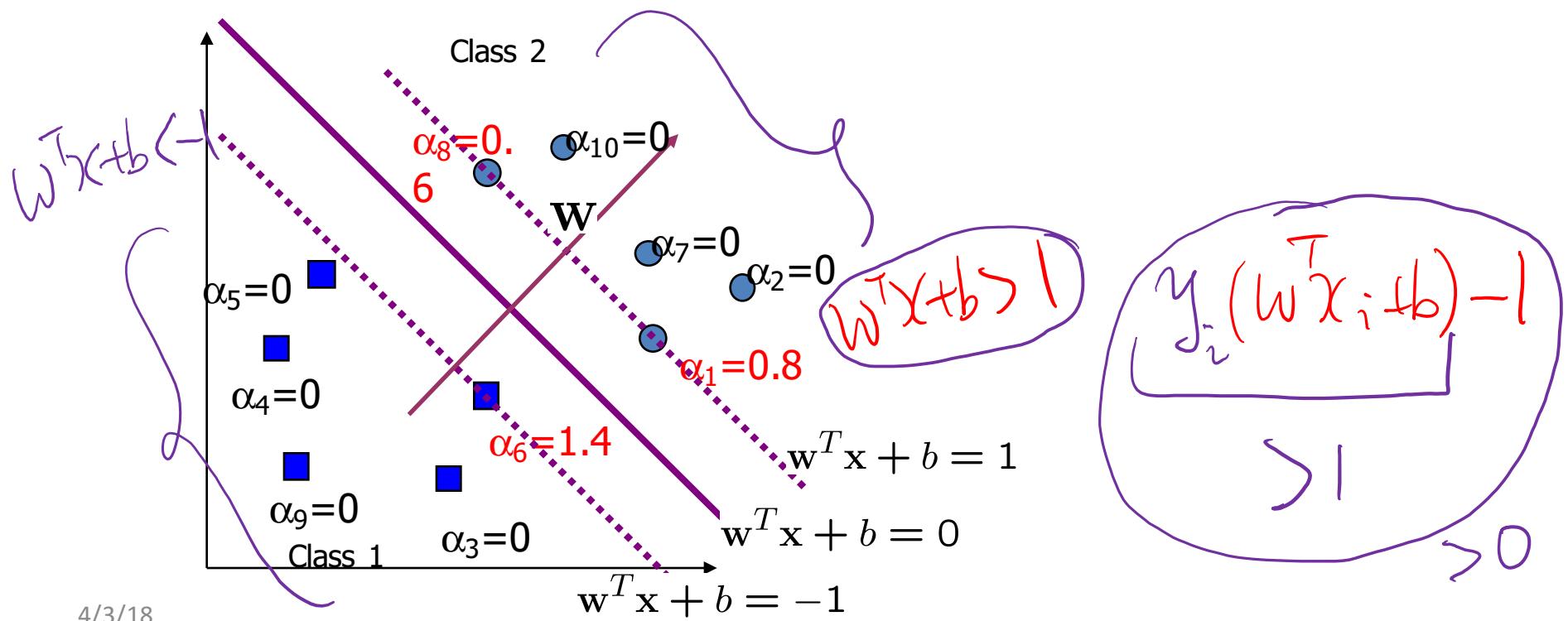
$$\alpha_i(y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1) = 0, \quad i=1,\dots,n$$

most $\alpha_i = 0$

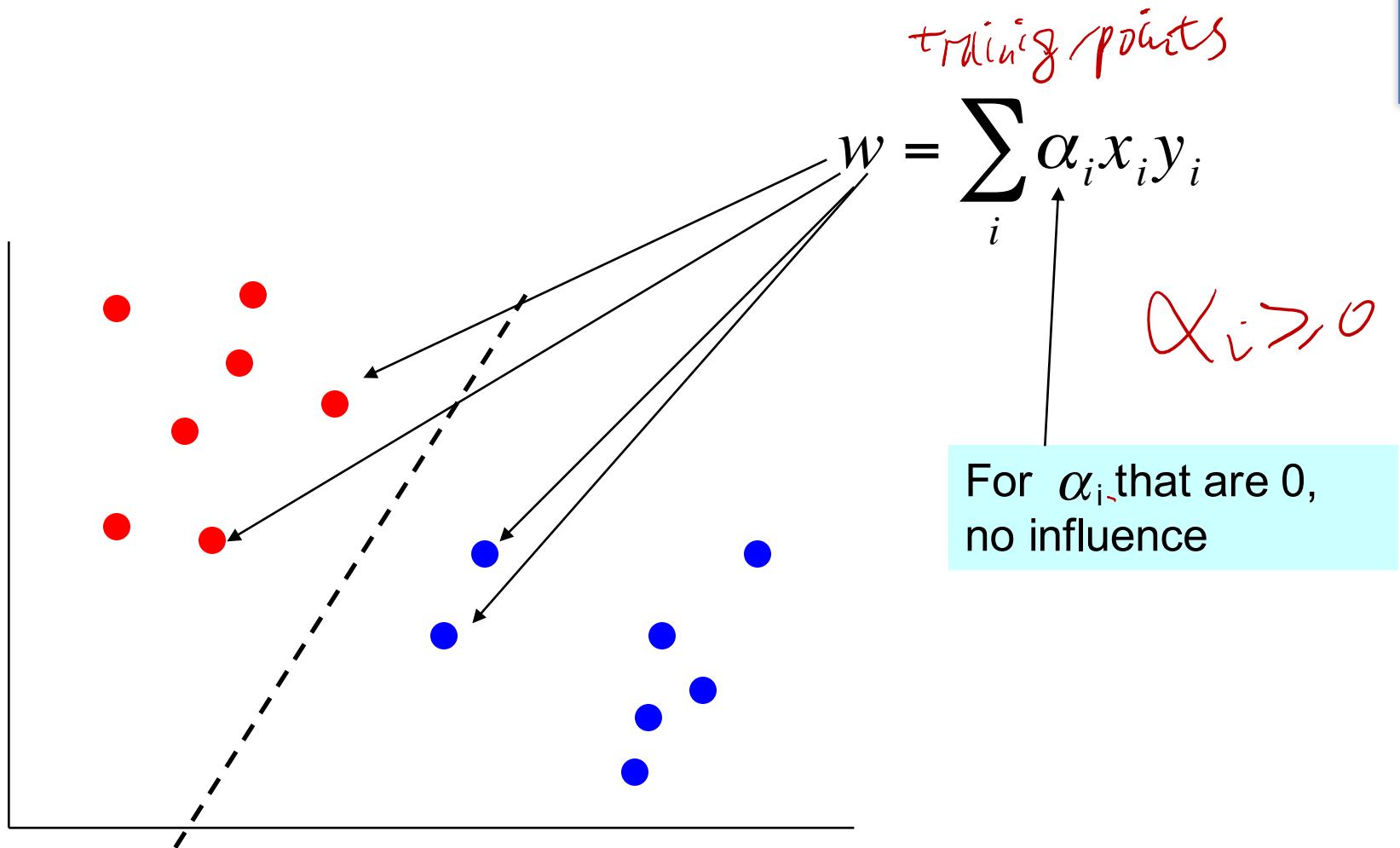


- only a few α_i can be nonzero!!

$$\alpha_i(y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1) = 0, \quad i=1,\dots,n \rightarrow \text{for most } \alpha_i \Rightarrow \alpha_i = 0$$



Dual SVM - interpretation



Dual SVM for linearly separable case –

Testing

To evaluate a new sample x_{ts}
we need to compute:

$$\hat{y}_{ts} = \text{sign}(w^T x_{ts} + b) = \text{sign}\left(\sum_{i=1..n} \alpha_i y_i x_i^T x_{ts} + b\right)$$

$$\hat{y}_{ts} = \text{sign}\left(\sum_{i \in SupportVectors} \alpha_i y_i (x_i^T x_{ts}) + b\right)$$

Dot product with (“all” ??)
training samples



For α_i that are 0,
no influence

Dual formulation for linearly non-separable case

Dual target function:

$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0$$

$$C > \alpha_i \geq 0, \forall i$$

Hyperparameter C
should be tuned
through k-folds CV

The only difference is
that the \alpha are now
bounded

This is very similar to the
optimization problem in the linear
separable case, except that there is
an upper bound C on α_i now

Once again, efficient algorithm exist
to find α_i

Dual formulation for linearly non-separable case

Dual target function:

$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0$$

$$C > \alpha_i \geq 0, \forall i$$

Hyperparameter C
should be tuned
through k-folds CV

The only difference is
that the \alpha are now
bounded

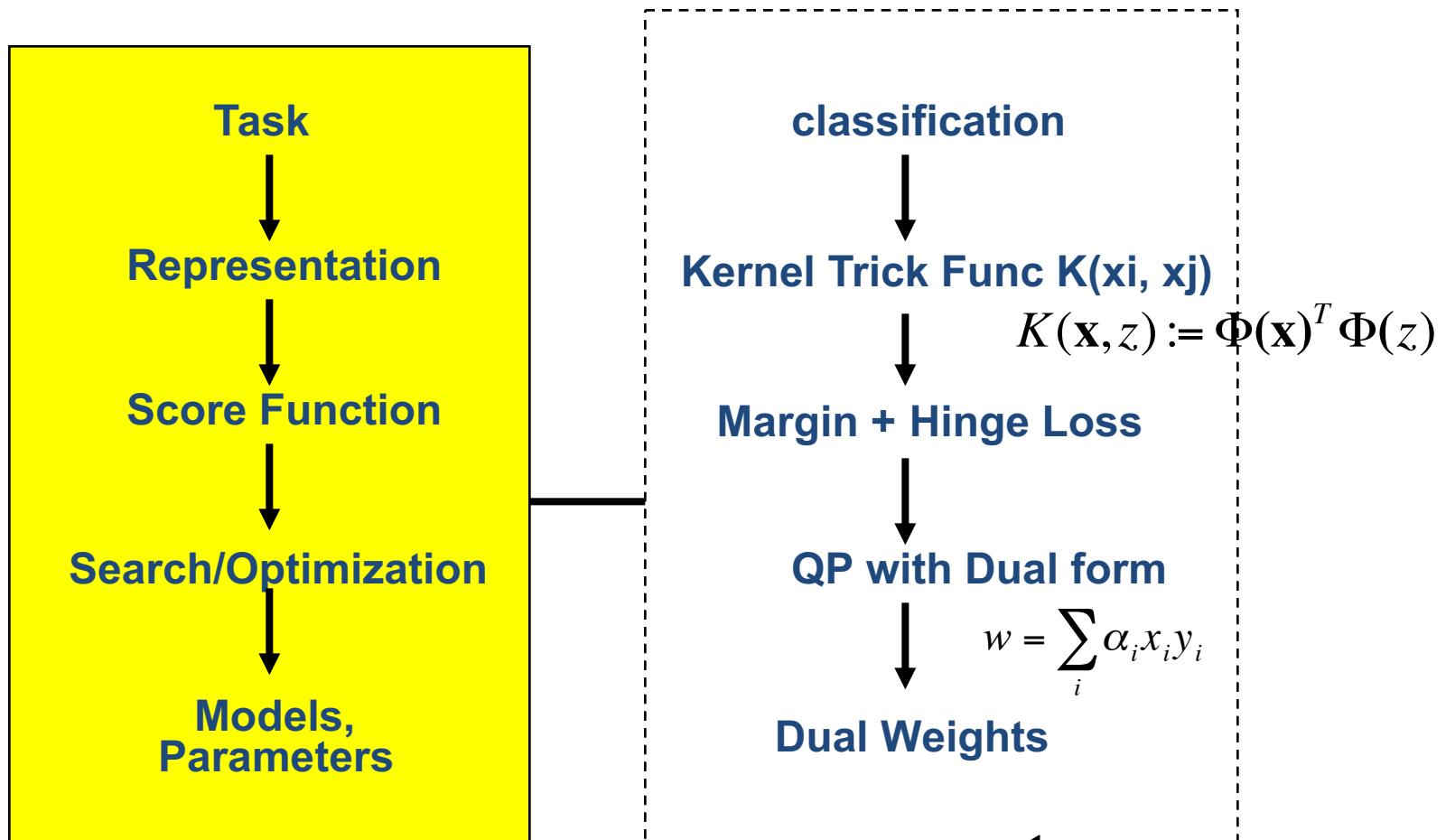
To evaluate a new sample \mathbf{x}_{ts}
we need to compute:

$$\mathbf{w}^T \mathbf{x}_{ts} + b = \sum_{i \in \text{supportV}} \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_{ts} + b$$

This is very similar to the
optimization problem in the linear
separable case, except that there is
an upper bound C on α_i now

Once again, efficient algorithm exist
to find α_i

Support Vector Machine



$$\underset{\mathbf{w}, b}{\operatorname{argmin}} \sum_{i=1}^p w_i^2 + C \sum_{i=1}^n \varepsilon_i$$

$$\text{subject to } \forall \mathbf{x}_i \in D_{\text{train}} : y_i (\mathbf{x}_i \cdot \mathbf{w} + b) \geq 1 - \varepsilon_i$$



$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0, \quad \alpha_i \geq 0 \quad \forall i$$

References

- Big thanks to Prof. Ziv Bar-Joseph and Prof. Eric Xing @ CMU for allowing me to reuse some of his slides
- Elements of Statistical Learning, by Hastie, Tibshirani and Friedman
- Prof. Andrew Moore @ CMU's slides
- Tutorial slides from Dr. Tie-Yan Liu, MSR Asia
- A Practical Guide to Support Vector Classification
Chih-Wei Hsu, Chih-Chung Chang, and Chih-Jen Lin,
2003-2010
- Tutorial slides from Stanford “Convex Optimization I –
Boyd & Vandenberghe