

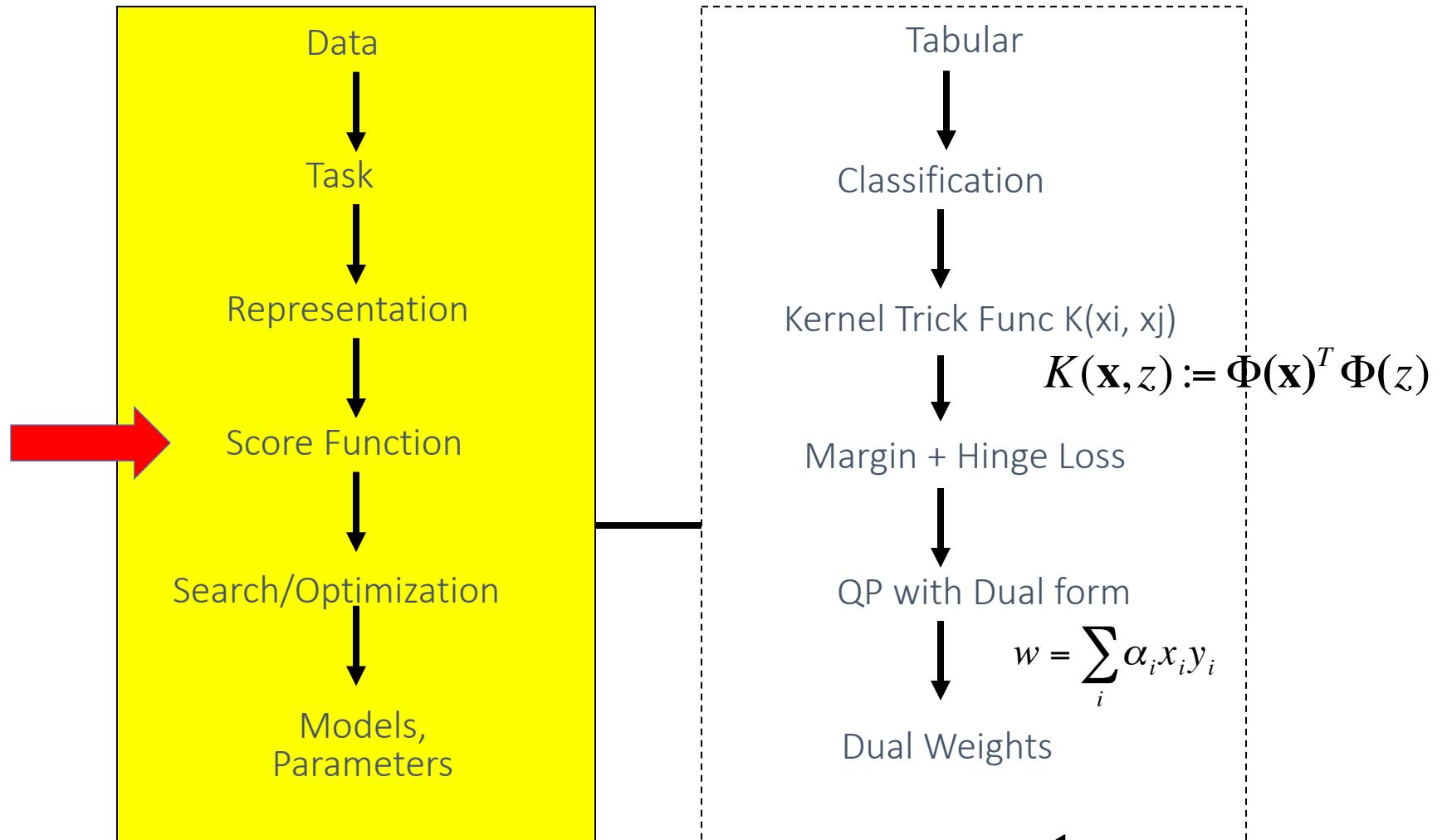
UVA CS 4774: Machine Learning

S4: Lecture 20: Support Vector Machine (Basics)

Dr. Yanjun Qi

University of Virginia
Department of Computer Science

Today: Basic Support Vector Machine



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

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

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

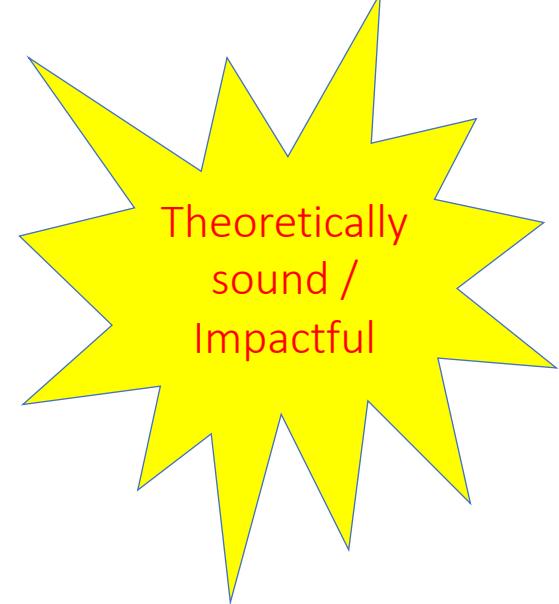
$$\sum_i \alpha_i y_i = 0, \quad \alpha_i \geq 0 \quad \forall i$$

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 digit 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.

Handwritten digit recognition

→ MNIST (sum)

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

Applications of SVMs

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

→ Lots of very successful applications!!!

Early History

- In 1950 English mathematician Alan Turing wrote a landmark paper titled “Computing Machinery and Intelligence” that asked the question: “Can machines think?”
- Further work came out of a 1956 workshop at Dartmouth sponsored by John McCarthy. In the proposal for that workshop, he coined the phrase a “study of artificial intelligence”
- 1950s
 - Samuel’s checker player : start of machine learning
 - Selfridge’s Pandemonium
- 1952-1969: Enthusiasm: Lots of work on neural networks
- 1970s: Expert systems, Knowledge bases to add on rule-based inference

Early History

- 1980s :
 - Advanced decision tree and rule learning
 - Valiant's PAC Learning Theory
- 1990s:
 - Reinforcement learning (RL)
 - Ensembles: Bagging, Boosting, and Stacking
 - Bayes Net learning
 - Convolutional neural network (CNN) and Recurrent neural network (RNN) were invented
 - **SVM**
- 2000s **~2010s**
 - Support vector machines (becoming popular and dominating)
 - Kernel methods
 - Graphical models
 - Statistical relational learning
 - Transfer learning
 - Sequence labeling
 - Collective classification and structured outputs

MIT Technology Review

10 Breakthrough Technologies 2013

Think of the most frustrating, intractable, or simply annoying problems you can imagine. Now think about what technology is doing to fix them. That's what we did in coming up with our annual list of 10 Breakthrough Technologies. We're looking for technologies that we believe will expand the scope of human possibilities.

Deep Learning

CNN/RNN
1990s

10 Breakthrough Technologies 2017

These technologies all have staying power. They will affect the economy and our politics, improve medicine, or influence our culture. Some are unfolding now; others will take a decade or more to develop. But you should know about all of them right now.

Deep
Reinforcement
Learning

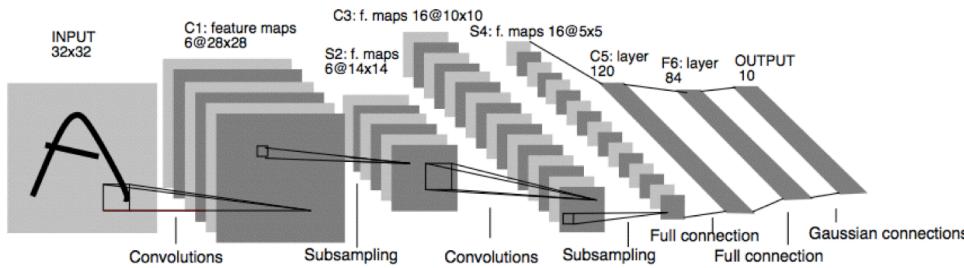
1990s~2010s



Generative
Adversarial
Network (GAN)

2010s

- 1952-1969 Enthusiasm: Lots of work on neural networks
- 1990s: Convolutional neural network (CNN) and Recurrent neural network (RNN) were invented



Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner,
Gradient-based learning applied to document
recognition, Proceedings of the IEEE 86(11):
2278–2324, 1998.

Reason of the recent breakthroughs of deep learning:

①

Low Bias
High Variable

Plenty of (Labeled) Data

Advanced Computer Architecture that fits DNNs

Powerful DNN platforms / Libraries

② XgBoost

Tree Based
reduce variance

X_1	X_2	X_3	Y

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

$\{-1, 1\}$

Output as Binary Class:
only two possibilities

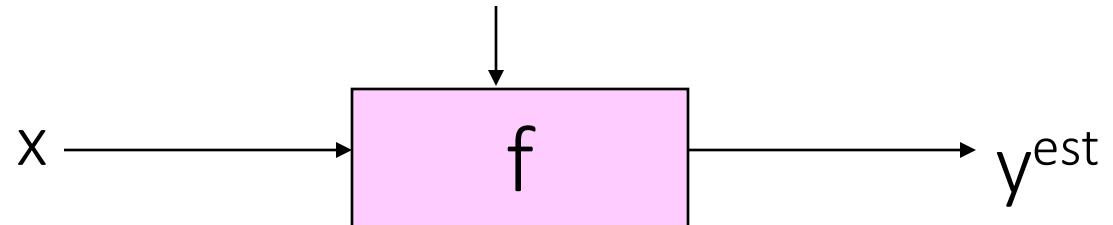
- Data/points/instances/examples/samples/records: [rows]
- Features/attributes/dimensions/independent variables/covariates/predictors/regressors: [columns, except the last]
- Target/outcome/response/label/dependent variable: special column to be predicted [last column]

Today

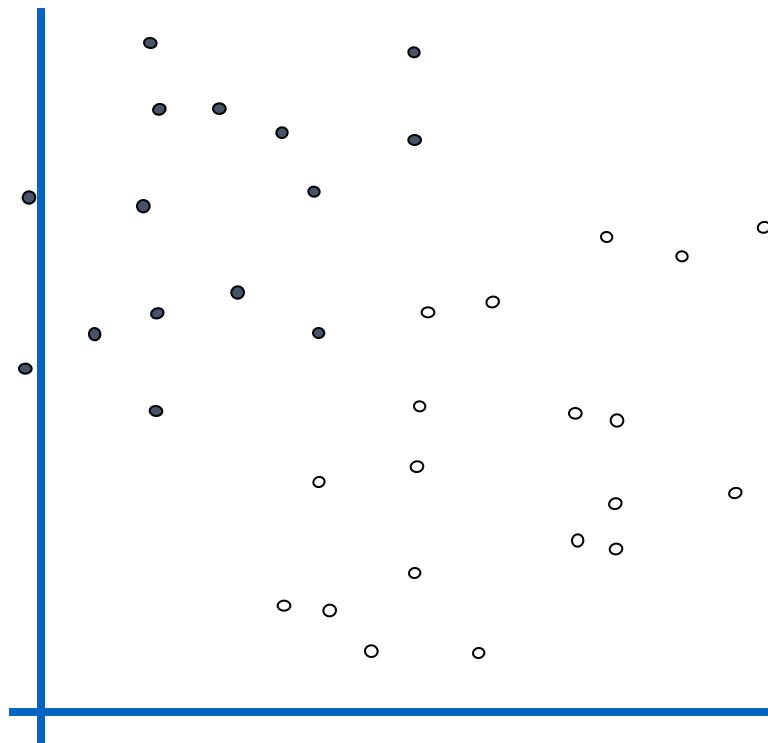
❑ Support Vector Machine (SVM)

- ✓ 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

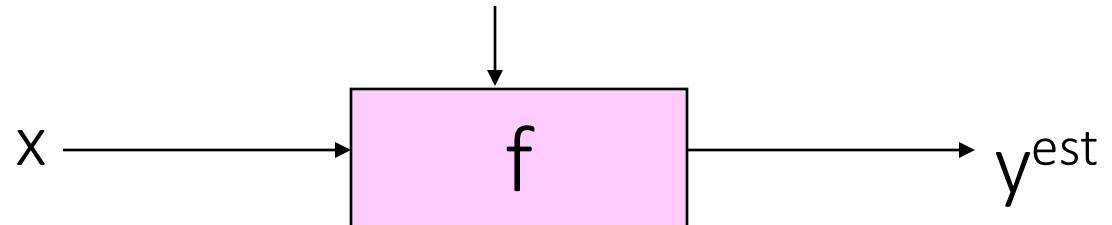


- denotes +1
- denotes -1



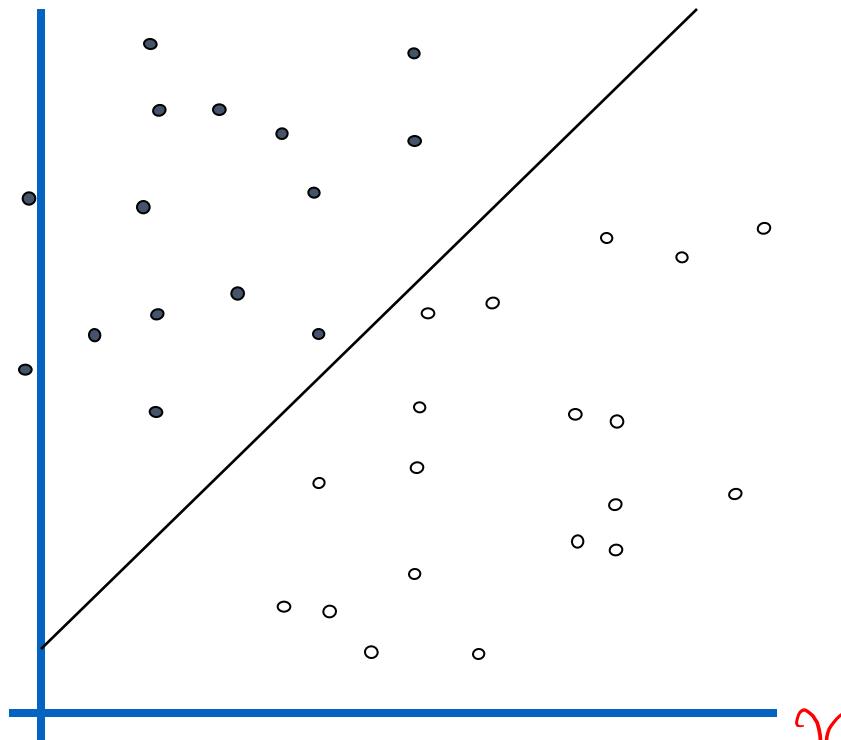
How would you
classify this data?

Linear Classifiers



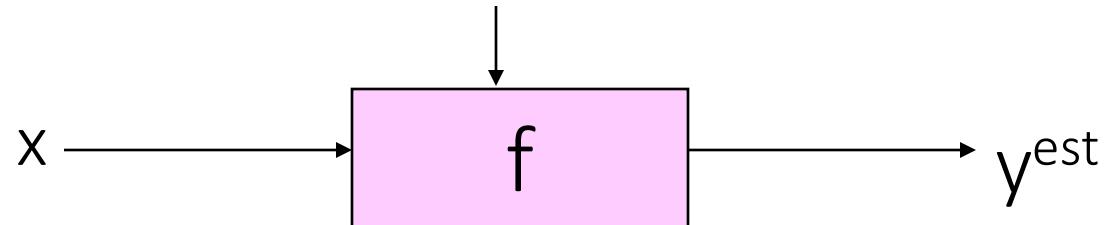
- denotes +1
- denotes -1

χ_2

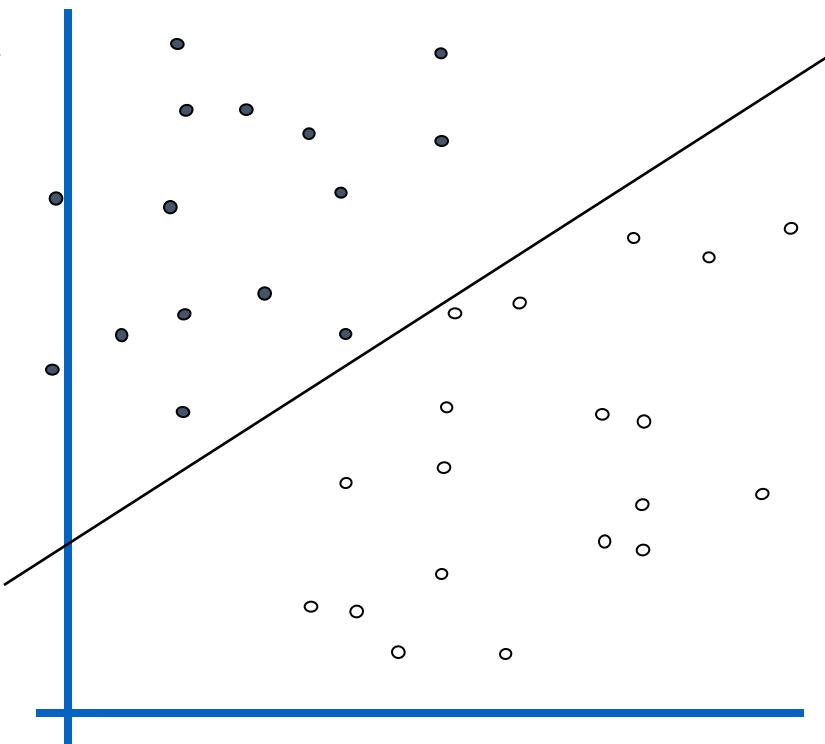


How would you
classify this data?

Linear Classifiers

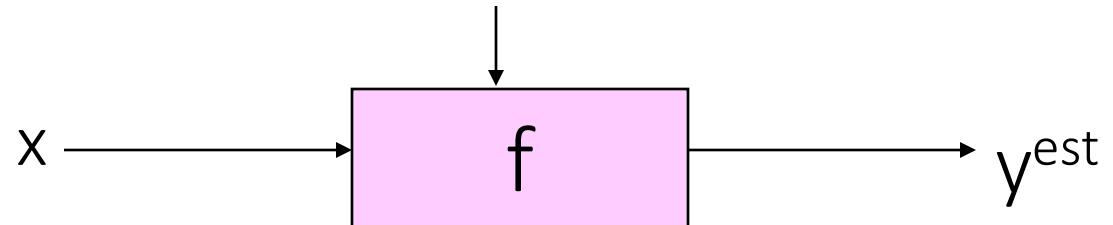


- denotes +1
- denotes -1

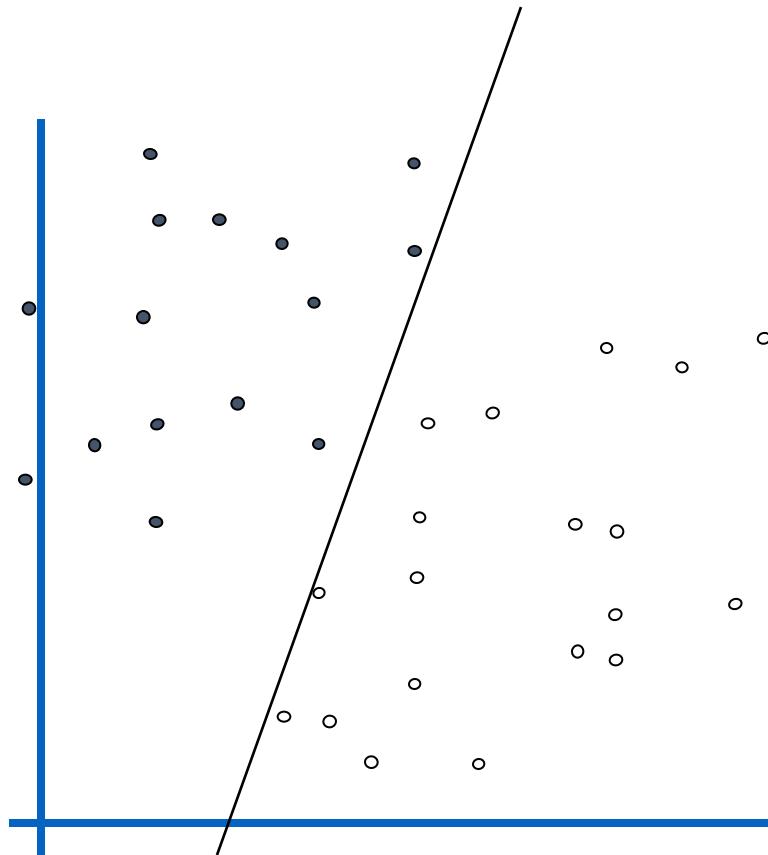


How would you
classify this data?

Linear Classifiers

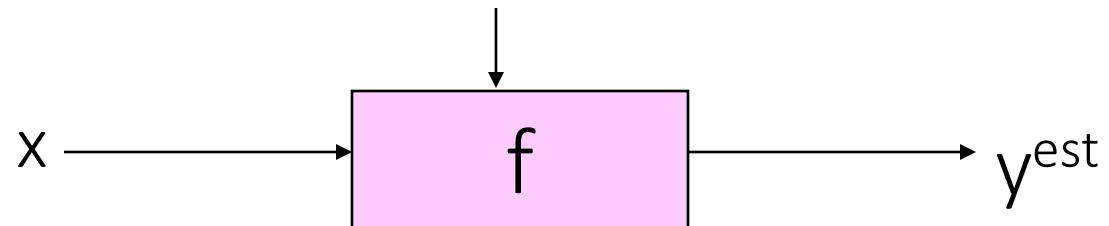


- denotes +1
- denotes -1

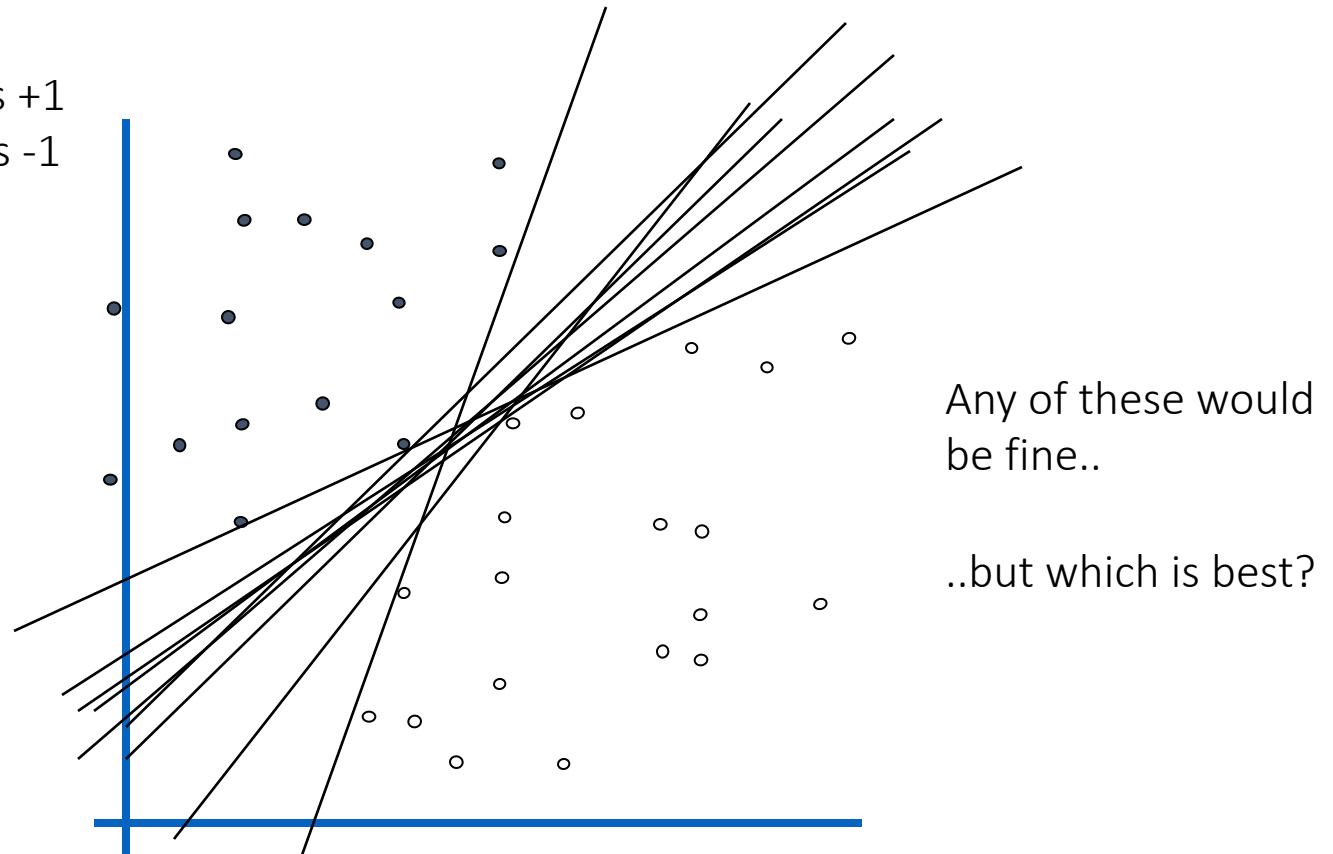


How would you
classify this data?

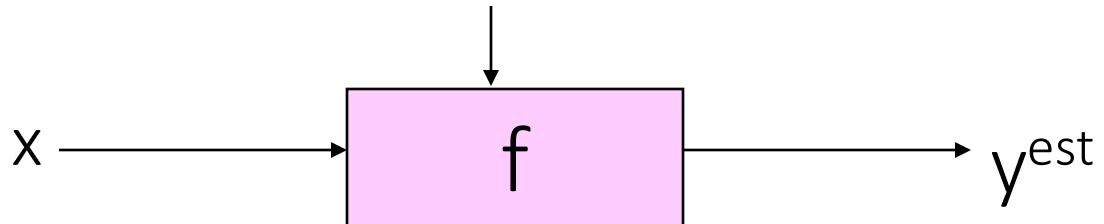
Linear Classifiers



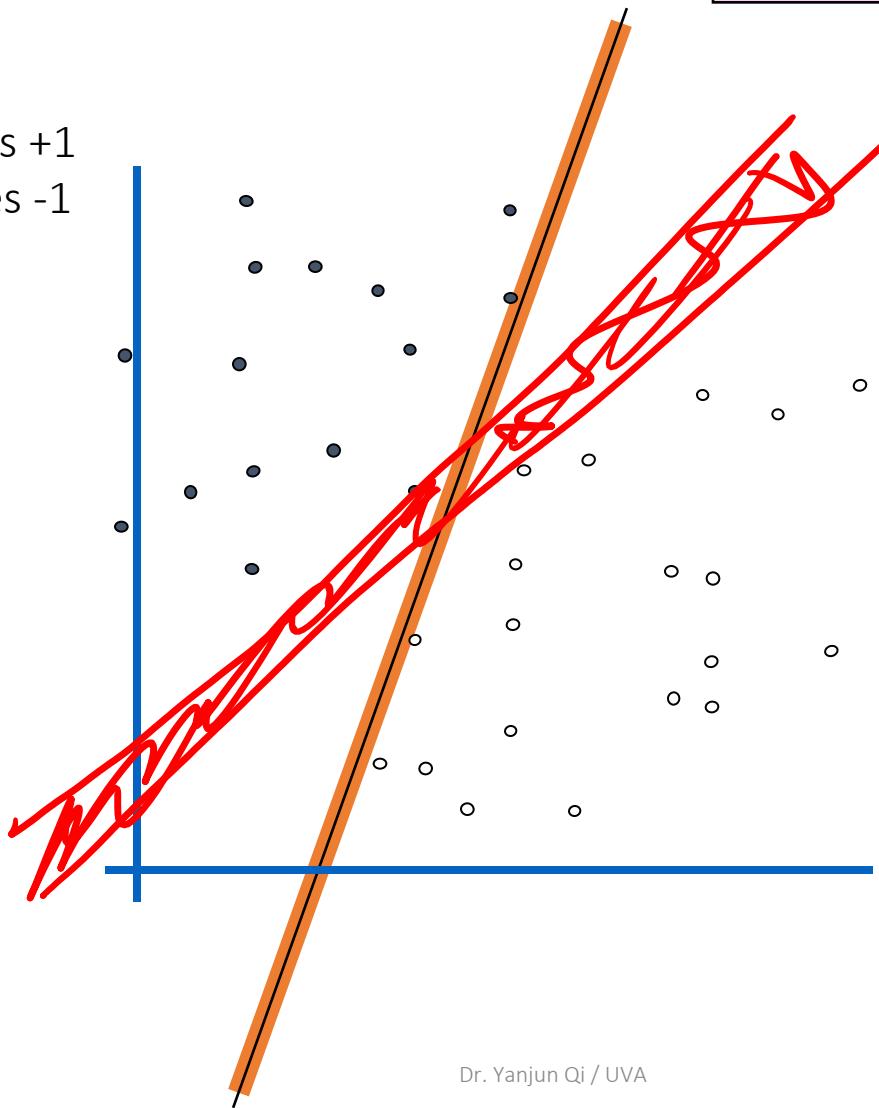
- denotes +1
- denotes -1



Classifier Margin

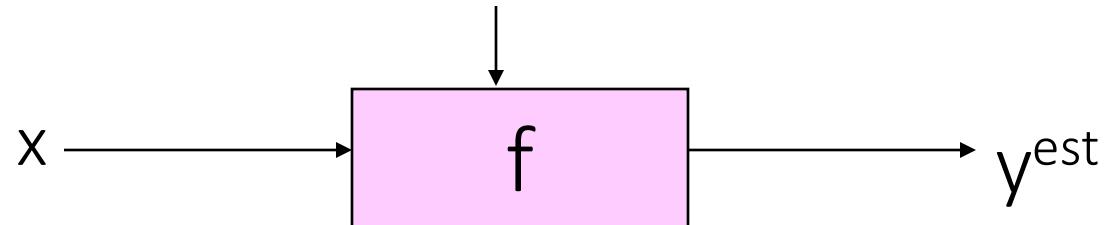


- denotes +1
- denotes -1

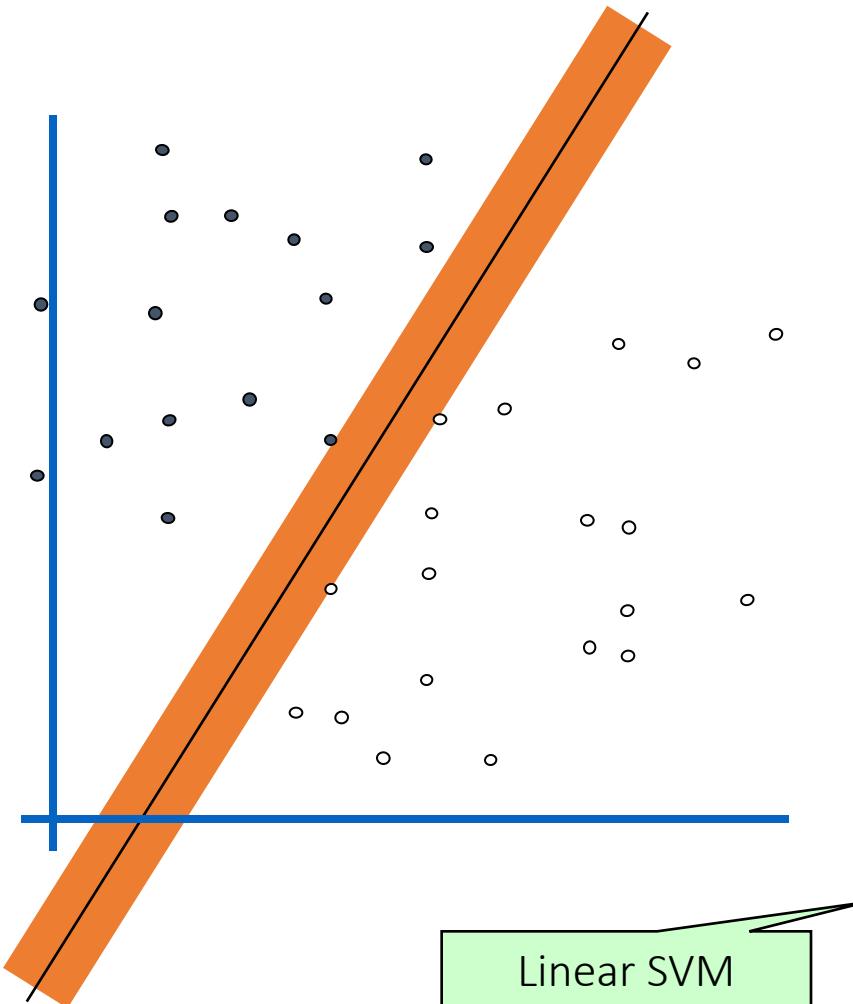


Define the **margin** of a linear classifier as the width that the **boundary could be increased by** before hitting a datapoint.

Maximum Margin



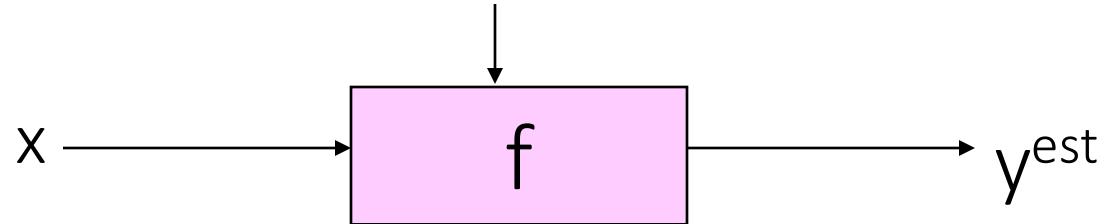
- denotes +1
- denotes -1



Linear SVM

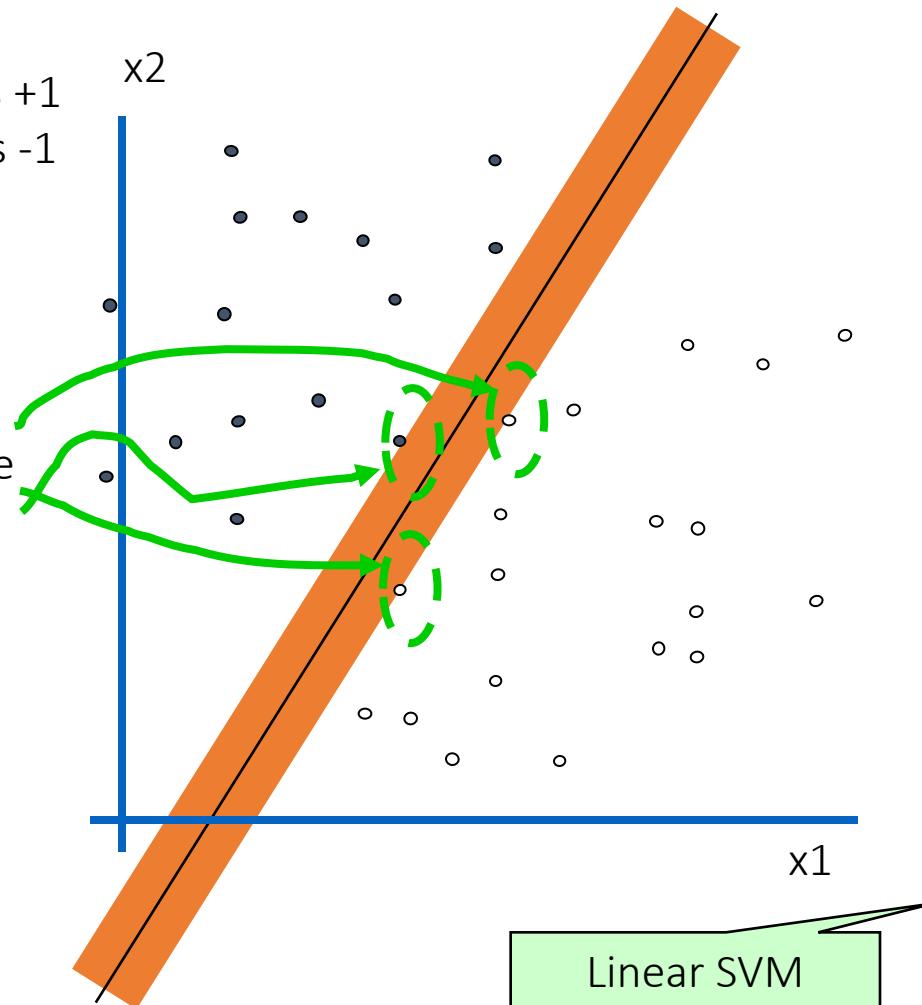
The **maximum margin linear classifier** is the linear classifier with the, **maximum** margin. This is the simplest kind of SVM (Called an LSVM)

Maximum Margin



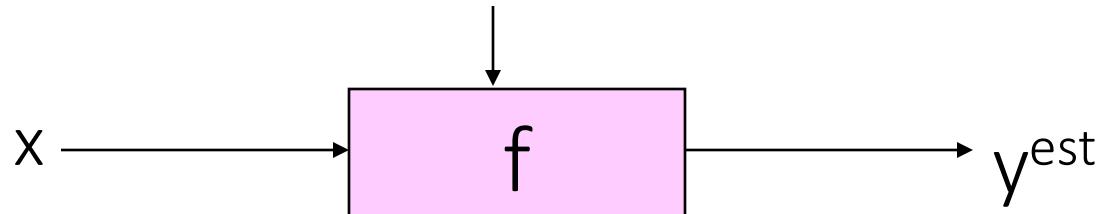
- denotes +1
- denotes -1

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)

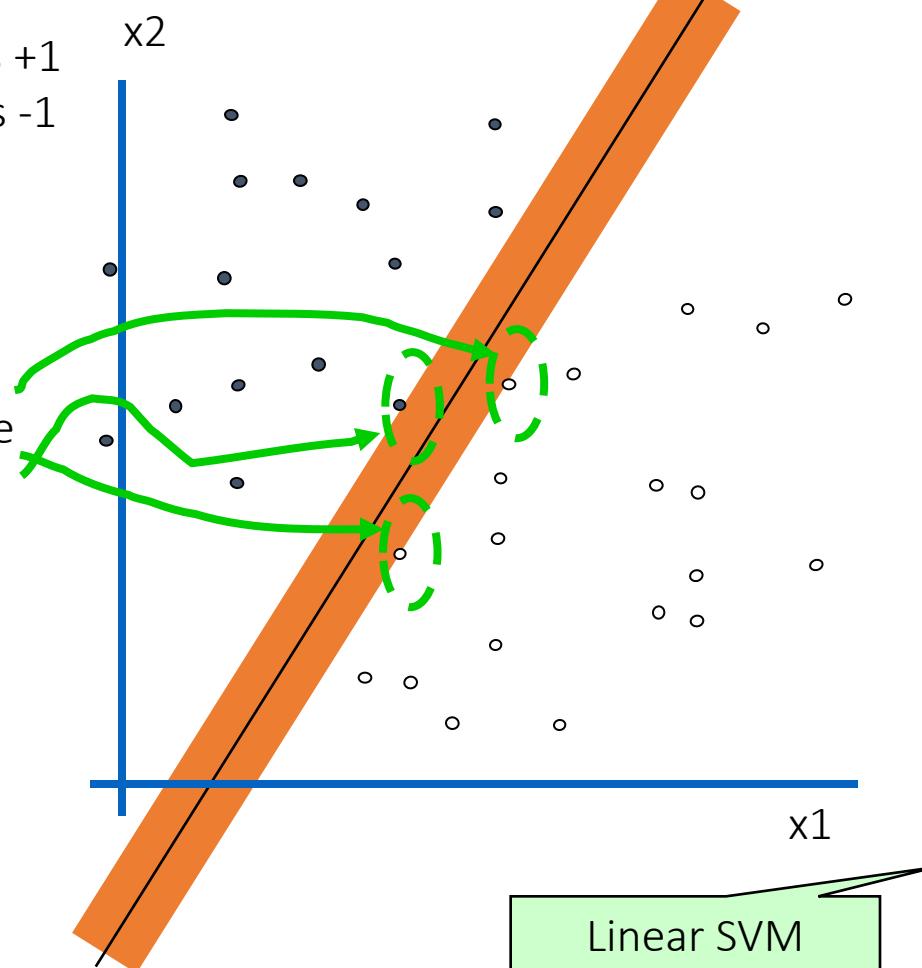
Maximum Margin



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

- denotes +1
- denotes -1

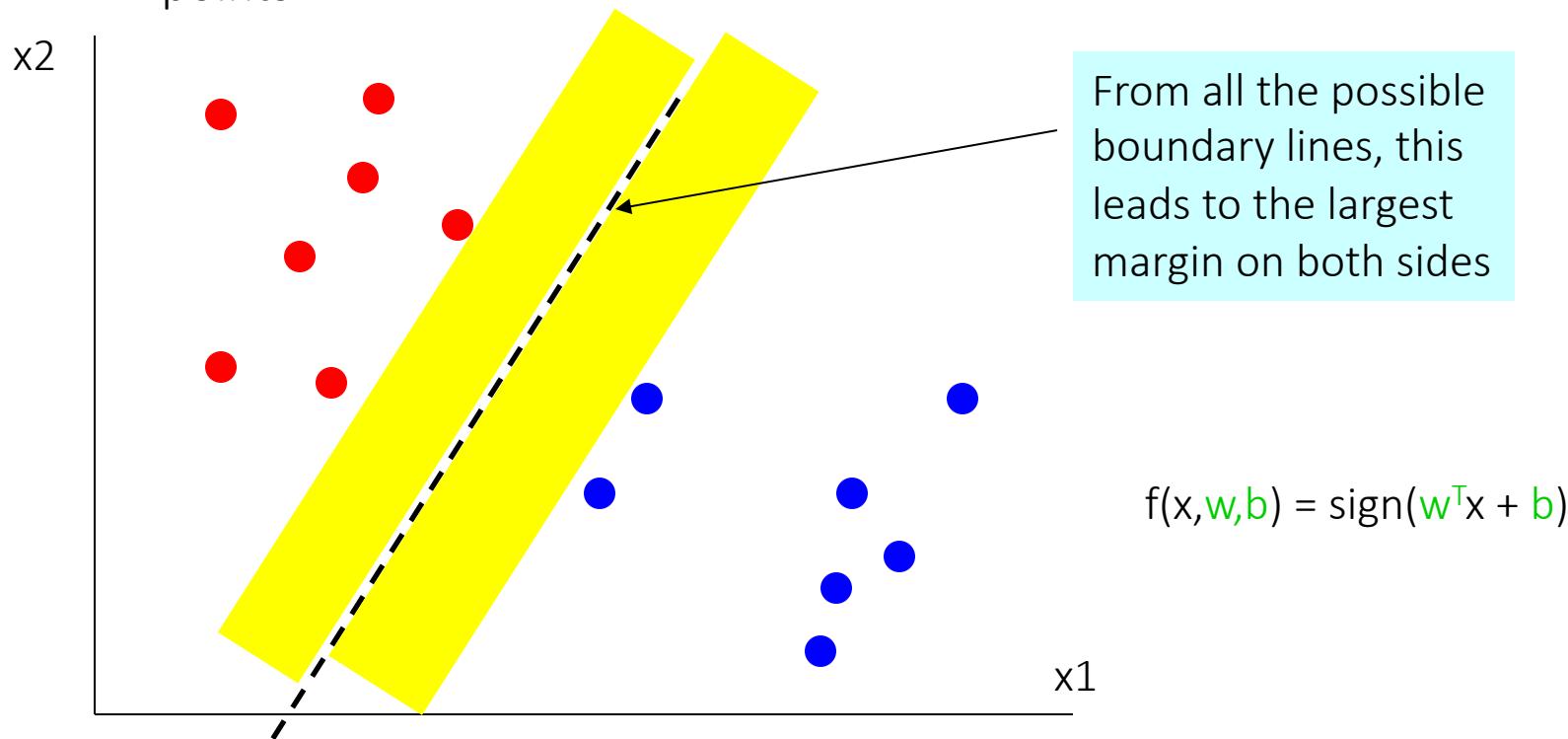
Support Vectors are those datapoints that the margin pushes up against



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

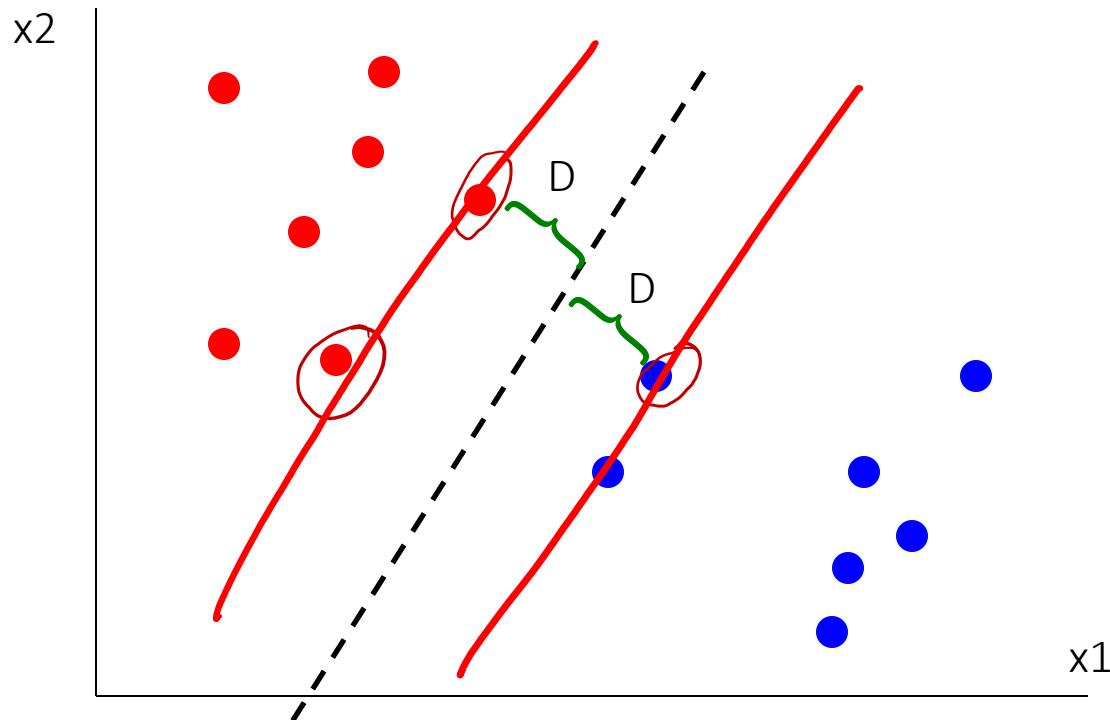
Max margin classifiers

- 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 (using VC dimension)
- Works well in practice



Thank You

Thank you

UVA CS 4774: Machine Learning

S4: Lecture 20: Support Vector Machine (Basics)

Module II

Dr. Yanjun Qi

University of Virginia
Department of Computer Science

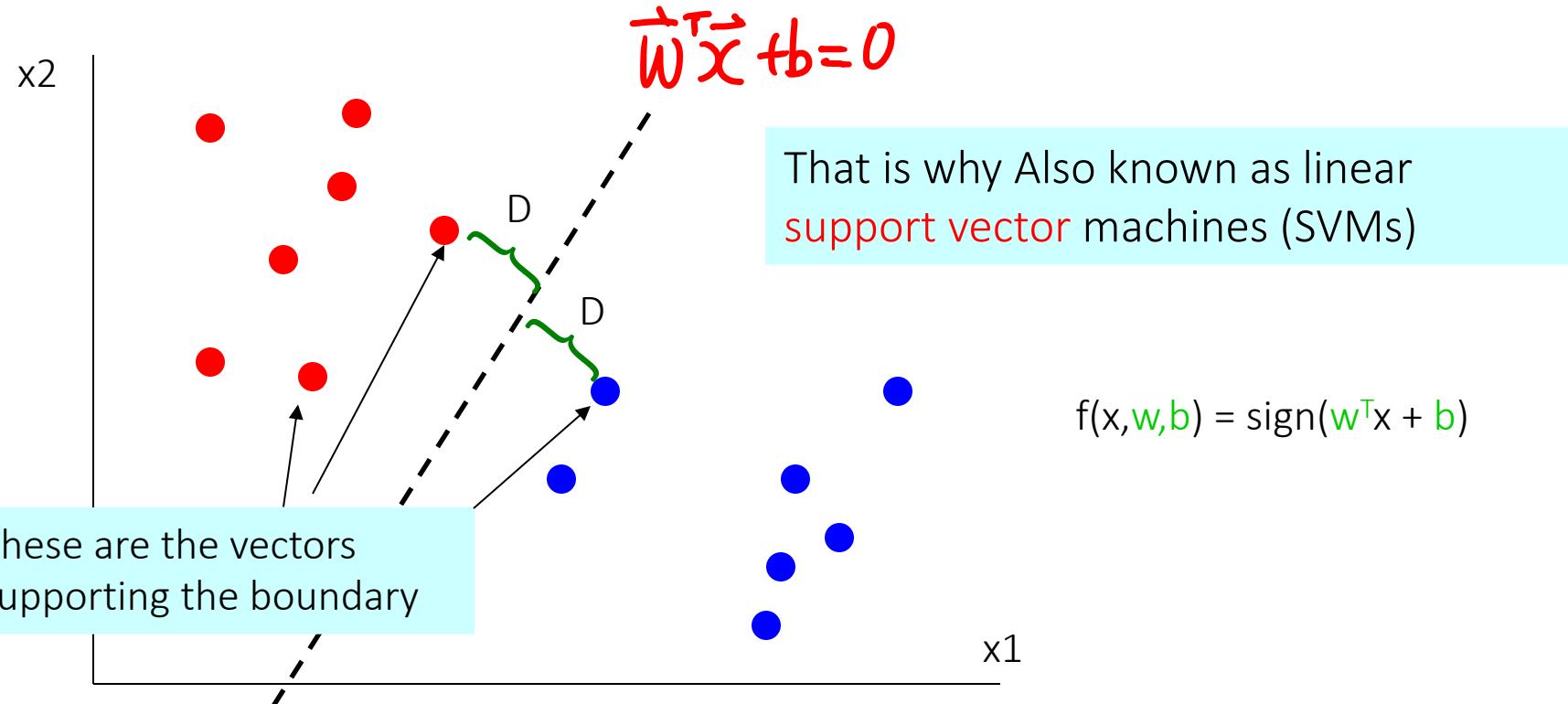
Today

- ❑ Supervised Classification
- ❑ Support Vector Machine (SVM)

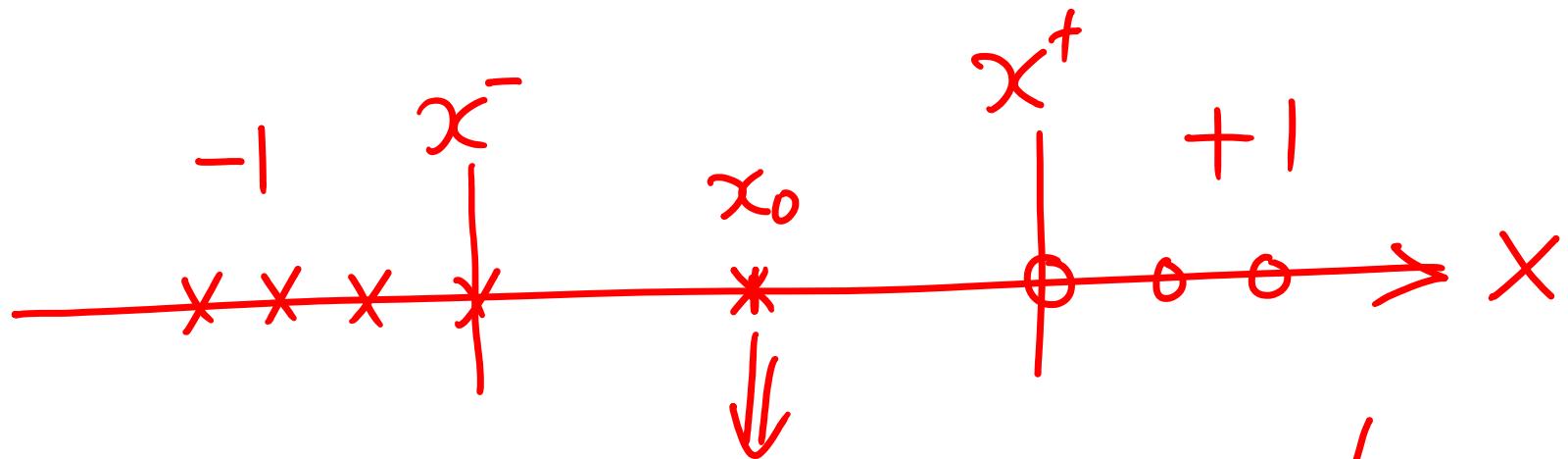
- ✓ History of SVM
- ✓ Large Margin Linear Classifier
- ✓ Define Margin (M) in terms of model parameter
- ✓ Optimization to learn model parameters (w, b)
- ✓ Linearly Non-separable case
- ✓ Optimization with dual form
- ✓ Nonlinear decision boundary
- ✓ Multiclass SVM

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

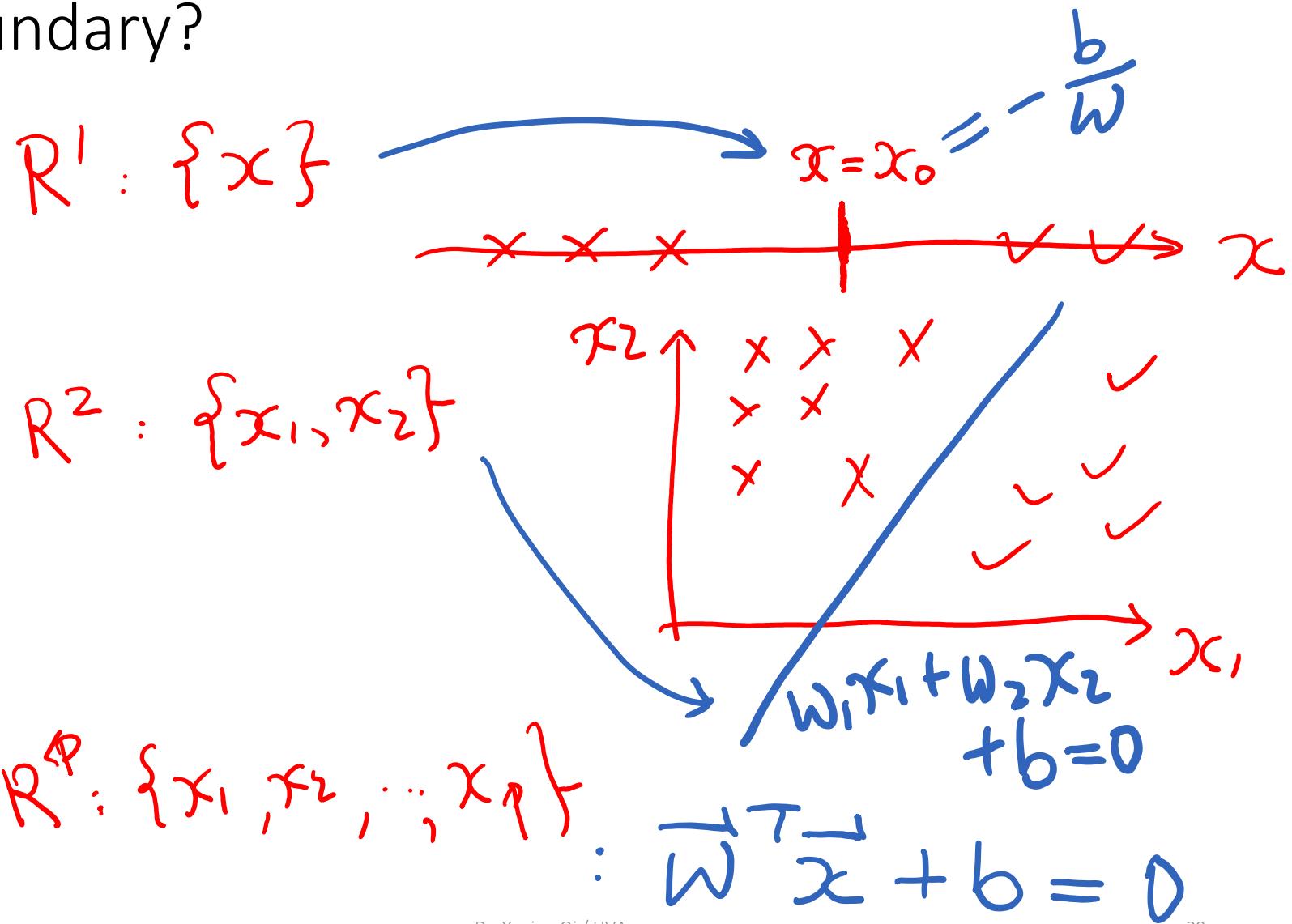


$$R^I : \{x\} \rightarrow \{+1, -1\}$$



$$\begin{cases} \vec{w}^T \vec{x}_0 + b = 0 \\ \vec{w}^T \vec{x}^- + b = -1 \\ \vec{w}^T \vec{x}^+ + b = 1 \end{cases} \quad f(x) = \text{Sign}(\vec{w}^T x + b)$$

How to represent a Linear Decision Boundary?



Review : 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

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

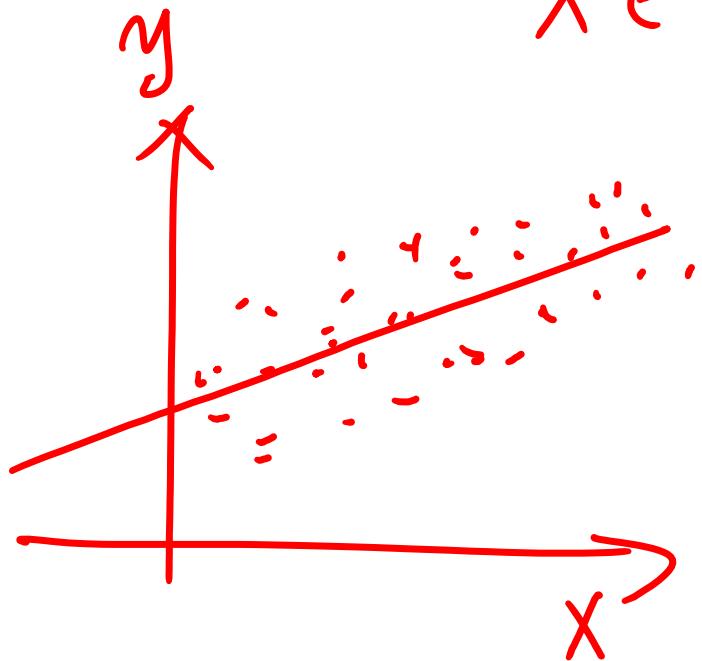
$$\begin{cases} \vec{x} \in \mathbb{R}^P \\ b \in \mathbb{R} \end{cases}$$

Q: How does this connect to linear regression?

[Regression]

$$y \in \mathbb{R}^1$$

$$\bar{x} \in \mathbb{R}^P, \beta=1 \text{ e.g.}$$



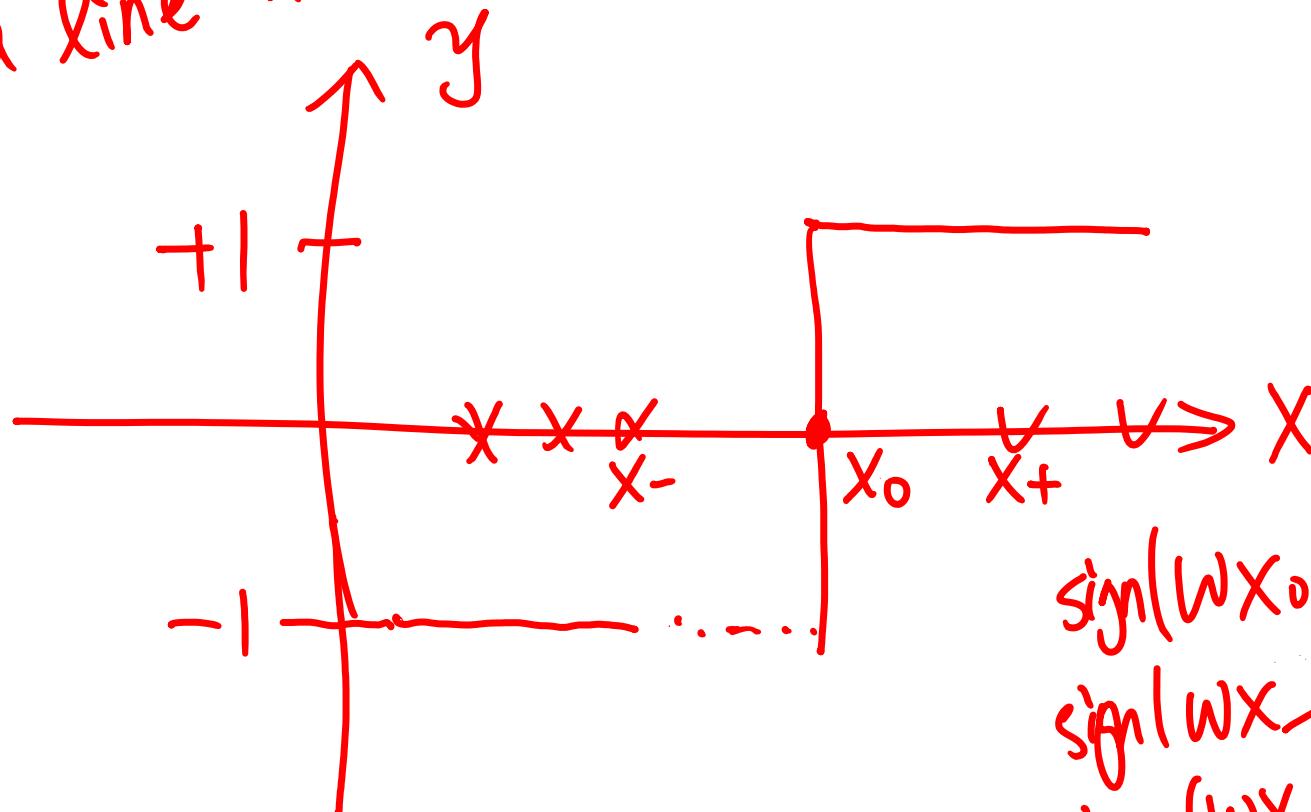
$$y = w^T x + w_0$$

$$-y + w^T x + w_0 = 0$$

a line in (x, y) space
 $P+1$

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

Decision Boundary is
a line in \mathbb{R}^P



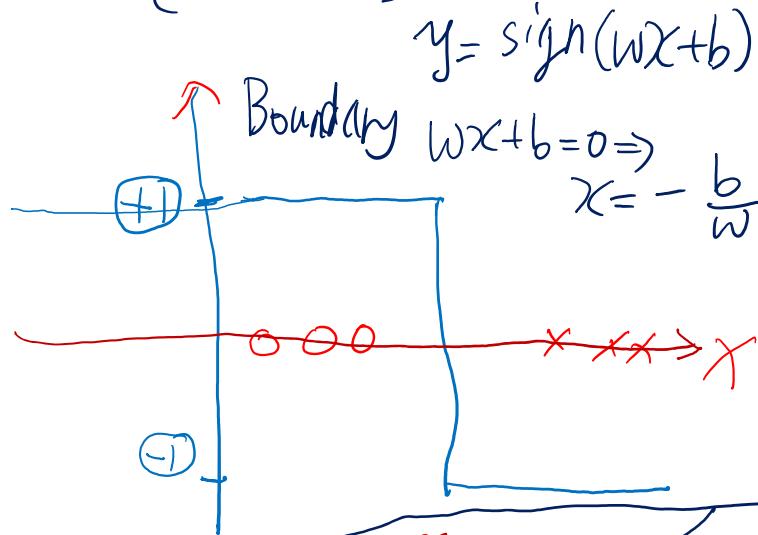
$X \in \mathbb{R}^P$, $P=1$ e.g.
 (x, y)

$$\begin{aligned}\text{sign}(w x_0 + b) &= 0 \\ \text{sign}(w x_- + b) &= -1 \\ \text{sign}(w x_+ + b) &= +1\end{aligned}$$

Binary Classification

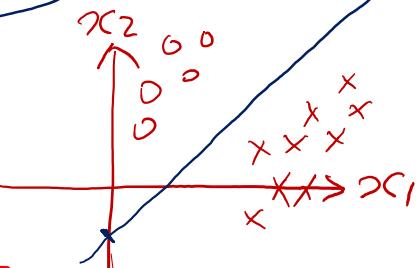
$$y \in \{-1, 1\}$$

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



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

Boundary $w^T x + b = 0$



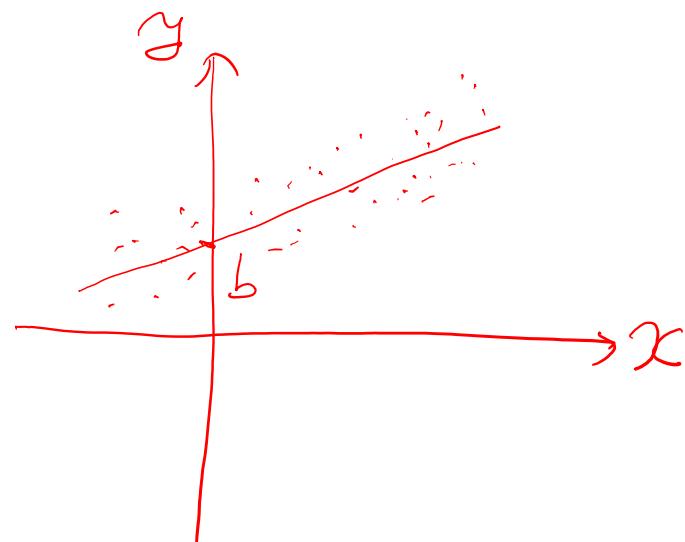
if $\Rightarrow p \text{ dim } [x \in \mathbb{R}^p]$

Boundary: hyperplane

Regression

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

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



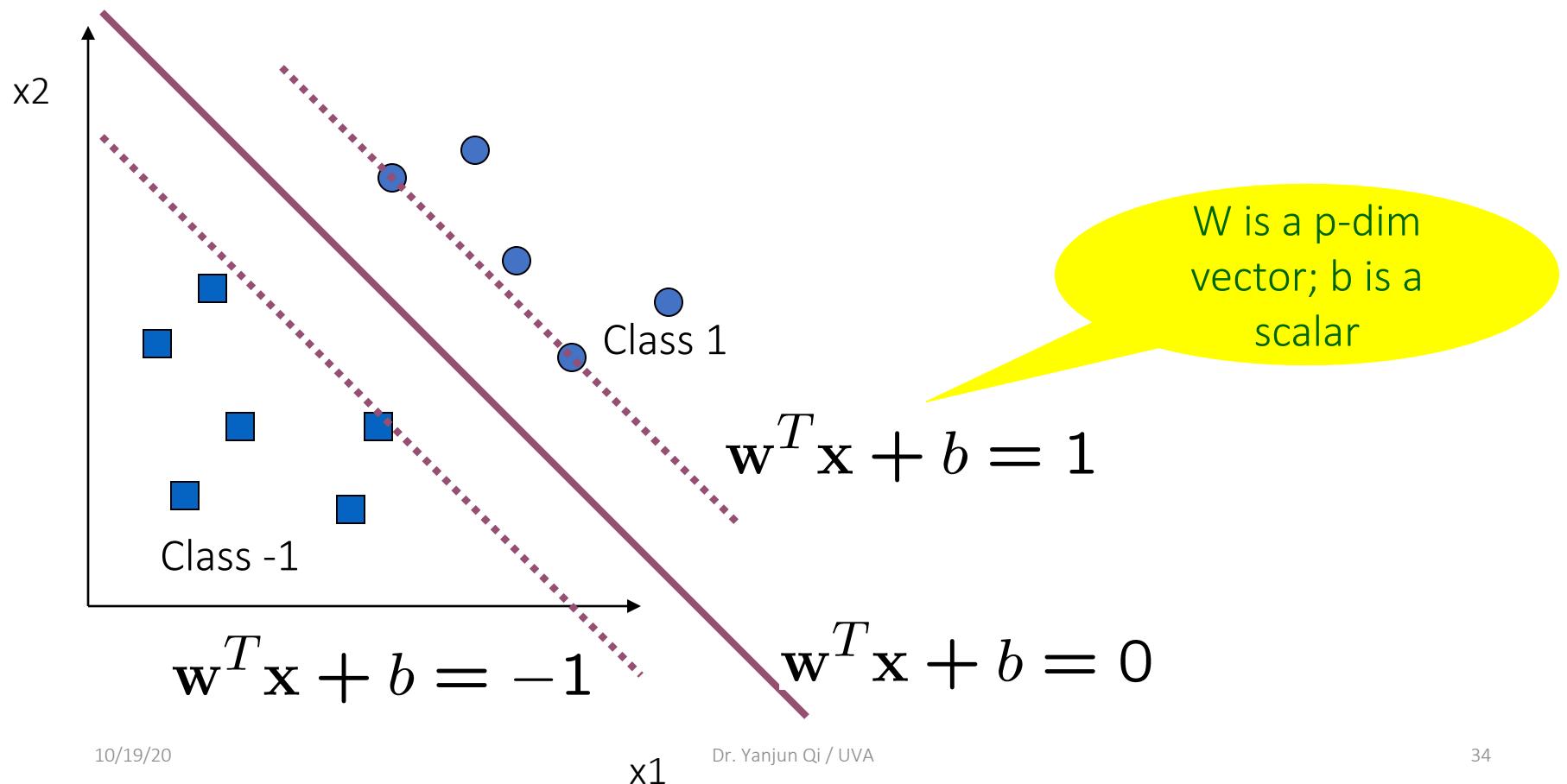
$$y = wx + b$$

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

$$y - wx - b = 0$$

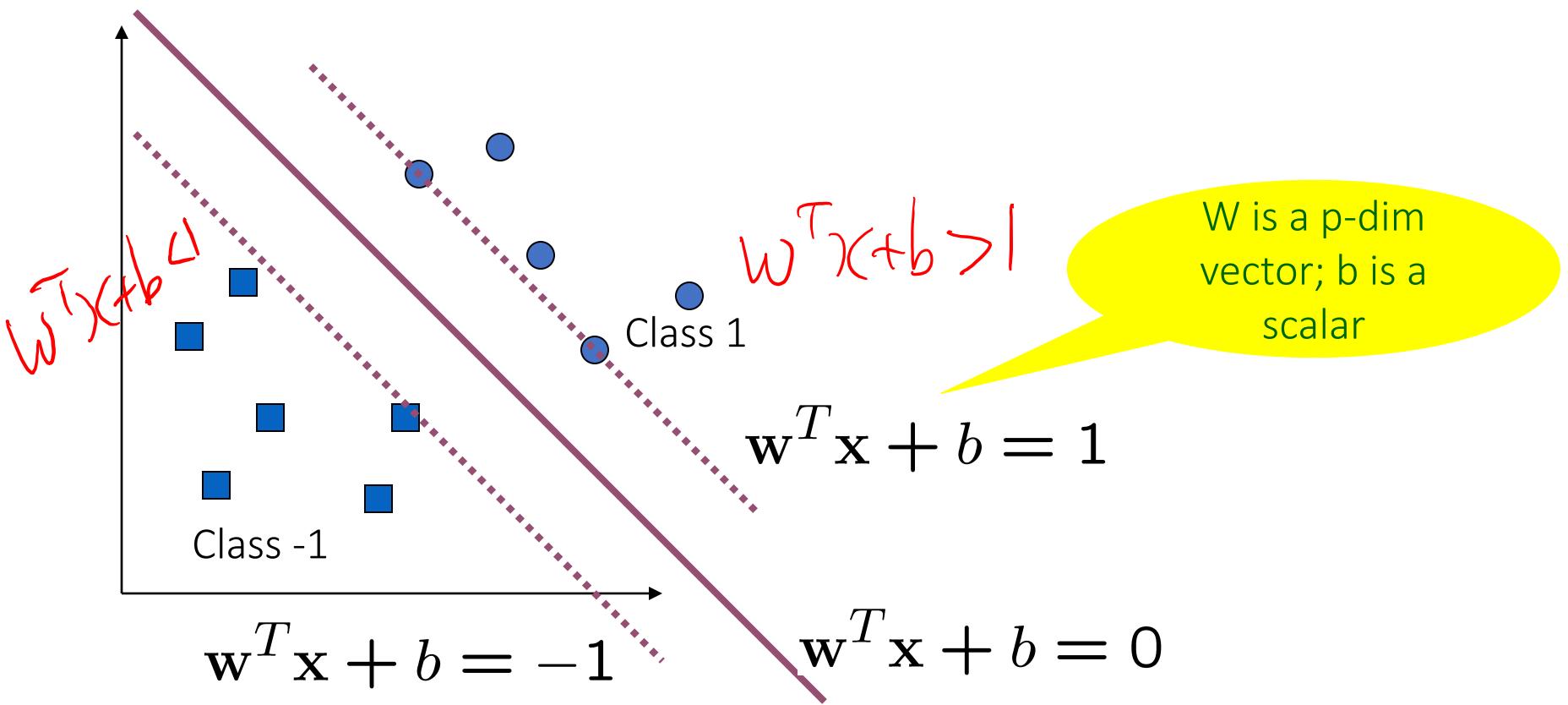
aka. linear line in 2-D space $\{x, y\}$

Max-margin & Decision Boundary



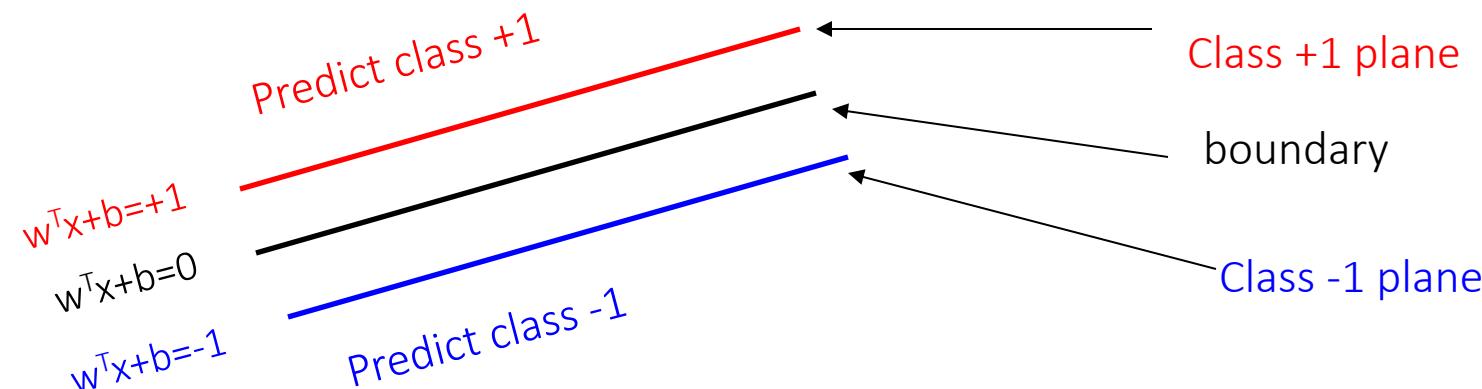
Max-margin & Decision Boundary $\mathbf{w}^T \mathbf{x} + b = 0$

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



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

Specifying a max margin classifier

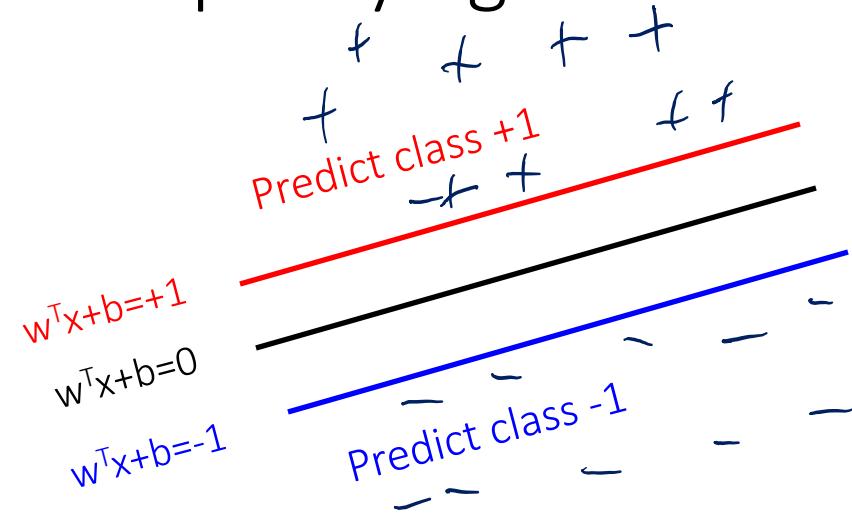


Classify as +1 if $w^T x + b \geq 1$

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

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

Classify as -1 if

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

Undefined if

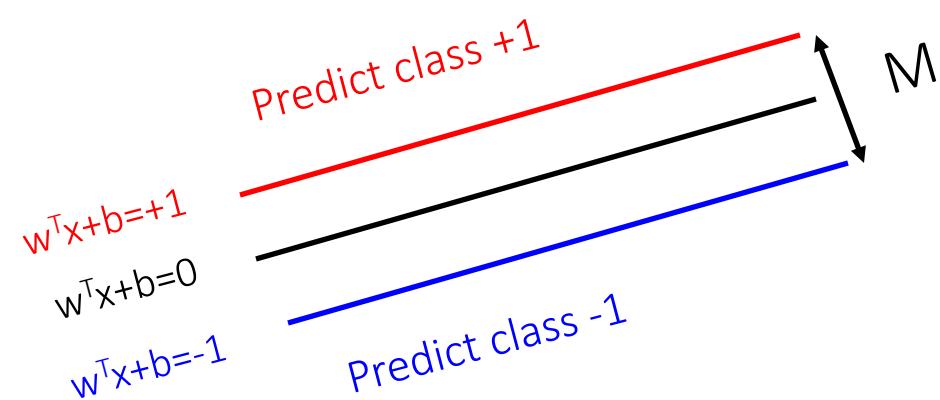
$$-1 < w^T x + b < 1$$

Is the linear separation assumption realistic?

We will deal with this shortly, but lets assume it for now

Now assuming such lines exist in our train

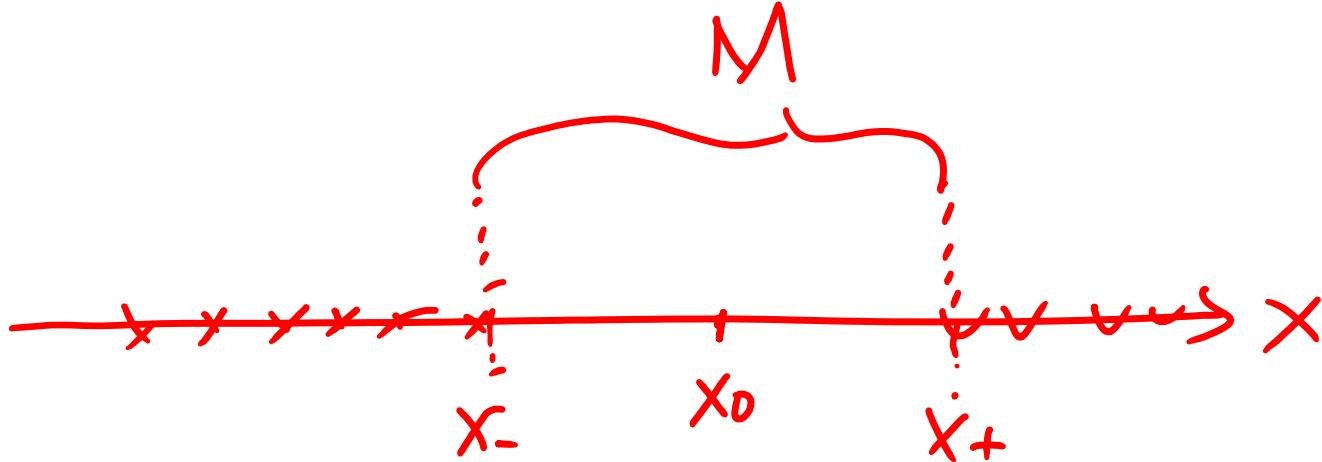
Maximizing the margin



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

- Let us 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 observations

See Concrete
derivations of M in
Extra slides



$$P=1$$

$$\Rightarrow M = |x_+ - x_-|$$

$$= \frac{2}{w}$$

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

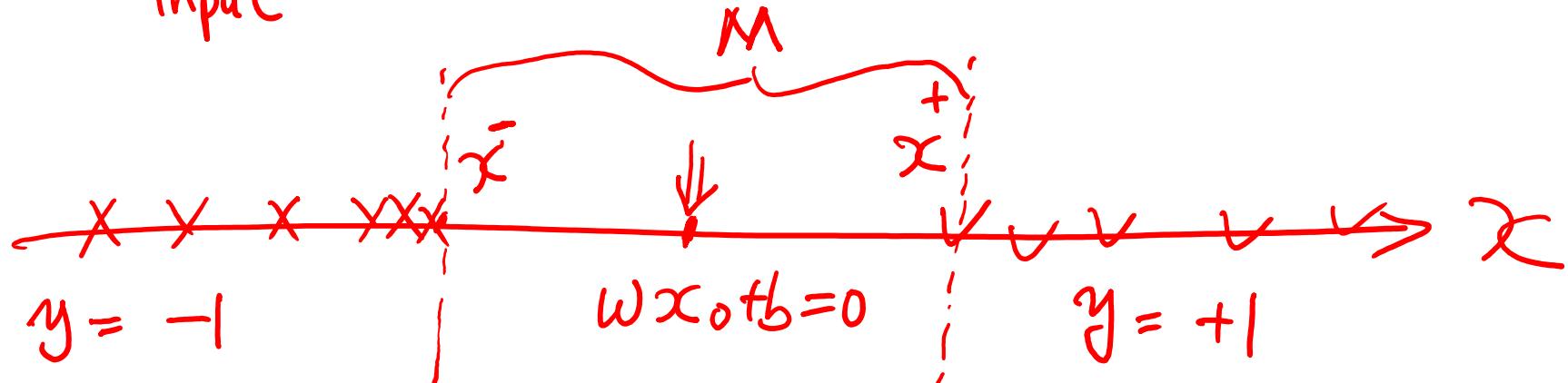
$$\begin{cases} w x_- + b = -1 \\ w x_+ + b = 1 \end{cases}$$



$$w(x_+ - x_-) = 2$$

$$\Rightarrow \max M \Rightarrow \min w^T w$$

If 1-D $\{x\} \rightarrow \{+1, -1\}$
Input

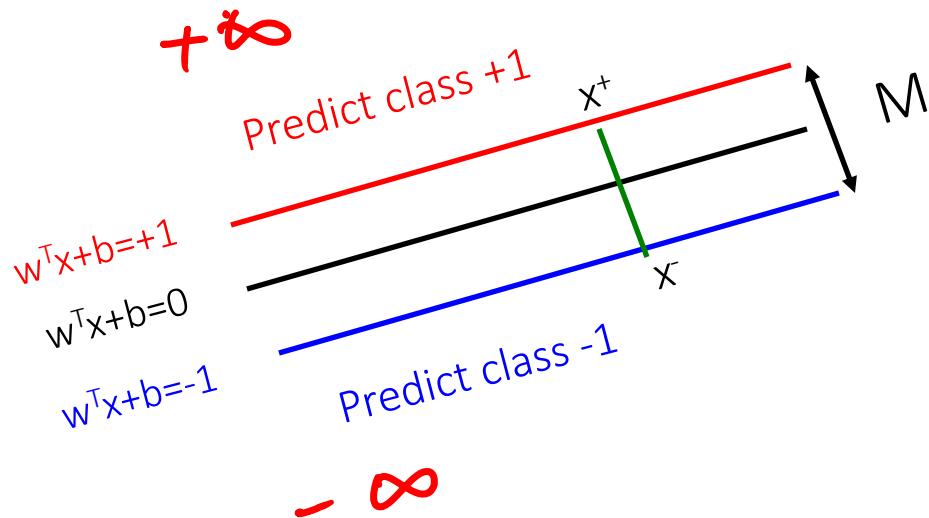


$$\Rightarrow \begin{cases} wx^- + b = -1 \\ wx^+ + b = 1 \end{cases} \Rightarrow w(x^+ - x^-) = 2$$

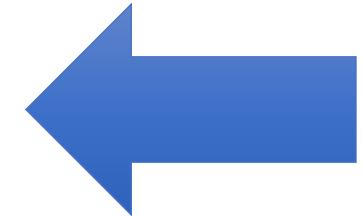
margin $M = |x^+ - x^-| = \frac{2}{|w|}$

(Interval when 1-D x)

Finding the optimal parameters



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

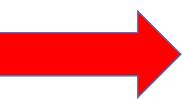


We can now search for the optimal parameters by finding a solution that:

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

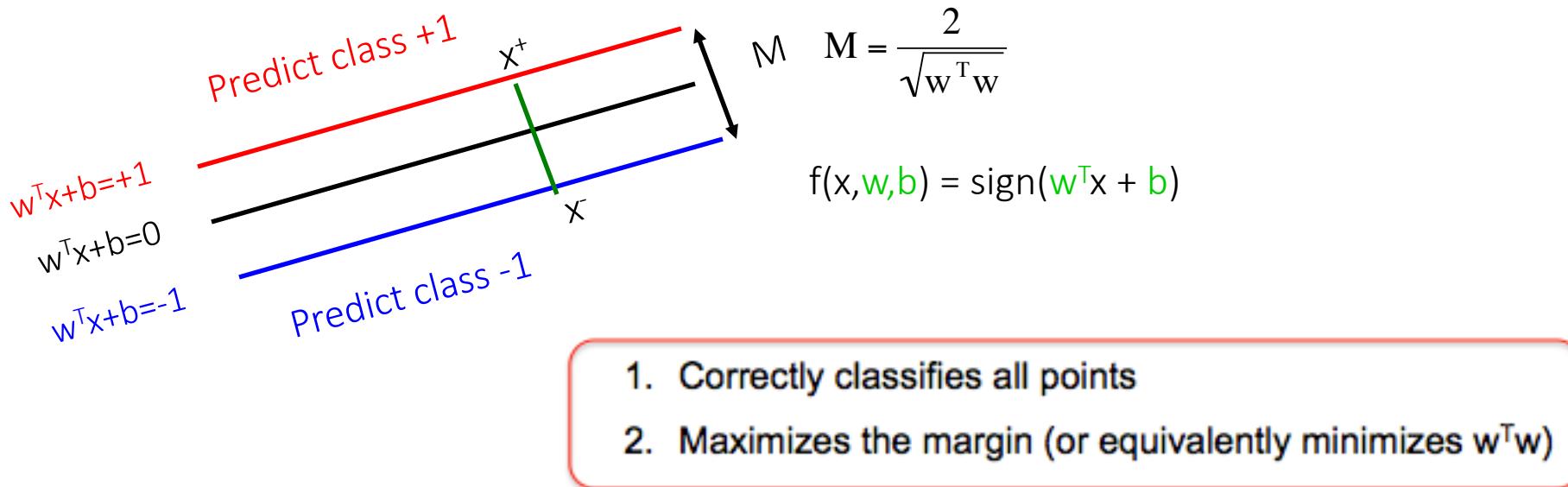
Several optimization methods can be used: Gradient descent, OR SMO (see extra slides)

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

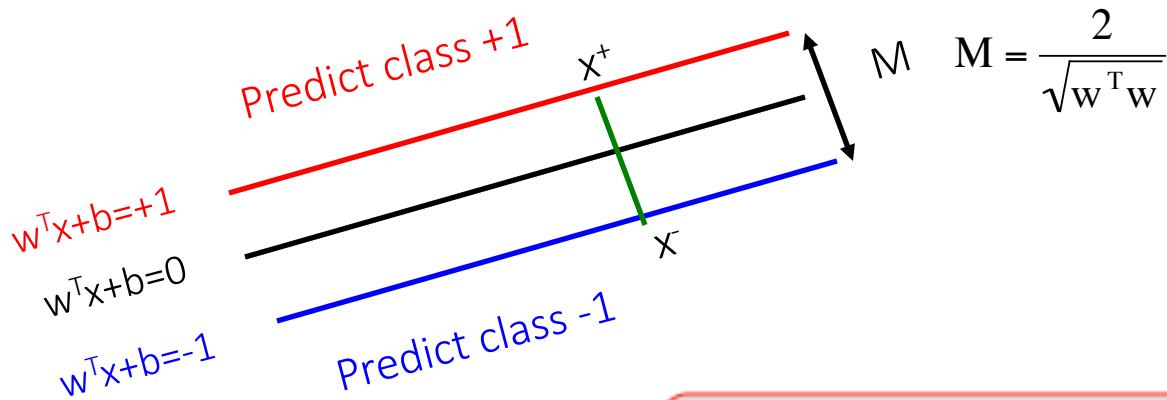
Optimization Step

i.e. learning optimal parameter for SVM



Optimization Step

i.e. learning optimal parameter for SVM



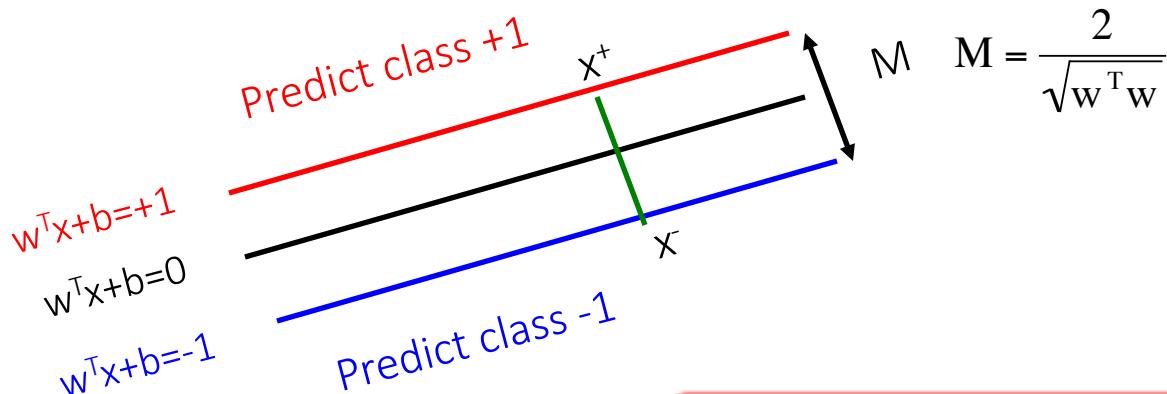
1. Correctly classifies all points
2. Maximizes the margin (or equivalently minimizes $w^T w$)

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

subject to the following constraints:

Optimization Step

i.e. learning optimal parameter for SVM



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

{ 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

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

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

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

$$y_i = 1$$

For all x in class -1

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

$$y_i = -1$$

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



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

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

$$\mathbf{x}_i^T \mathbf{w} = \\ \mathbf{w}^T \mathbf{x}_i$$

Optimization Reformulation

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

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

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

For all x in class -1

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

}

A total of n constraints if we have n input samples



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

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

$$y_i \in \{+1, -1\}$$

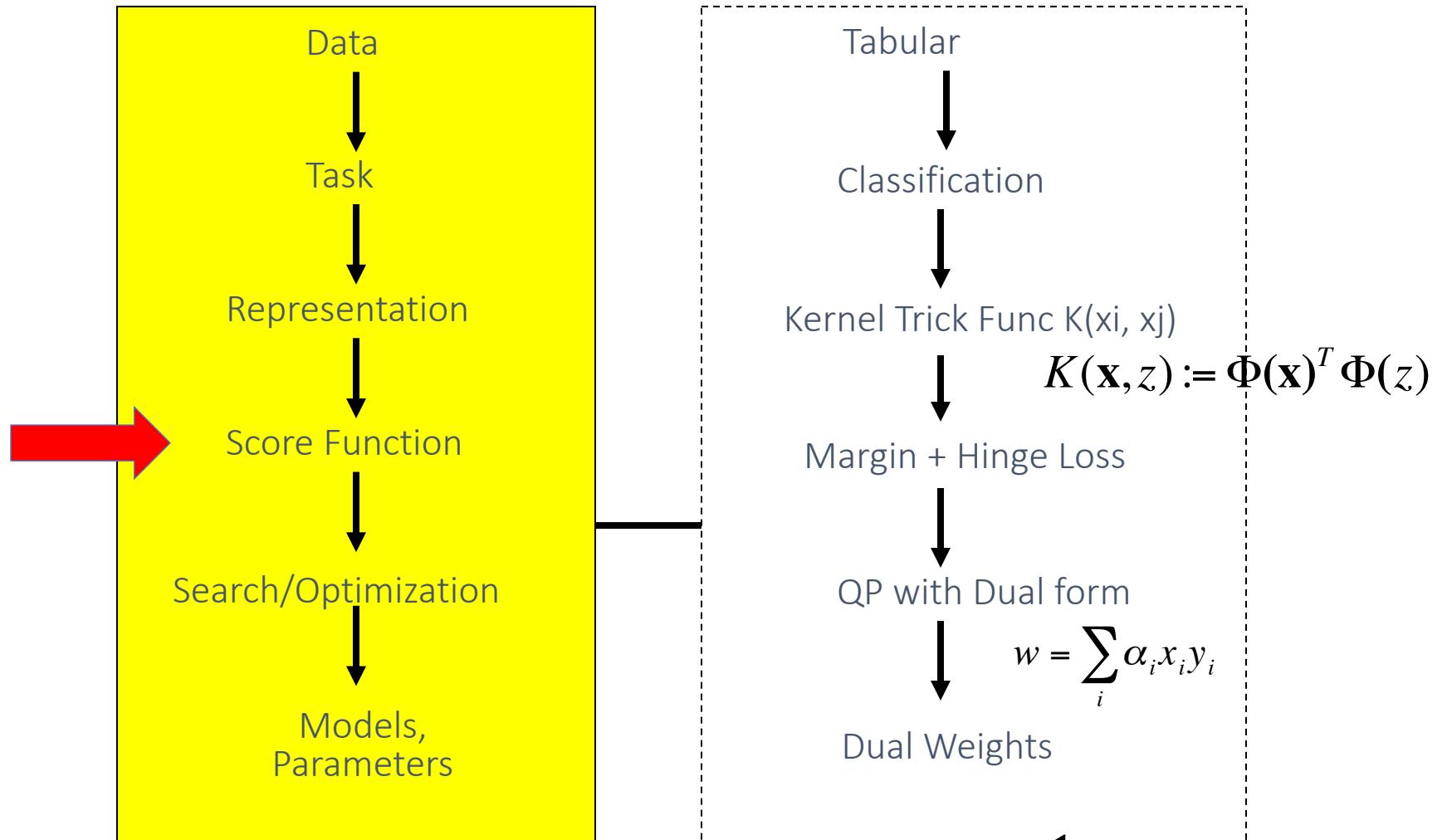
Quadratic Objective

$$\underbrace{w^T}_{|x|} \underbrace{x}_{|x|} \underbrace{p}_{|x|} \underbrace{p}_{|x|} \underbrace{x}_{|x|}$$

Quadratic programming i.e.,

- Quadratic objective
- Linear constraints

Today: Basic Support Vector Machine



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

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

$$\begin{aligned}
 & \max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\
 & \sum_i \alpha_i y_i = 0, \quad \alpha_i \geq 0 \quad \forall i
 \end{aligned}$$

What Next?

- ❑ 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)
 - ✓ Linerly Non-separable case (soft SVM)
 - ✓ Optimization with dual form
 - ✓ Nonlinear decision boundary
 - ✓ Practical Guide



Thank You

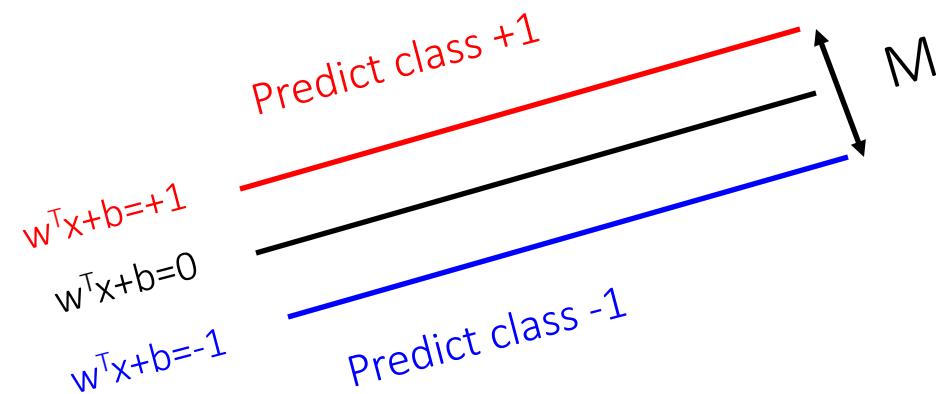
Thank you

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

EXTRA

How to define the width of the margin by M (EXTRA)

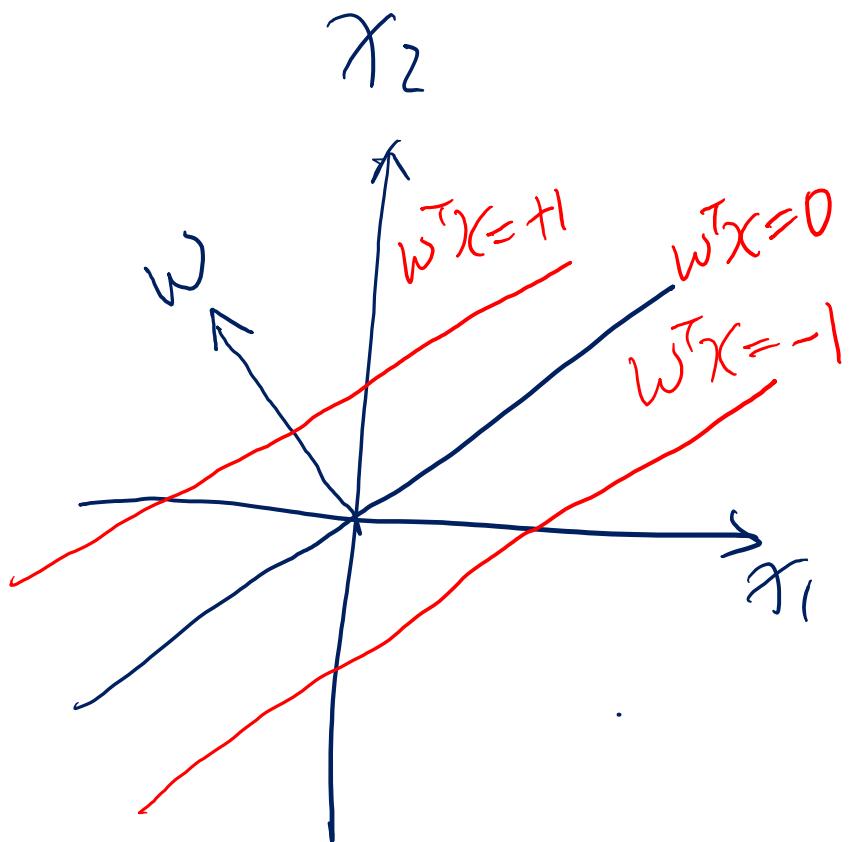


Classify as +1 if $w^T x + b \geq 1$
Classify as -1 if $w^T x + b \leq -1$
Undefined if $-1 < w^T x + b < 1$

- 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

Concrete derivations of M see Extra

[Classification]

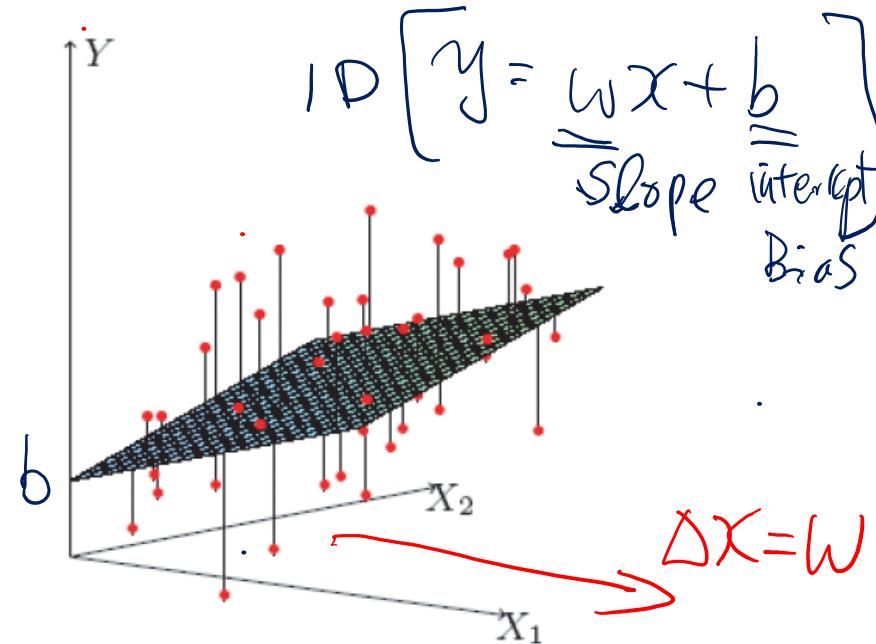


[Regression]

$$y = \vec{w}^T \vec{x} + b$$

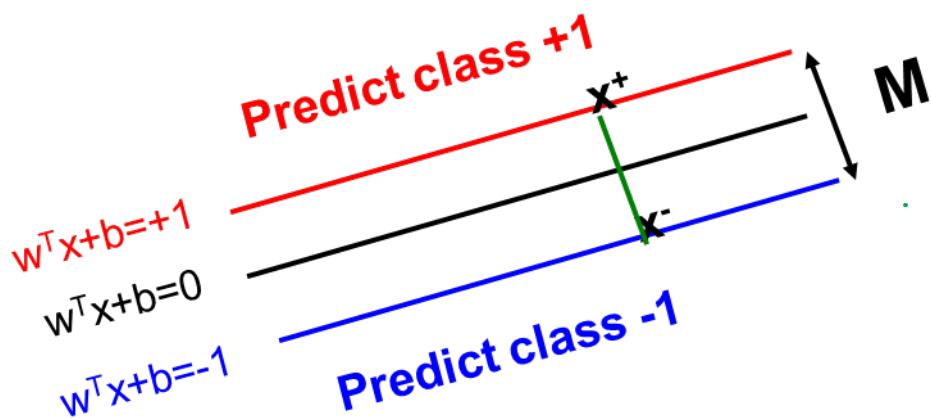
$$\Delta x = \frac{\partial y}{\partial x} = \vec{w}$$

slope



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

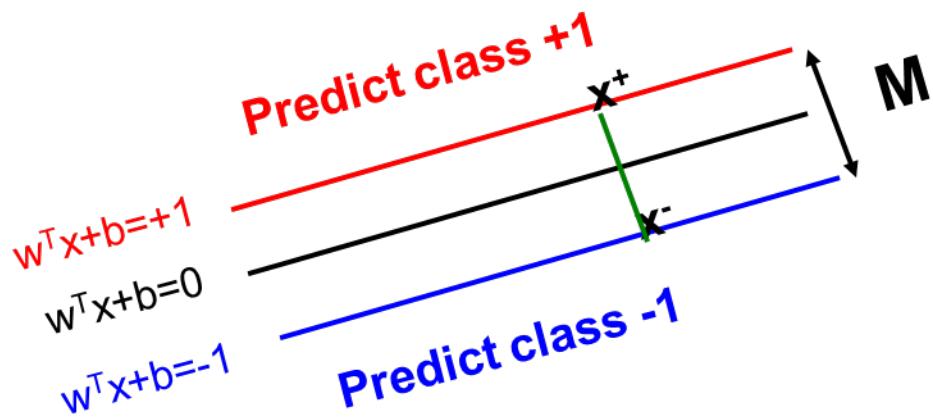
Margin M



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

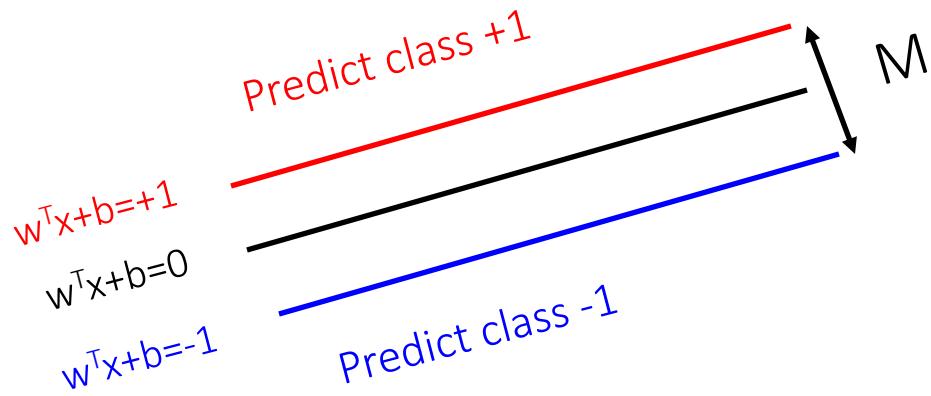
$$M = |x^+ - x^-| \quad \begin{array}{l} \text{length of} \\ \text{Vector } (x^+ - x^-) \end{array}$$

⇒ How to represent $(x^+ - x^-)$???



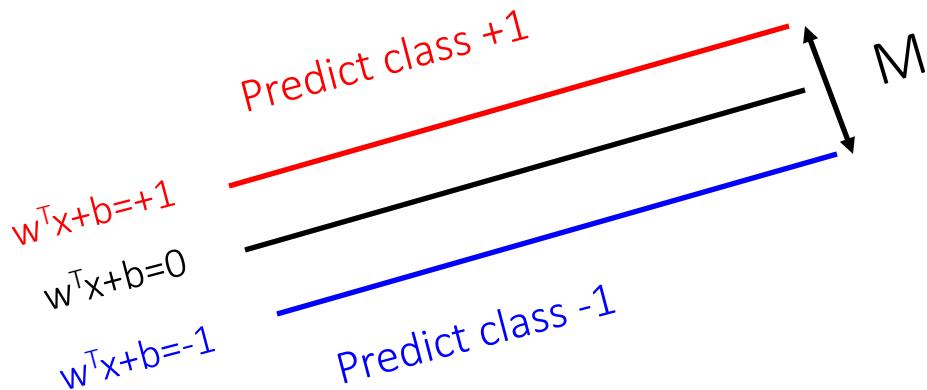
- $w^T x^+ + b = +1$
- $w^T x^- + b = -1$
- $M = |x^+ - x^-| = ?$

Maximizing the margin: observation-1



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

Maximizing the margin: observation-1

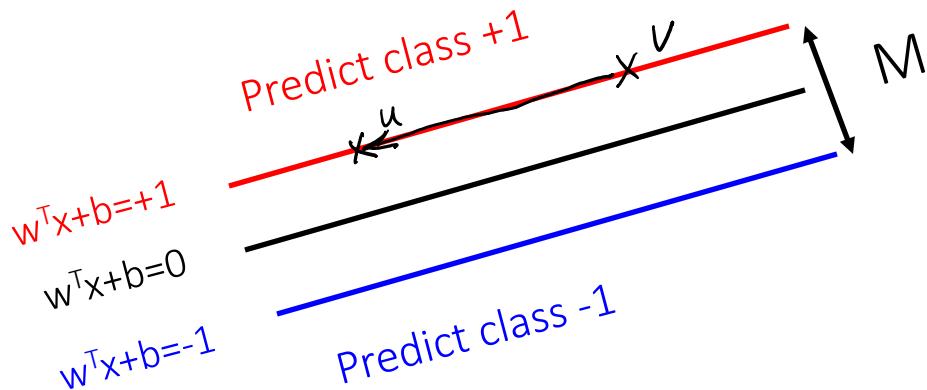


Classify as +1 if $w^T x + b \geq 1$
Classify as -1 if $w^T x + b \leq -1$
Undefined if $-1 < w^T x + b < 1$

- 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

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

$\Rightarrow w$ 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 $w^T(u-v) = 0$

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

$$x^T y = 0$$

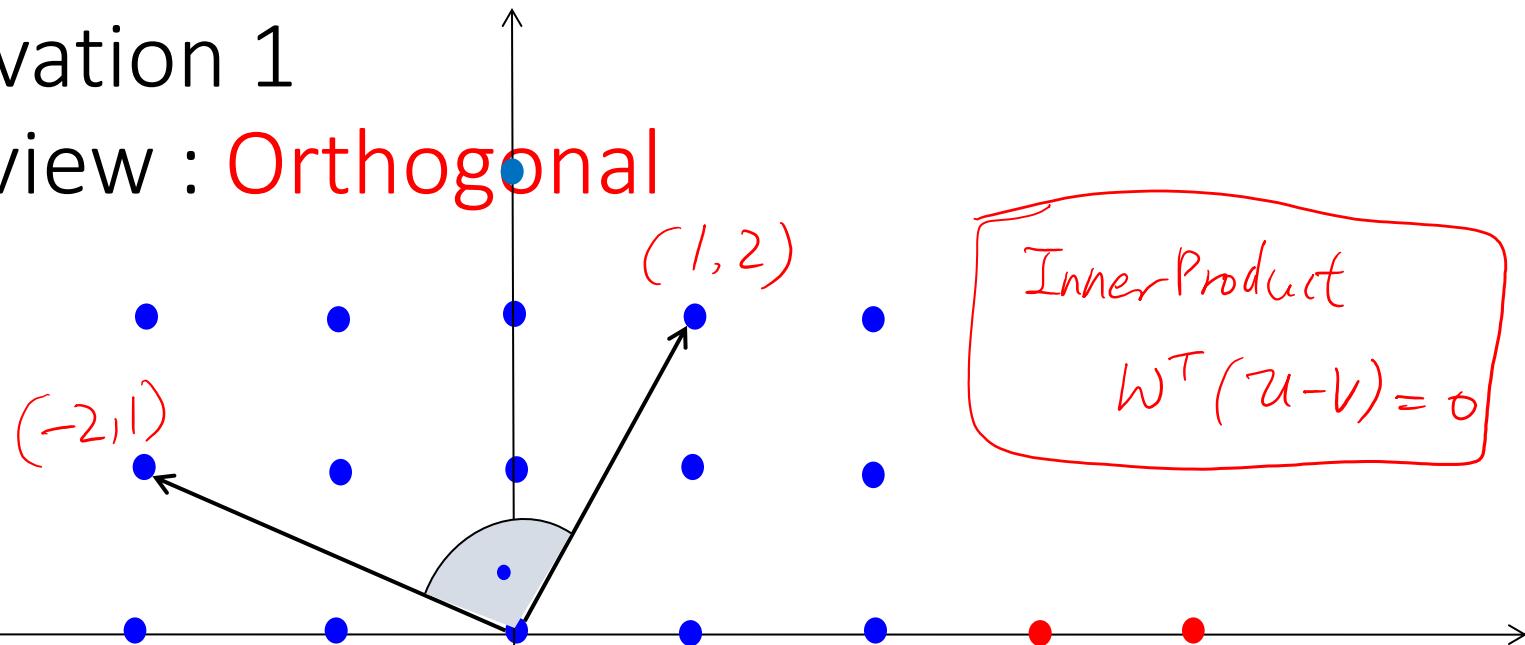
The square root of the product of a vector with itself,

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

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

Observation 1

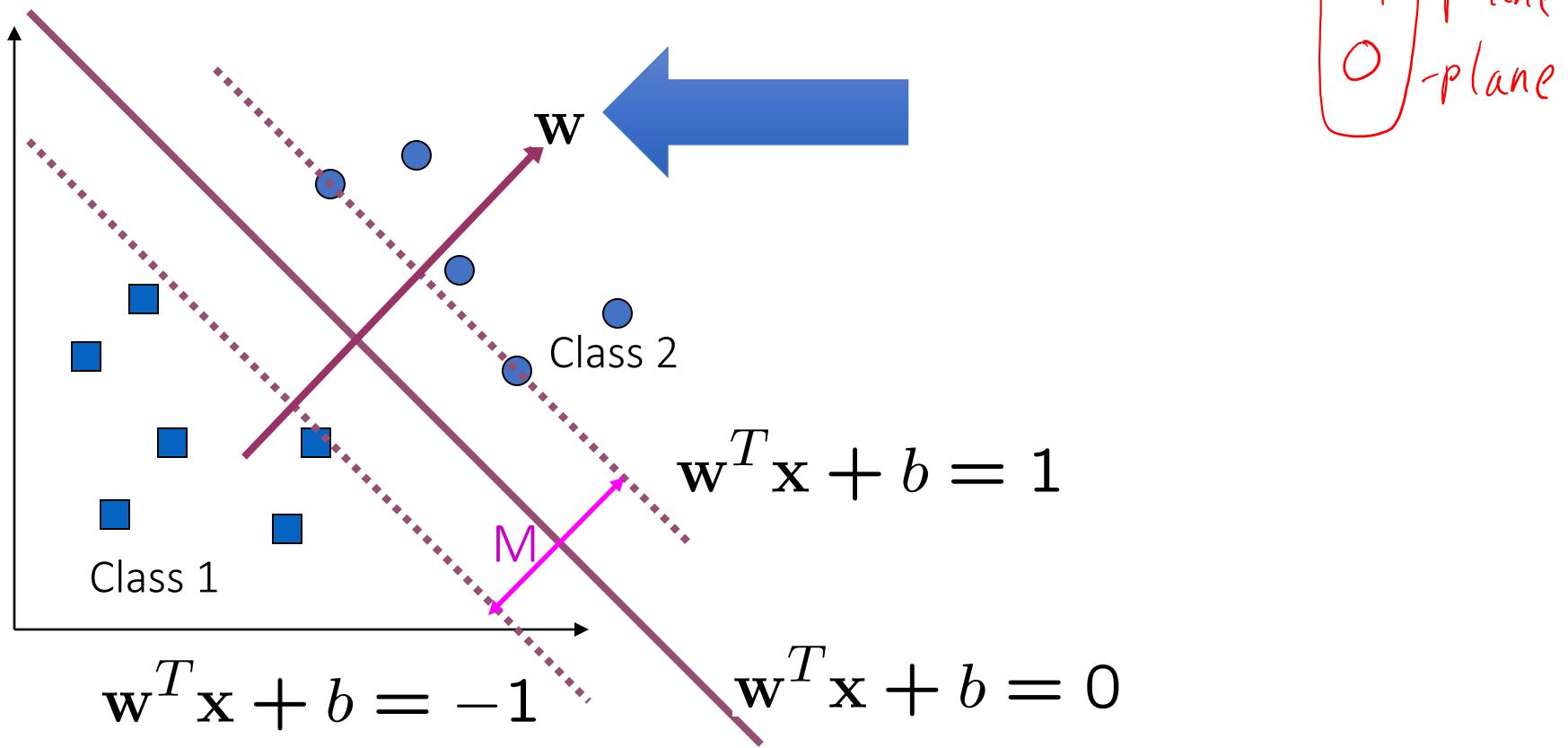
→ Review : Orthogonal



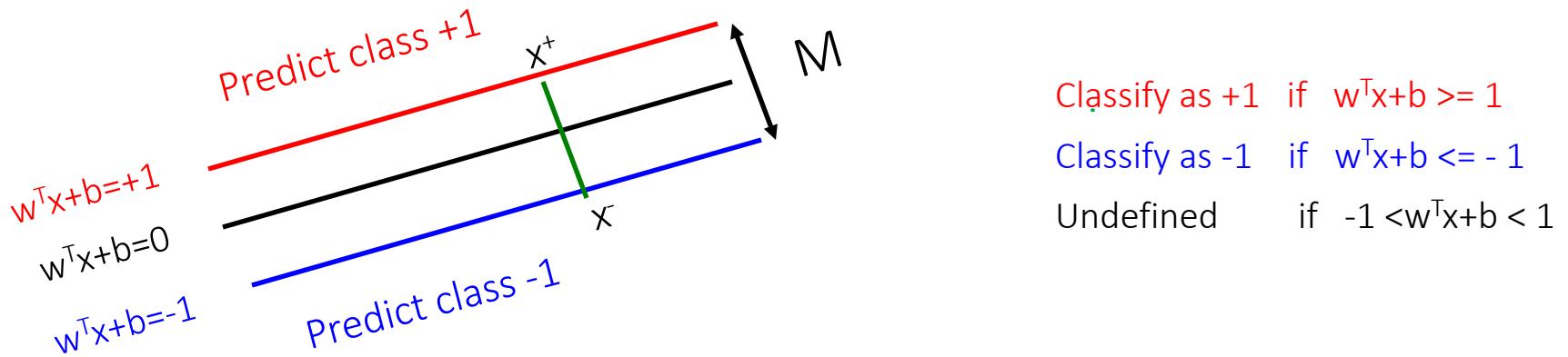
$$\begin{pmatrix} -2 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = 0$$

Maximizing the margin: observation-1

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



Maximizing the margin: observation-2

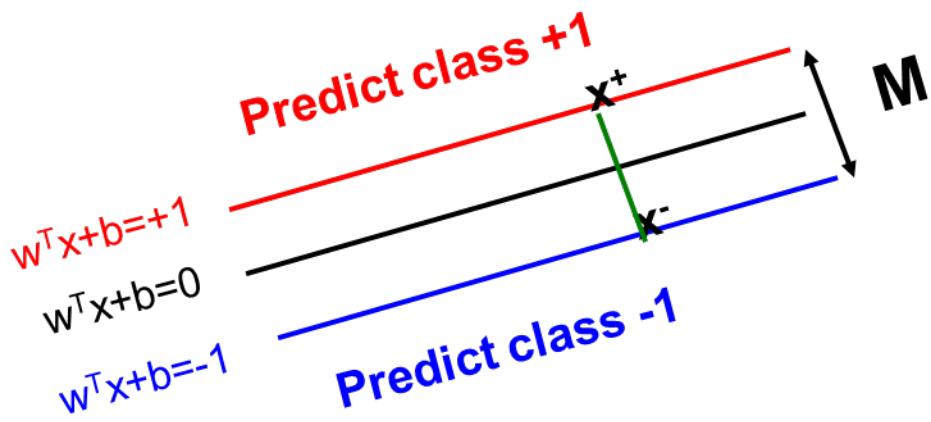


- Observation 1: the vector w is orthogonal to the +1 and -1 planes
- Observation 2: if x^+ is a point on the +1 plane and x^- is the closest point to x^+ on the -1 plane then

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

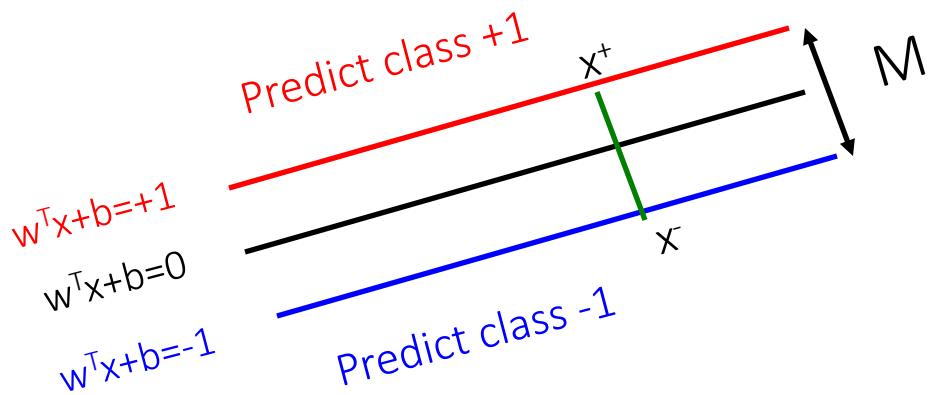
Since w is orthogonal to both planes we need to ‘travel’ some distance along w to get from x^+ to x^-

Putting it together



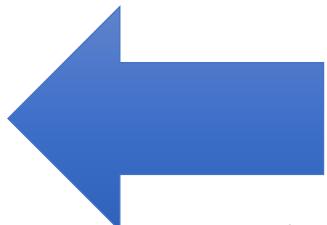
- $w^T x^+ + b = +1$
- $w^T x^- + b = -1$
- $M = | x^+ - x^- | = ?$
- $x^+ = \lambda w + x^-$

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



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

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

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

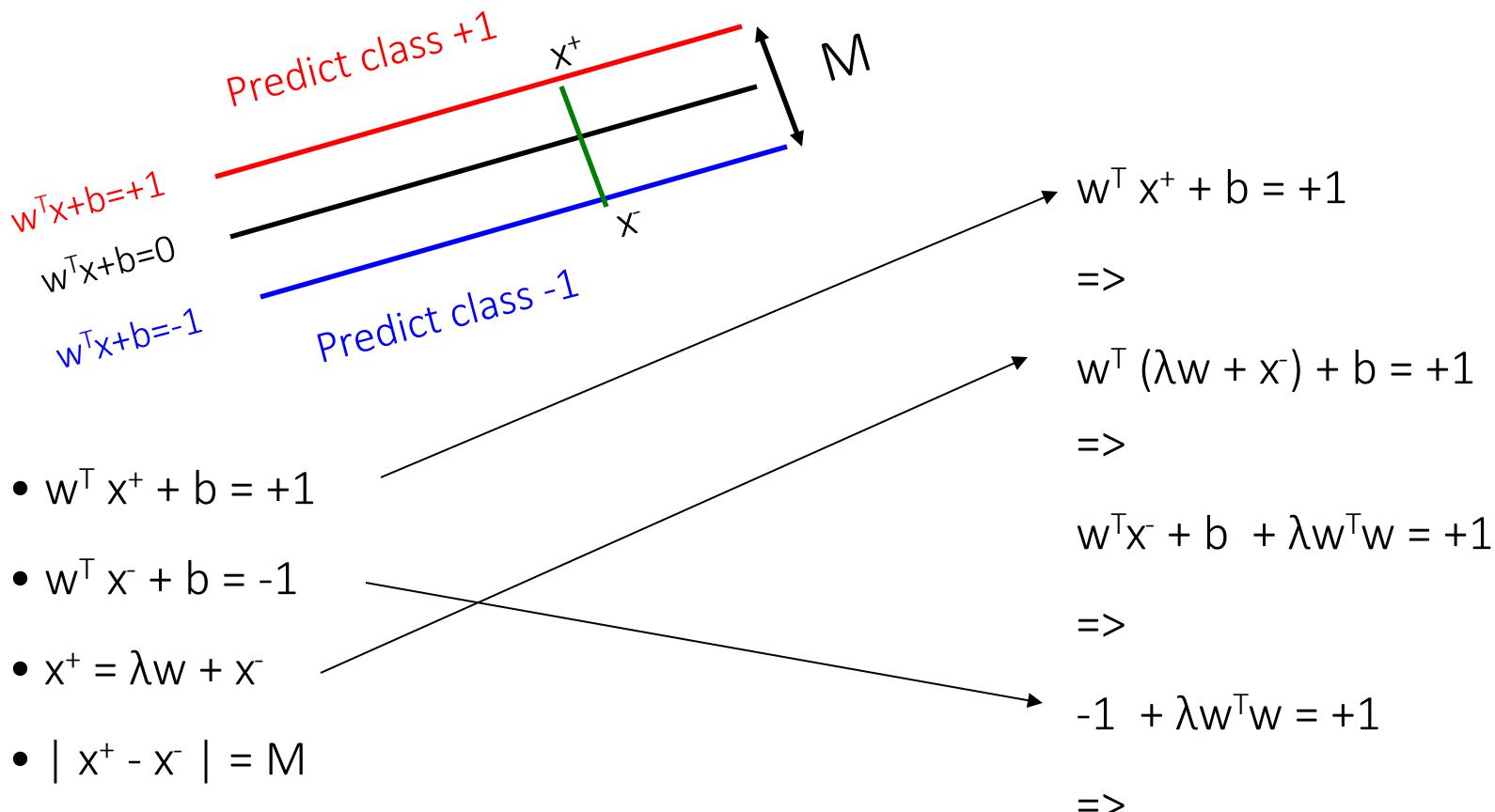
$$\lambda w^T w + w^T x^- + b = 1$$

$\underbrace{}_{\Rightarrow} -1$

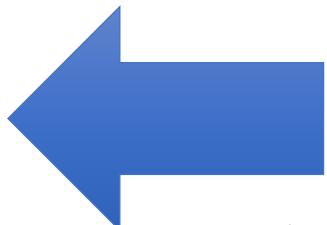
$$\lambda w^T w = 2$$

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

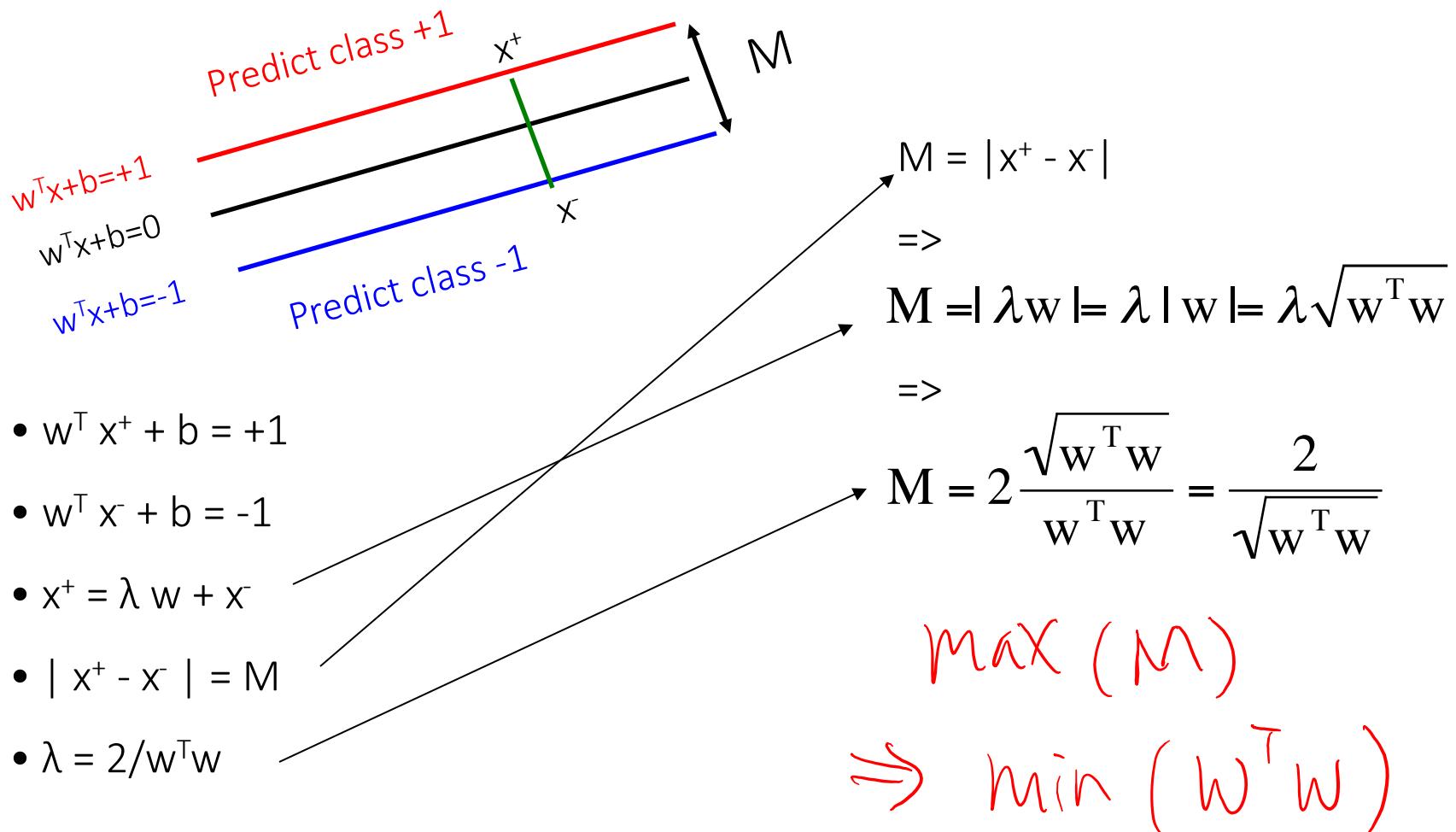
Putting it together



We can now define M in terms of w and b



Putting it together



We can now define M in terms of w and b