# Support Vector Machine

## Foundations of Data Analysis
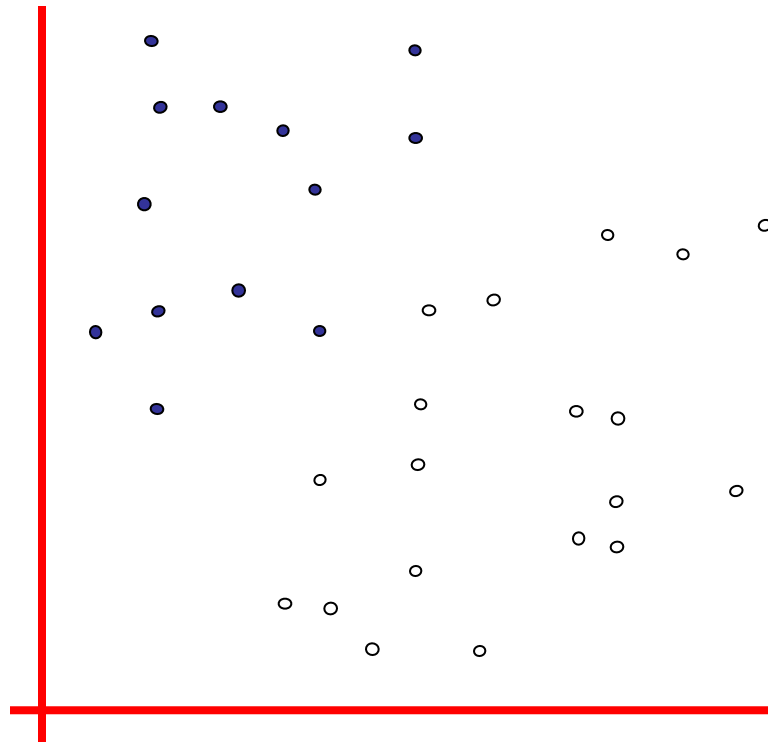## 04/14/2020

Slide credits to Andrew Moore: http://www.cs.cmu.edu/~awm/tutorials

# Linear Classifiers

$$f(x, w) = \text{sign}(w \cdot x)$$

• denotes +1

○ denotes -1

How would you classify this data?

# Linear Classifiers

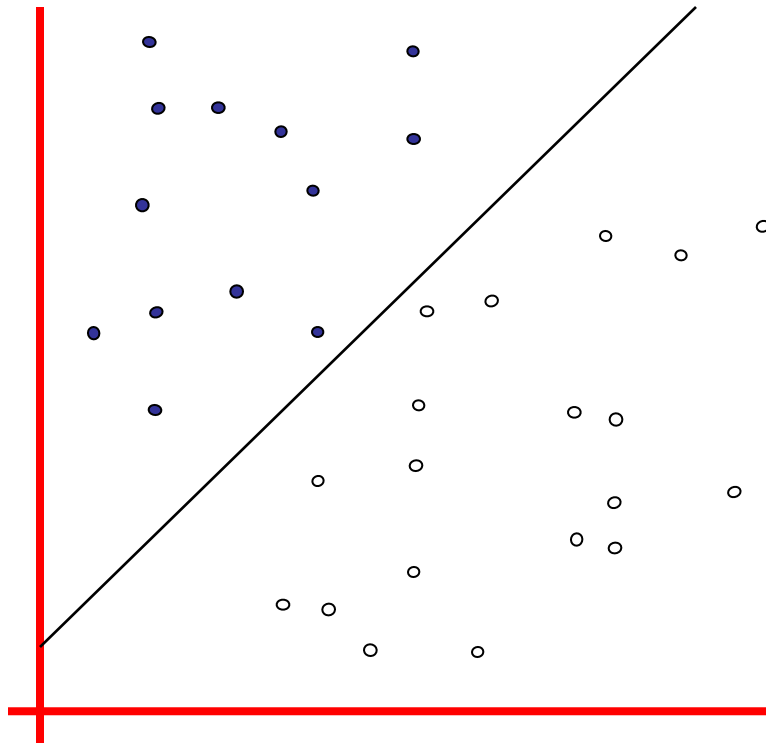$$f(x, w) = \text{sign}(w \cdot x)$$

• denotes +1

◦ denotes -1

How would you classify this data?

# Linear Classifiers

$$f(x, w) = \text{sign}(w \cdot x)$$

- denotes +1
- denotes -1

How would you classify this data?

# Linear Classifiers

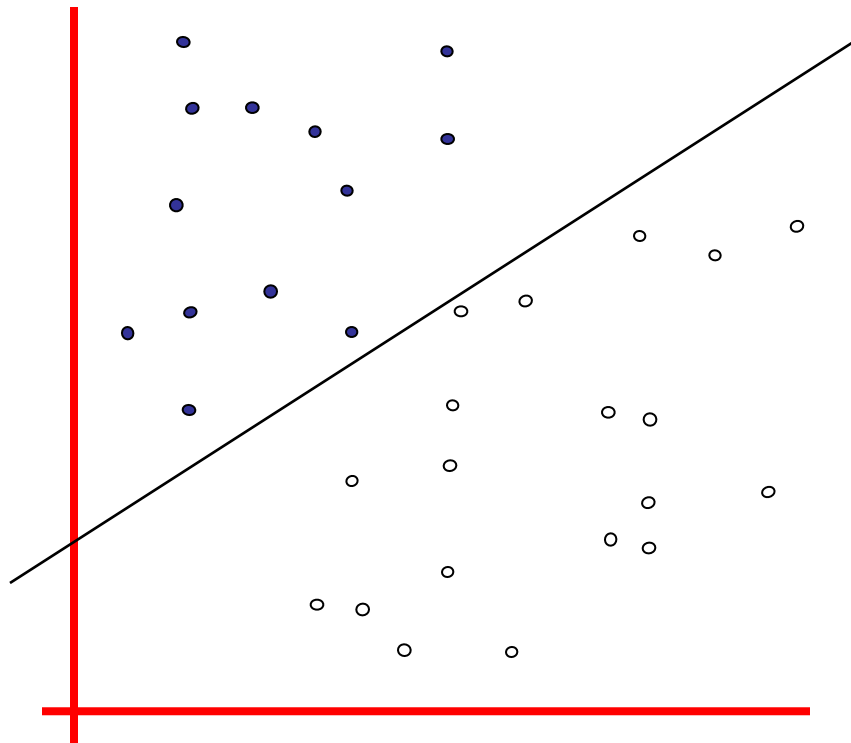$$f(x, w) = \text{sign}(w \cdot x)$$

- denotes +1
- denotes -1

How would you classify this data?

# Linear Classifiers

$$f(x, w) = \text{sign}(w \cdot x)$$

• denotes +1

○  denotes -1

How would you classify this data?

# Linear Classifiers

$$f(x, w) = \text{sign}(w \cdot x)$$
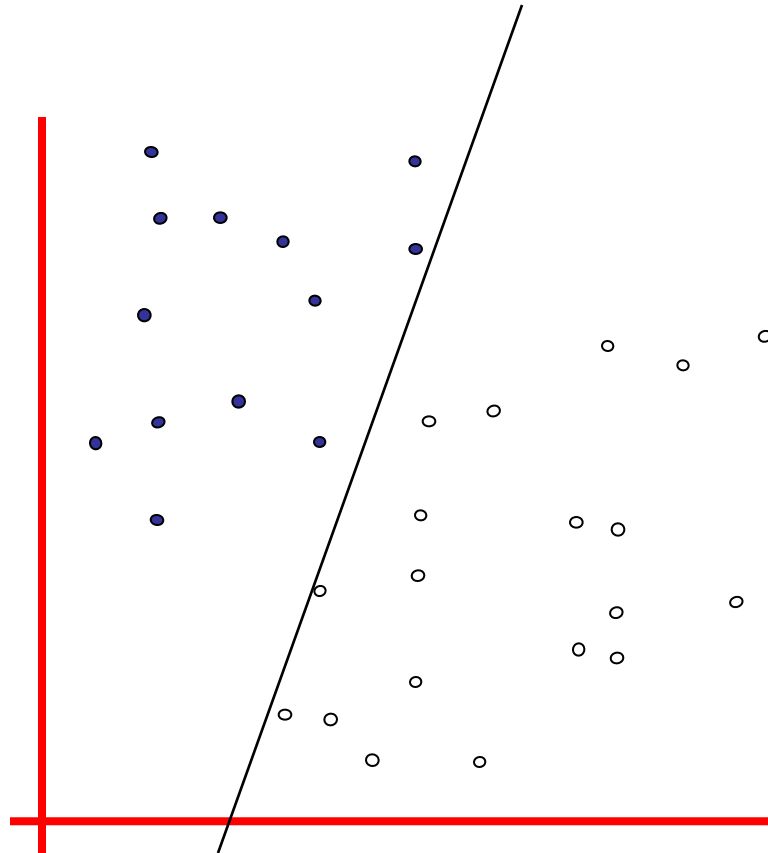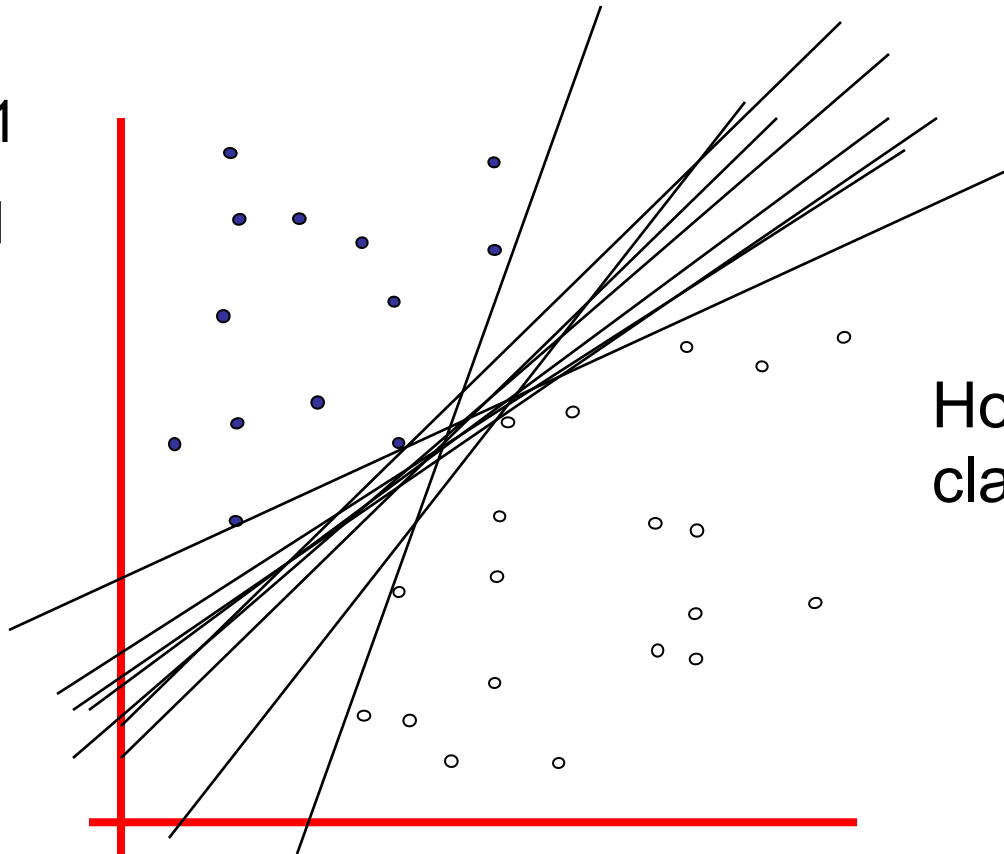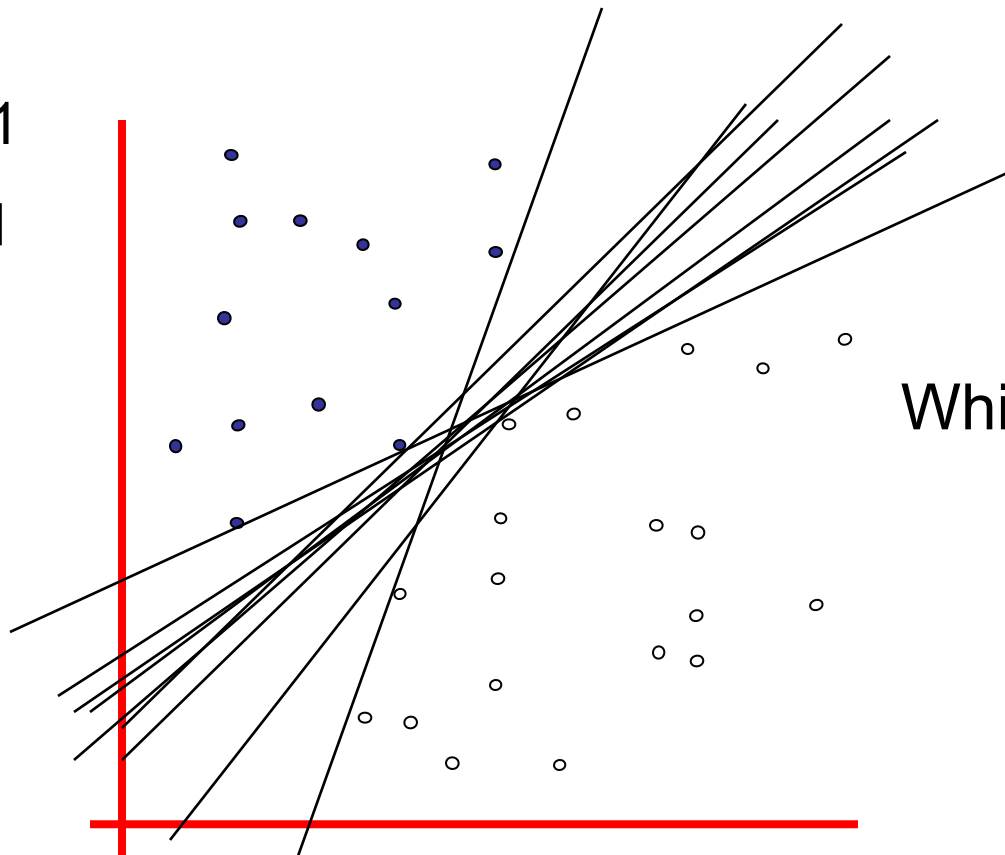
- denotes +1
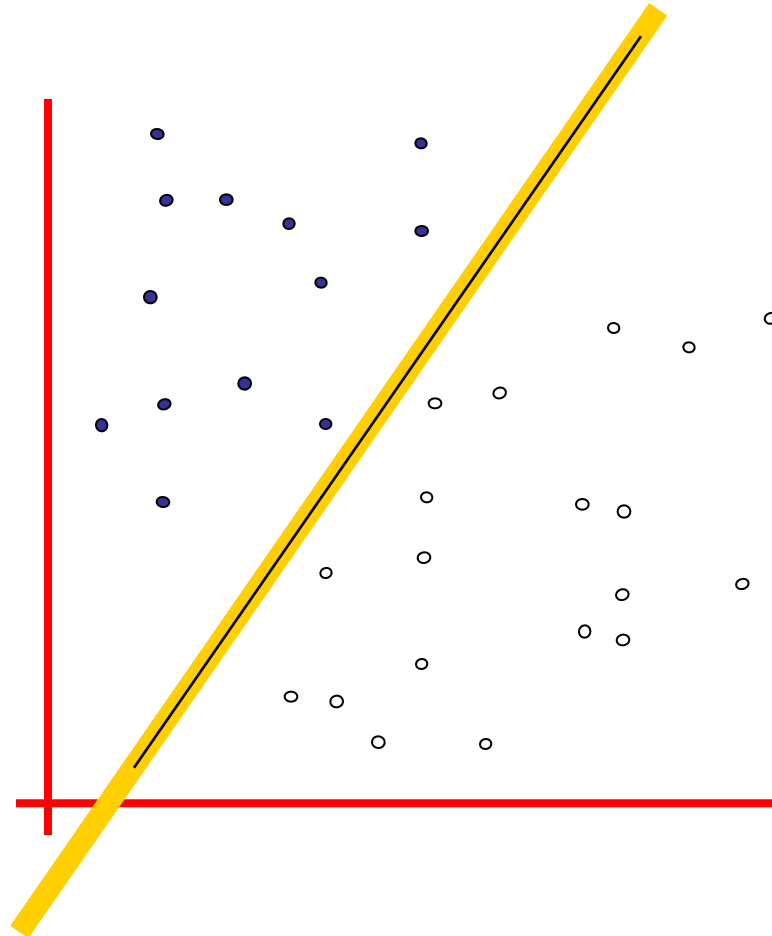- denotes -1



Which is the best?

# Linear Classifiers

$$f(x, w) = \text{sign}(w \cdot x)$$

• denotes +1

∘ denotes -1
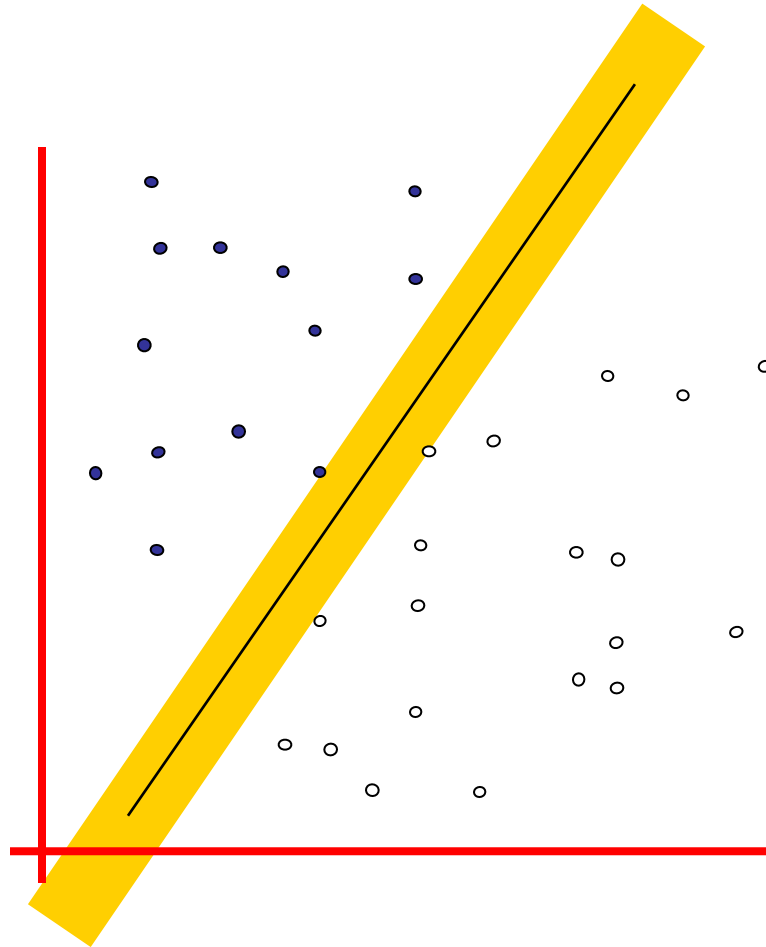
Define the margin of a linear classifier as the width that the boundary could be increased before hitting a datapoint.

# Maximum Margin

$$f(x, w) = \text{sign}(w \cdot x)$$

- denotes +1

- denotes -1

Maximum margin: the widest margin that maximally separates two data groups.
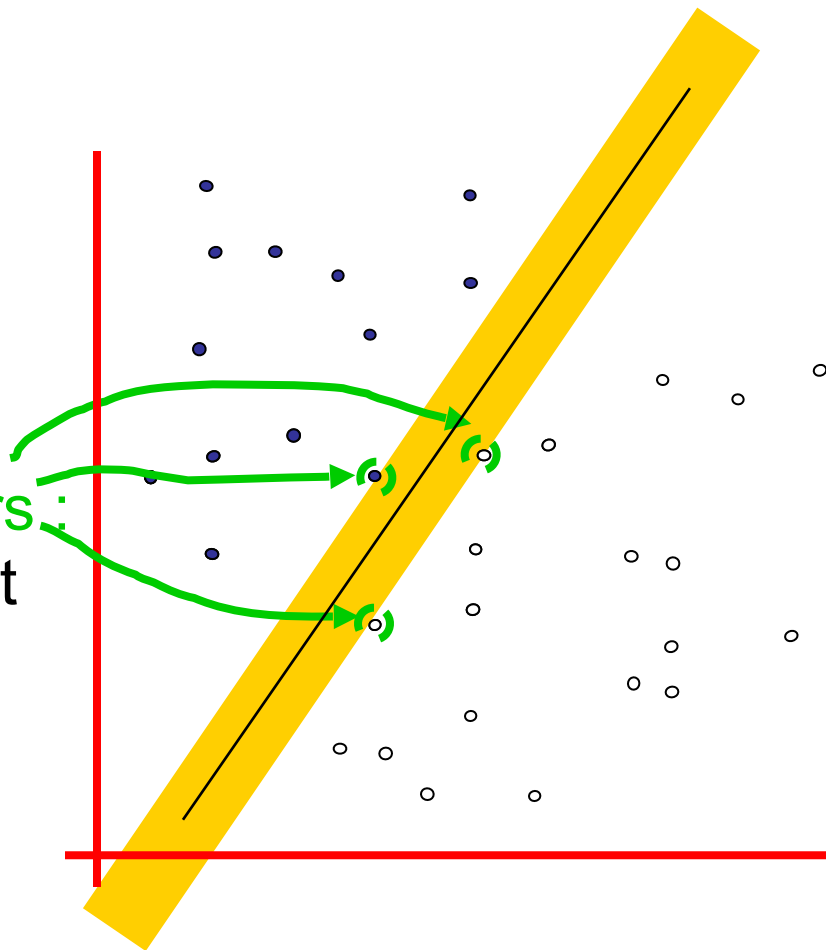
# Support Vector Machines

$$f(x, w) = \text{sign}(w \cdot x)$$

• denotes +1

○ denotes -1

Support Vectors:
data points that
the margin
pushes up
against.

Maximum margin
linear classifier is
the simplest kind of
SVM (LinearSVM)

# Specifying a line and margin

"Predict Class = +1" zone

Plus-Plane

Classifier Boundary

Minus-Plane

"Predict Class = -1" zone

- How do we represent this mathematically?
- …in m input dimensions?

# Specifying a line and margin

"Predict Class = +1" zone

Plus-Plane

Classifier Boundary

Minus-Plane

"Predict Class = -1" zone

| Class | +1 | if | $w \cdot x >= 1$ |
|---|---|---|---|
| | -1 | if | $w \cdot x <= -1$ |
| | embarrassing points | if | $-1 < w \cdot x < 1$ |

# Computing the Margin Width

"Predict Class = +1" zone

wx=1

wx=0

wx=-1

"Predict Class = -1" zone

M = Margin Width

How to compute M?

# Computing the Margin Width
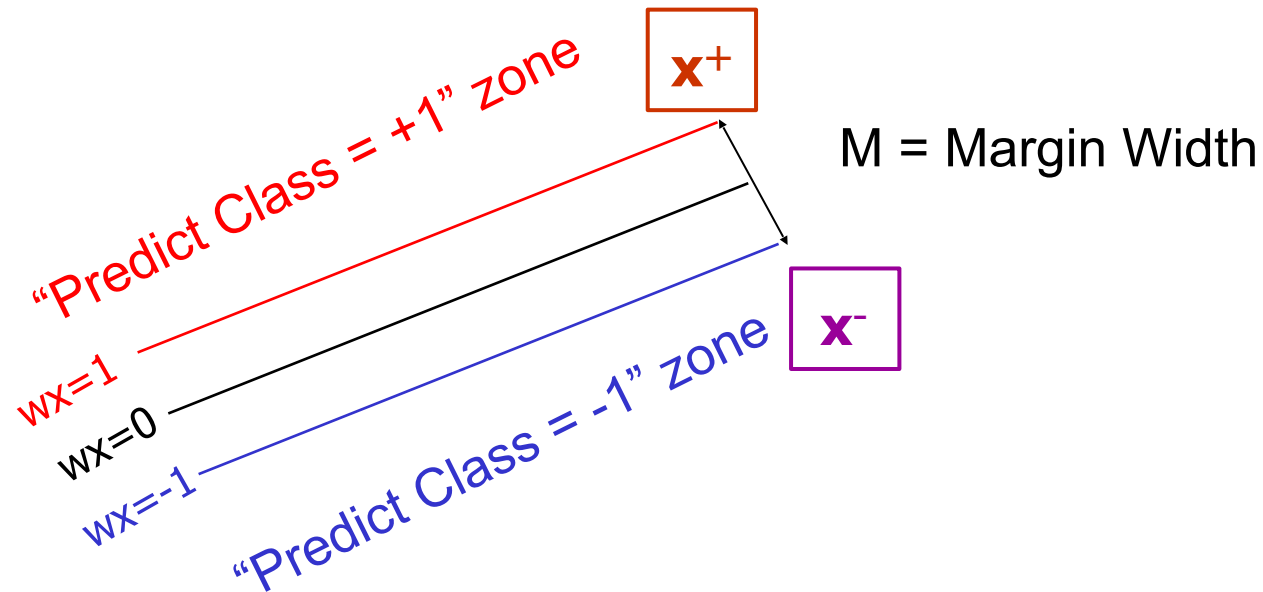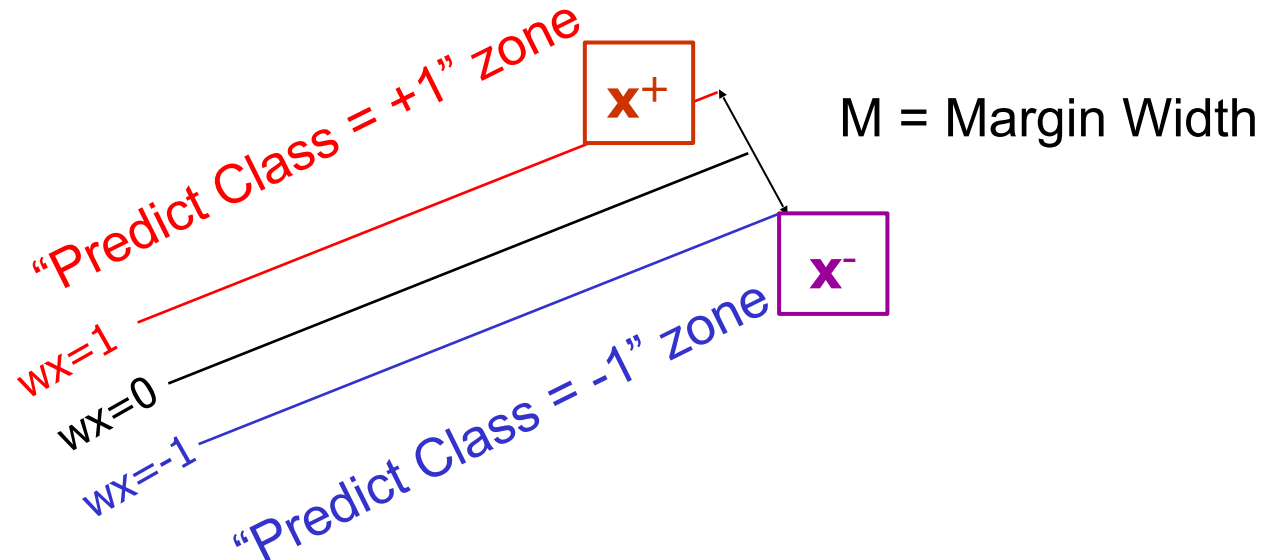


- Let **x⁻** be any point on the minus plane
- Let **x⁺** be the closest plus-plane-point to **x⁻**
- Claim: **x⁺** = **x⁻** + λ **w**  for some value of λ. Why?

# Computing the Margin Width



"Predict Class = +1" zone

x$^+$

M = Margin Width

wx=1

wx=0

wx=-1

x$^-$

"Predict Class = -1" zone

- Let **x**$^-$ be any point on the minus plane
- Let **x**$^+$ be the closest plus-plane-point to **x**$^-$
- Claim: **x**$^+$ = **x**$^-$ + λ **w**  for some value of λ. Why?

  The line from **x**$^-$ to **x**$^+$ is perpendicular to the planes. So to get from **x**$^-$ to **x**$^+$ travel some distance in direction w.

# Computing the Margin Width



**What we know:**

- $\mathbf{w} \cdot \mathbf{x}^+ = +1$
- $\mathbf{w} \cdot \mathbf{x}^- = -1$
- $\mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$
- $|\mathbf{x}^+ - \mathbf{x}^-| = M$

# Computing the Margin Width

"Predict Class = +1" zone

$\mathbf{x}^+$

M = Margin Width

wx=1

wx=0

wx=-1

"Predict Class = -1" zone

$\mathbf{x}^-$

**What we know:**

- $\mathbf{w} \cdot \mathbf{x}^+ = +1$
- $\mathbf{w} \cdot \mathbf{x}^- = -1$
- $\mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$
- $|\mathbf{x}^+ - \mathbf{x}^-| = M$

$\mathbf{w} \cdot (\mathbf{x}^- + \lambda \mathbf{w}) + b = 1$

=>

$\mathbf{w} \cdot \mathbf{x}^- + b + \lambda \mathbf{w} \cdot \mathbf{w} = 1$

=>

$-1 + \lambda \mathbf{w} \cdot \mathbf{w} = 1$

=> $\lambda = \dfrac{2}{\mathbf{w} \cdot \mathbf{w}}$

# Computing the Margin Width

"Predict Class = +1" zone

$\mathbf{x}^+$

M = Margin Width

wx=1

wx=0

wx=-1

$\mathbf{x}^-$

"Predict Class = -1" zone

$$M = |\mathbf{x}^+ - \mathbf{x}^-| = |\lambda \mathbf{w}| =$$

**What we know:**

- $\mathbf{w} \cdot \mathbf{x}^+ = +1$
- $\mathbf{w} \cdot \mathbf{x}^- = -1$
- $\mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$
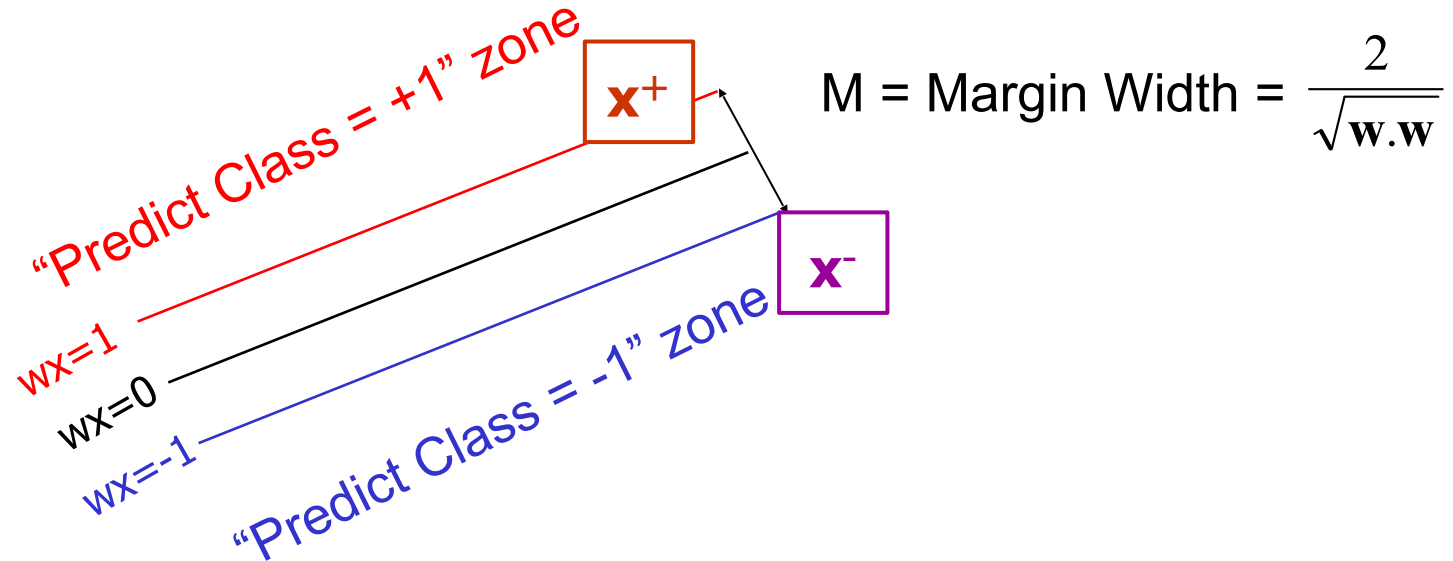- $|\mathbf{x}^+ - \mathbf{x}^-| = M$
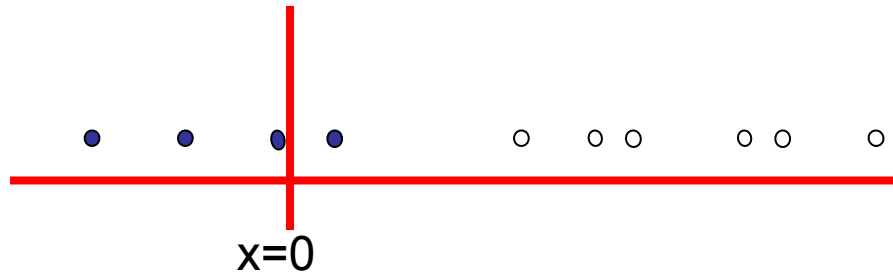
$$\lambda = \frac{2}{\mathbf{w}.\mathbf{w}}$$

$$= \lambda |\mathbf{w}| = \lambda \sqrt{\mathbf{w}.\mathbf{w}}$$

$$= \frac{2\sqrt{\mathbf{w}.\mathbf{w}}}{\mathbf{w}.\mathbf{w}} = \frac{2}{\sqrt{\mathbf{w}.\mathbf{w}}}$$

# Learning the Maximum Margin Classifier



M = Margin Width = $\dfrac{2}{\sqrt{\mathbf{w}.\mathbf{w}}}$

"Predict Class = +1" zone

$\mathbf{x}^+$

$\mathbf{x}^-$

wx=1

wx=0

wx=-1

"Predict Class = -1" zone

Use optimization to search the space of W to find the widest margin that matches all the data points.

# Simple 1-D Example

What would SVMs do?


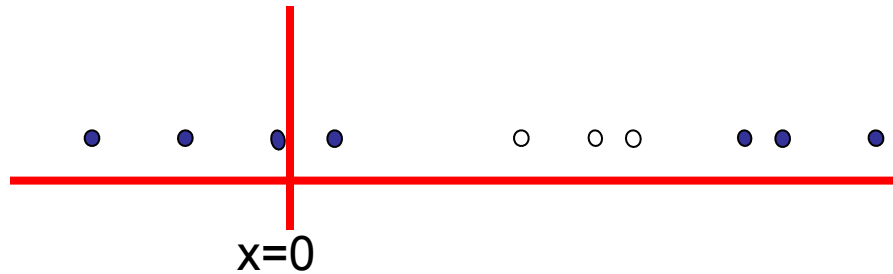
x=0

# Simple 1-D Example

Not a big surprise

x=0

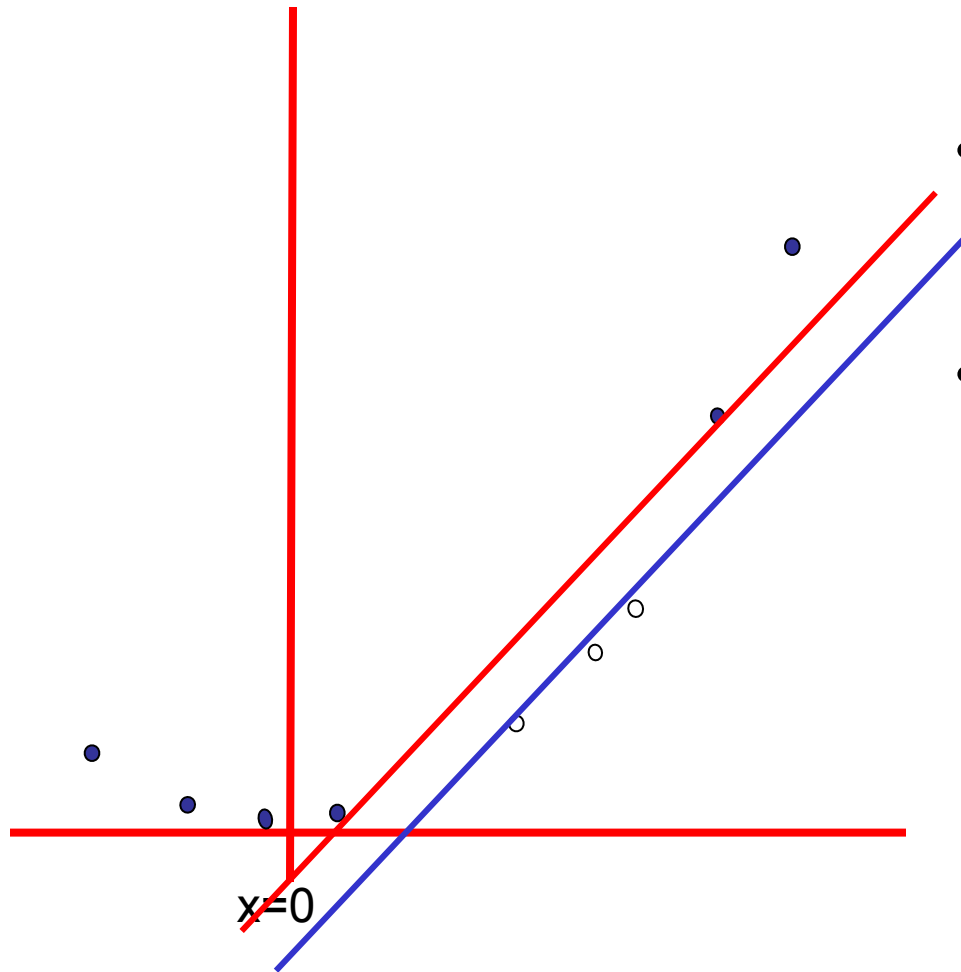Positive "plane"

Negative "plane"

# Harder 1-D Example

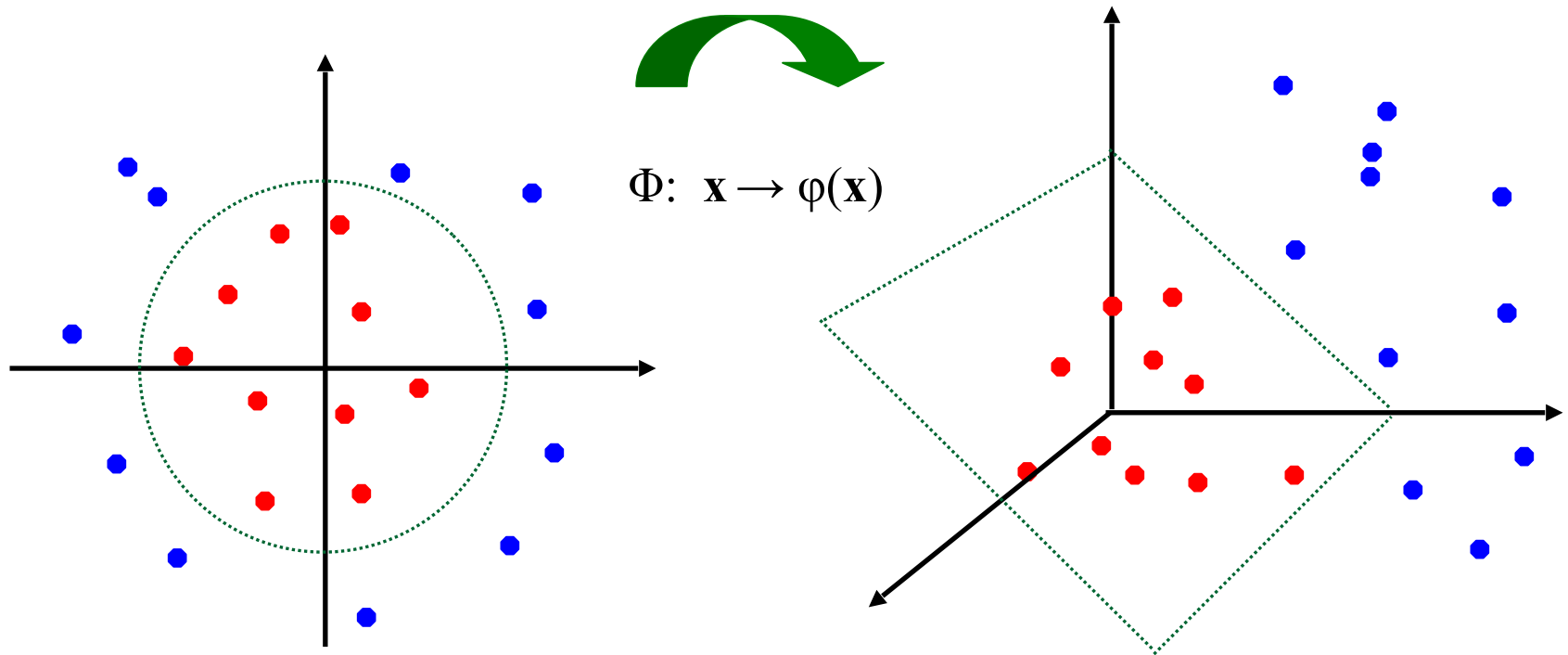What can SVM do about this?



x=0

# Harder 1-D Example

- Permit non-linear basis functions made linear regression
- Let's permit them here too

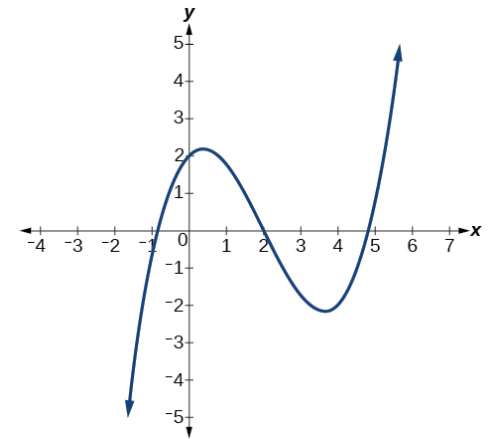x=0

$$\mathbf{z}_k = (x_k, x_k^2)$$

# Nonlinear SVMs: Feature Space

- General idea: the original input space can always be mapped to some higher-dimensional feature space where the training set is separable:
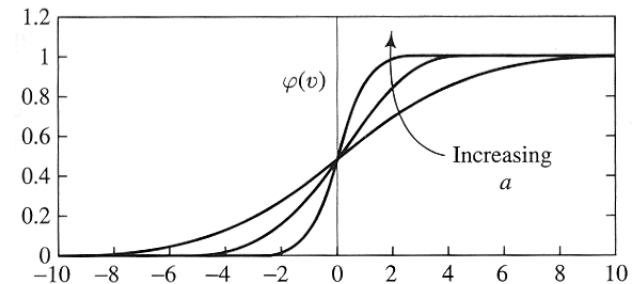
$$\Phi:\ \mathbf{x} \rightarrow \varphi(\mathbf{x})$$

# Nonlinear SVM Basis Functions

$\Phi(x_k)$ = ( polynomial terms of $\mathbf{x}_k$ of degree 1 to q )

$\Phi(x_k)$ = ( sigmoid functions of $\mathbf{x}_k$ )

$\Phi(x_k)$ = ( Gaussian radial basis functions of $\mathbf{x}_k$ )