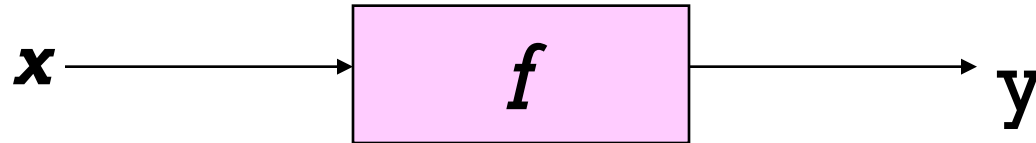


# LINEAR CLASSIFIER

# FIND A FUNCTION TO CLASSIFY HIGH VALUE CUSTOMERS



**High Value Customers**

Salary	Nb Orders
150	70
300	100
200	80
120	100

**Low Value Customers**

Salary	Nb Orders
40	80
220	20
100	20
175	10

Task: Find  $\alpha_1, \alpha_2, \alpha_3$  :

High value customer  $\alpha_1 \cdot salary + \alpha_2 \cdot orders - \alpha_3 > 0$

Low value customer  $\alpha_1 \cdot salary + \alpha_2 \cdot orders - \alpha_3 < 0$

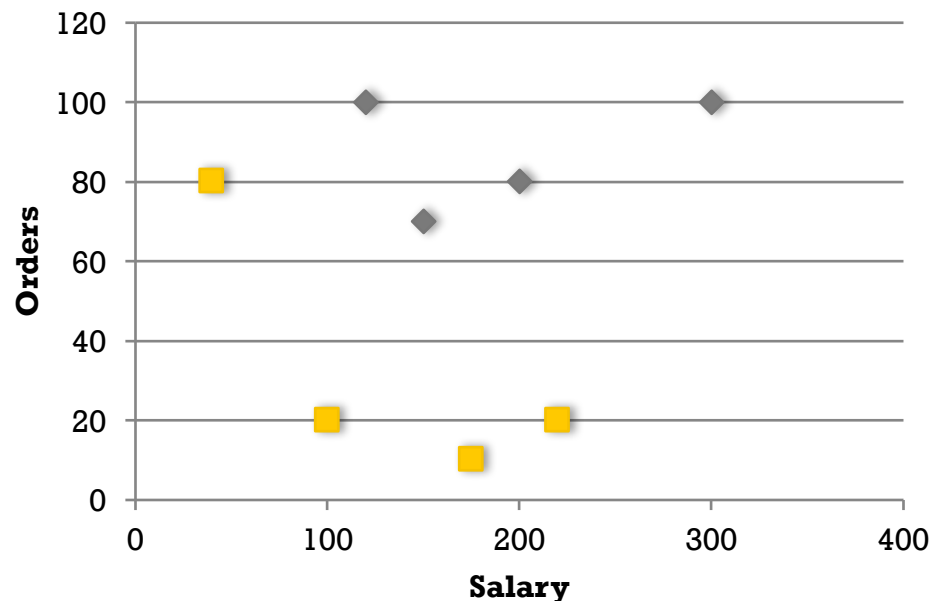
# FIND A FUNCTION TO CLASSIFY HIGH VALUE CUSTOMERS

## High Value Customers

Salary	Orders
150	70
300	100
200	80
120	100

## Low Value Customers

Salary	Orders
40	80
220	20
100	20
175	10



Task: Find  $\alpha_1, \alpha_2, \alpha_3$  :

High value customer  $\alpha_1 \cdot salary + \alpha_2 \cdot orders - \alpha_3 > 0$

Low value customer  $\alpha_1 \cdot salary + \alpha_2 \cdot orders - \alpha_3 < 0$

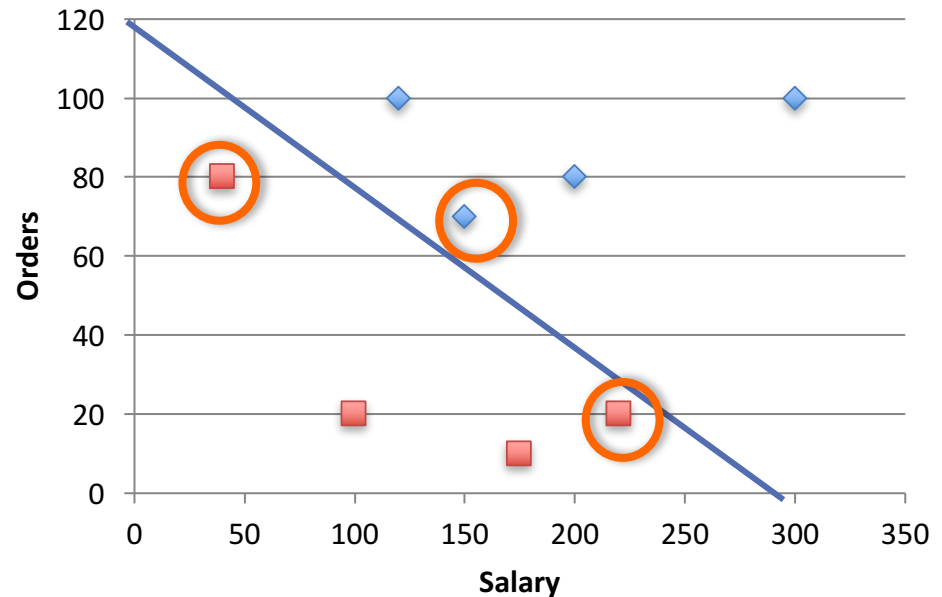
# FIND A FUNCTION TO CLASSIFY HIGH VALUE CUSTOMERS

## High Value Customers

Salary	Orders
150	70
300	100
200	80
120	100

## Low Value Customers

Salary	Orders
40	80
220	20
100	20
175	10

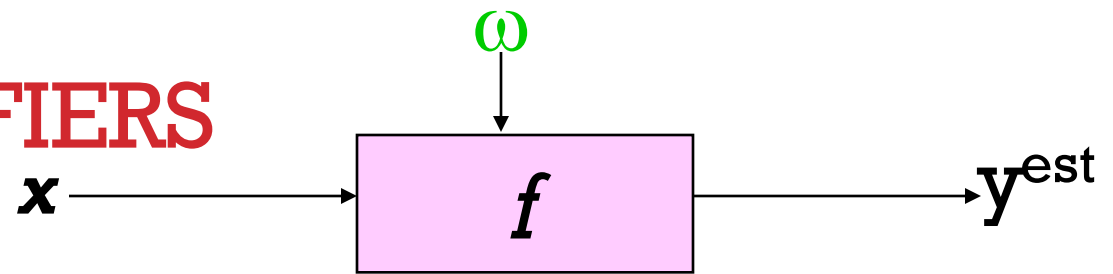


Task: Find  $\alpha_1, \alpha_2, \alpha_3$  :

High value customer  $\alpha_1 \cdot salary + \alpha_2 \cdot orders - \alpha_3 > 0$

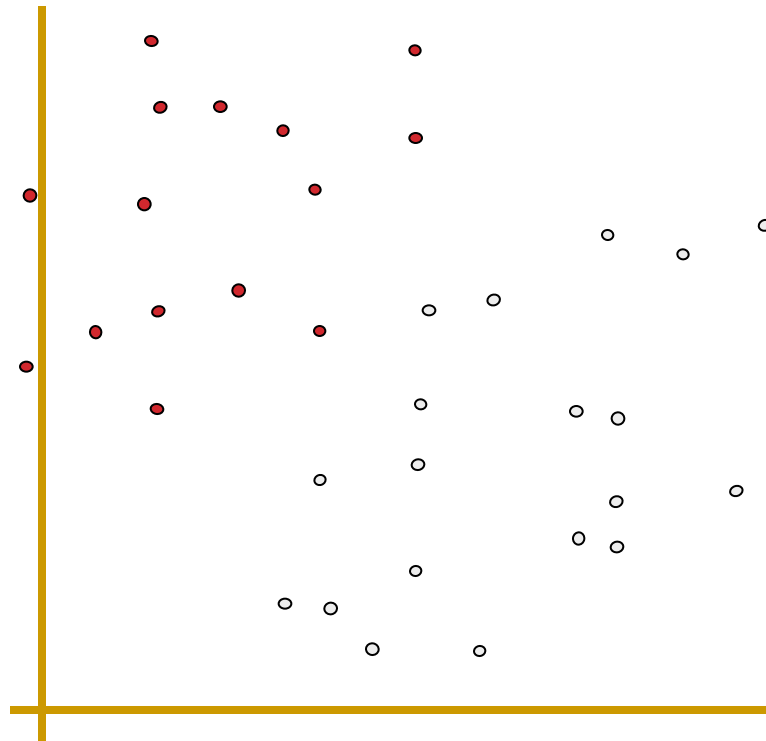
Low value customer  $\alpha_1 \cdot salary + \alpha_2 \cdot orders - \alpha_3 < 0$

# LINEAR CLASSIFIERS



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

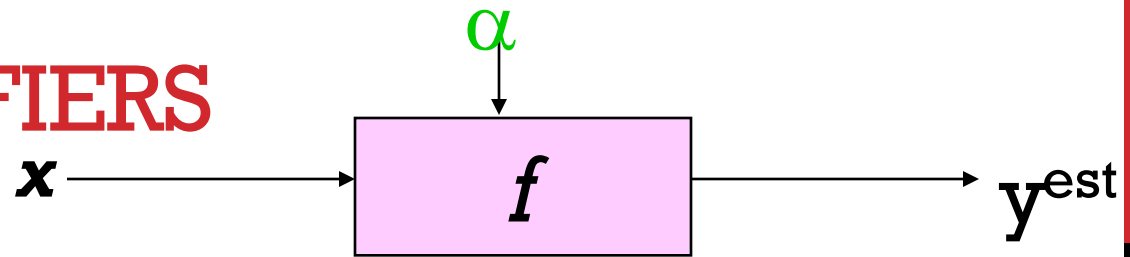
- denotes +1
- denotes -1



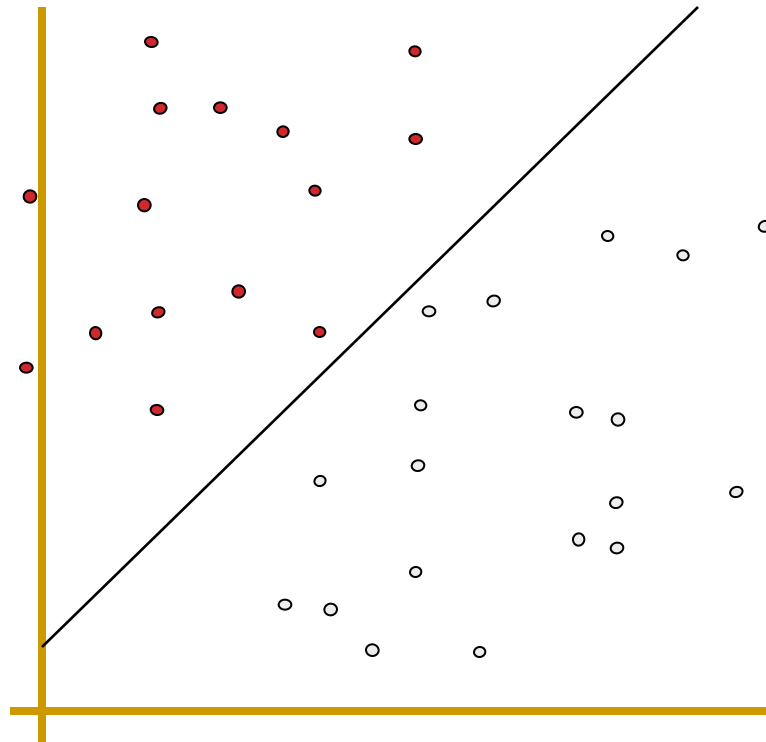
How would you classify this data?

$\mathbf{x}$  vector of feature values,  $\mathbf{w}$  vector of weights.  $\mathbf{x} \cdot \mathbf{w} = b$  is a line (hyperplane)

# LINEAR CLASSIFIERS



- denotes +1
- denotes -1

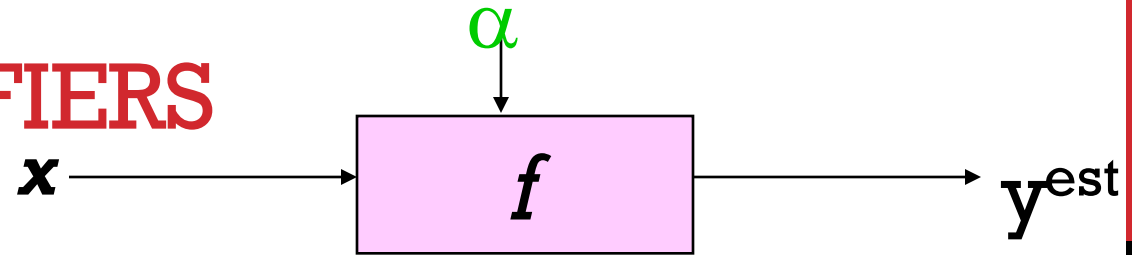


$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

How would you classify this data?

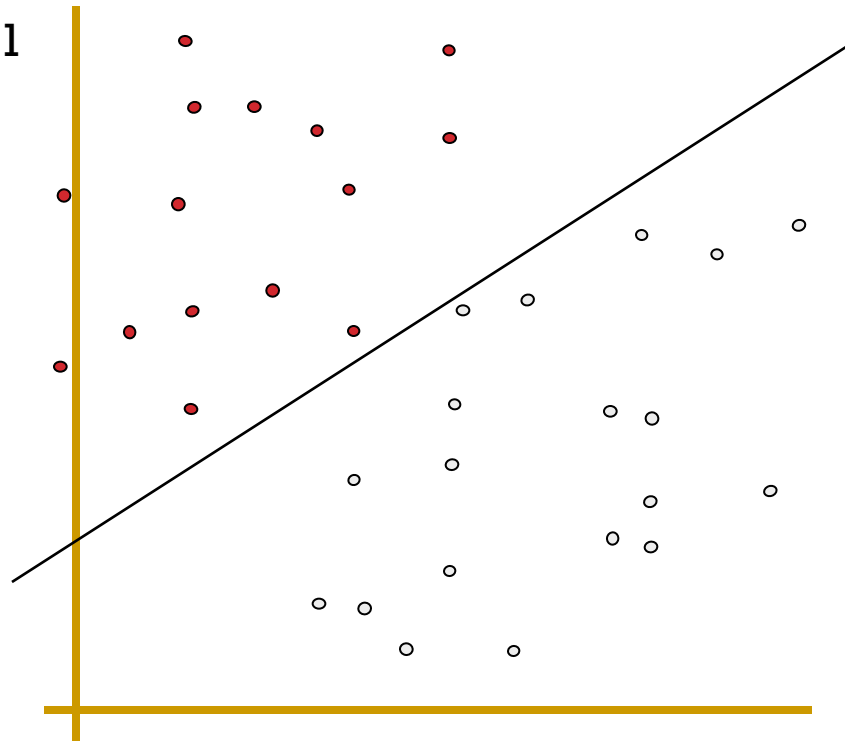
$\mathbf{x} \cdot \mathbf{w} = b$  or  $w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + \dots = b$  is an hyperplane.

# LINEAR CLASSIFIERS



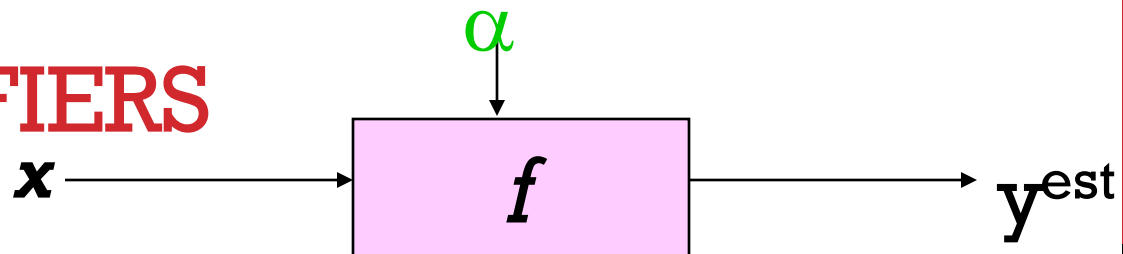
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1



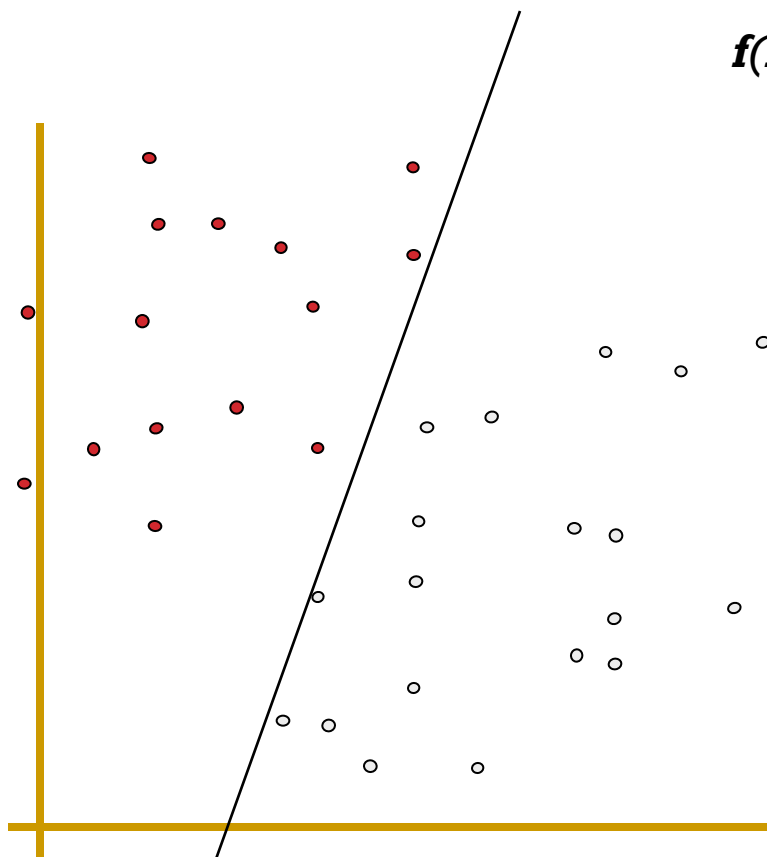
How would you classify this data?

# LINEAR CLASSIFIERS



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

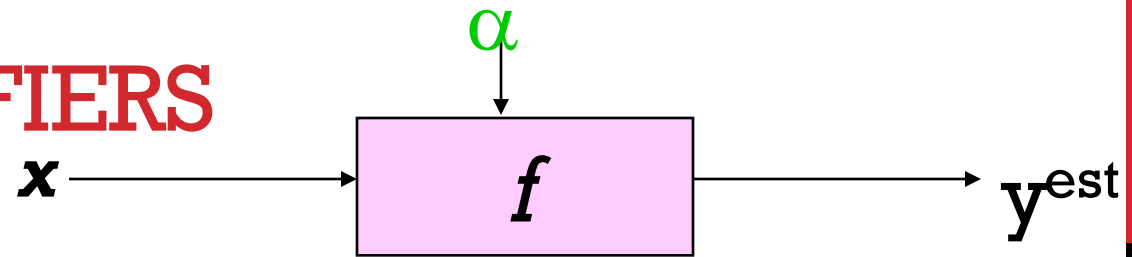
- denotes +1
- denotes -1



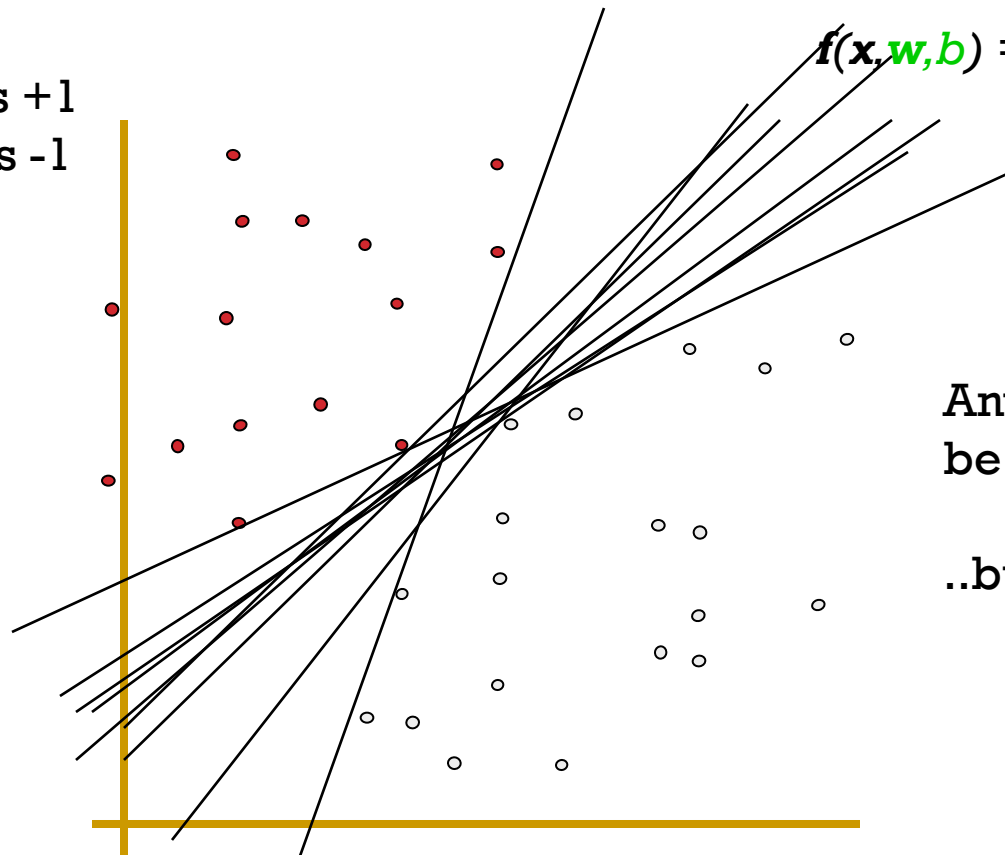
How would you  
classify this data?



# LINEAR CLASSIFIERS



- denotes +1
- denotes -1

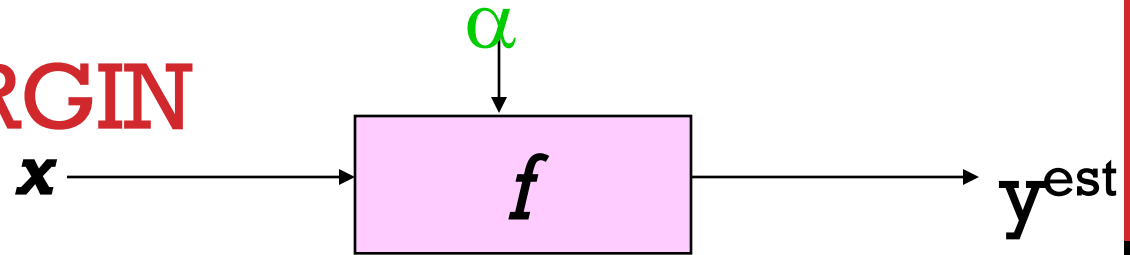


$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

Any of these would be fine..

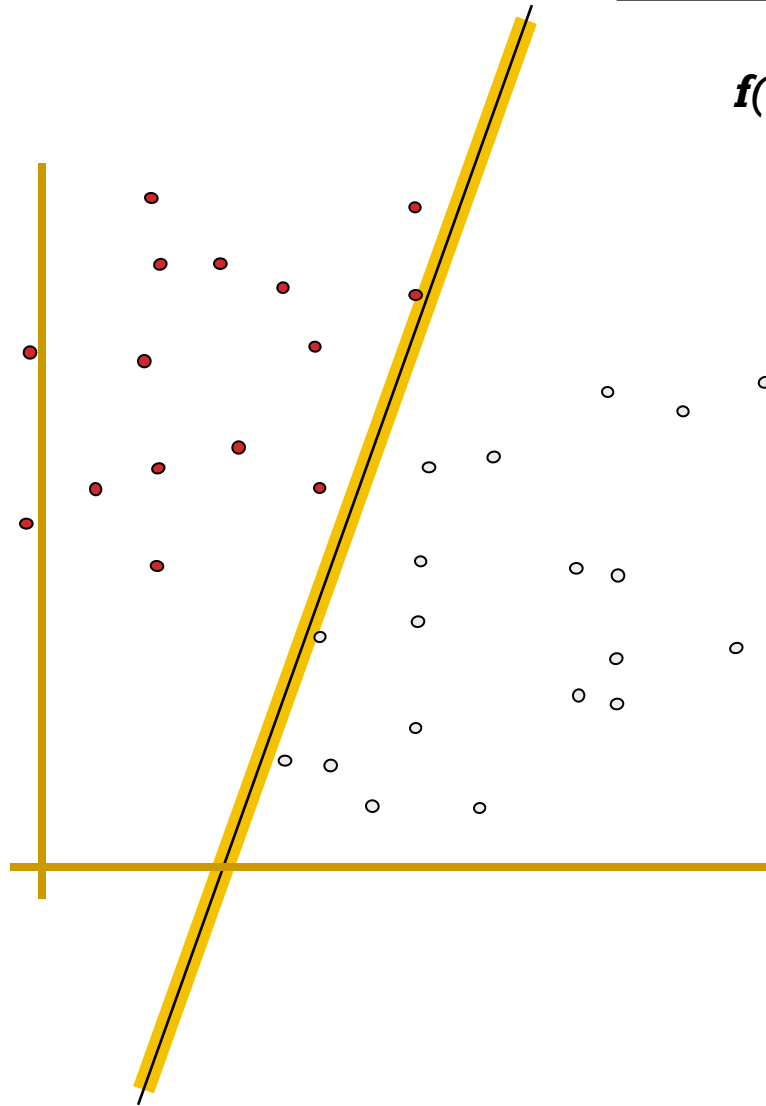
..but which is best?

# CLASSIFIER MARGIN



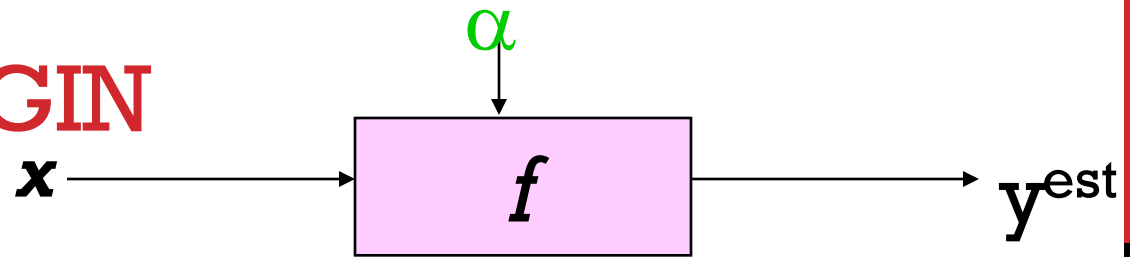
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1

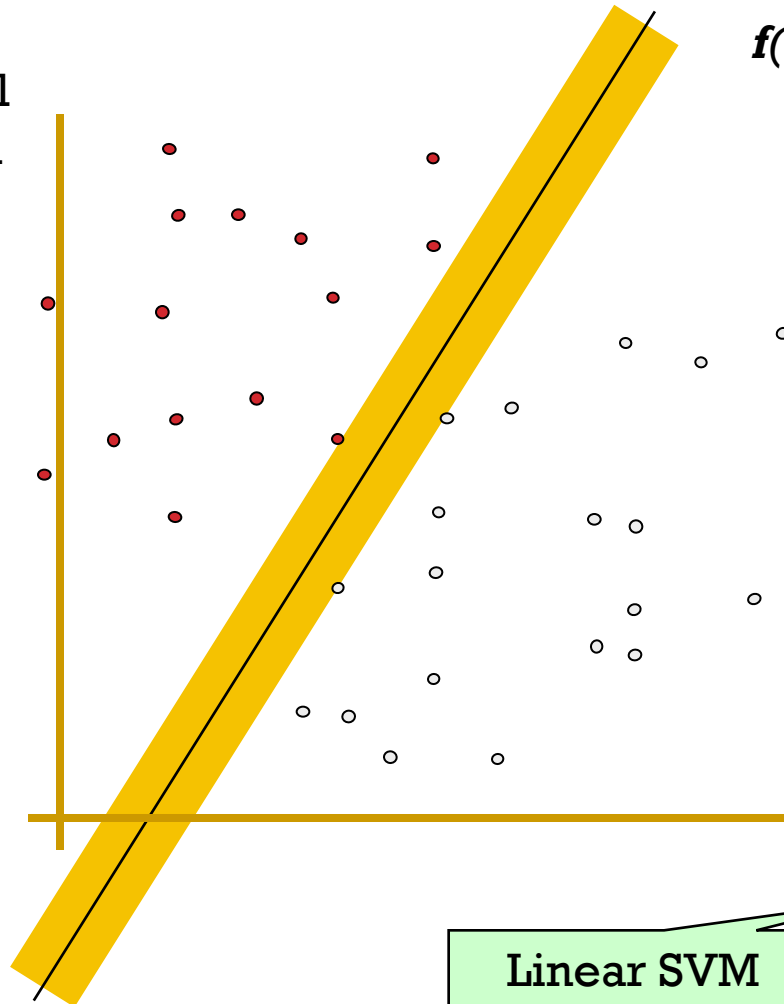


Define the **margin** of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

# MAXIMUM MARGIN



- denotes +1
- denotes -1

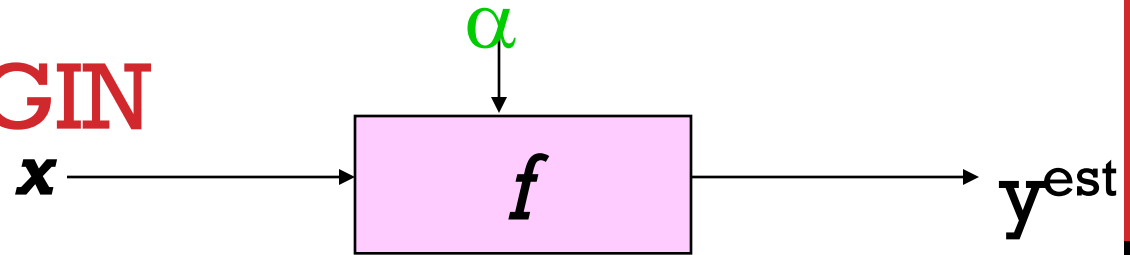


$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

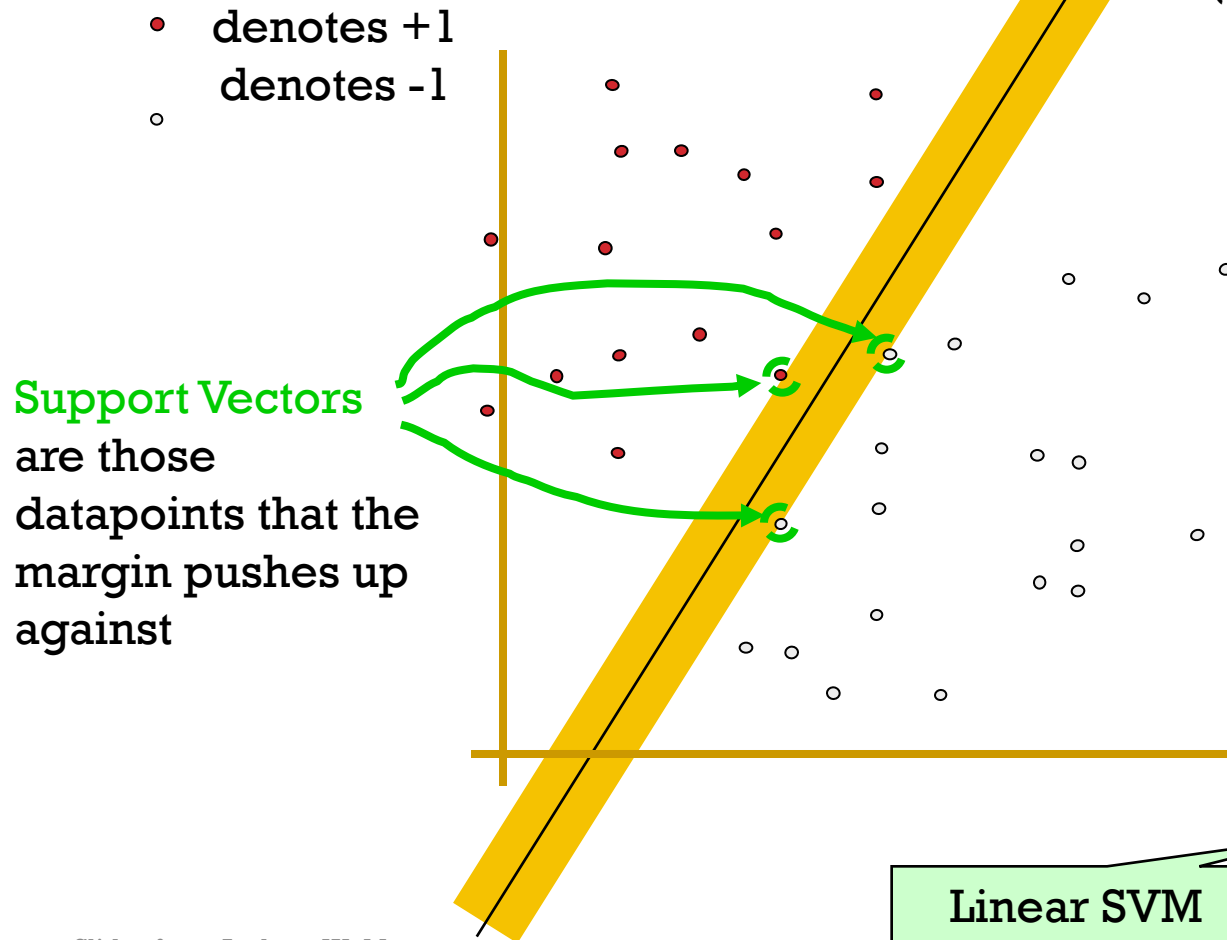
The **maximum margin linear classifier** is the linear classifier with the, um, maximum margin. This is the simplest kind of SVM (Called an LSVM)

Linear SVM

# MAXIMUM MARGIN

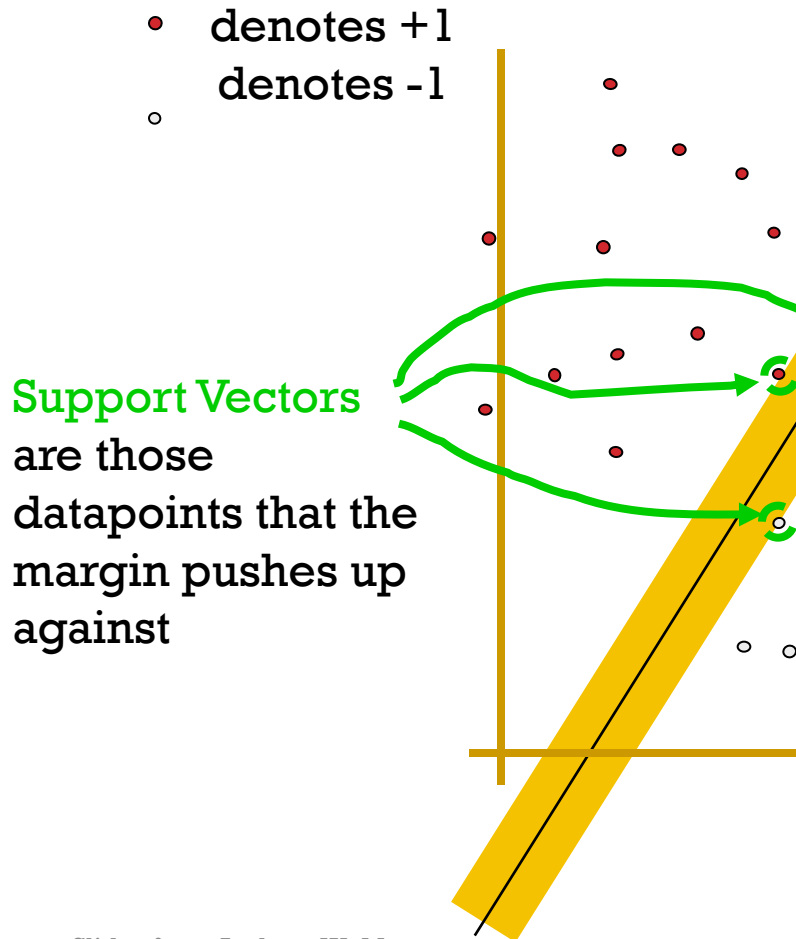


$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$



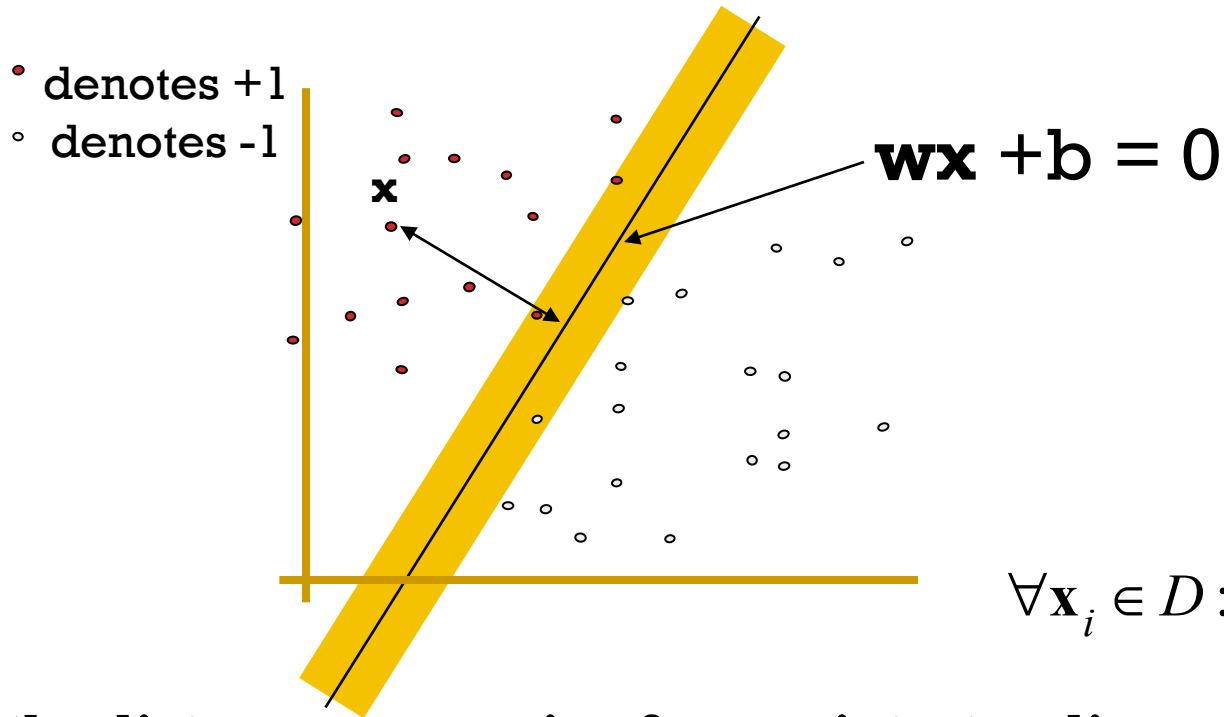
The **maximum margin linear classifier** is the linear classifier with the, um, maximum margin. This is the simplest kind of SVM (Called an LSVM)

# WHY MAXIMUM MARGIN?



1. Intuitively this feels safest.
2. If we've made a small error in the location of the boundary (it's been jolted in its perpendicular direction) this gives us least chance of causing a misclassification.
3. Leave-one-out-cross-validation (LOOCV) is easy since the model is immune to removal of any non-support-vector datapoints.
4. There's some theory (using VC dimension) that is related to (but not the same as) the proposition that this is a good thing.
5. Empirically it works very very well.

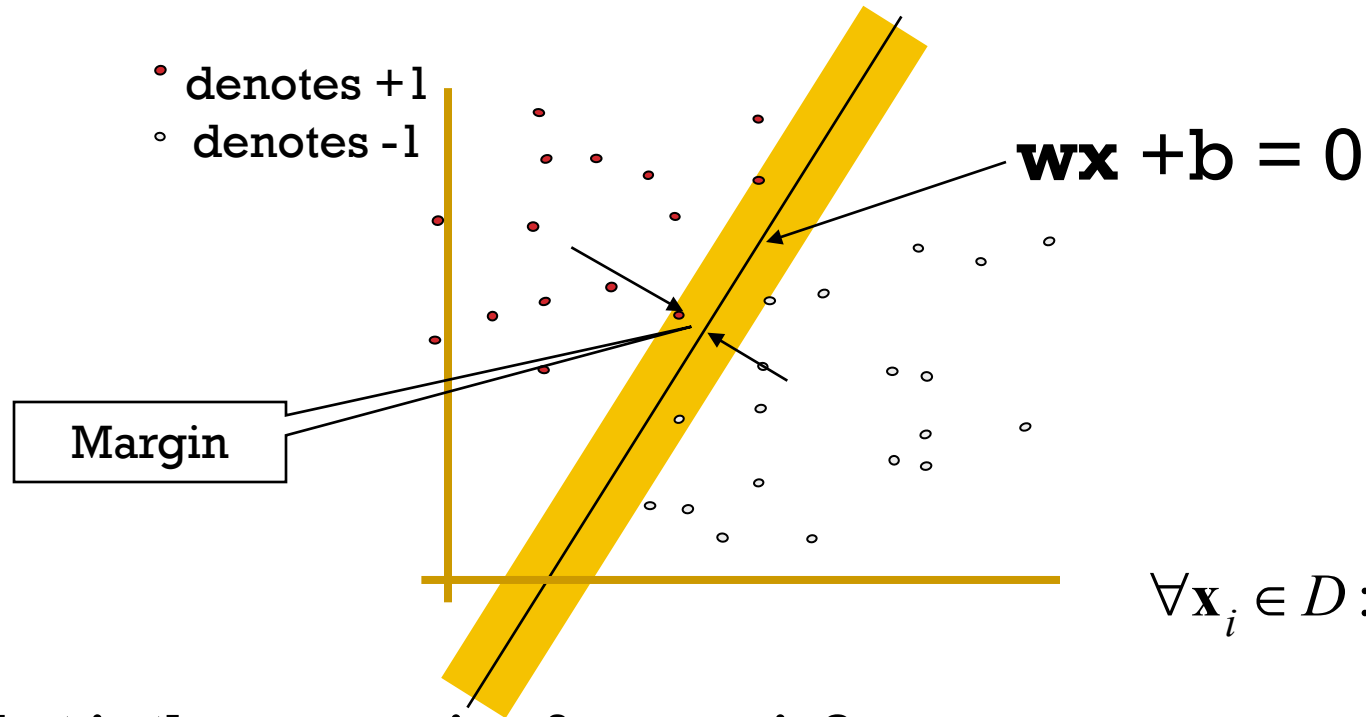
# ESTIMATE THE MARGIN



**What is the distance expression for a point  $\mathbf{x}$  to a line  $\mathbf{w}\mathbf{x} + b = 0$ ?**

$$d(\mathbf{x}) = \frac{|\mathbf{x} \cdot \mathbf{w} + b|}{\|\mathbf{w}\|_2} = \frac{|\mathbf{x} \cdot \mathbf{w} + b|}{\sqrt{\sum_{i=1}^d w_i^2}}$$

# ESTIMATE THE MARGIN

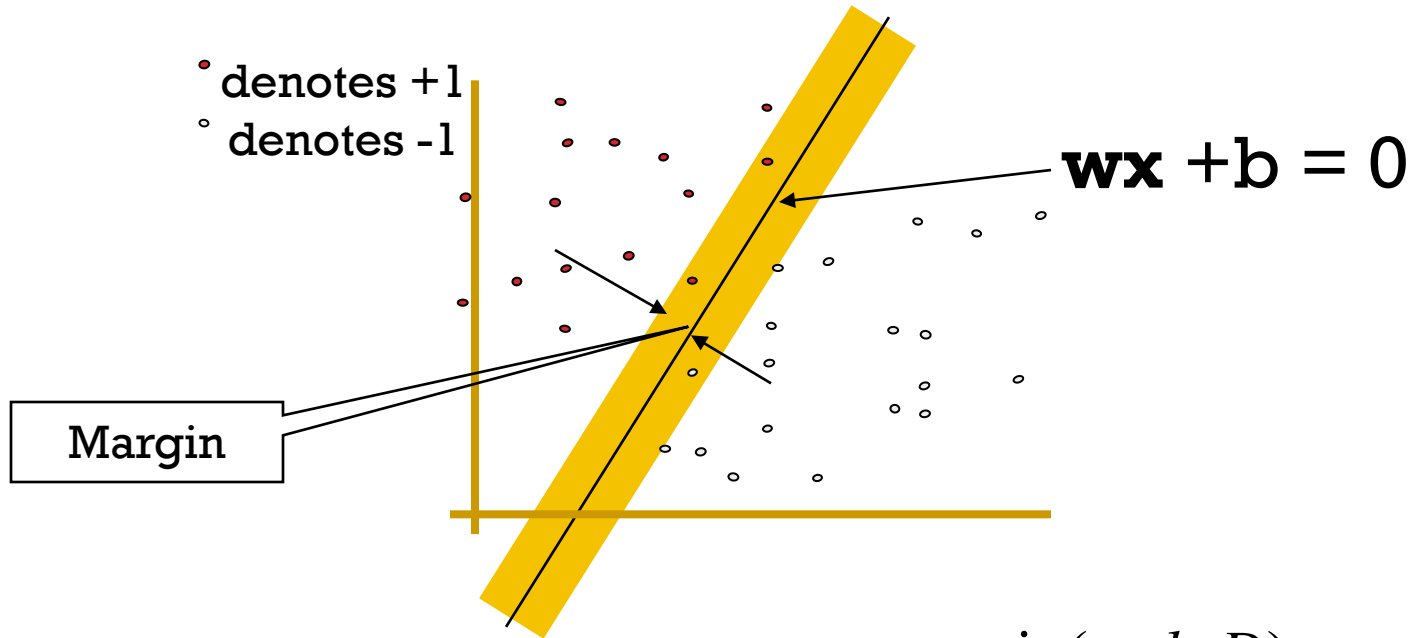


$$\forall \mathbf{x}_i \in D : y_i (\mathbf{x}_i \cdot \mathbf{w} + b) > 0$$

**What is the expression for margin?**

$$\text{margin} \equiv \min_{\mathbf{x} \in D} d(\mathbf{x}) = \min_{\mathbf{x} \in D} \frac{|\mathbf{x} \cdot \mathbf{w} + b|}{\sqrt{\sum_{i=1}^d w_i^2}}$$

# MAXIMIZE MARGIN



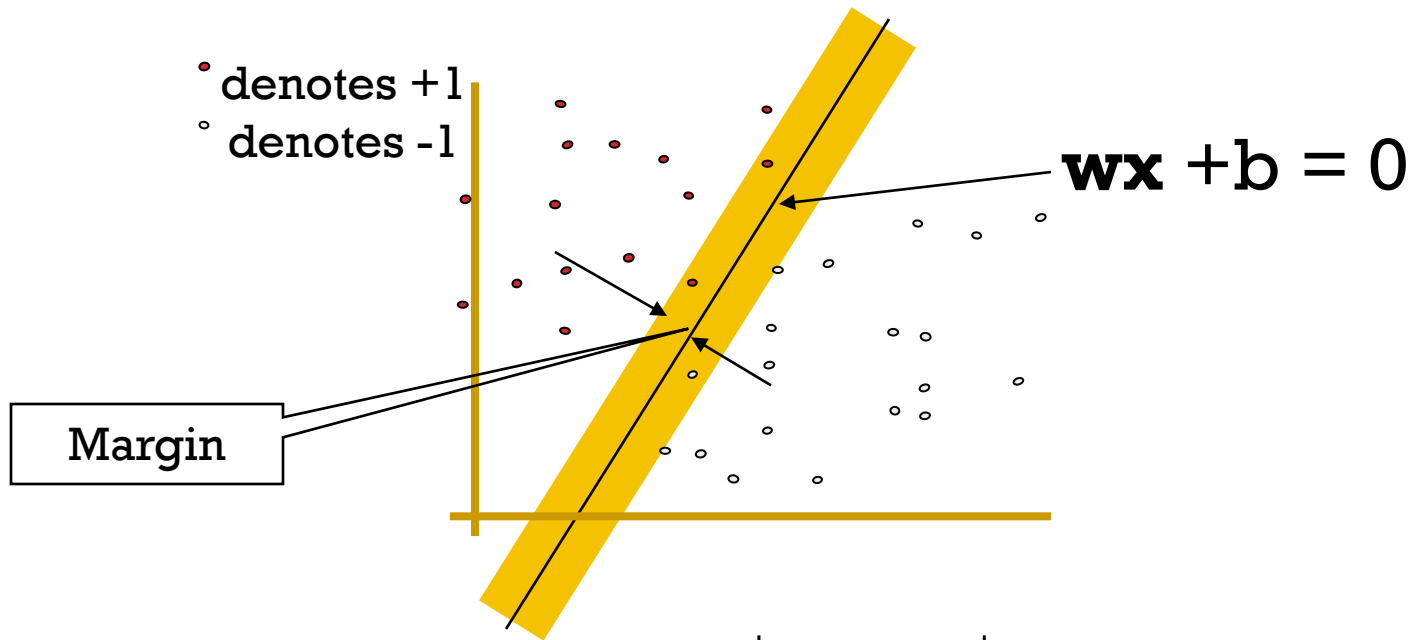
$$\operatorname{argmax}_{\mathbf{w}, b} \operatorname{margin}(\mathbf{w}, b, D)$$

$$= \operatorname{argmax}_{\mathbf{w}, b} \min_{\mathbf{x}_i \in D} d(\mathbf{x}_i)$$

$$= \operatorname{argmax}_{\mathbf{w}, b} \min_{\mathbf{x}_i \in D} \frac{|b + \mathbf{x}_i \cdot \mathbf{w}|}{\sqrt{\sum_{i=1}^d w_i^2}}$$



# MAXIMIZE MARGIN

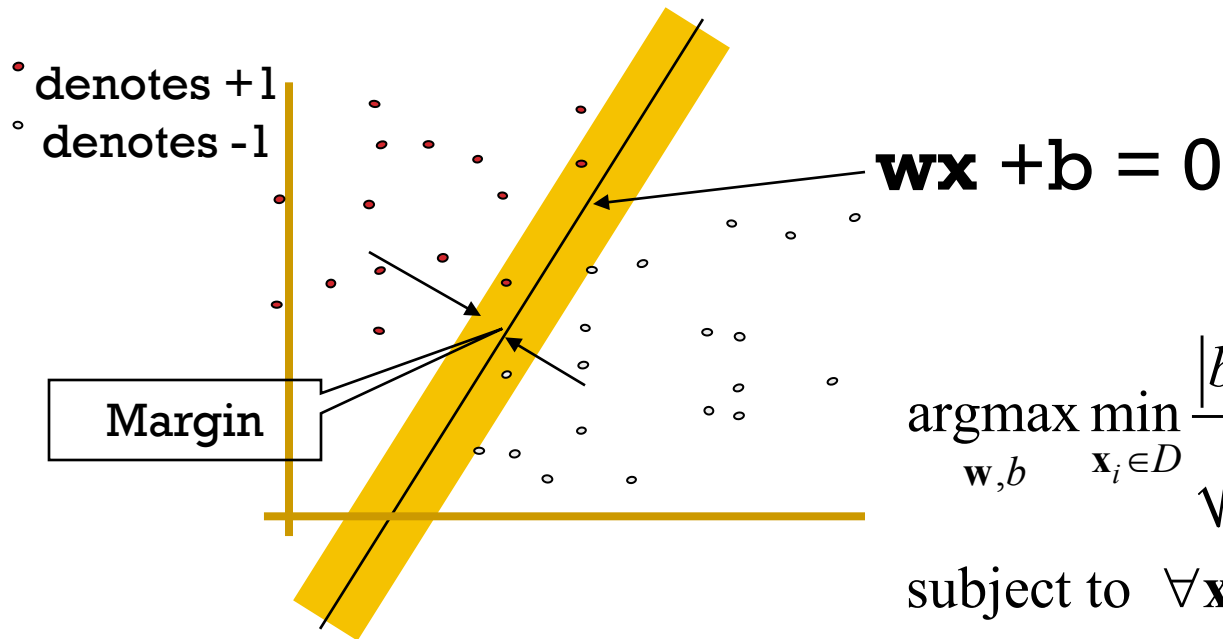


$$\operatorname{argmax}_{\mathbf{w}, b} \min_{\mathbf{x}_i \in D} \frac{|b + \mathbf{x}_i \cdot \mathbf{w}|}{\sqrt{\sum_{i=1}^d w_i^2}}$$

$$\text{subject to } \forall \mathbf{x}_i \in D : y_i (\mathbf{x}_i \cdot \mathbf{w} + b) > 0$$

**Min-max problem  $\rightarrow$  game problem**

# MAXIMIZE MARGIN



$$\operatorname{argmax}_{\mathbf{w}, b} \min_{\mathbf{x}_i \in D} \frac{|b + \mathbf{x}_i \cdot \mathbf{w}|}{\sqrt{\sum_{i=1}^d w_i^2}}$$

$$\text{subject to } \forall \mathbf{x}_i \in D : y_i (\mathbf{x}_i \cdot \mathbf{w} + b) \geq 0$$



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

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

**Strategy:**

$$\forall \mathbf{x}_i \in D : |b + \mathbf{x}_i \cdot \mathbf{w}| \geq 1$$

If you want to learn why this holds I highly recommend the following video  
<https://www.youtube.com/watch?v=PwhiWxHK8o> (not necessary to know for the

# MAXIMUM MARGIN LINEAR CLASSIFIER

$$\{\vec{w}^*, b^*\} = \operatorname{argmin}_{\vec{w}, b} \sum_{k=1}^d w_k^2$$

subject to

$$y_1 (\vec{w} \cdot \vec{x}_1 + b) \geq 1$$

$$y_2 (\vec{w} \cdot \vec{x}_2 + b) \geq 1$$

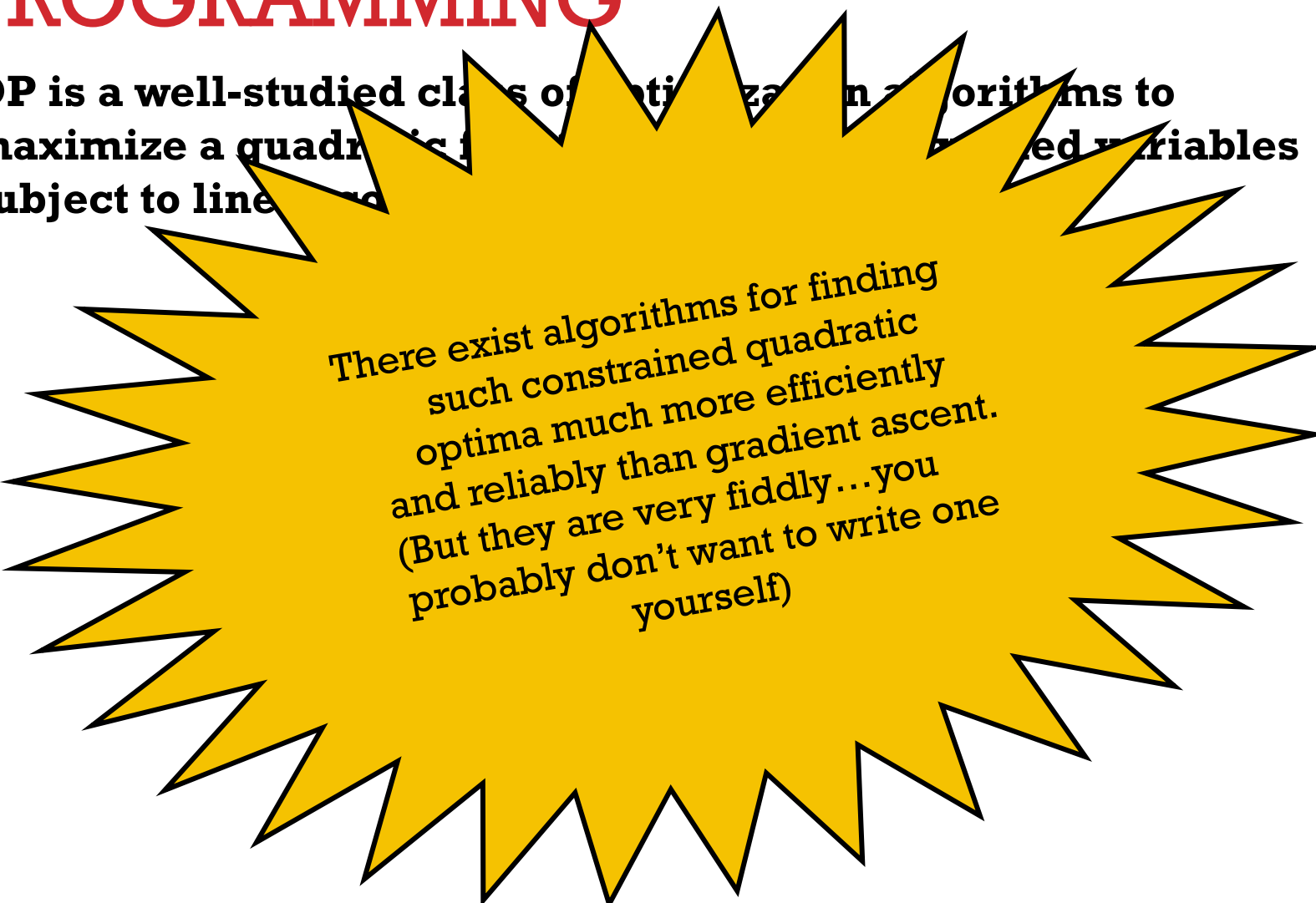
....

$$y_N (\vec{w} \cdot \vec{x}_N + b) \geq 1$$

**How to solve it?**

# LEARNING VIA QUADRATIC PROGRAMMING

QP is a well-studied class of optimization algorithms to maximize a quadratic function of  $n$  variables subject to linear constraints.



There exist algorithms for finding such constrained quadratic optima much more efficiently and reliably than gradient ascent. (But they are very fiddly...you probably don't want to write one yourself)

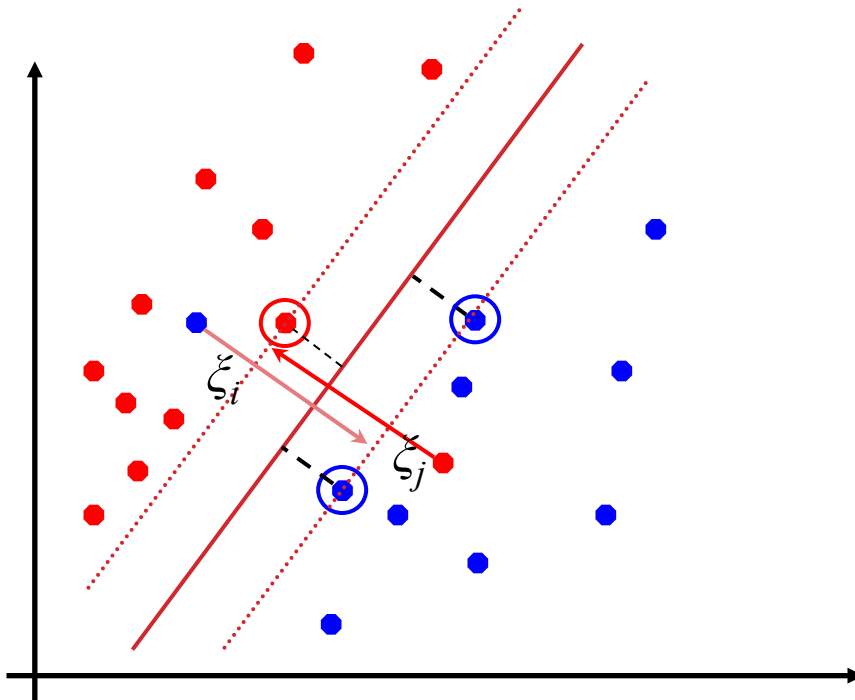
# SOFT MARGIN CLASSIFICATION

If the training data is not linearly separable, *slack variables*  $\xi_i$  can be added to allow misclassification of difficult or noisy examples.

Allow some errors

- Let some points be moved to where they belong, at a cost

Still, try to minimize training set errors, and to place hyperplane “far” from each class (large margin)



# REGULARIZATION

## SOFT MARGIN CLASSIFICATION

### MATHEMATICALLY

Find  $\mathbf{w}$  and  $b$  such that

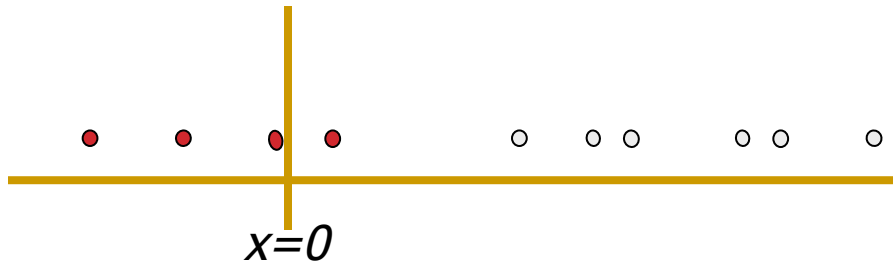
$$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} \text{ is minimized and for all } \{(\mathbf{x}_i, y_i)\}$$
$$y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

Find  $\mathbf{w}$  and  $b$  such that

$$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum \xi_i \text{ is minimized and for all } \{(\mathbf{x}_i, y_i)\}$$
$$y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \text{ and } \xi_i \geq 0 \text{ for all } i$$

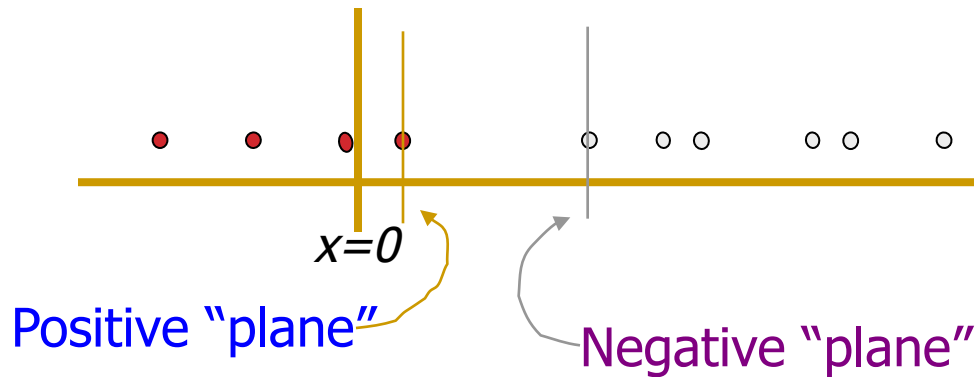
# SUPPOSE WE'RE IN 1-DIMENSION

What would  
SVMs do with  
this data?



# SUPPOSE WE'RE IN 1-DIMENSION

Not a big surprise

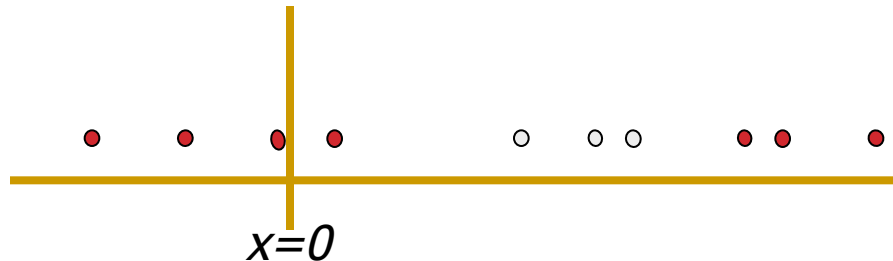




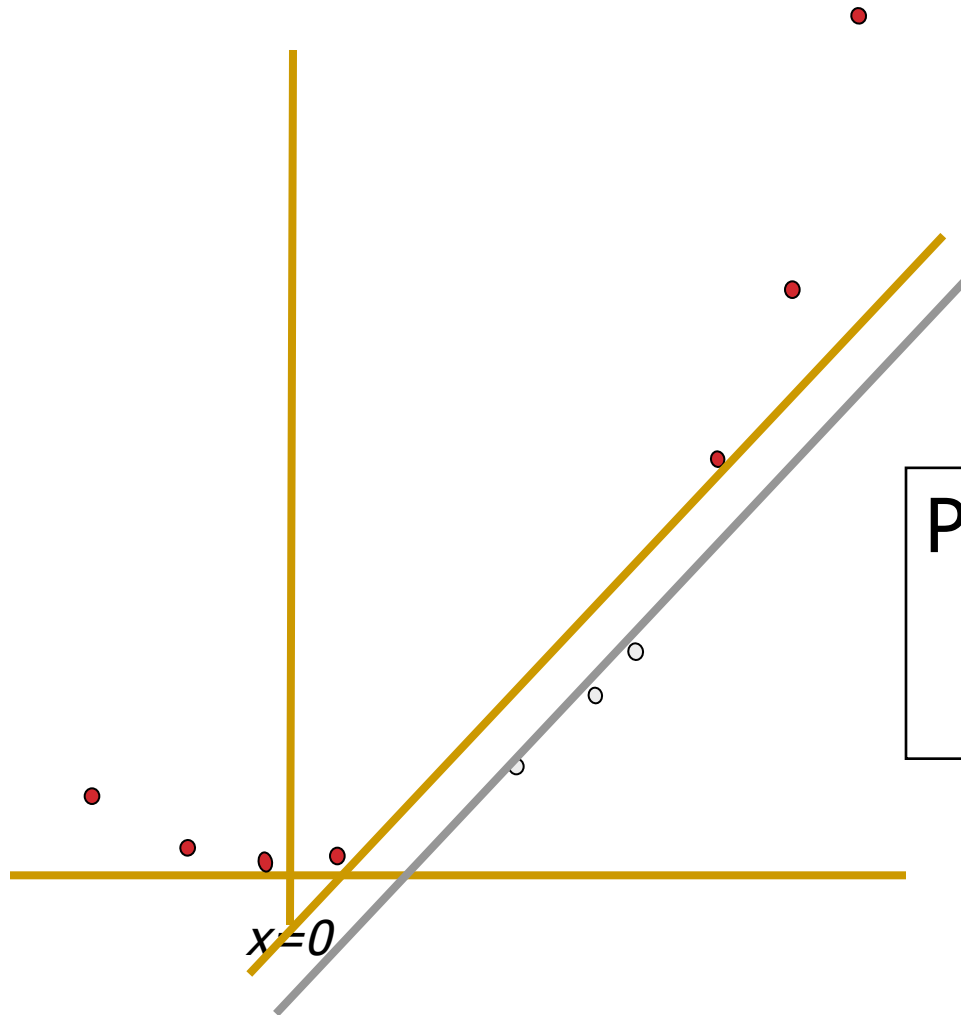
# HARDER 1-DIMENSIONAL DATASET

That's wiped the smirk off SVM's face.

What can be done about this?



# HARDER 1-DIMENSIONAL DATASET



Permitting non-linear basis functions

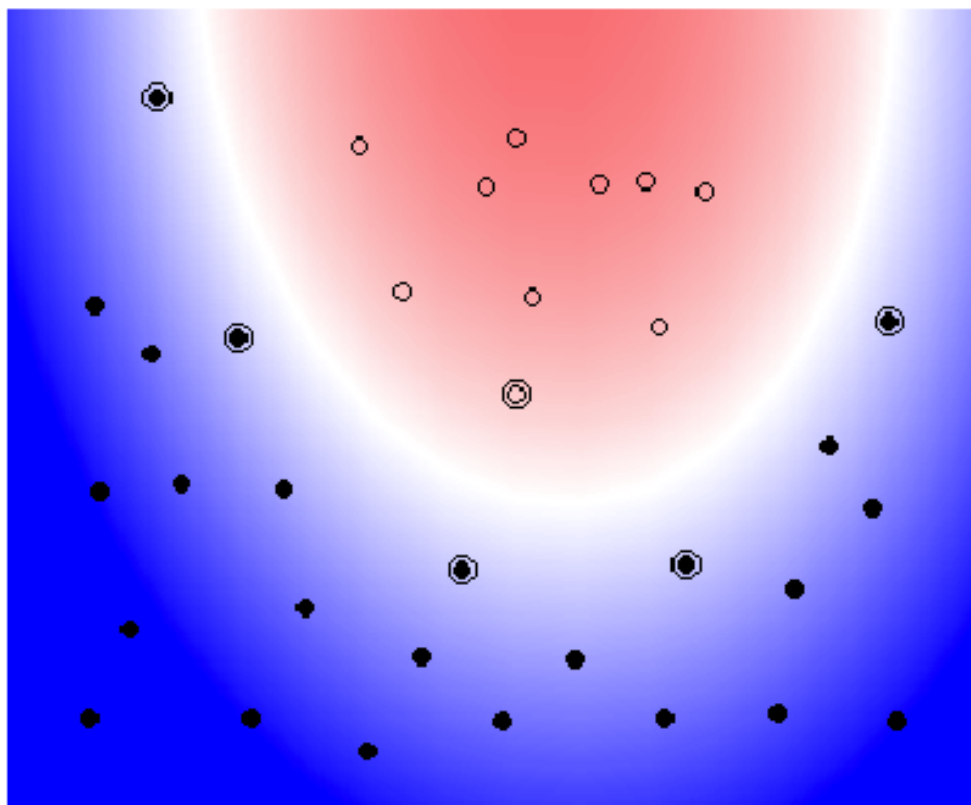
$$\mathbf{z}_k = (x_k, x_k^2)$$

# NONLINEAR KERNEL (I)

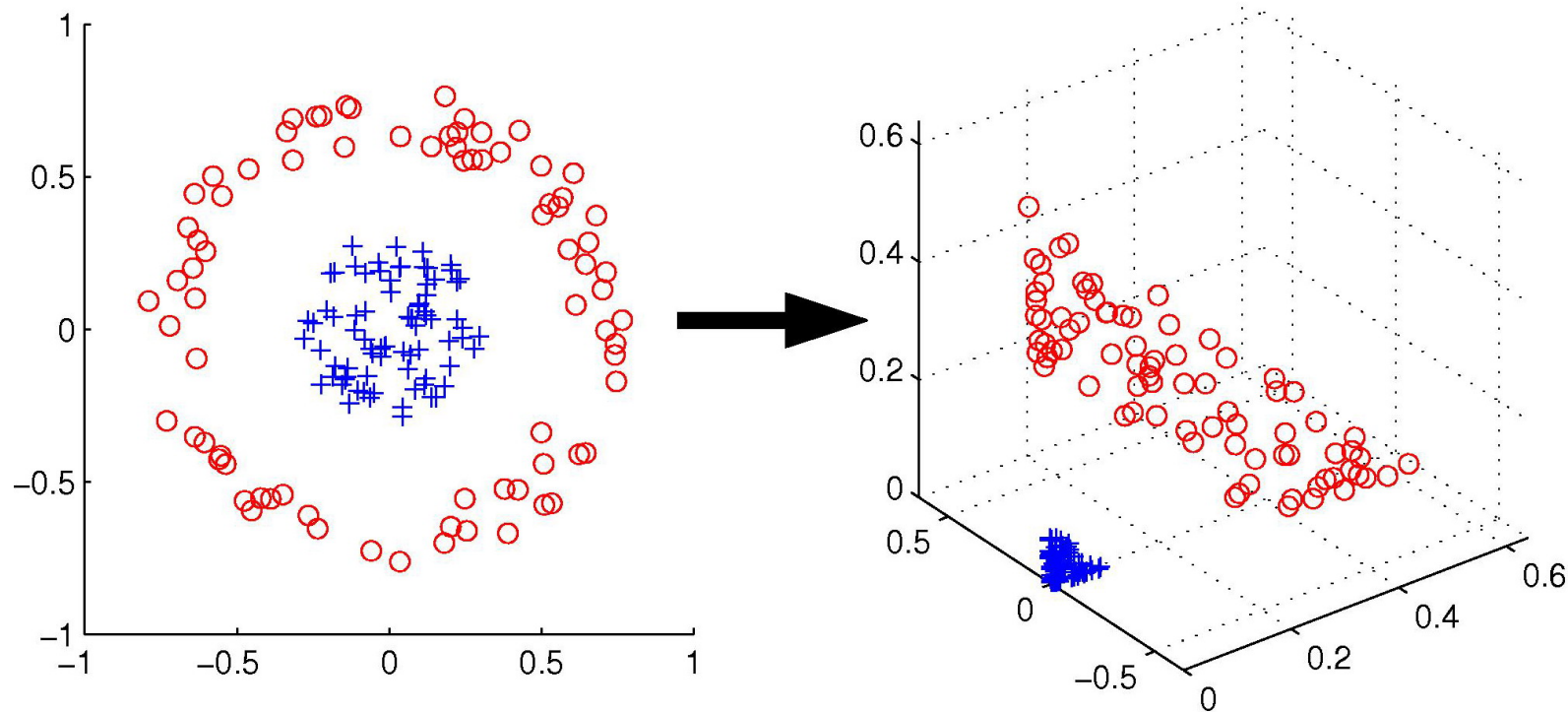
## Example: SVM with Polynomial of Degree 2

Kernel:  $K(\vec{x}_i, \vec{x}_j) = [\vec{x}_i \cdot \vec{x}_j + 1]^2$

plot by Bell SVM applet



# NON-LINEAR KERNEL



$$\begin{aligned}\phi: \quad \mathbb{R}^2 &\longrightarrow \mathbb{R}^3 \\ (x_1, x_2) &\longmapsto (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)\end{aligned}$$

# SVM with a polynomial Kernel visualization

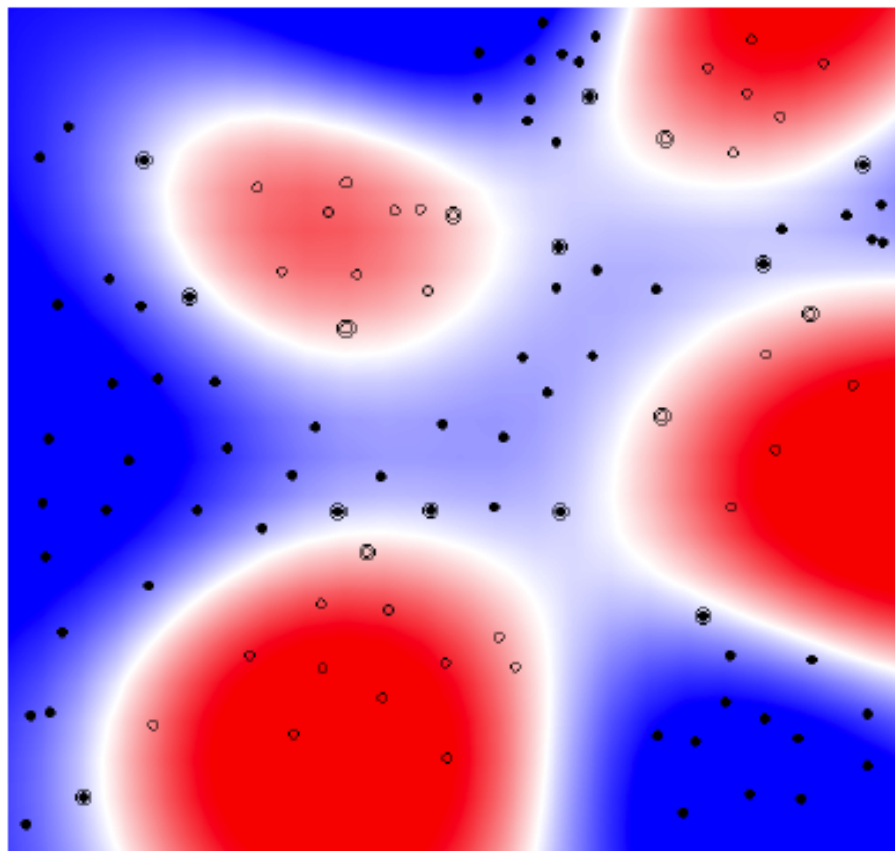
Created by:  
Udi Aharoni

# NONLINEAR KERNEL (II)

## Example: SVM with RBF-Kernel

Kernel:  $K(\vec{x}_i, \vec{x}_j) = \exp(-|\vec{x}_i - \vec{x}_j|^2 / \sigma^2)$

plot by Bell SVM applet



Nice video lecture from CalTech on RBF kernels:  
<https://www.youtube.com/watch?v=O8CfrnOPTLc>

# DEMO

**<http://cs.stanford.edu/people/karpathy/svmjs/demo/>**

# KERNEL TRICKS

## **Pro**

- Introducing nonlinearity into the model
- Computational cheap

## **Con**

- Still have potential overfitting problems



# SVM PERFORMANCE

**Anecdotally they work very very well.**

**Example: They are currently the best-known classifier on a well-studied hand-written-character recognition benchmark**

**Anecdotally reliable people doing practical real-world work who claim that SVMs have saved them when their other favorite classifiers did poorly.**

**There is a lot of excitement and religious fervor about SVMs**

**Despite this, some practitioners are a little skeptical.**

# REFERENCES

## **An excellent tutorial on VC-dimension and Support Vector Machines:**

C.J.C. Burges. A tutorial on support vector machines for pattern recognition. Data Mining and Knowledge Discovery, 2(2):955-974, 1998.

## **The VC/SRM/SVM Bible:**

Statistical Learning Theory by Vladimir Vapnik, Wiley-Interscience; 1998

**Software: SVM-light, <http://svmlight.joachims.org/>, free download**