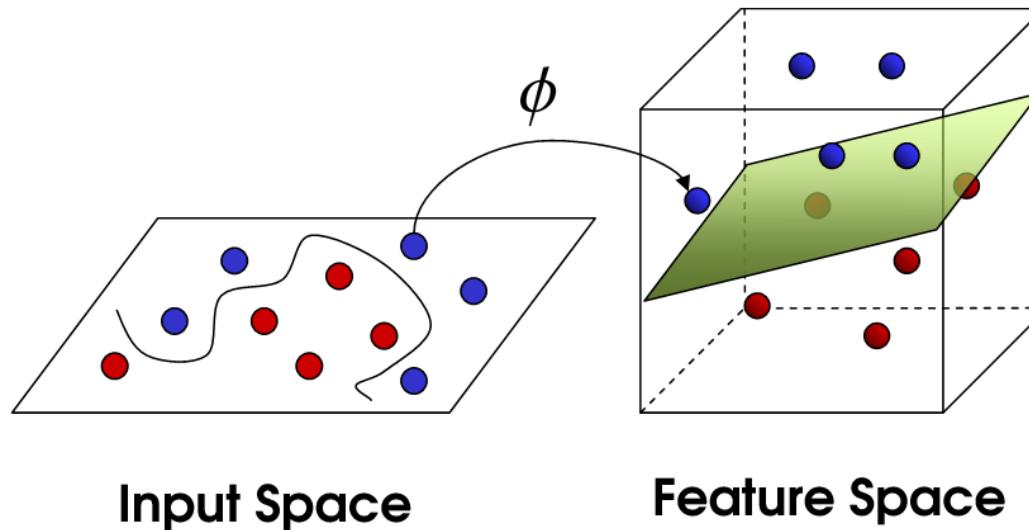


# CS109 – Data Science

## SVM, Performance evaluation

Joe Blitzstein, Hanspeter Pfister, Verena Kaynig-Fittkau

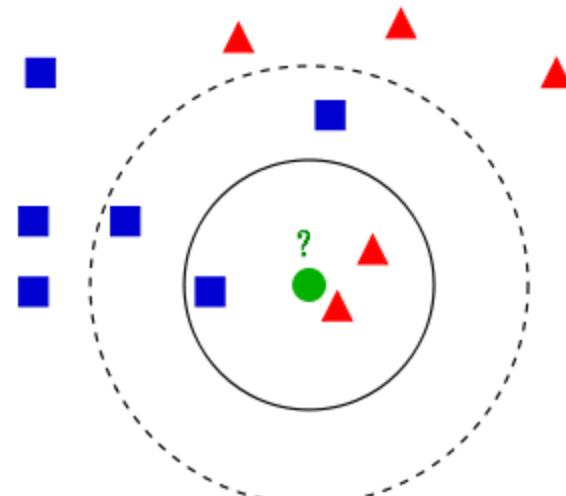


# Announcements

- HW1 grades went out yesterday
- They are looking really good, well done everyone!
- HW2 is due this Thursday!
- You should submit an executed notebook
- But please without pages of test output

# Recap K-NN

- Keeps all training data
  - Training is fast
  - Prediction is slow
- Go through all data points → slow prediction
- related.*



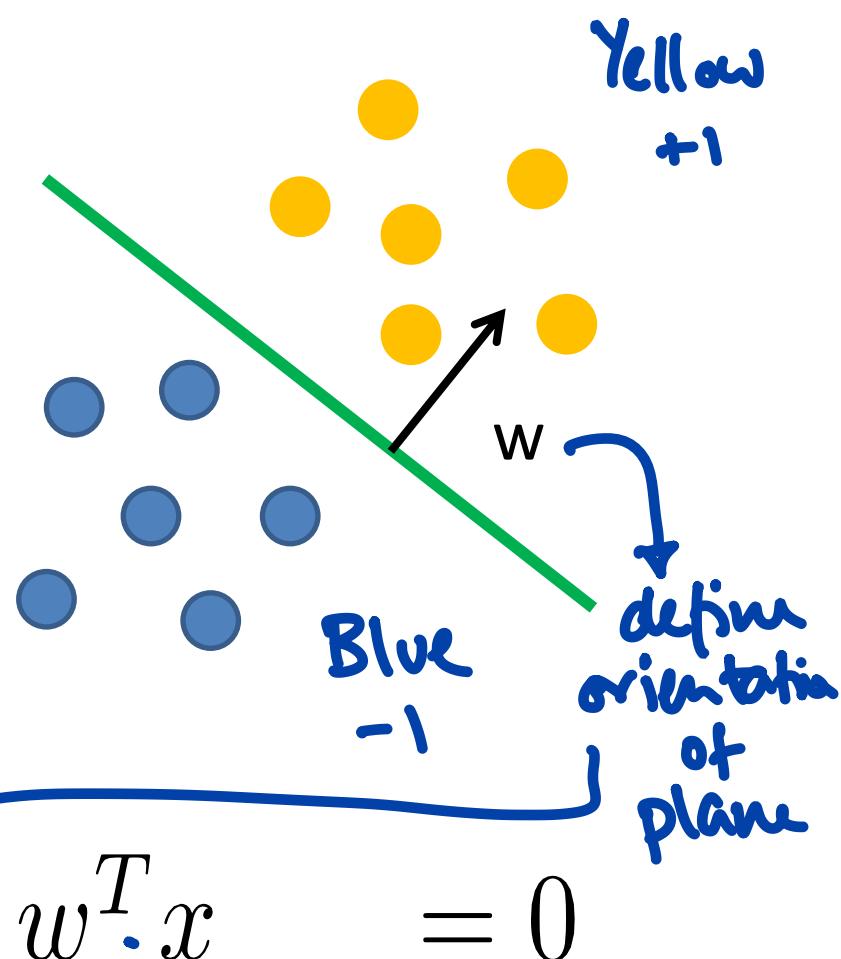
# Separating Hyperplane

image  $\rightarrow$  cat / dog

- $x$ : data point
- $y$ : label  $\in \{-1, +1\}$
- $w$ : weight vector

Some  
use  
0/1

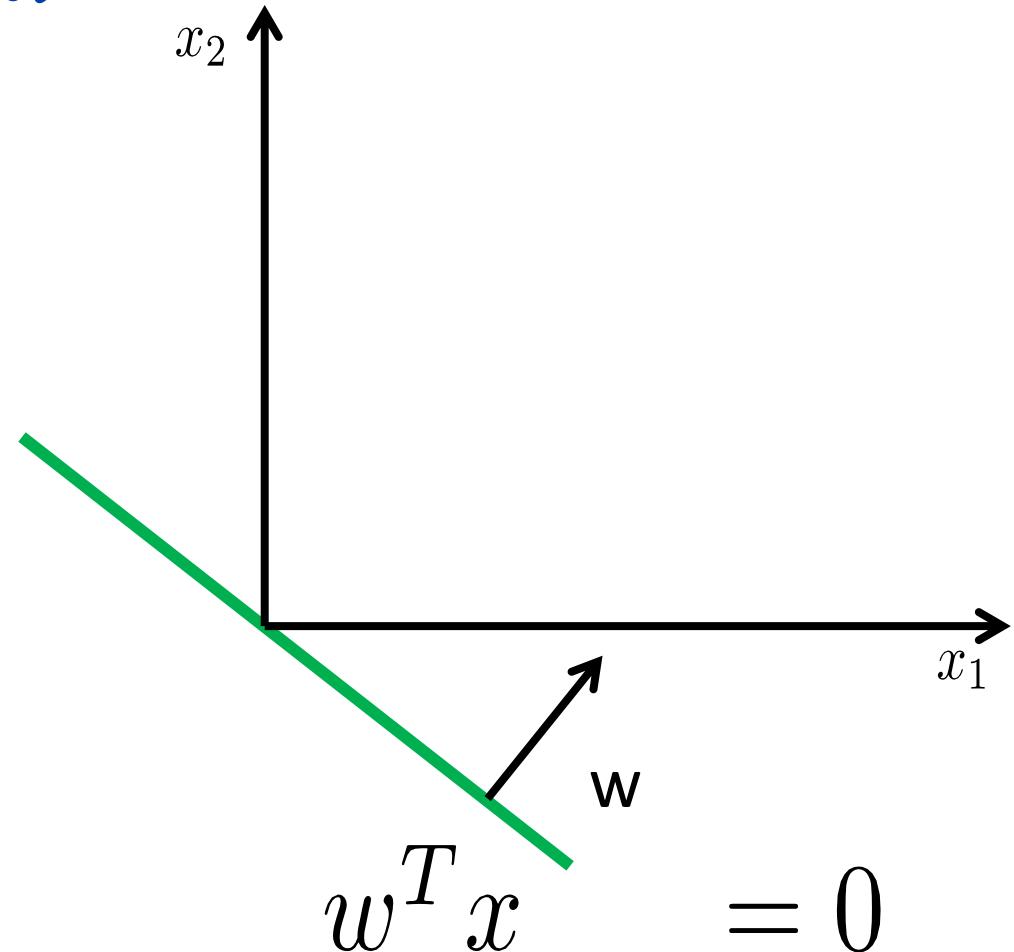
1 Goal is  
to estimate  
 $w$



# Separating Hyperplane

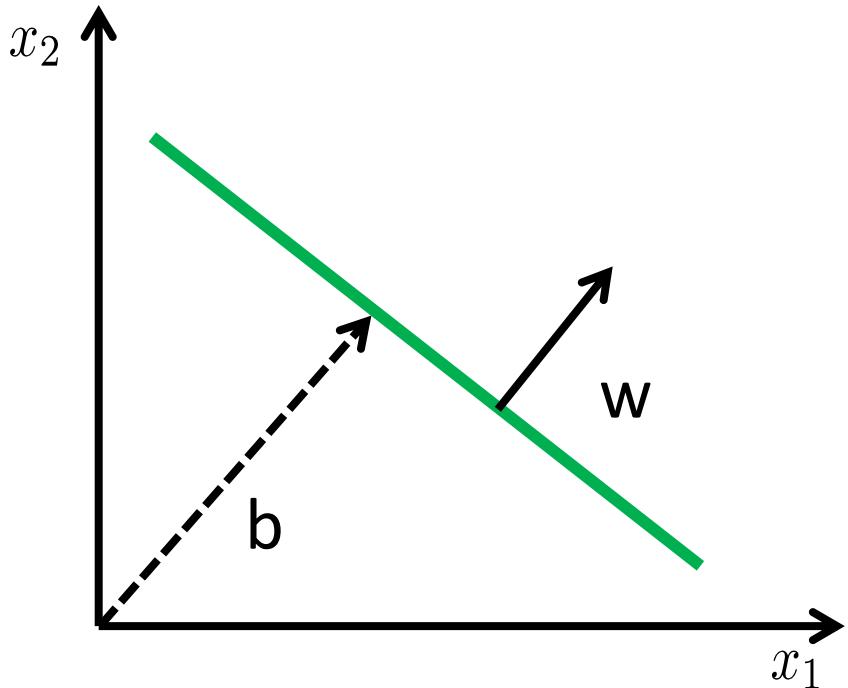
has 2 dimensions)

- $x$ : data point
- $y$ : label  $\in \{-1, +1\}$
- $w$ : weight vector



# Separating Hyperplane

- $x$ : data point
- $y$ : label  $\in \{-1, +1\}$
- $w$ : weight vector
- $b$ : bias



$$w^T x + b = 0$$

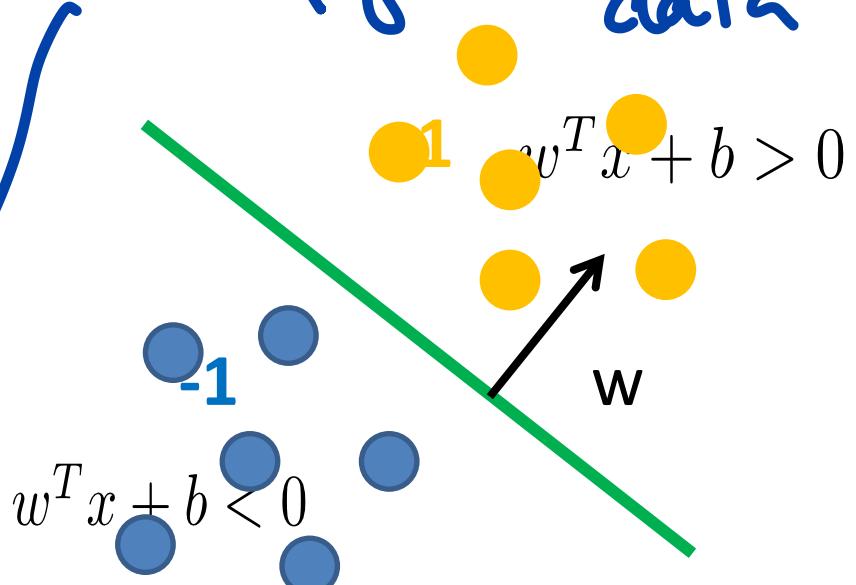
*↳ orient    ↳ shift*

# Separating Hyperplane

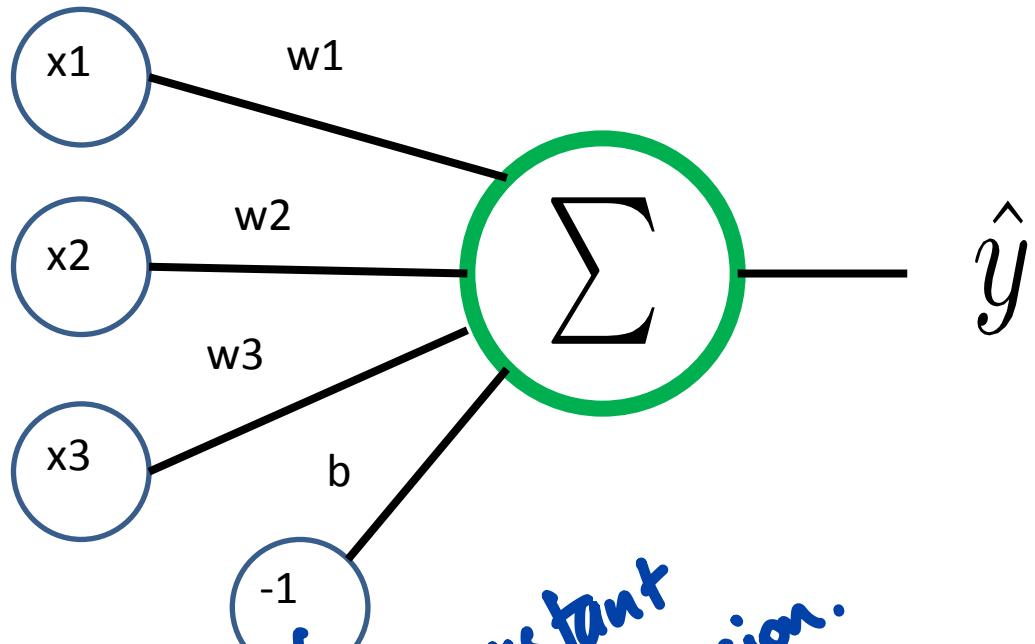
- $x$ : data point
- $y$ : label  $\in \{-1, +1\}$
- $w$ : weight vector
- $b$ : bias

Fast prediction.

Once you have the parameters  $w + b$ , can forget about the data



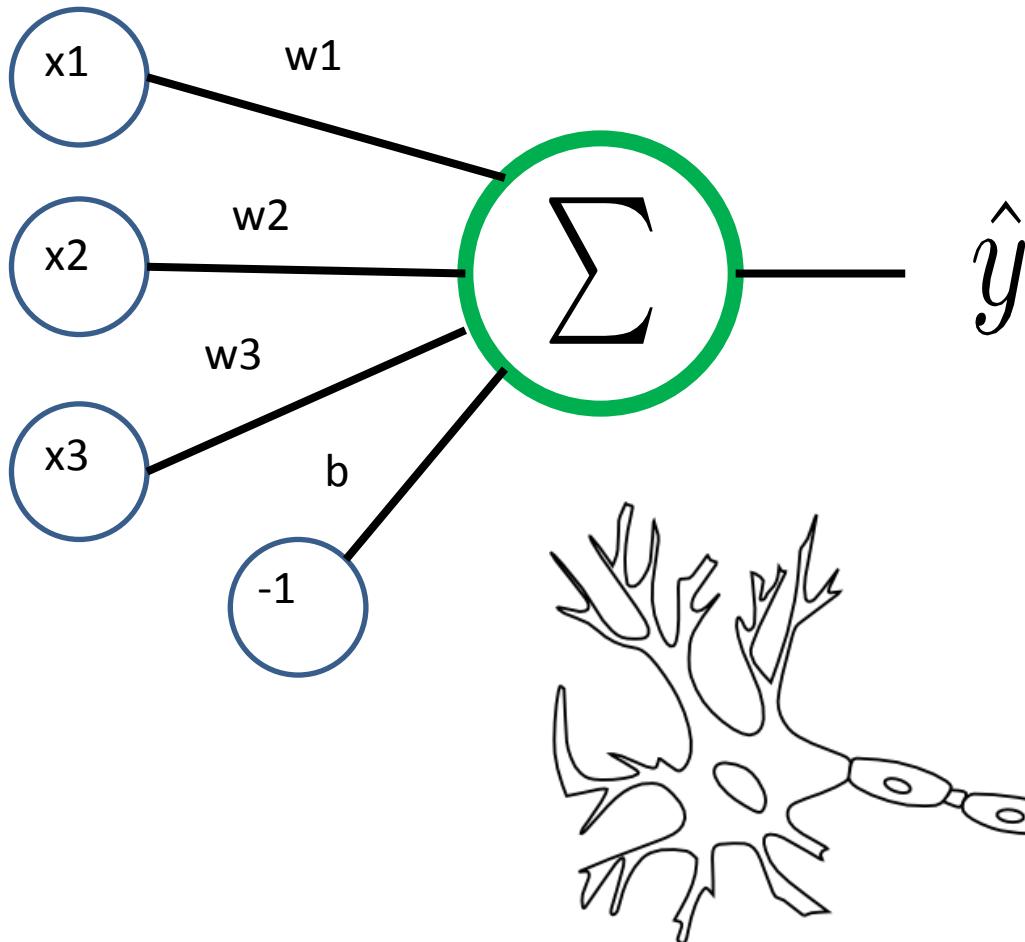
# Perceptron



↳ Constant dimension.

$$w^T x + b = 0$$

# Perceptron



- Simple model of a neuron

Mathematical description of a neuron

# Perceptron History

- invented 1957
- by Frank Rosenblatt
- the embryo of an electronic computer that [the Navy] expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence. (NYT 1958)

(<http://en.wikipedia.org/wiki/Perceptron>)

Thought they  
found a  
model of the  
brain

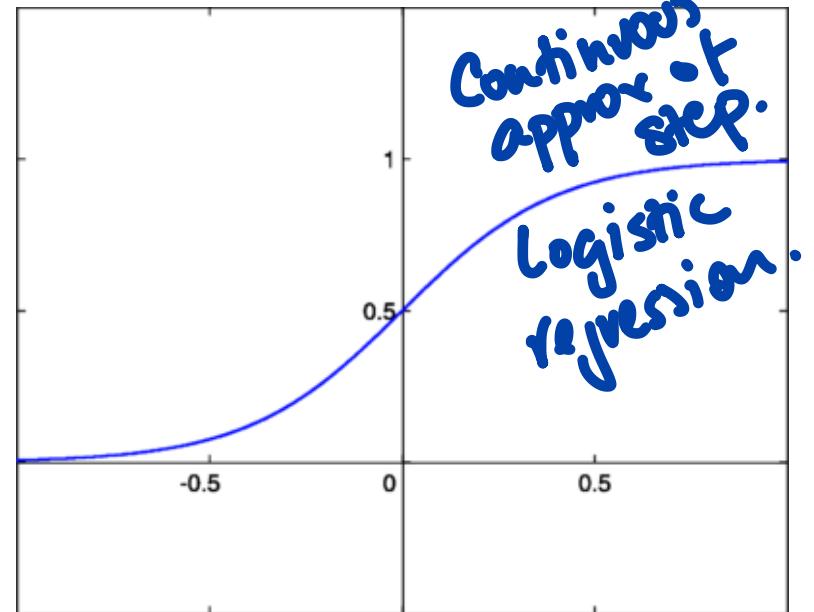
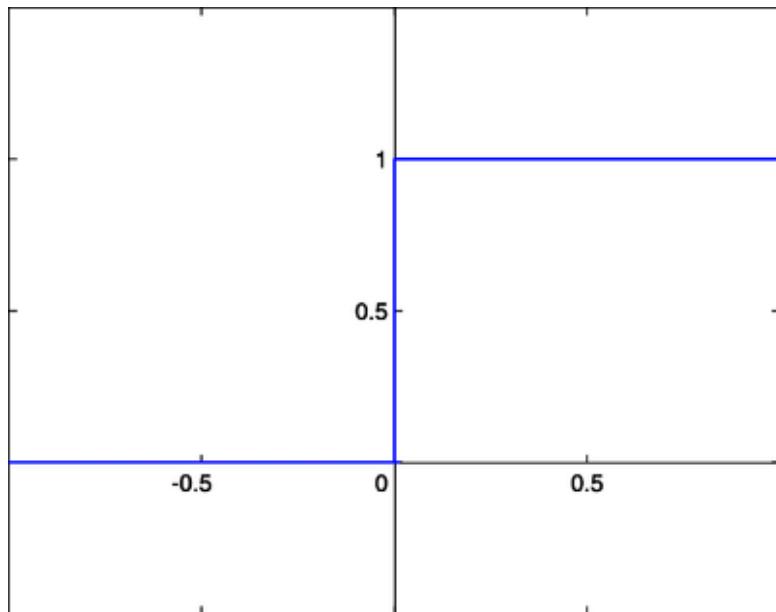
Are now the  
basis of  
deep learning.



Perceptron.mp4

[https://www.youtube.com/watch?v=cNxadbrN\\_ai&list=PLdVOMWcqwwIaygvb9ZteZ1r4Br6kRuBO](https://www.youtube.com/watch?v=cNxadbrN_ai&list=PLdVOMWcqwwIaygvb9ZteZ1r4Br6kRuBO)

# Side Note: Step vs Sigmoid Activation

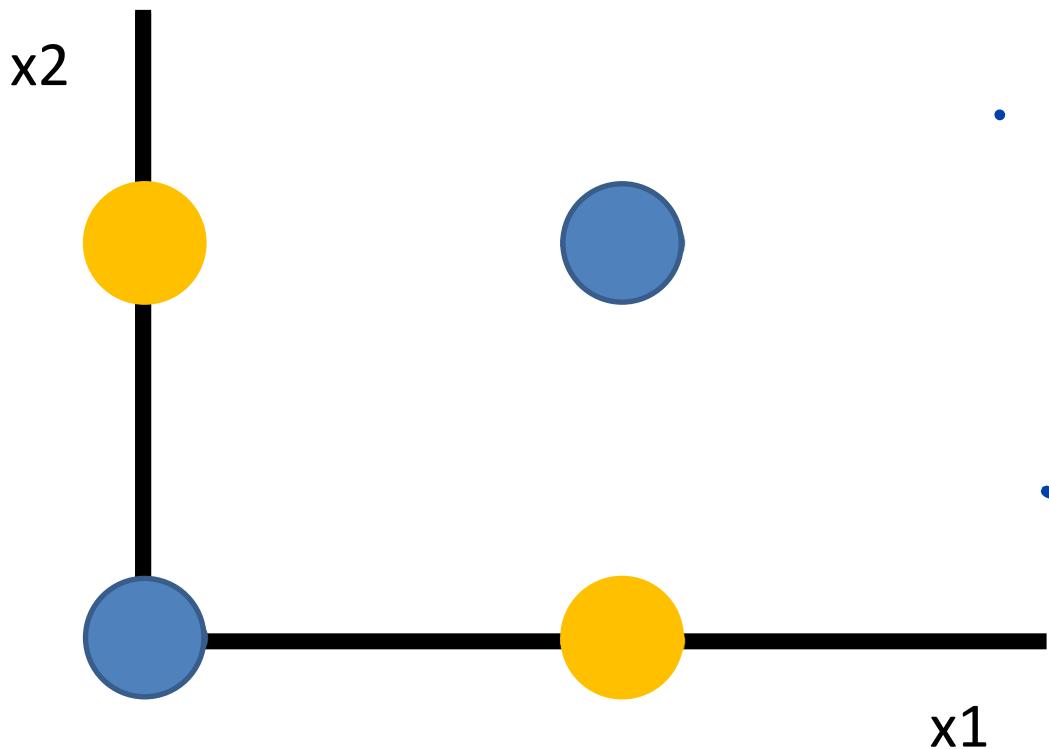


$$s(x) = \frac{1}{1 + e^{-cx}}$$

# The Critics

- 1969: Minsky and Papert publish their book “Perceptrons”  
*↳ included weakness of the method.*
- Very controversial book, some blame the book for causing the whole research area to stagnate.

# The XOR Problem



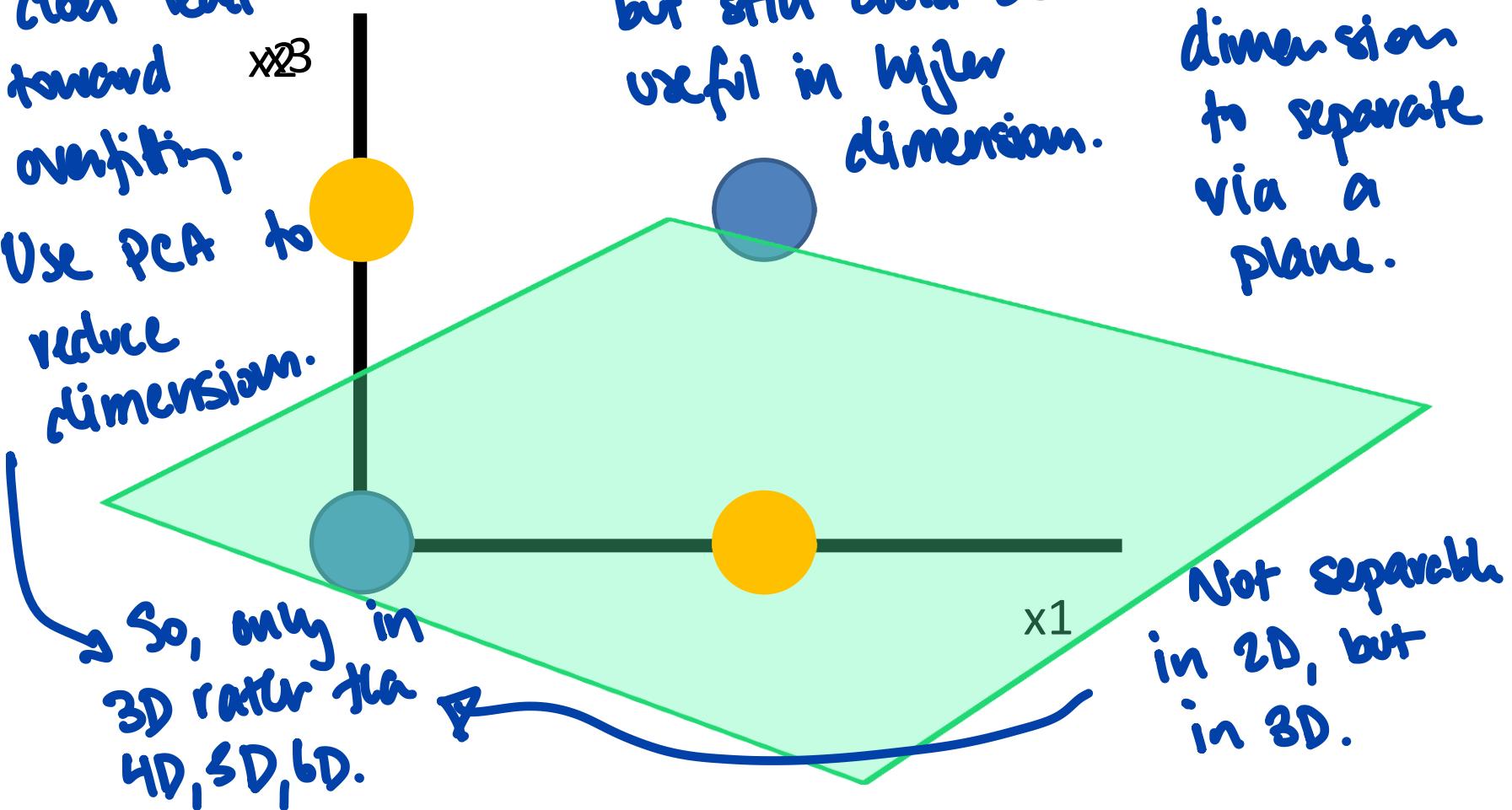
- Can't separate w/ a hyperplane
- Fundamental limitation.

# The XOR Problem

- Adding dimension can lead toward overfitting.

$x_2^3$

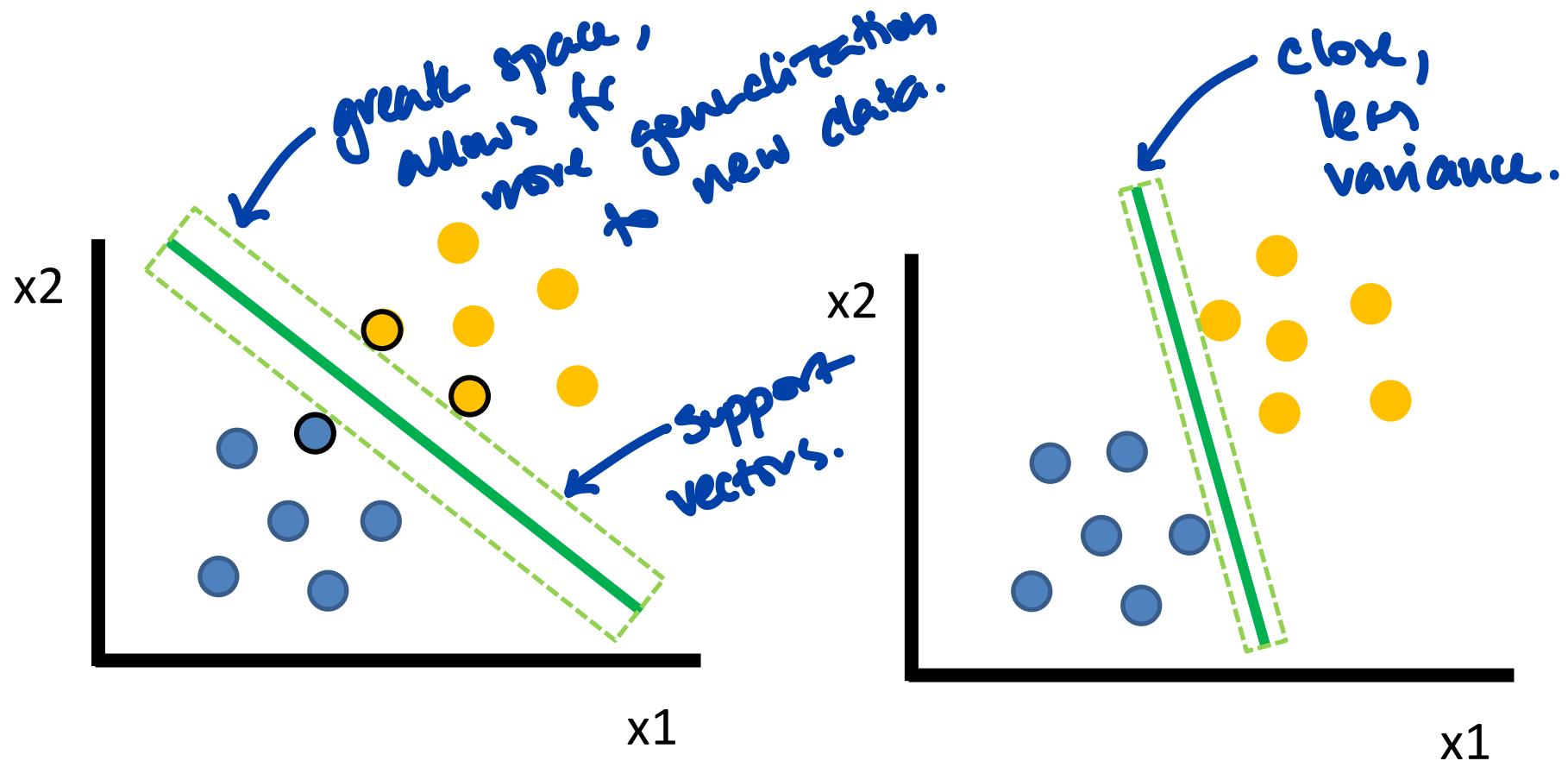
- Use PCA to reduce dimension.



# Support Vector Machine

- Widely used for all sorts of classification problems
- Some people say it is the best off the shelf classifier out there  
*Does not really require parameter tuning*

# Maximum Margin Classification



Solution depends only on the support vectors!

↳ only need these, no dath.

# Maximum Margin Classification

sci-kit learn have "C" and "gamma"

want  $x^{(i)}$   
blc  $x^{(i)}$  we know  
 $\gamma^{(i)}$

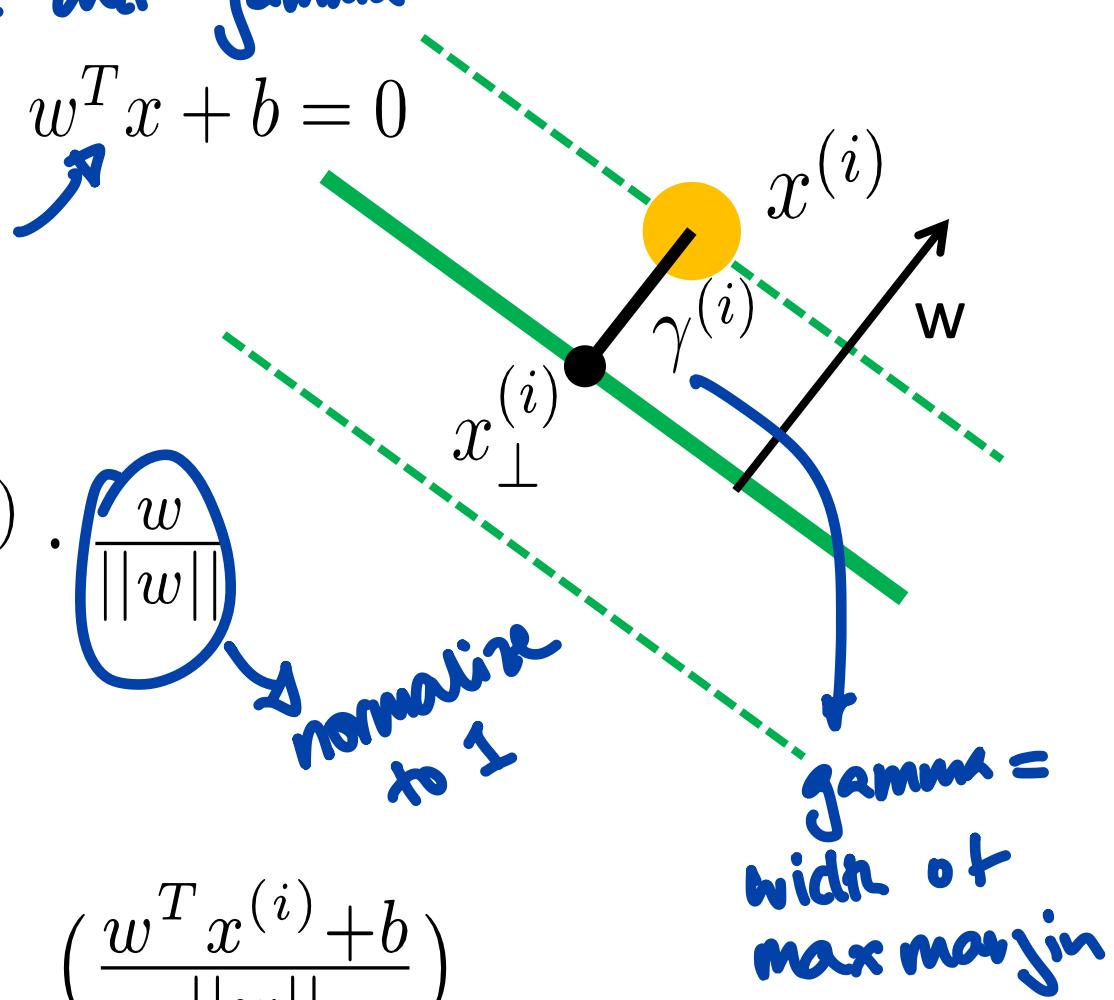
margin:

$$x_{\perp}^{(i)} = x^{(i)} - \gamma^{(i)} \cdot \frac{w}{\|w\|}$$

$$w^T x_{\perp}^{(i)} + b = 0$$

$$\gamma^{(i)} = \left( \frac{w^T x^{(i)} + b}{\|w\|} \right)$$

Combo of eqns above.



# Maximum Margin Classification

$$\gamma^{(i)} = y^{(i)}(w^T x + b)$$

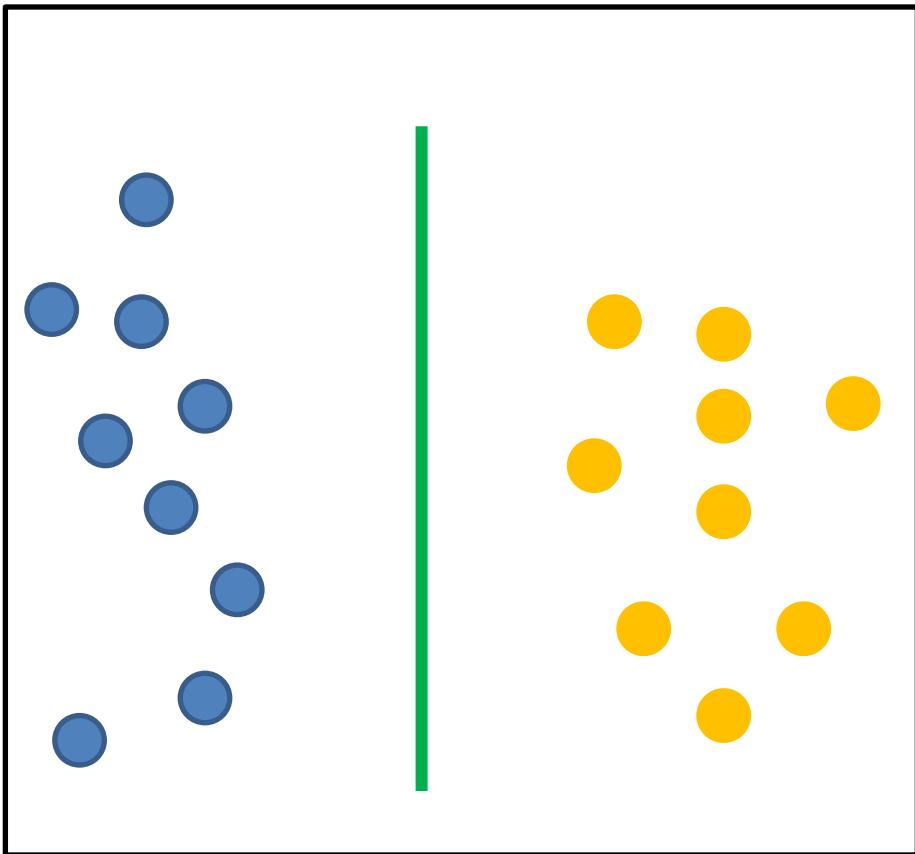
*tells pos & neg away from hyperplane.*

$$\begin{aligned} & \max_{\gamma, w, b} \quad \gamma \\ \text{s.t.} \quad & y^{(i)}(w^T x^{(i)} + b) \geq \gamma, \quad i = 1, \dots, m \\ & \|w\| = 1. \end{aligned}$$

*Software does this;  
we don't write our  
own optimizers.*

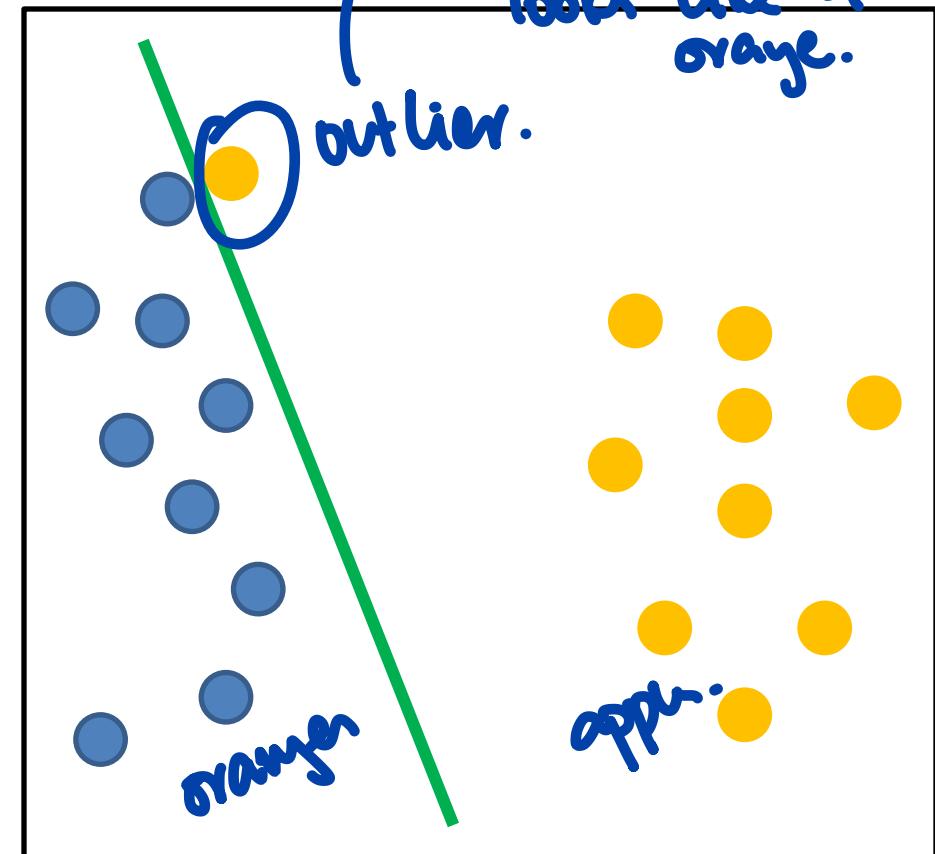
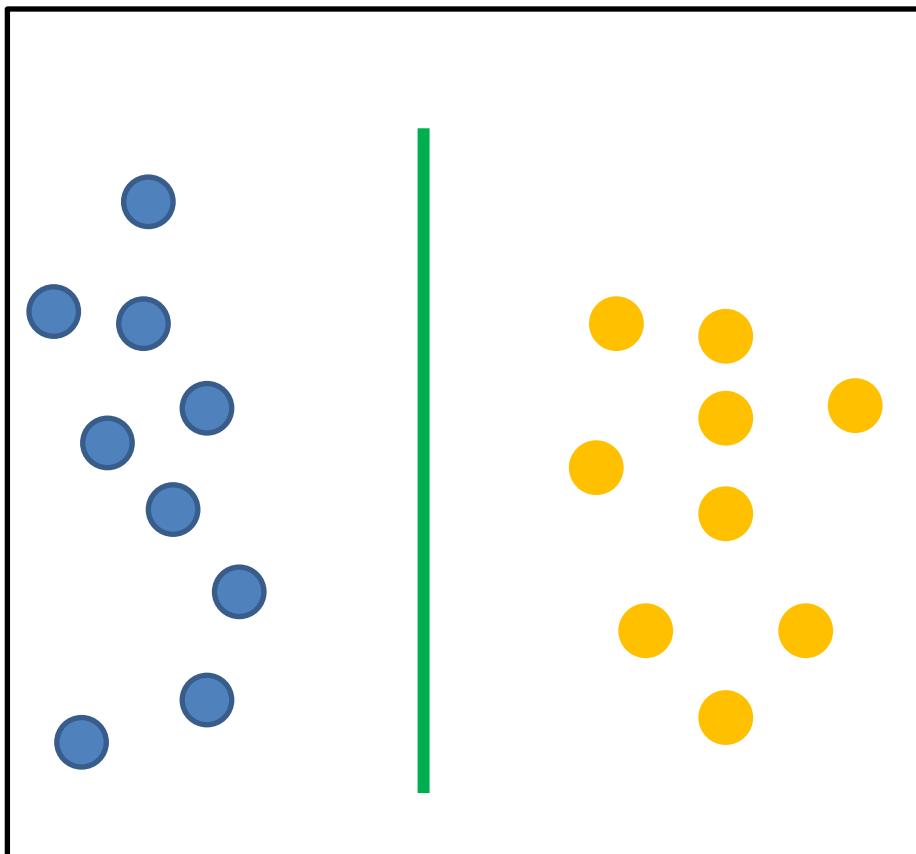
**non-convex**

# This Is Kind of Odd



- Which data points do we care the most about? **Rare**
- What would those samples look like?

# Two Very Similar Problems



# What about outliers?

Don't want a perfect classification on training data.

$\xi_i$ : slack variables

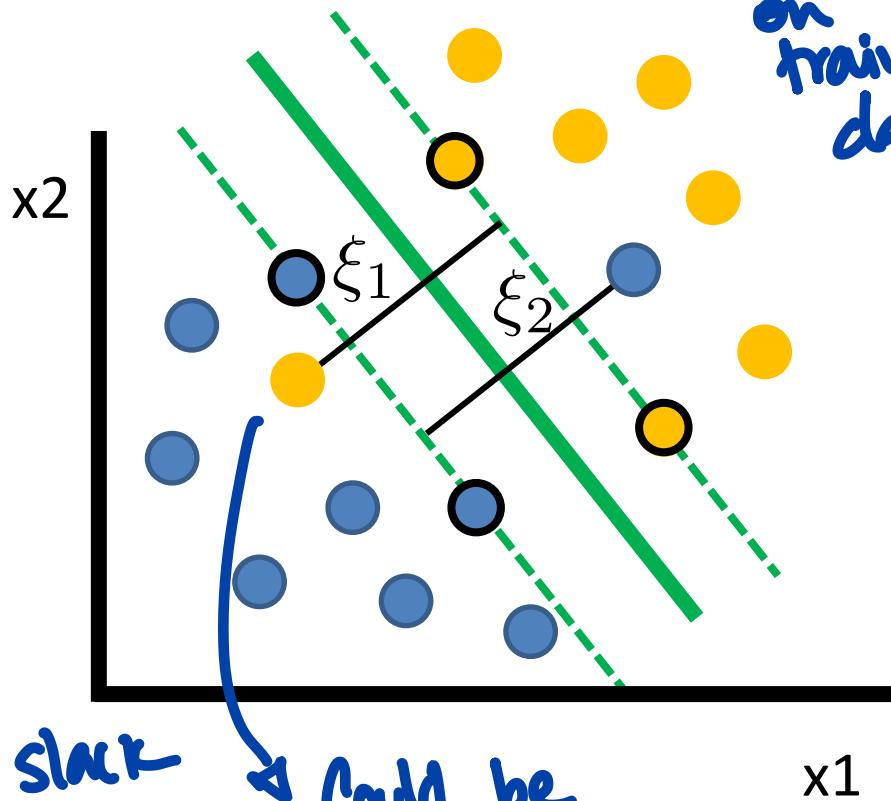
$$\min_{w,b,\xi} \frac{1}{2} \|w\|^2 + C \sum_i \xi_i$$

subject to:

$$y^{(i)}(w^T x^{(i)} + b) \geq 1 - \xi_i$$

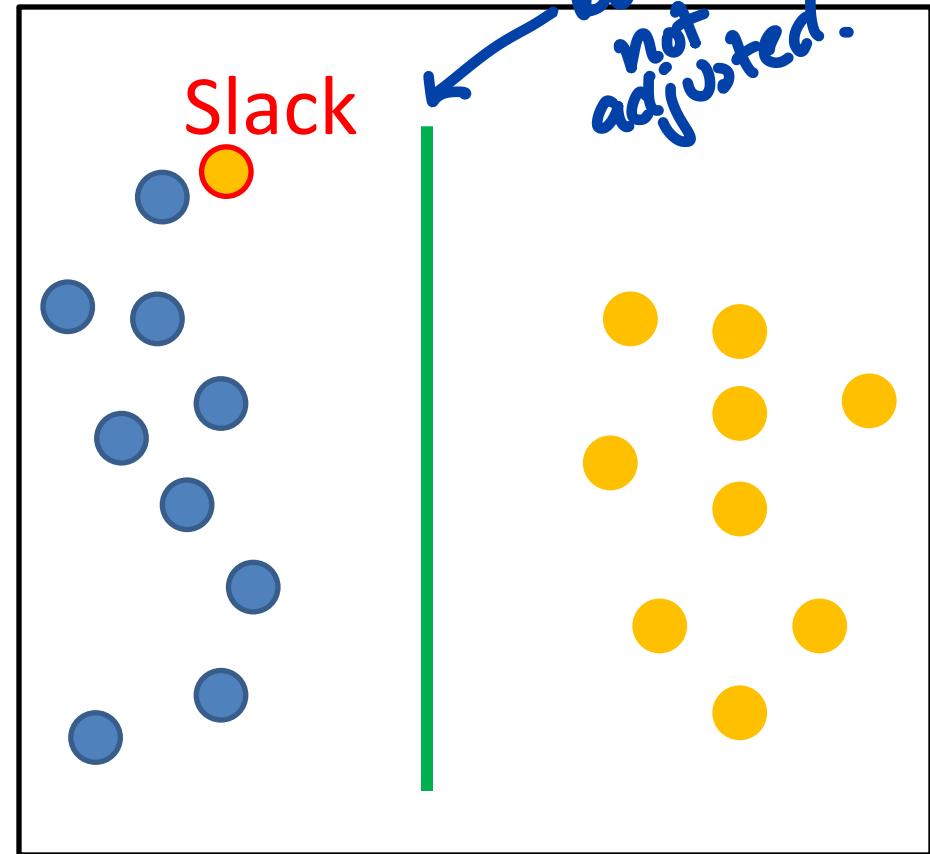
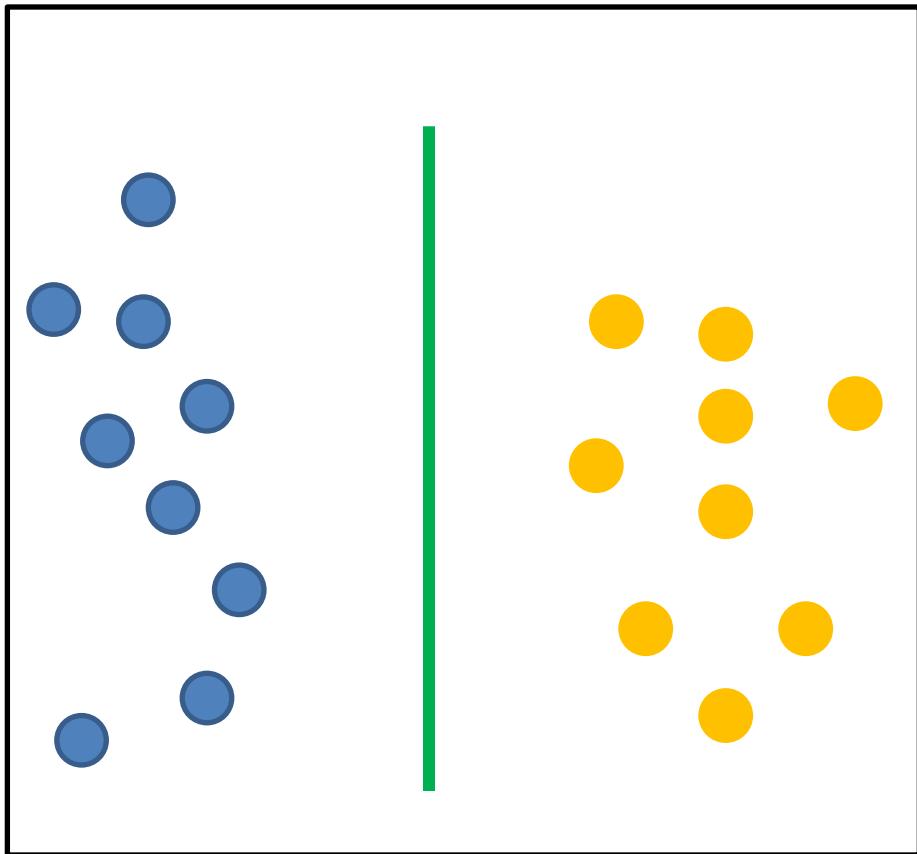
$$(i = 1, \dots, n)$$

Large  $C \rightarrow$  not much slack  
Small  $C \rightarrow$  more room (slack)



Could be mislabelled.

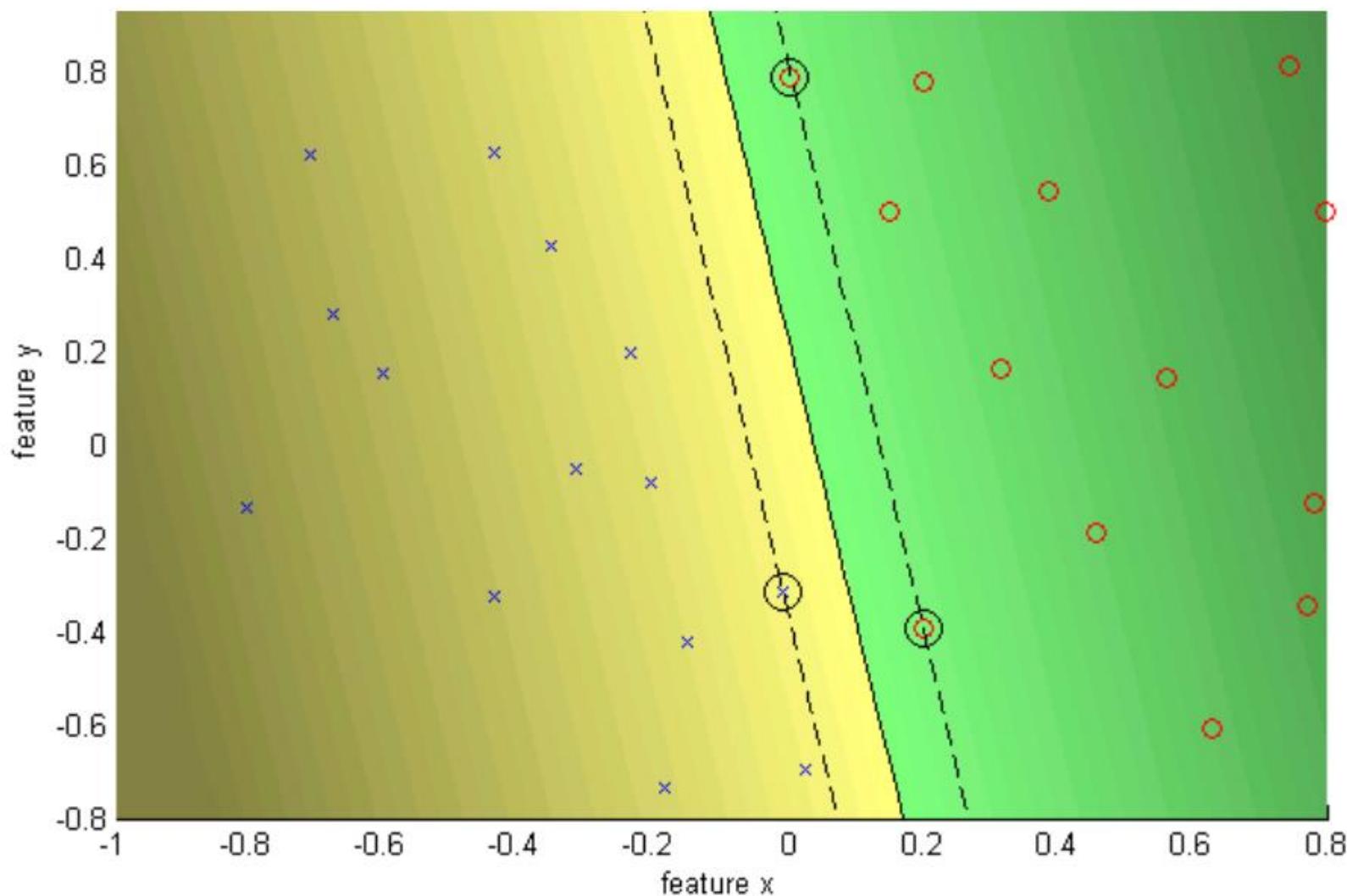
# Two Very Similar Problems



*boundary  
not  
adjusted.*

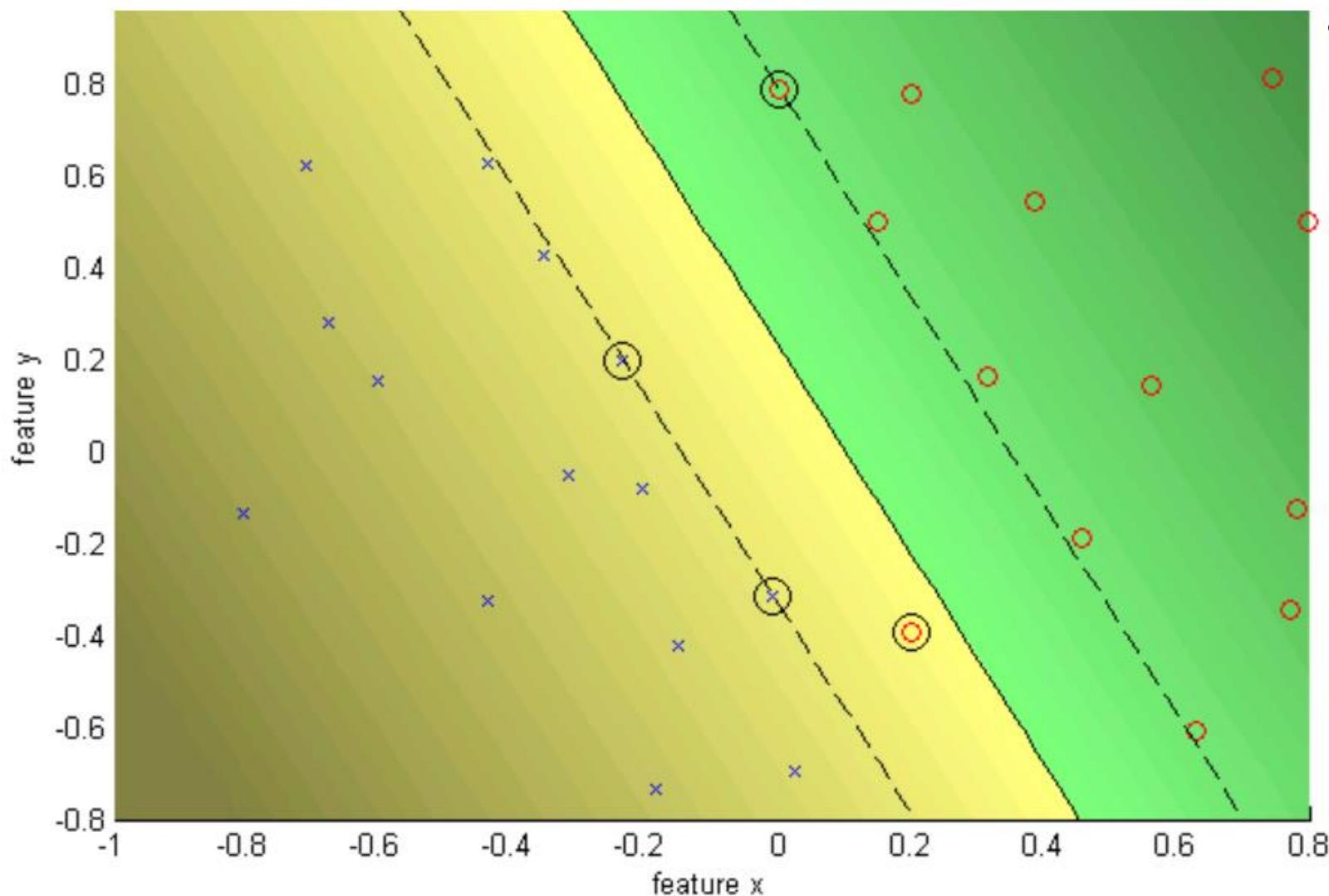
# Hard Margin ( $C = \text{Infinity}$ )

no mistake;  
margin  
small.

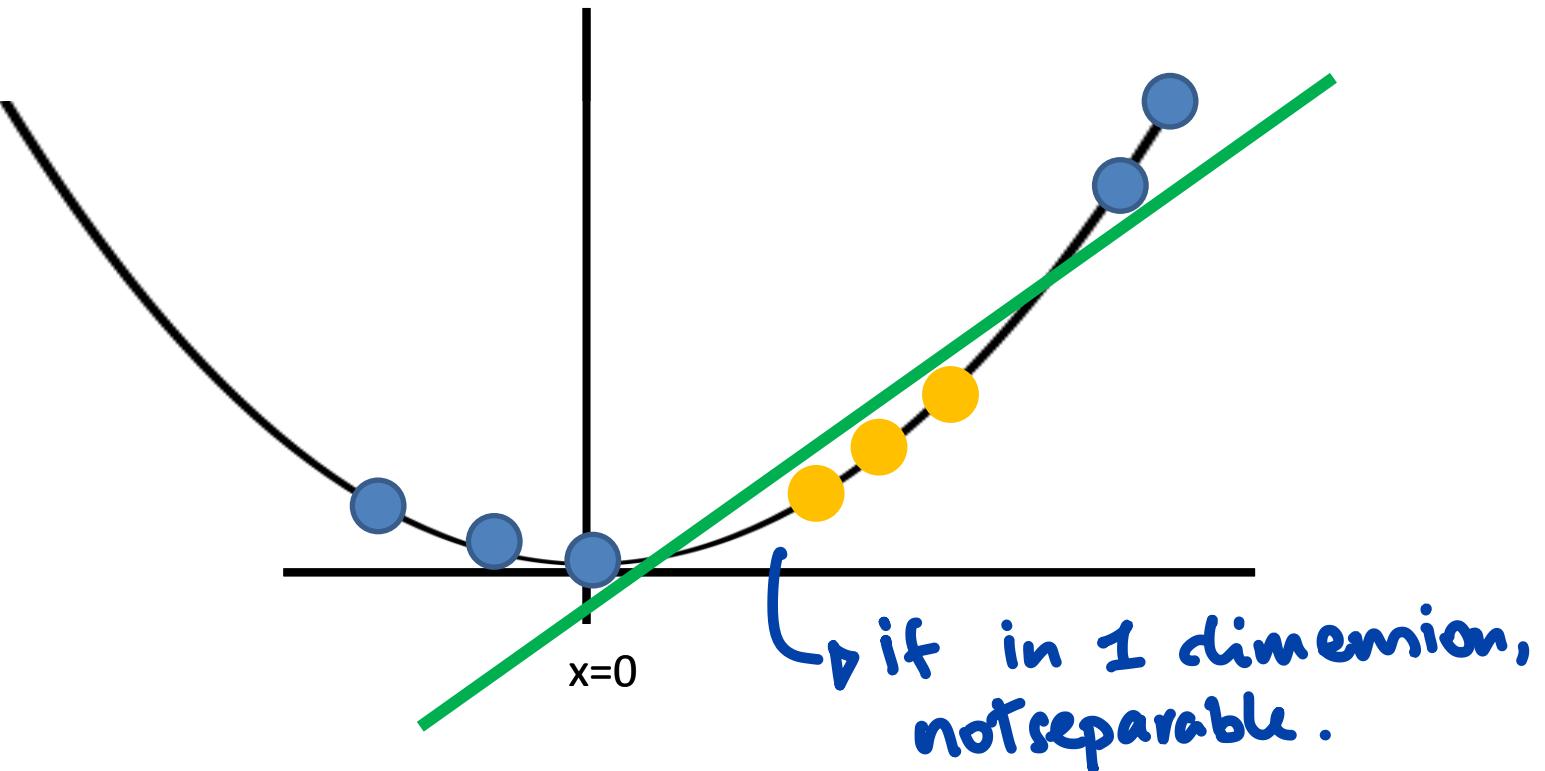


# Soft Margin ( $C = 10$ )

Not perfect  
but great  
slack.



# XOR problem revised



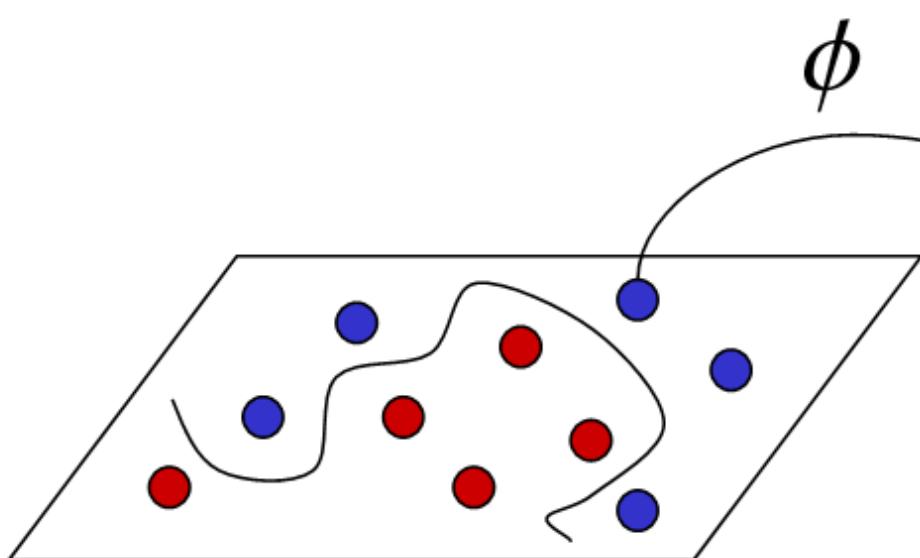
Did we add information to make the problem separable?

↳ No. Just dimension A.

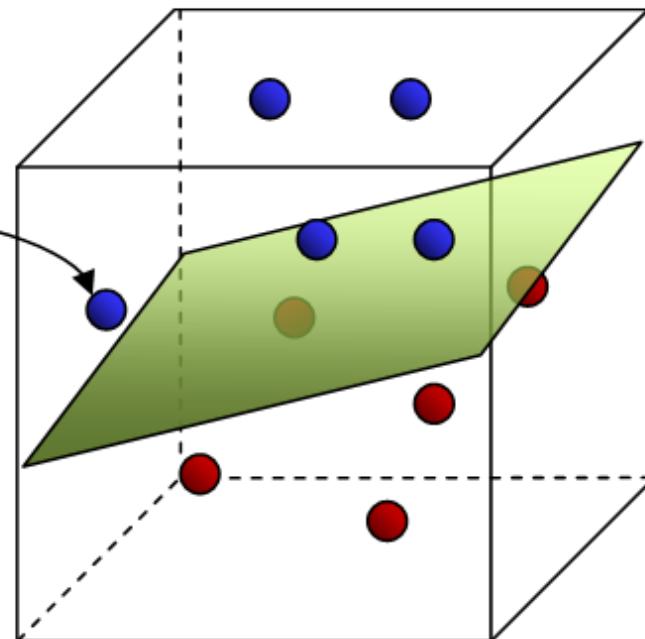
Cross-validation to determine  $m$  # of dimension.

# Non-Linear Decision Boundary

↳ Adding dimension makes the hyperplane non-linear.



**Input Space**



**Feature Space**

# SVM with a polynomial Kernel visualization

Created by:  
Udi Aharoni

Quadratic  
Expansion  
of 2D  
Point

# Quadratic Kernel

$$x = (x_1, x_2)$$

$$\Phi(x) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

$$\Phi(x) \cdot \Phi(z) = 1 + 2 \sum_{i=1}^d x_i z_i$$

$$\checkmark \quad \text{need } m \text{ dot product; high computational cost.}$$
$$+ \sum_{i=1}^d x_i^2 z_i^2 + 2 \sum_{i=1}^d \sum_{j=i+1}^d x_i x_j z_i z_j$$

$$= (1 + x \cdot z)^2$$

Kernel trick  
→ Just use this;  
no  $\Sigma$

# Kernel Functions

$$K(x, z) = \Phi(x) \cdot \Phi(z)$$

- Polynomial:

$$K(x, z) = (1 + x \cdot z)^s$$

can do any dimension 2+  
longer to tune;  
might be "d" in  
sklearn

- Radial basis function (RBF):  
*(very popular)*

$$K(x, z) = \exp(-\gamma(x - z)^2)$$

work b/c don't need th actual point, just th dot product

# So what is the excitement?

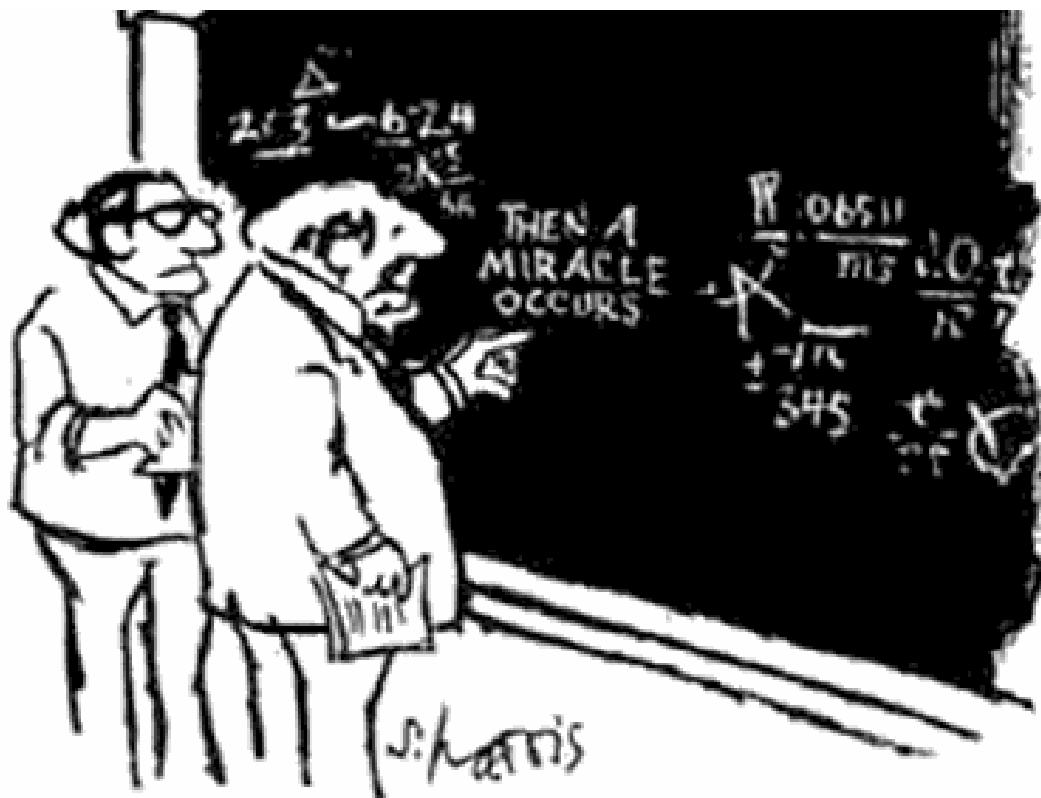
$$\max_{\alpha} \sum$$

$$\text{s.t. } \alpha_i$$

$$\sum$$

$$\arg \Gamma$$

$$\text{s.t. } y$$



$$(i)^T x^{(j)}$$

Andrew  
Ng  
for this  
explanation

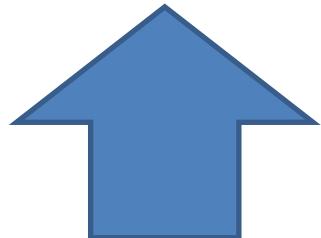
"I THINK YOU SHOULD BE MORE EXPLICIT  
HERE IN STEP TWO."

# So what is the excitement?

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j x^{(i)T} x^{(j)}$$

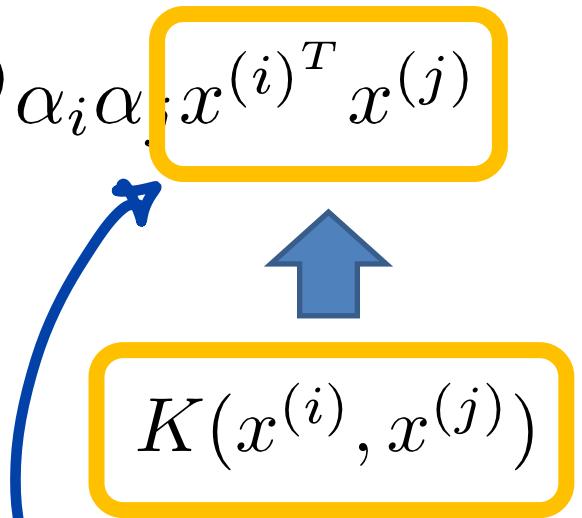
$$\text{s.t. } \alpha_i \geq 0, \quad i = 1, \dots, m$$

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0$$



$$\arg \min_{w,b} \frac{1}{2} \|w\|^2$$

$$\text{s.t. } y^{(i)} (w^T x^{(i)} + b) \geq 1$$



# Prediction

$$w^T x + b = \sum_{i=1}^m \alpha_i y^{(i)} \langle x^{(i)}, x \rangle + b.$$

- Again we can use the kernel trick!
- Prediction speed depends on number of support vectors

# The Miracle Explained

- Andrew Ng does this really well
- [http://cs229.stanford.edu/notes/cs229-  
notes3.pdf](http://cs229.stanford.edu/notes/cs229-notes3.pdf)
- Course is also on Youtube, ItunesU, etc.

# Kernel Trick for SVMs

- Arbitrary many dimensions
- Little computational cost
- Maximal margin helps with curse of dimensionality

PCA → SVM.

# Face Recognition

pred: Colin\_Powell  
true: Colin\_Powell



pred: George\_W\_Bush  
true: George\_W\_Bush



pred: Colin\_Powell  
true: Colin\_Powell



pred: Tony\_Blair  
true: Tony\_Blair



pred: George\_W\_Bush  
true: George\_W\_Bush



pred: Colin\_Powell  
true: Colin\_Powell



pred: George\_W\_Bush  
true: George\_W\_Bush



pred: George\_W\_Bush  
true: George\_W\_Bush



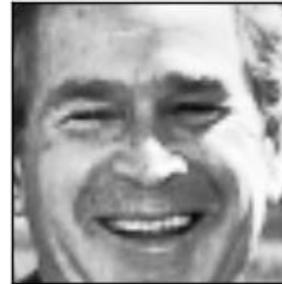
pred: Tony\_Blair  
true: Tony\_Blair



pred: Colin\_Powell  
true: Colin\_Powell



pred: George\_W\_Bush  
true: George\_W\_Bush



pred: Donald\_Rumsfeld  
true: Donald\_Rumsfeld



# Face Recognition

- Load image data
- Put your test data aside
- Extract Eigenfaces (PCA on images)
- Train SVM
- Evaluate performance
- Red are cross validation steps

Cross validation.  
Do these steps for all



SVM\_sign\_language.mp4

Jhon Gonzalez

[https://www.youtube.com/watch?v=cxHMgl2\\_5zg](https://www.youtube.com/watch?v=cxHMgl2_5zg)

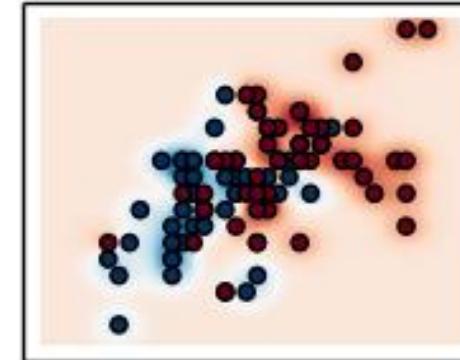
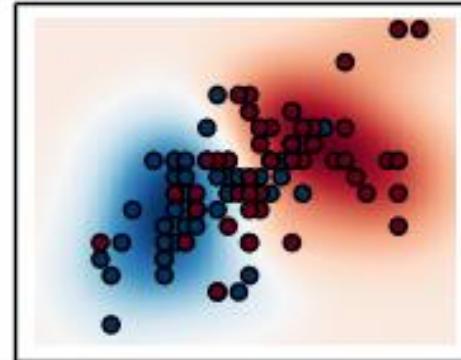
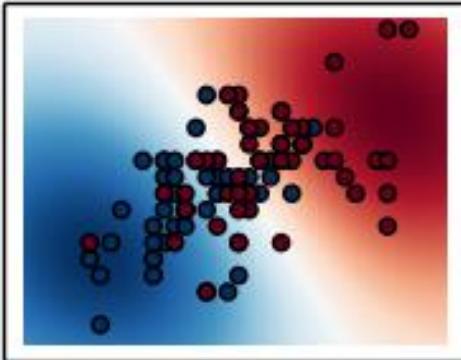
→ tune parameters in exponential steps.

gamma

gamma=10<sup>-1</sup>, C=10<sup>-2</sup> gamma=10<sup>0</sup>, C=10<sup>-2</sup> gamma=10<sup>1</sup>, C=10<sup>-2</sup>

larger

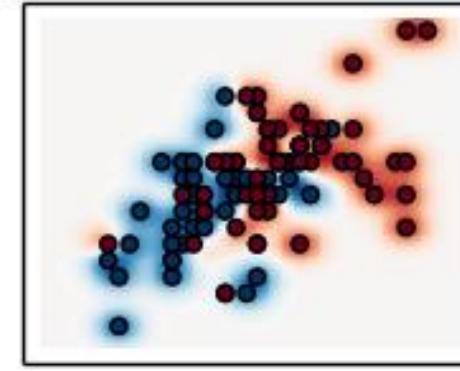
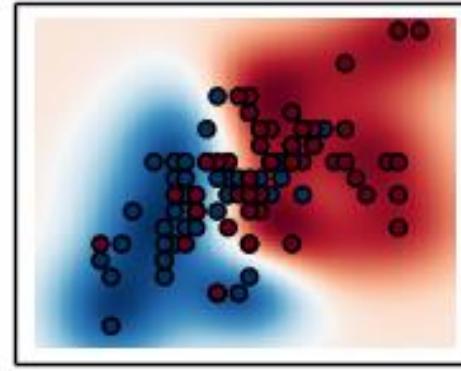
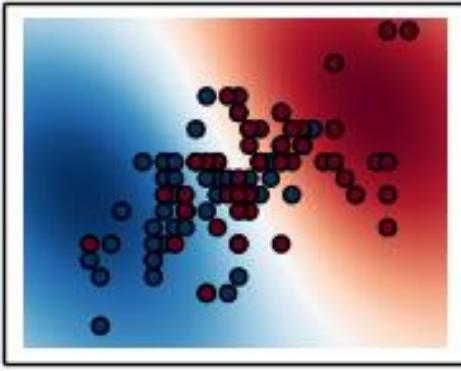
C



gamma=10<sup>-1</sup>, C=10<sup>0</sup>

gamma=10<sup>0</sup>, C=10<sup>0</sup>

gamma=10<sup>1</sup>, C=10<sup>0</sup>

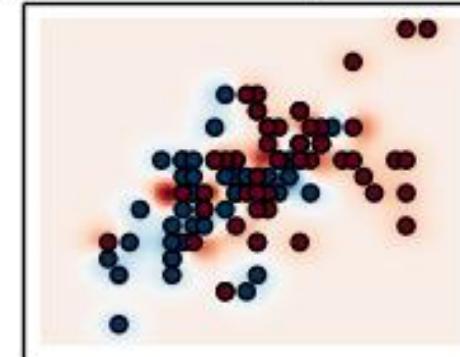
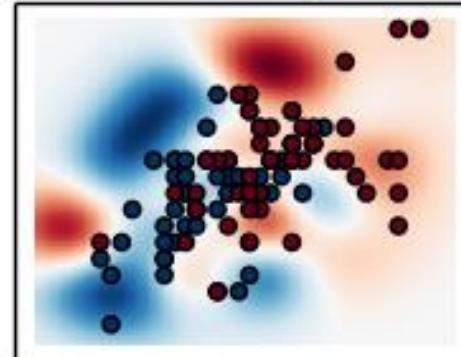
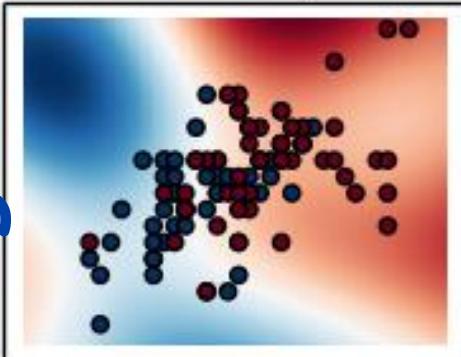


gamma=10<sup>-1</sup>, C=10<sup>2</sup>

gamma=10<sup>0</sup>, C=10<sup>2</sup>

gamma=10<sup>1</sup>, C=10<sup>2</sup>

large  
(err  
size)



Convol  
cloudie  
+  
outfit  
to  
train

# Tips and Tricks

- SVMs are not scale invariant (*mile vs. nanometer*).
- Check if your library normalizes by default
- Normalize your data
  - mean: 0 , std: 1
  - map to [0,1] or [-1,1]
- Normalize test set in same way!

*↳ likely not.*

# Tips and Tricks

- RBF kernel is a good default
- For parameters try exponential sequences
- Read:

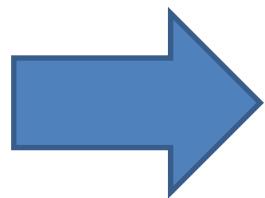
Chih-Wei Hsu et al., “A Practical Guide to Support Vector Classification”,  
Bioinformatics (2010)

# SVM vs KNN

- What are the main key differences?

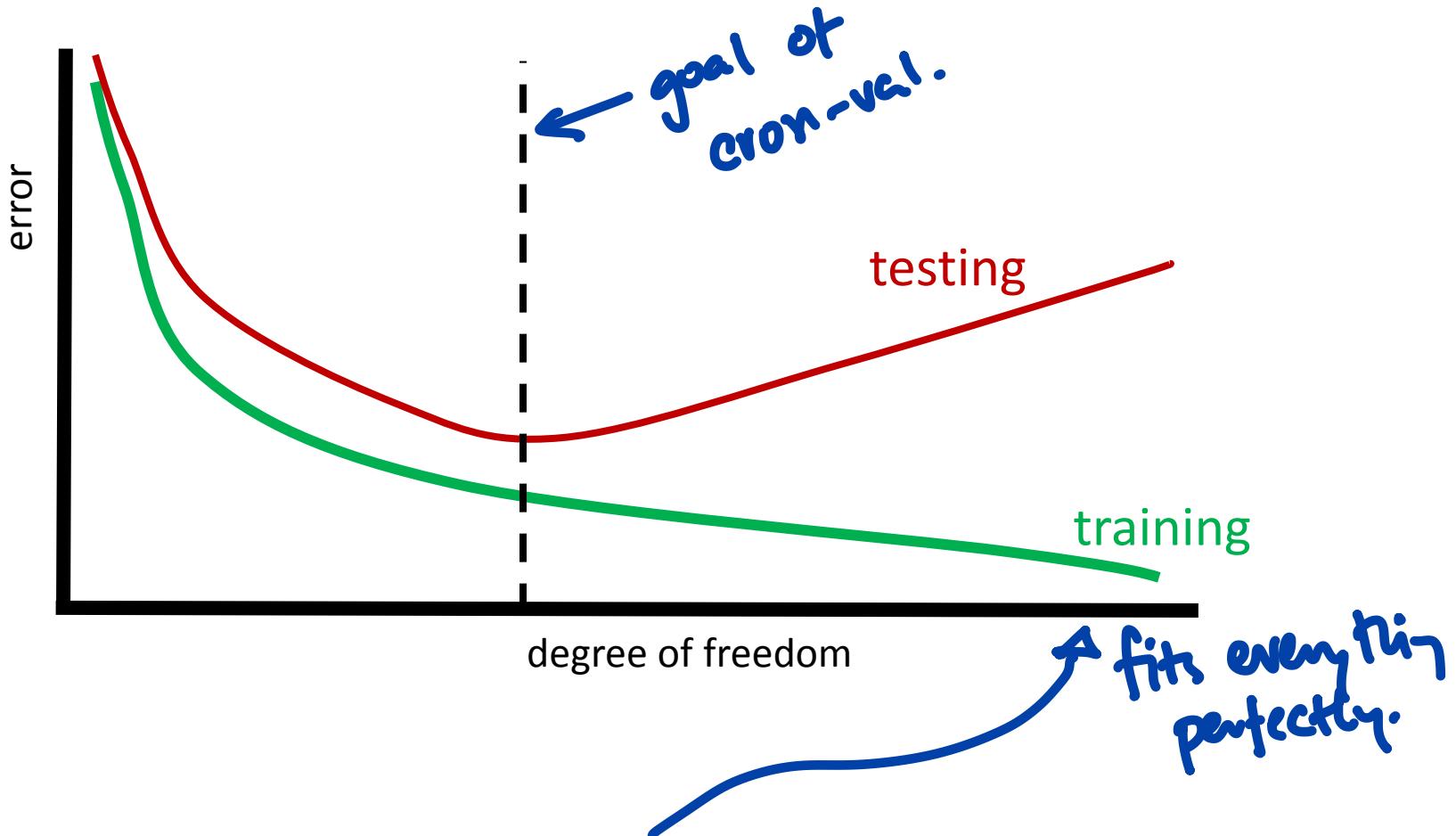
# Parameter Tuning

- Given a classification task
- Which kernel ?
- Which kernel parameter values?
- Which value for C?



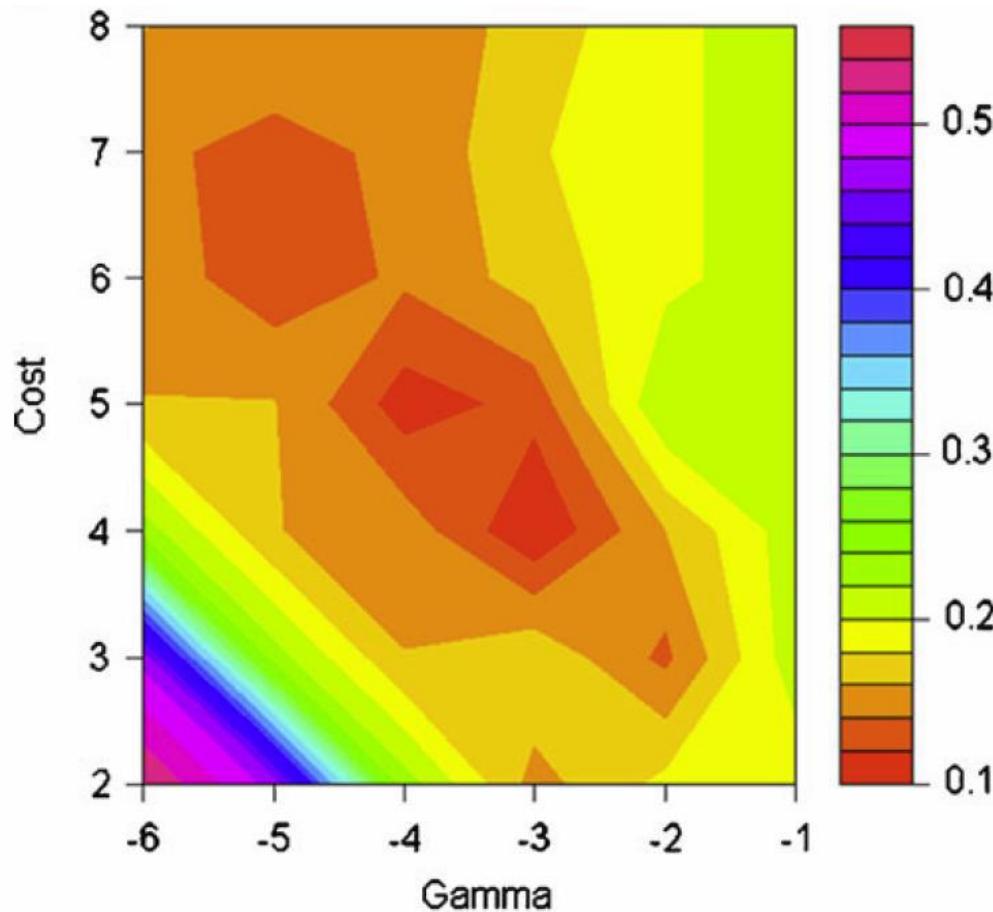
Try different combinations  
and take the **best**.

# Train vs. Test Error



Where is KNN on this graph for  $K=1$ , or for  $K=\infty$ ?

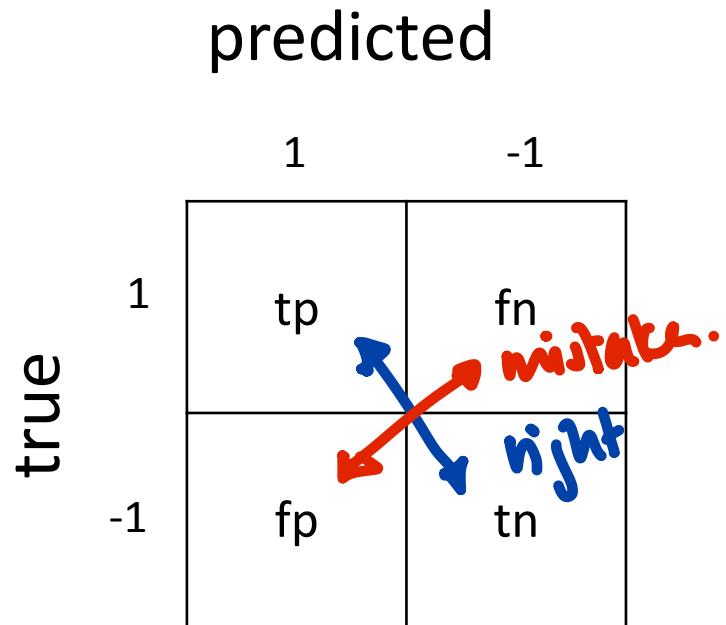
# Grid Search



Zang et al., "Identification of heparin samples that contain impurities or contaminants by chemometric pattern recognition analysis of proton NMR spectral data", Anal Bioanal Chem (2011)

# Error Measures

- True positive (tp)
- True negative (tn)
- False positive (fp)
- False negative (fn)



# TPR and FPR

- True Positive Rate:

$$\frac{tp}{tp + fn}$$

↳ all positive examples.

- False Positive Rate:

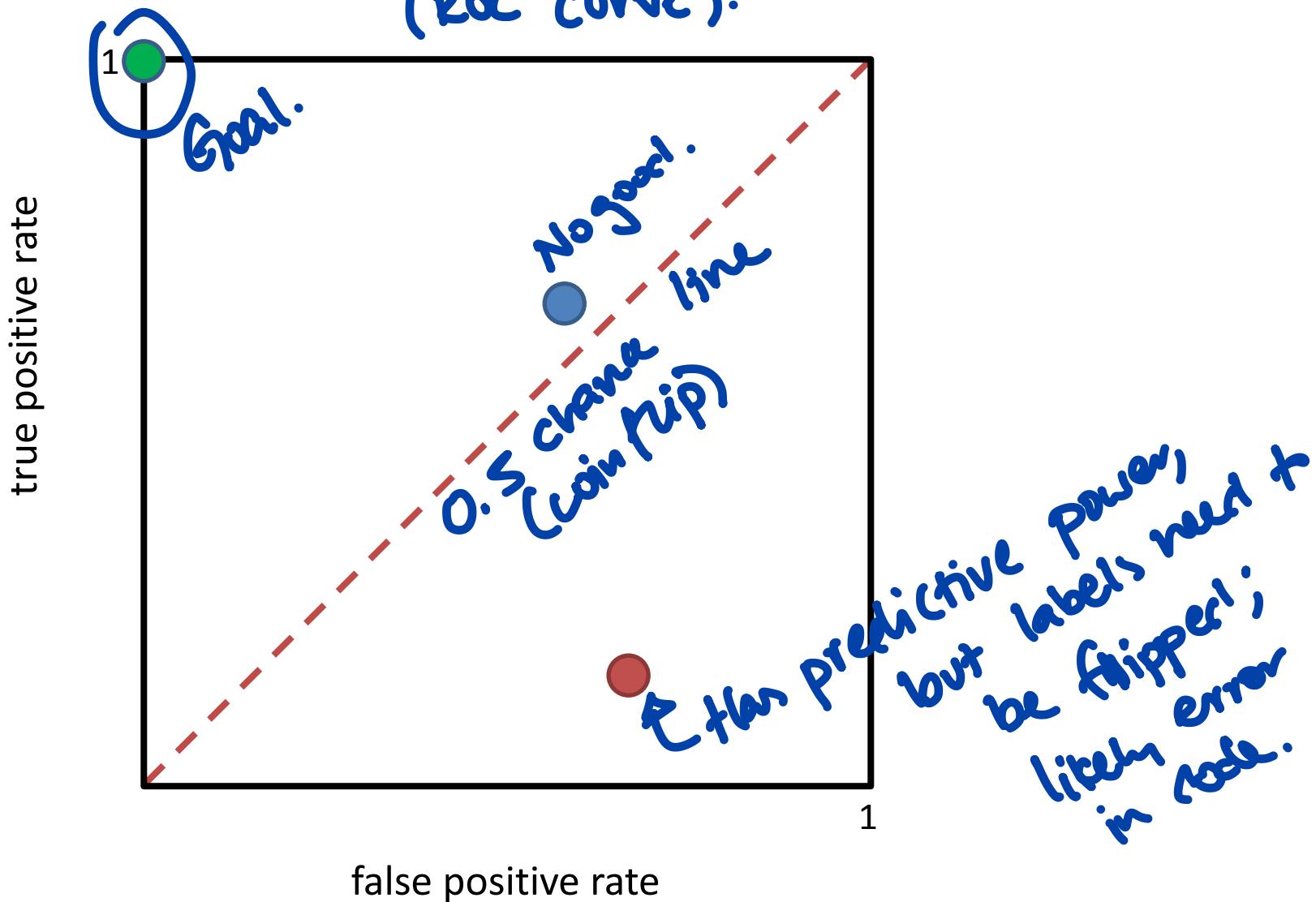
$$\frac{fp}{fp + tn}$$

↳ all negative examples.

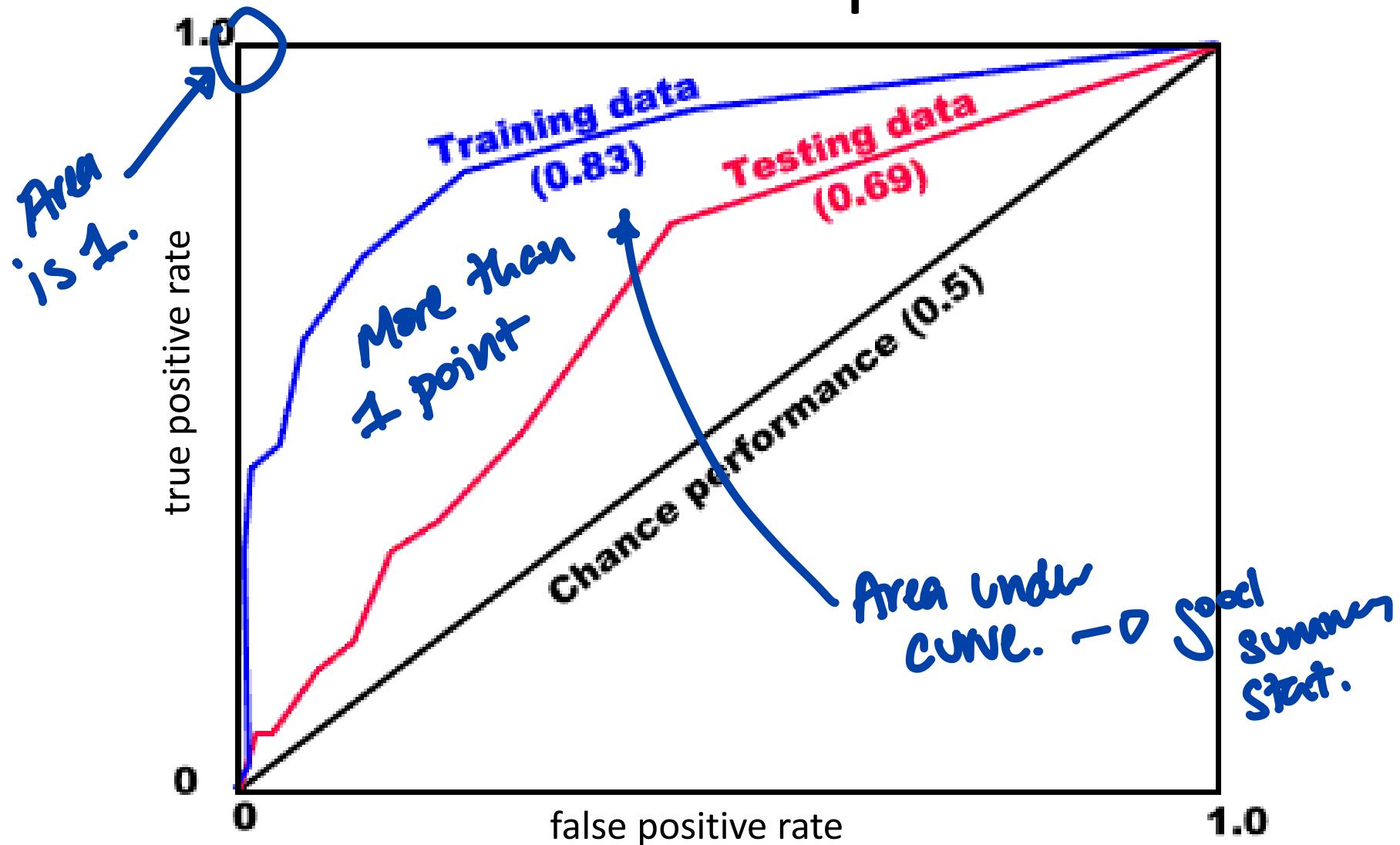
predicted

|   |    | 1  | -1 |
|---|----|----|----|
| 1 | 1  | tp | fn |
|   | -1 | fp | tn |

# Reciever Operating Characteristic (ROC curve).



# ROC Example



# Precision Recall

Good for imbalanced data sets, like churn.  
predicted

- Recall: 
$$\frac{tp}{tp + fn}$$

(TPR)  
(> same;  
diff name)

- Precision: 
$$\frac{tp}{tp + fp}$$

↳ not the same as F1R

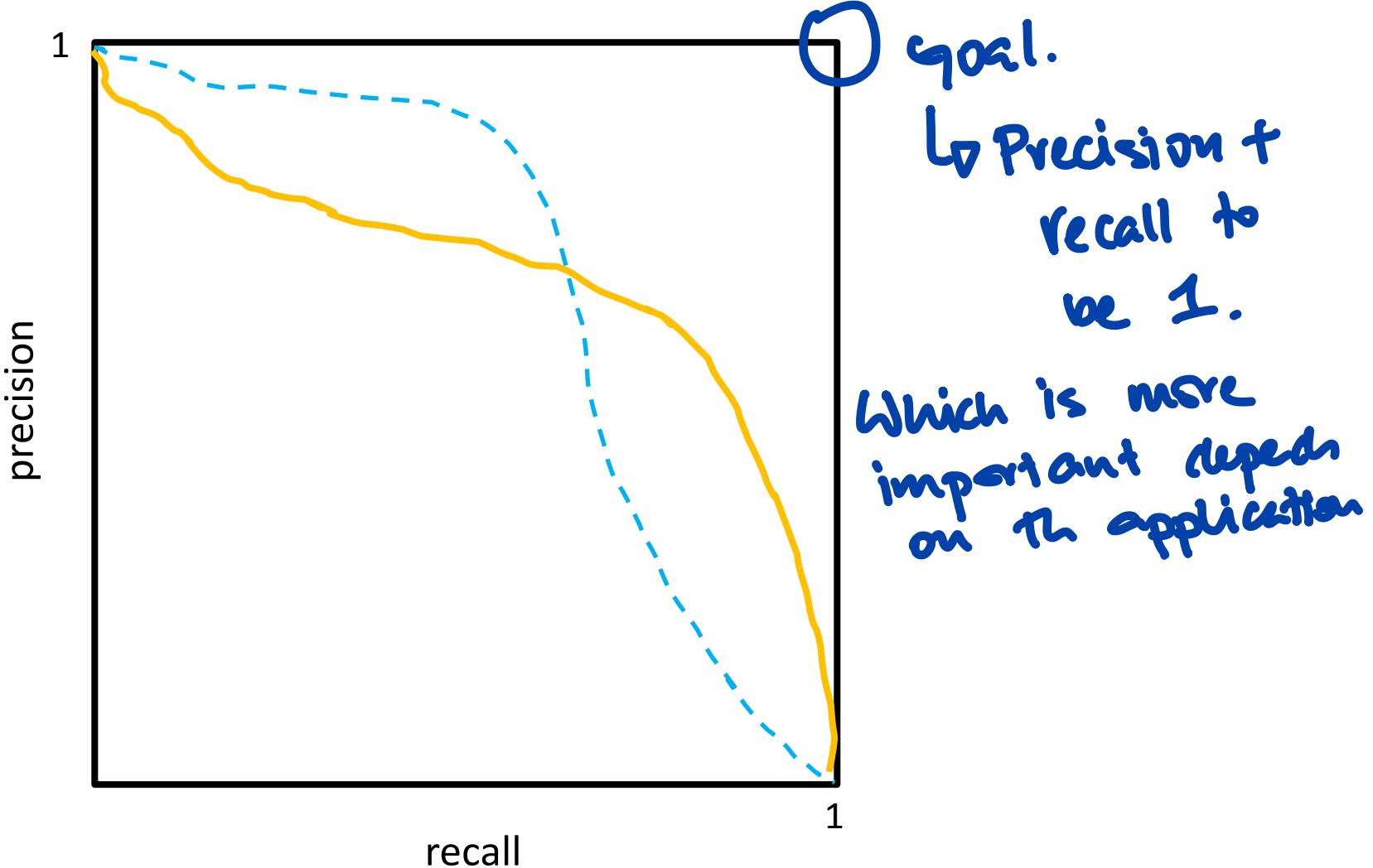
true

| 1  | -1 |
|----|----|
| 1  | fn |
| -1 | tn |
| tp | fp |

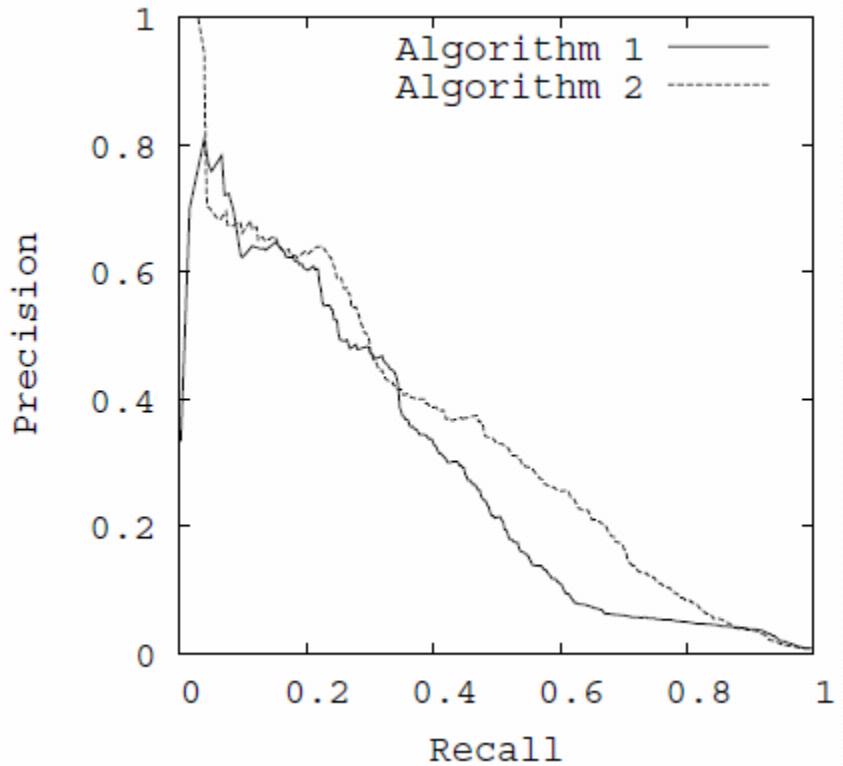
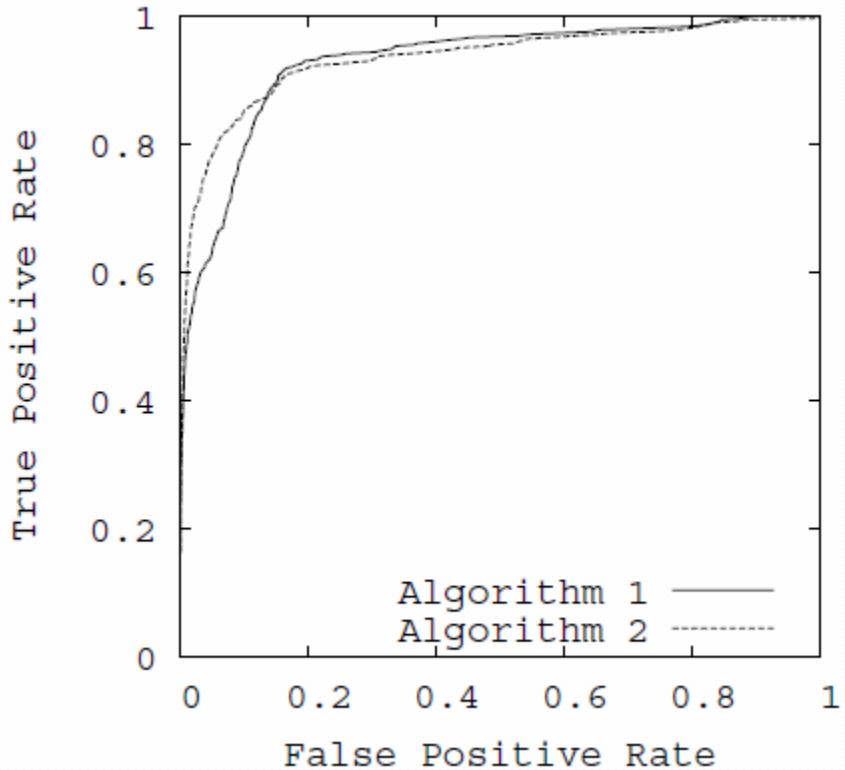
# Precision Recall

- **Recall:** If I pick a random positive example, what is the probability of making the right prediction?
- **Precision:** If I take a positive prediction example, what is the probability that it is indeed a positive example?

# Precision Recall Curve



# Comparison



J. Davis & M. Goadrich,  
“The Relationship Between Precision-Recall and ROC Curves.”,  
*ICML (2006)*

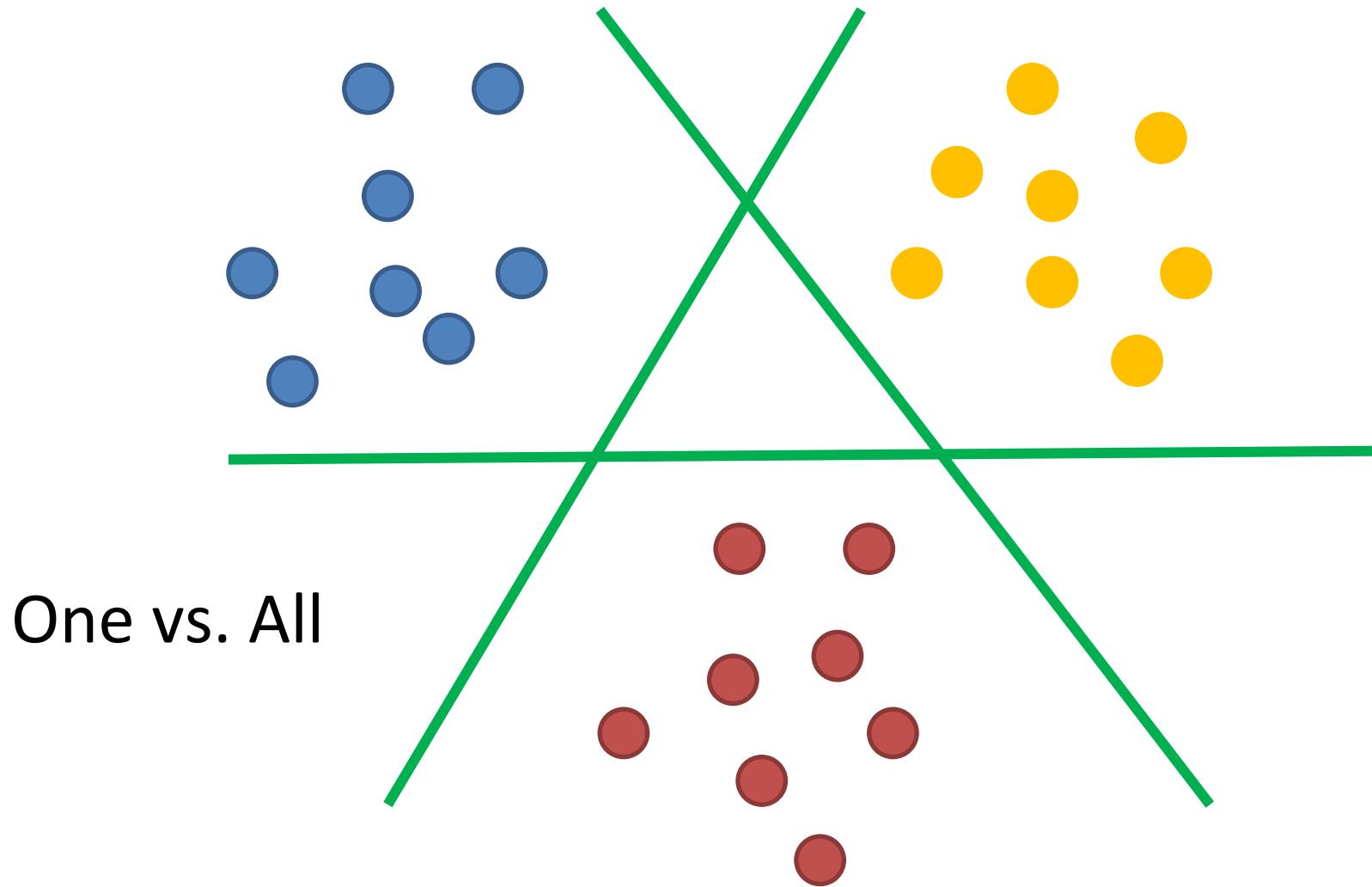
# F-measure

- Weighted average of precision and recall

$$F_{\beta} = \frac{(\beta^2 + 1) \cdot P \cdot R}{\beta^2 \cdot P + R}$$

- Usual case:  $\beta = 1$
- Increasing  $\beta$  allocates weight to recall

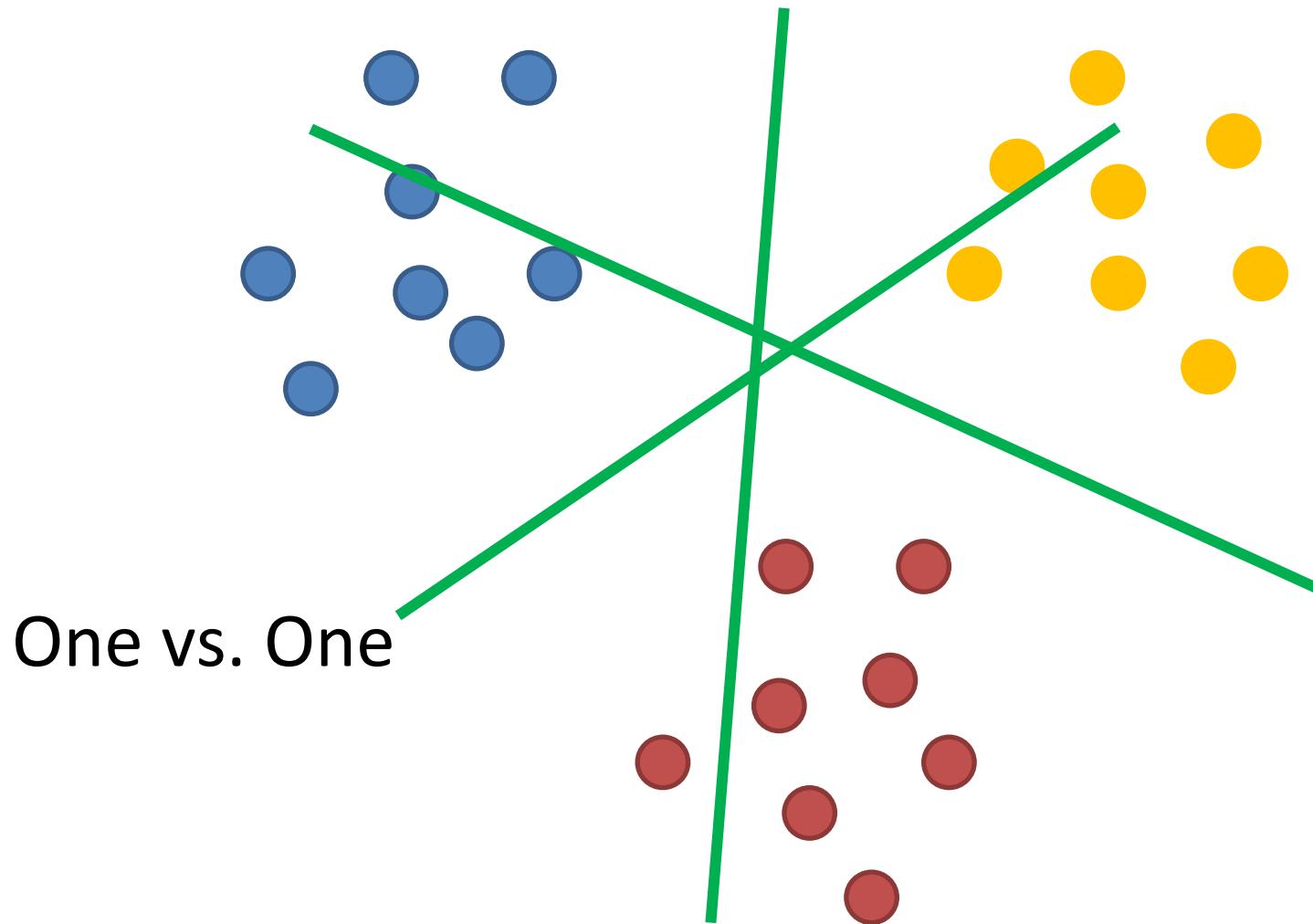
# Multi Class



# One vs All

- Train n classifier for n classes
- Take classification with greatest margin
- Slow training

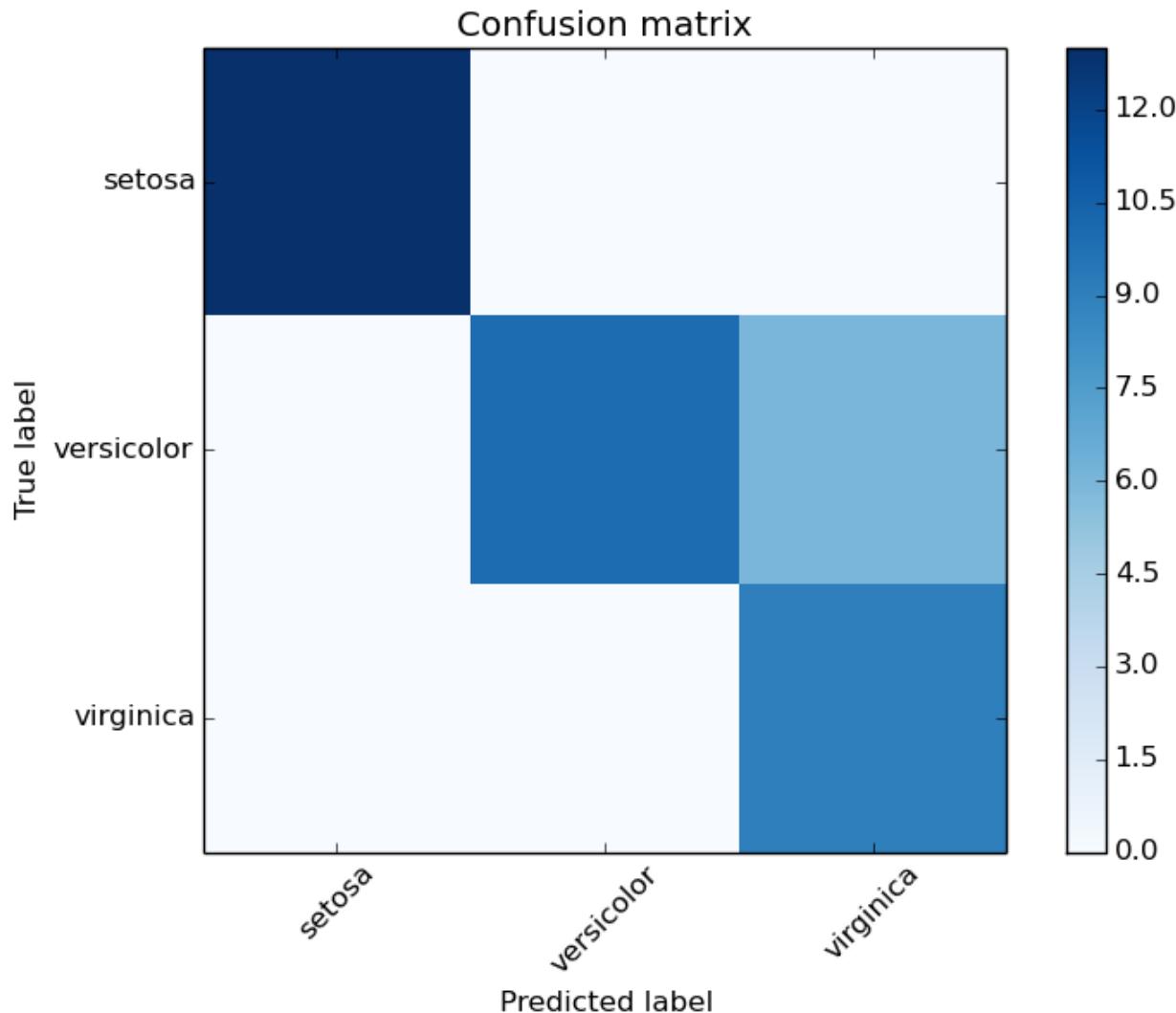
# Multi Class



# One vs One

- Train  $n(n-1)/2$  classifiers
- Take majority vote
- Fast training

# Confusion Matrix



# Recap

- Perceptrons are great
- But really just a separating hyperplane
- So is SVM
- Kernels are neat
- Evaluation metrics are important