

# Lecture 9:

# Support Vector Machines

Olexandr Isayev

Department of Chemistry, CMU

[olexandr@cmu.edu](mailto:olexandr@cmu.edu)

# The Ridge regression

The solution to the ordinary least squares fitting procedure is the vector  $\hat{\beta} = (\hat{\beta}_0, \dots, \hat{\beta}_p)$  that minimizes the Residual Sum of Squares

$$RSS = \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2.$$

**Ridge regression**, similarly, seeks the vector  $\hat{\beta}^{\text{Ridge}}$  that minimizes the *penalized or regularized* RSS

$$\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 = RSS + \lambda \sum_{j=1}^p \beta_j^2$$

where  $\lambda \geq 0$  is a **complexity parameter**.

# The LASSO Regression

“LASSO” stands for **L**east **A**bsolute **S**hrinkage and **S**election **O**perator

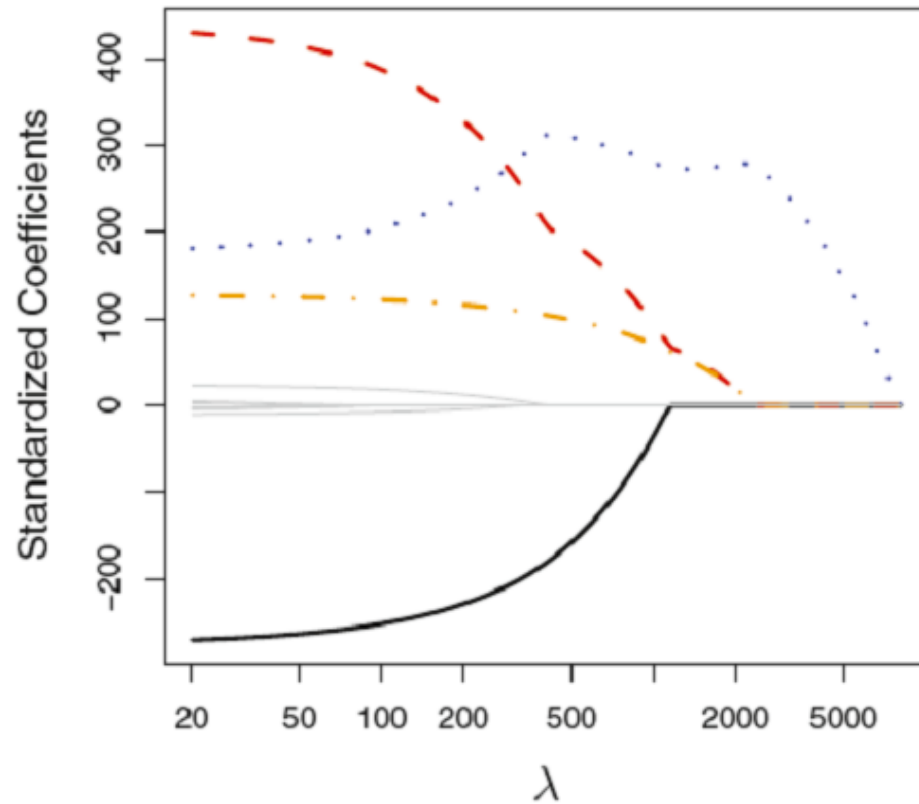
The idea of constraining the size of the OLS estimates can be extended to consider different kinds of penalizations.

While the Ridge penalty encompasses the  $L2$  norm of the estimates vector, the **LASSO** makes use of the  $L1$  norm:

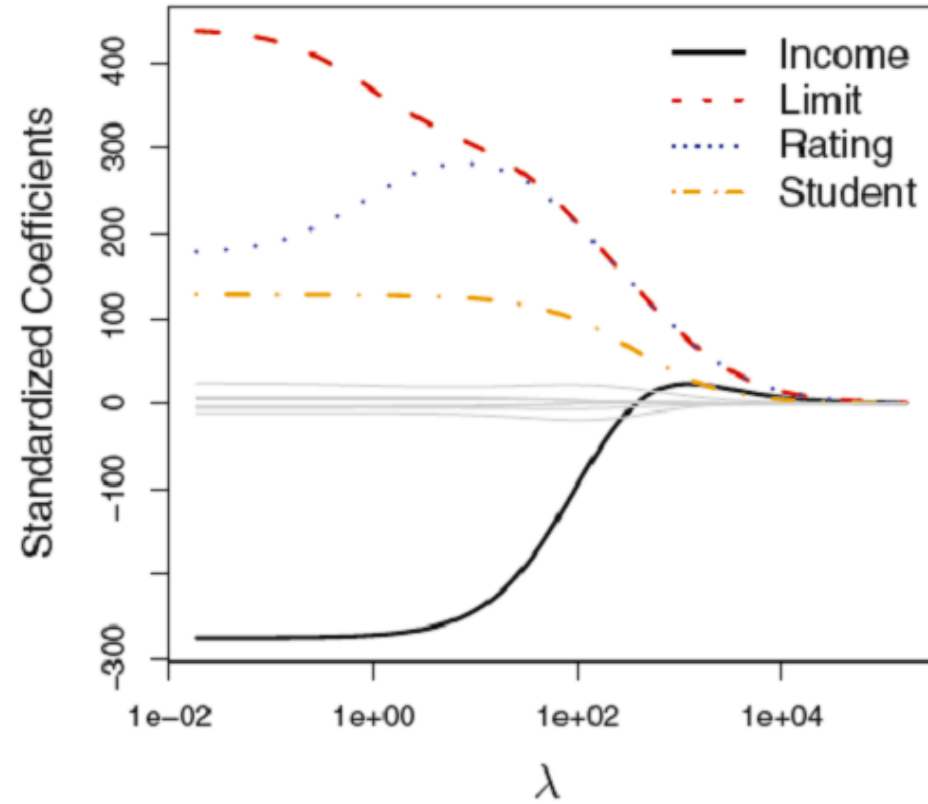
$$\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j| = RSS + \lambda \sum_{j=1}^p |\beta_j|$$

# Comparing LASSO and Ridge

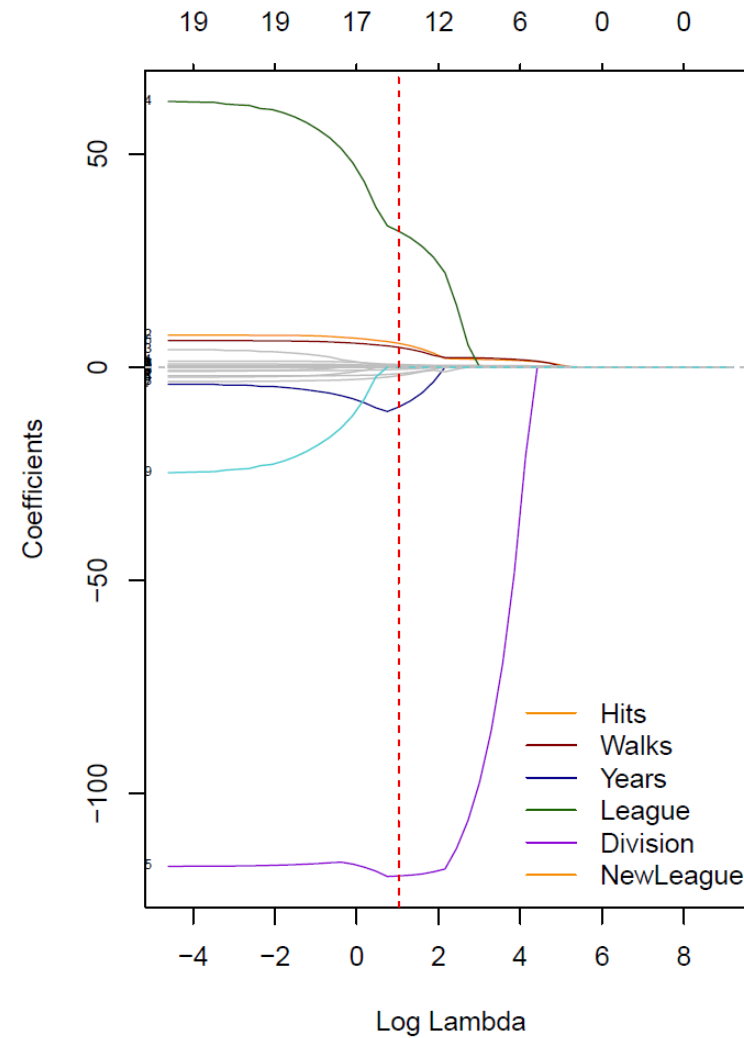
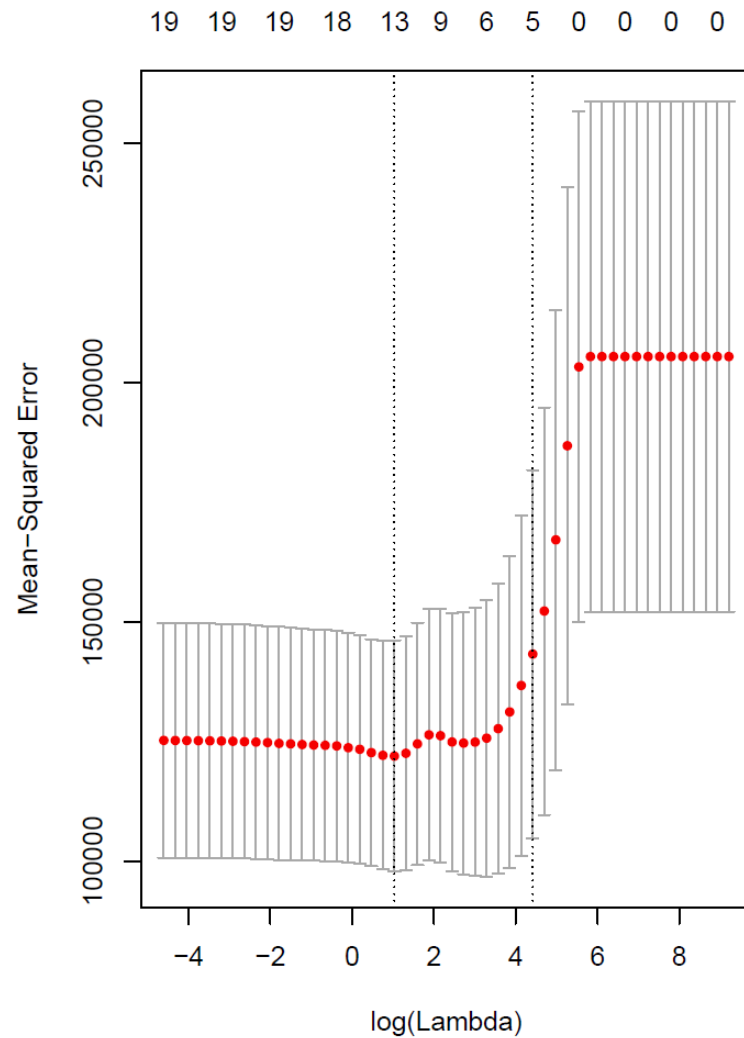
Lasso



Ridge



# Choosing $\lambda$ : cross-validation



# Difference between L1 and L2 regularization

## L1 Regularization

- L1 penalizes sum of absolute value of weights.
- L1 has a sparse solution
- L1 has multiple solutions
- L1 has built in feature selection
- L1 is robust to outliers
- L1 generates model that are simple and interpretable but cannot learn complex patterns

## L2 Regularization

- L2 regularization penalizes sum of square weights.
- L2 has a non sparse solution
- L2 has one solution
- L2 has no feature selection
- L2 is not robust to outliers
- L2 gives better prediction when output variable is a function of all input features
- L2 regularization is able to learn complex data patterns

# A hybrid penalty: the Elastic Nets

The LASSO sometimes does not perform well with highly correlated variables, and often performs worse than Ridge in prediction.

To overcome this limitations, a penalty that combines the  $L1$  and  $L2$  constraints has been developed.

An **elastic net** is a regularization and variable selection procedure that makes use of the penalty

$$\lambda \left[ \frac{1}{2}(1 - \alpha) \sum_{j=1}^p \beta_j^2 + \alpha \sum_{j=1}^p |\beta_j| \right]$$

where  $\alpha \in [0, 1]$  is called the **mixing** parameter and  $\lambda$  has the usual interpretation. LASSO and Ridge are special cases, respectively for  $\alpha = 1$  and  $\alpha = 0$ .

# Elastic nets Summary

Elastic net regression is a regularization and variable selection procedure that overcomes some of the limitations of the LASSO by borrowing strength from the Ridge. Specifically, it

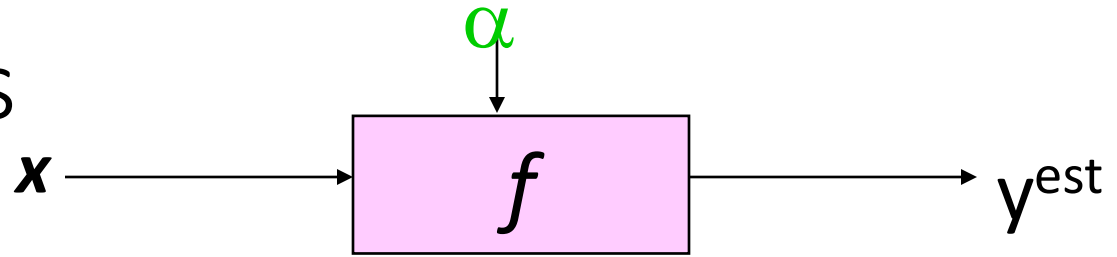
- allows to select more than  $n$  variables
- tends to jointly select or leave out groups of highly correlated variables
- improves the predictive performance w.r.t. LASSO
- is readily extendable to use with more general methods, such as GLM.

Elastic nets are especially useful when a sparse solution is either necessary or desirable (such as in  $p \gg n$  problems) and small groups of highly correlated predictors are present.



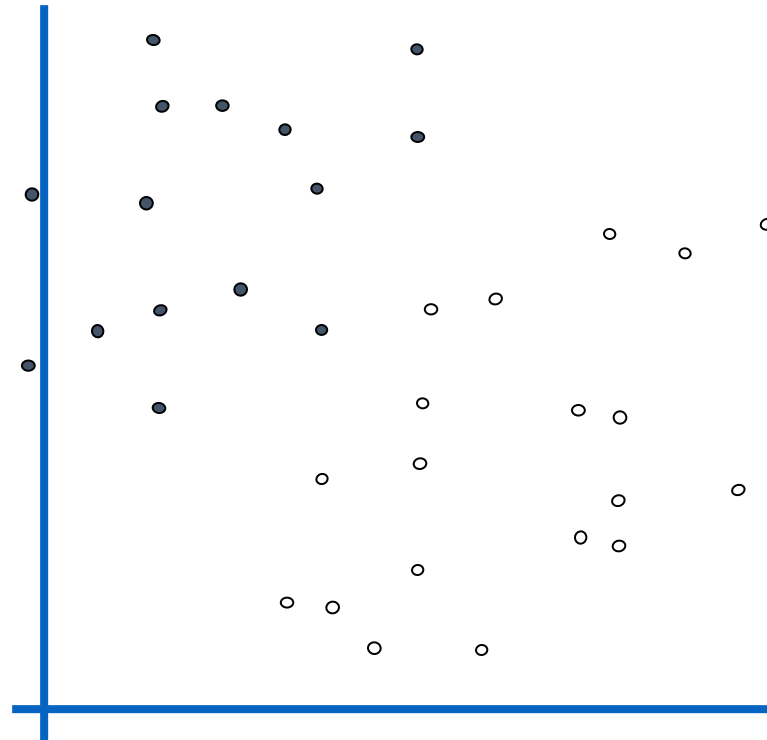
# Classification

# Linear Classifiers



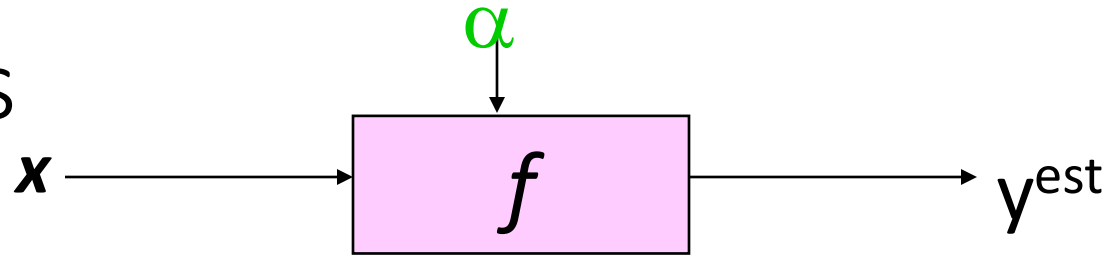
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1

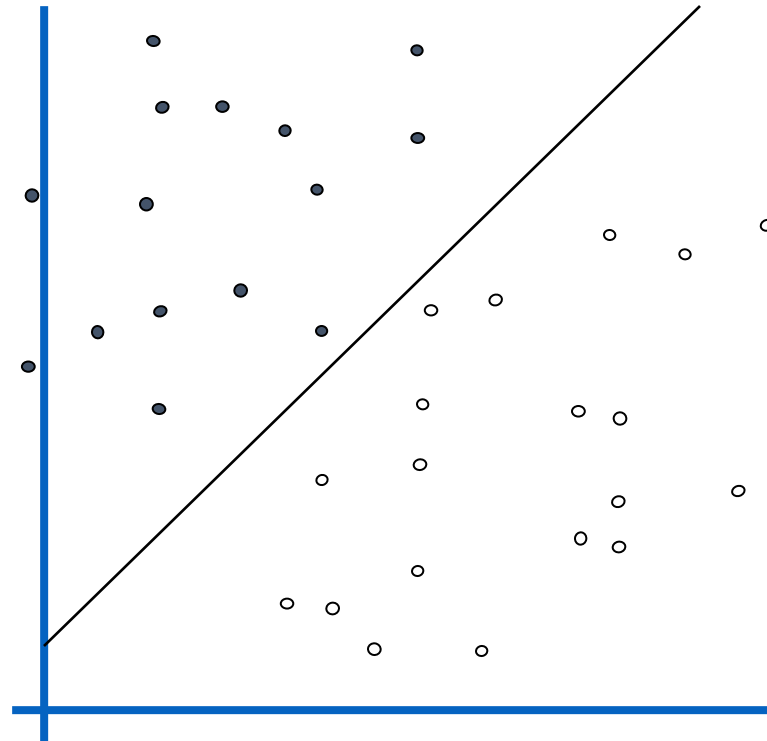


How would you  
classify this data?

# Linear Classifiers



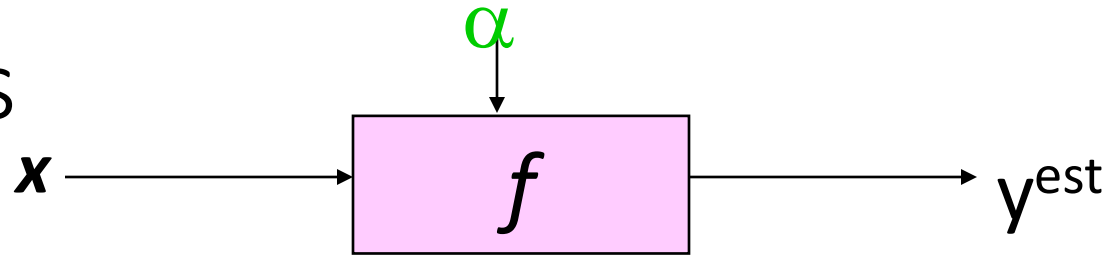
- denotes +1
- denotes -1



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

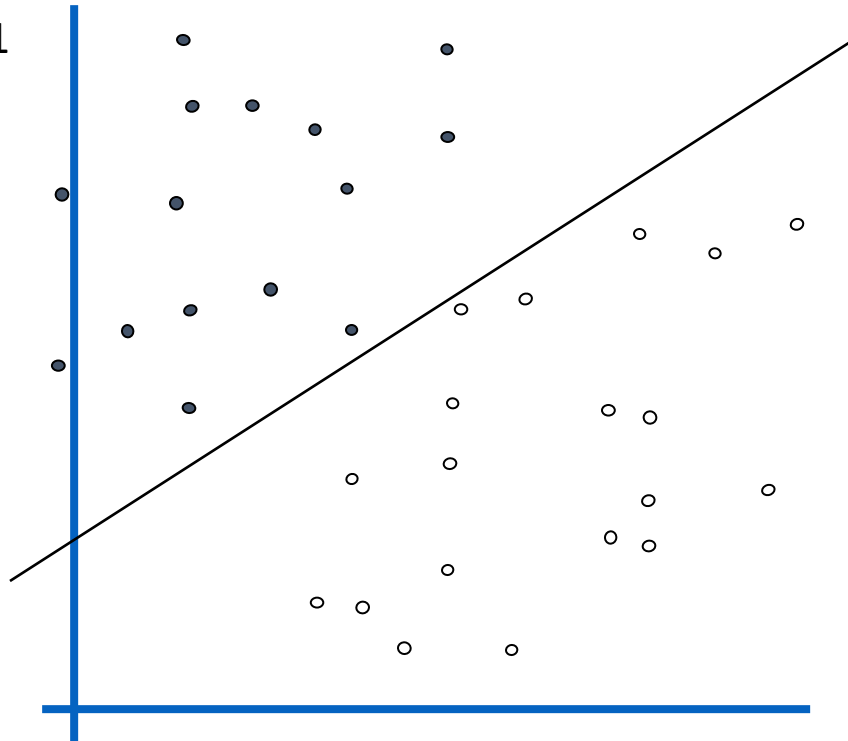
How would you  
classify this data?

# Linear Classifiers



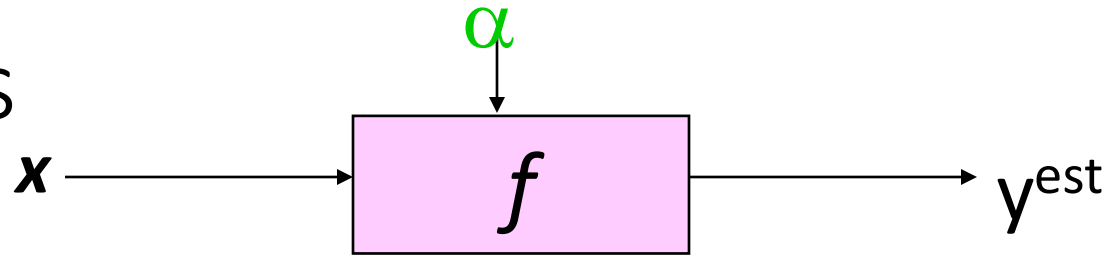
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1



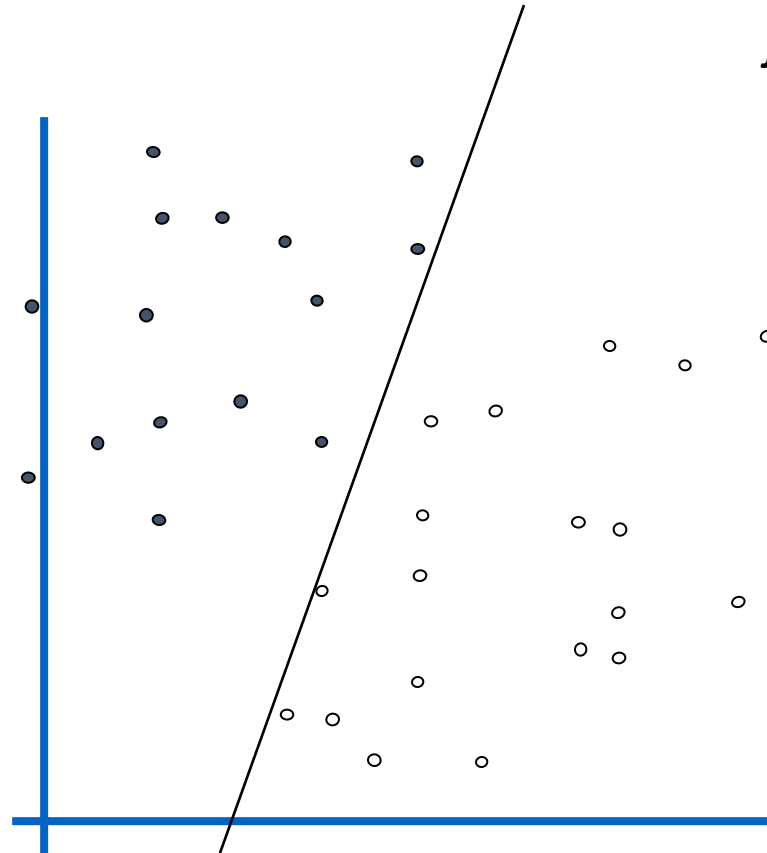
How would you  
classify this data?

# Linear Classifiers



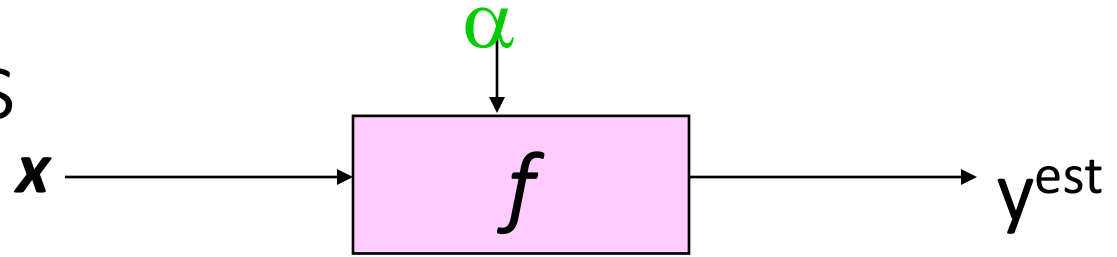
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1

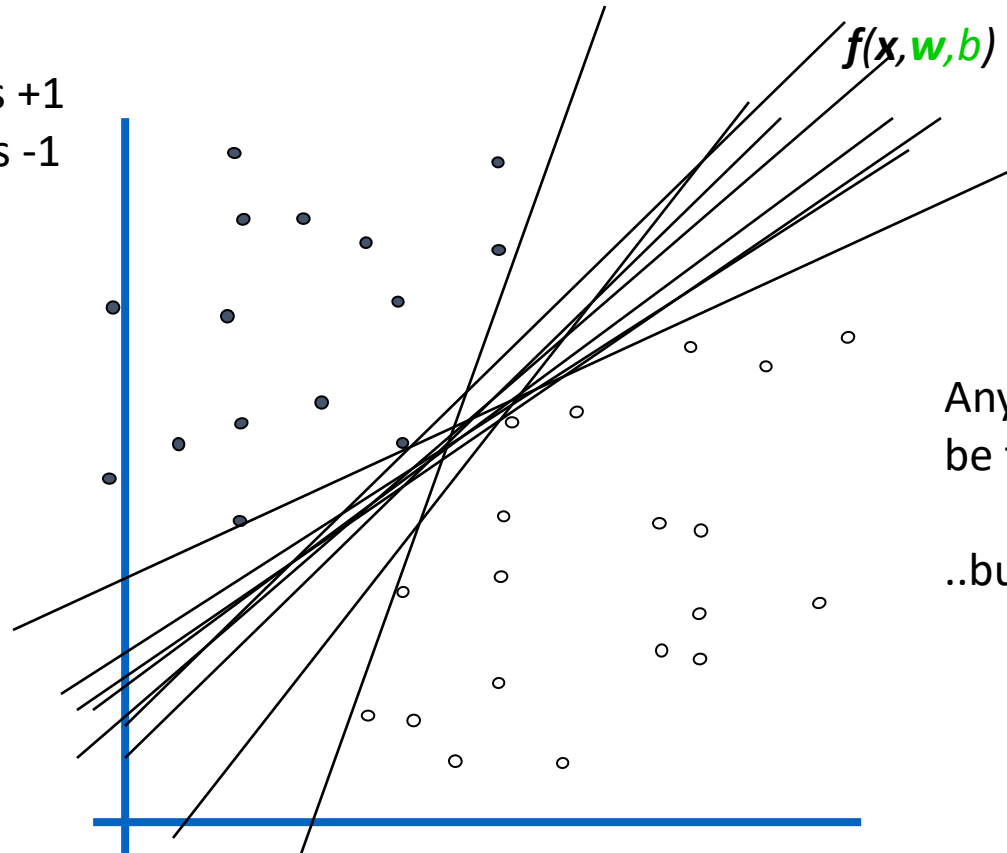


How would you  
classify this data?

# Linear Classifiers



- denotes +1
- denotes -1

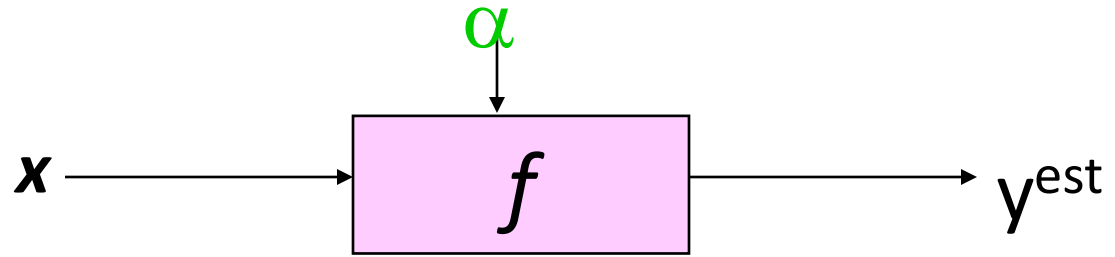


$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

Any of these would  
be fine..

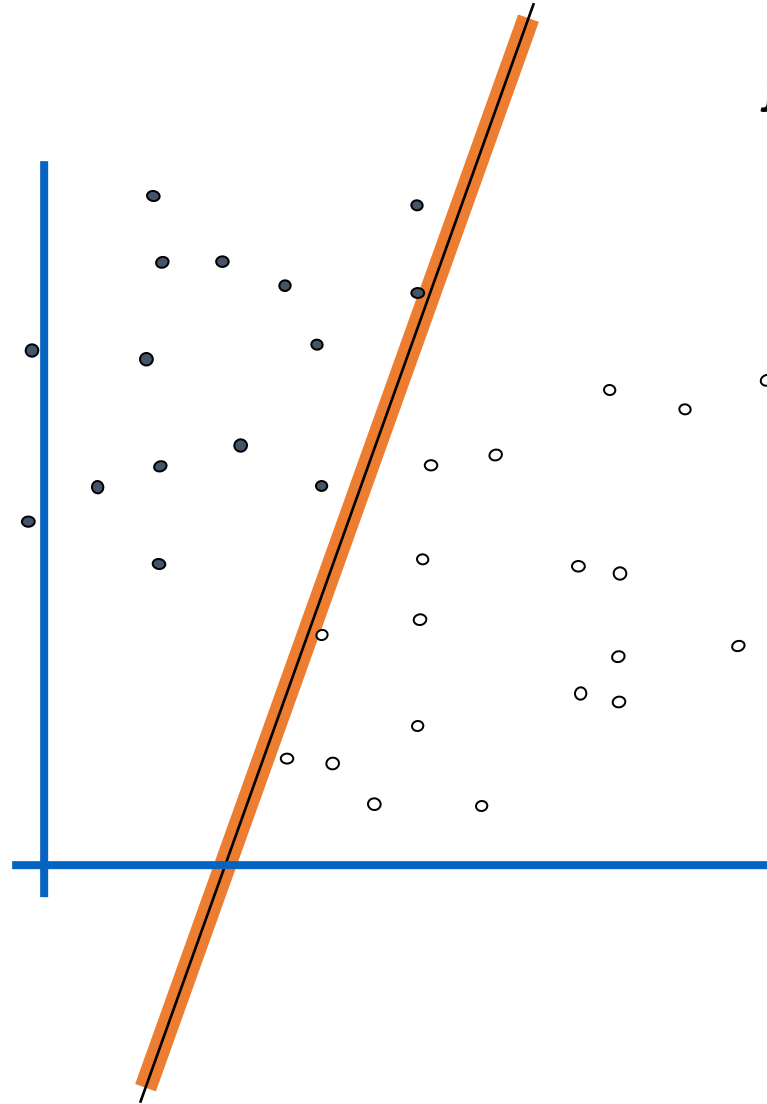
..but which is best?

# Classifier Margin



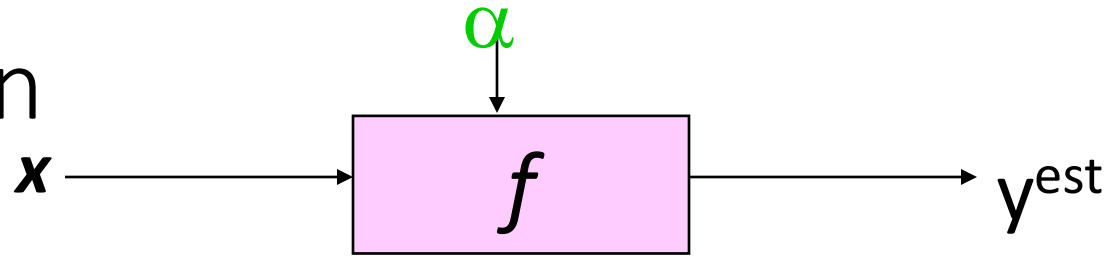
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

- denotes +1
- denotes -1

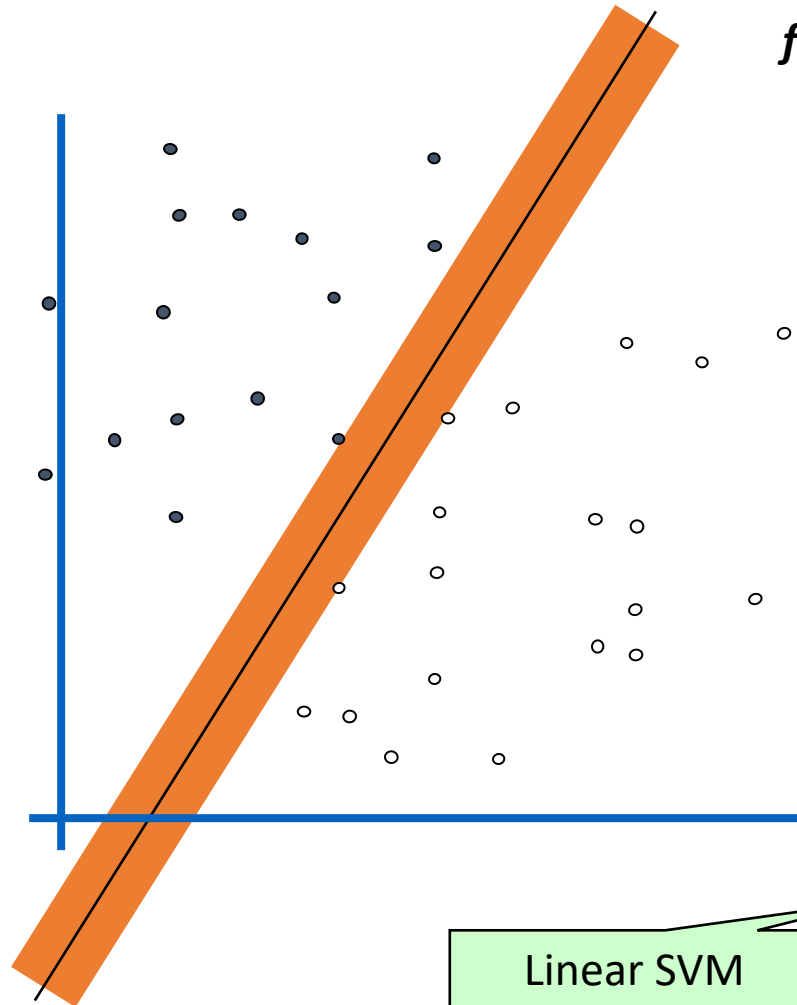


Define the **margin** of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

# Maximum Margin



- denotes +1
- denotes -1



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

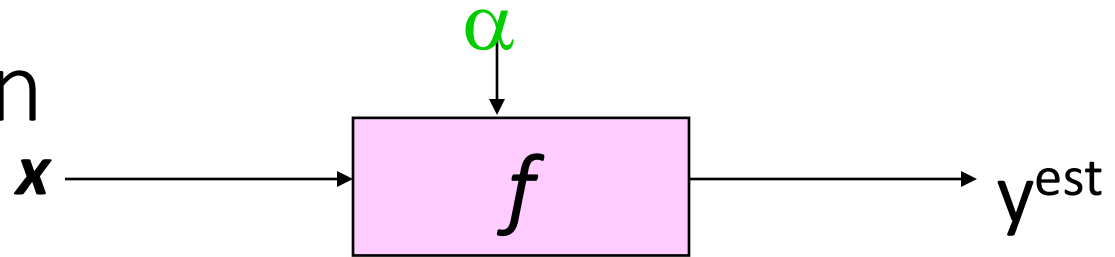
The **maximum margin linear classifier** is the linear classifier with the, um, maximum margin.

This is the simplest kind of SVM (Called an LSVM)

Linear SVM



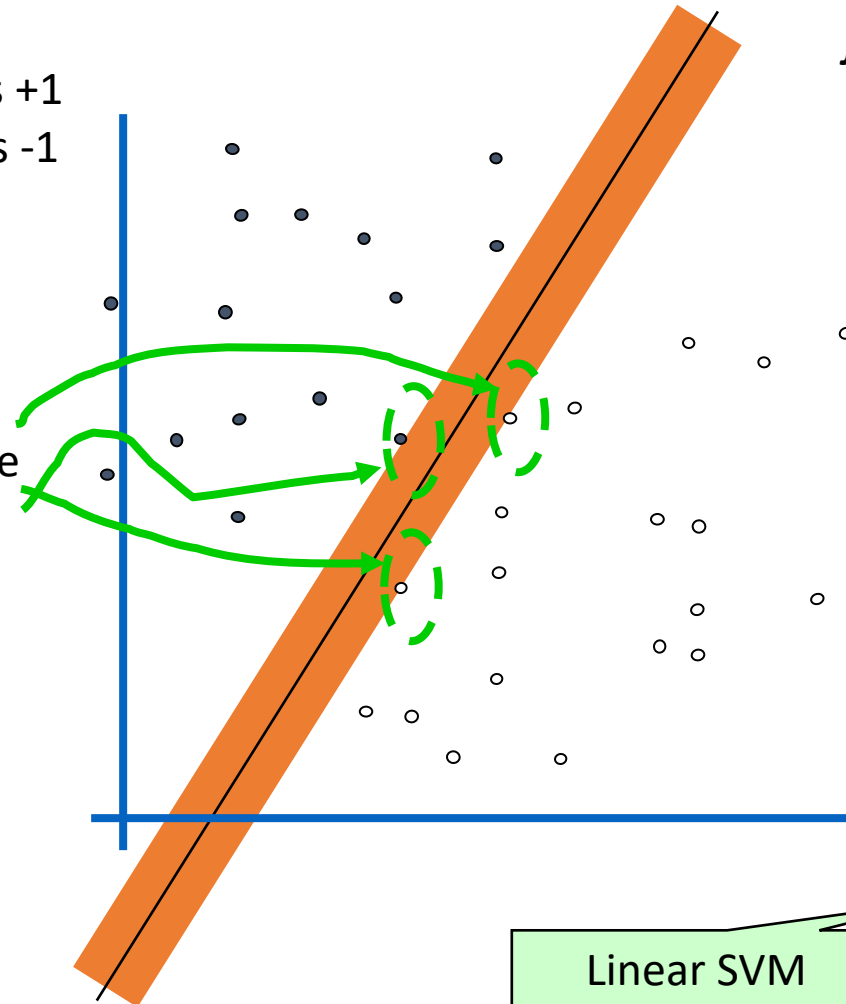
# Maximum Margin



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

- denotes +1
- denotes -1

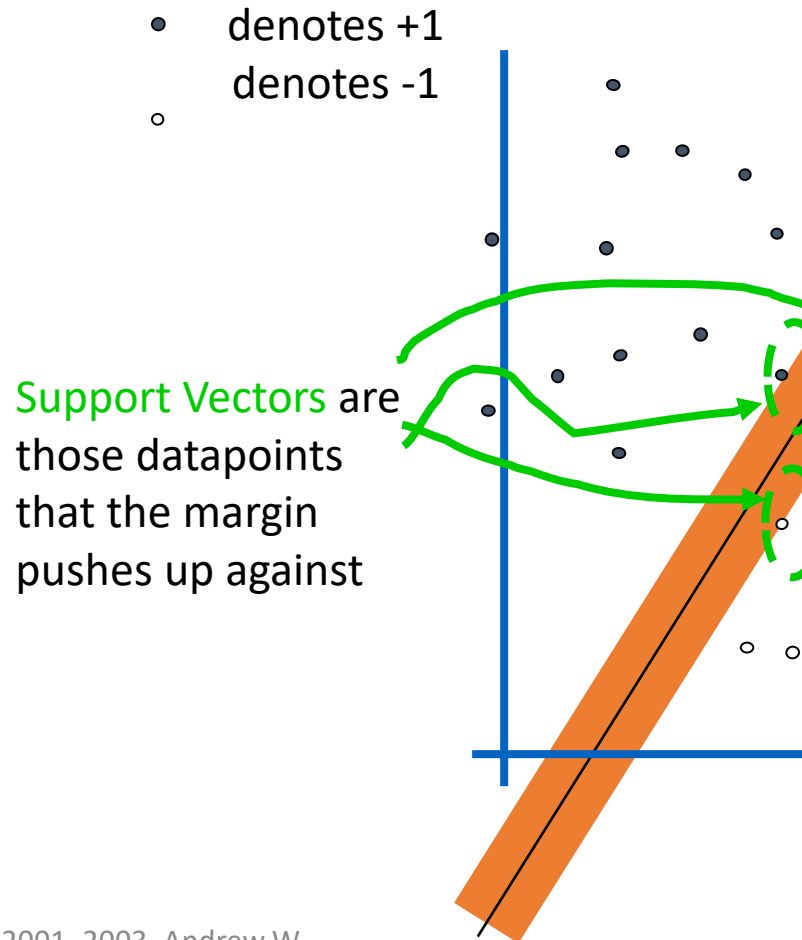
Support Vectors are those datapoints that the margin pushes up against



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

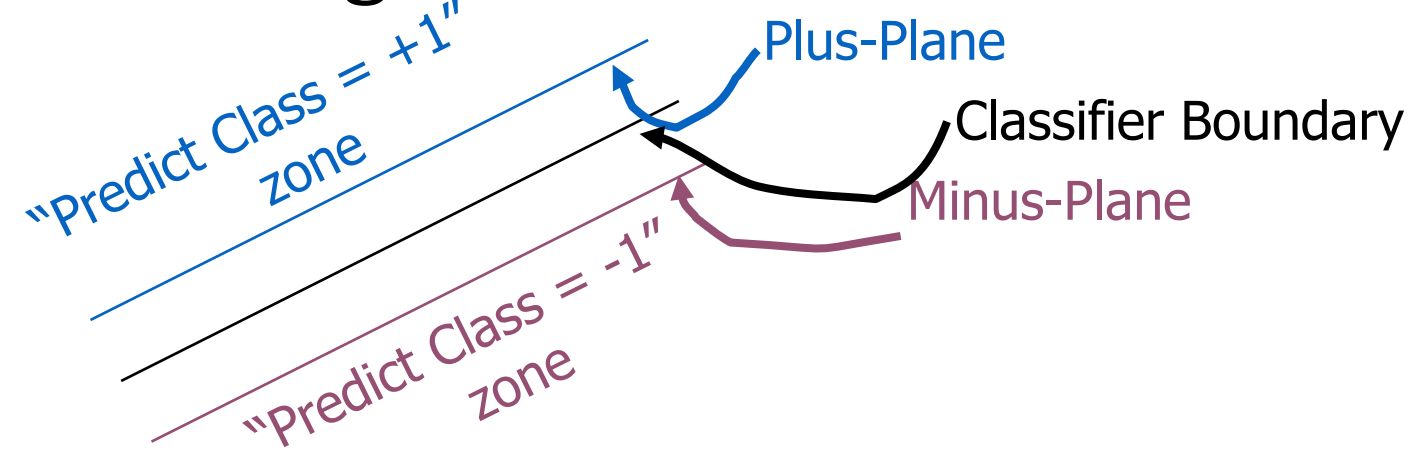
Linear SVM

# 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. 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.

# Specifying a line and margin



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

# Distances to the Boundary

- The decision boundary consists of all points  $\mathbf{x}$  that are solutions to equation:  $\mathbf{w}^T \mathbf{x} + b = 0$ .
  - $\mathbf{w}$  is a column vector of parameters (weights).
  - $\mathbf{x}$  is an input vector.
  - $b$  is a scalar value (a real number).
- If  $\mathbf{x}_n$  is a training point, its distance to the boundary is computed using this equation:

$$D(\mathbf{x}_n, \mathbf{w}) = \left| \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|} \right|$$

# Distances to the Boundary

- If  $\mathbf{x}_n$  is a training point, its distance to the boundary is computed using this equation:

$$D(\mathbf{x}_n, \mathbf{w}) = \left| \frac{\mathbf{w}^T \mathbf{x}_n + b}{\|\mathbf{w}\|} \right|$$

- Since the training data are linearly separable, the data from each class should fall on opposite sides of the boundary.
- Suppose that  $t_n = -1$  for points of one class, and  $t_n = +1$  for points of the other class.
- Then, we can rewrite the distance as:

$$D(\mathbf{x}_n, \mathbf{w}) = \frac{t_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|}$$

# Distances to the Boundary

- So, given a decision boundary defined  $\mathbf{w}$  and  $b$ , and given a training input  $\mathbf{x}_n$ , the distance of  $\mathbf{x}_n$  to the boundary is:

$$D(\mathbf{x}_n, \mathbf{w}) = \frac{t_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|}$$

- If  $t_n = -1$ , then:
  - $\mathbf{w}^T \mathbf{x}_n + b < 0$ .
  - $t_n(\mathbf{w}^T \mathbf{x}_n + b) > 0$ .
- If  $t_n = 1$ , then:
  - $\mathbf{w}^T \mathbf{x}_n + b > 0$ .
  - $t_n(\mathbf{w}^T \mathbf{x}_n + b) > 0$ .
- So, in all cases,  $t_n(\mathbf{w}^T \mathbf{x}_n + b)$  is positive.

# Optimization Criterion

- If  $\mathbf{x}_n$  is a training point, its distance to the boundary is computed using this equation:

$$D(\mathbf{x}_n, \mathbf{w}) = \frac{t_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|}$$

- Therefore, the optimal boundary  $\mathbf{w}_{\text{opt}}$  is defined as:

$$(\mathbf{w}_{\text{opt}}, b_{\text{opt}}) = \operatorname{argmax}_{\mathbf{w}, b} \left\{ \min_n \left[ \frac{t_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|} \right] \right\}$$

- In words: find the  $\mathbf{w}$  and  $b$  that maximize the minimum distance of any training input from the boundary.

# Constrained Optimization

- Summarizing the previous slides, we want to find:

$$\mathbf{w}_{\text{opt}} = \operatorname{argmin}_{\mathbf{w}} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 \right\}$$

subject to the following constraints:

$$\forall n \in \{1, \dots, N\}, t_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1$$

- This is a different optimization problem than what we have seen before.
- We need to minimize a quantity **while satisfying a set of inequalities**.
- This type of problem is a **constrained optimization problem**.

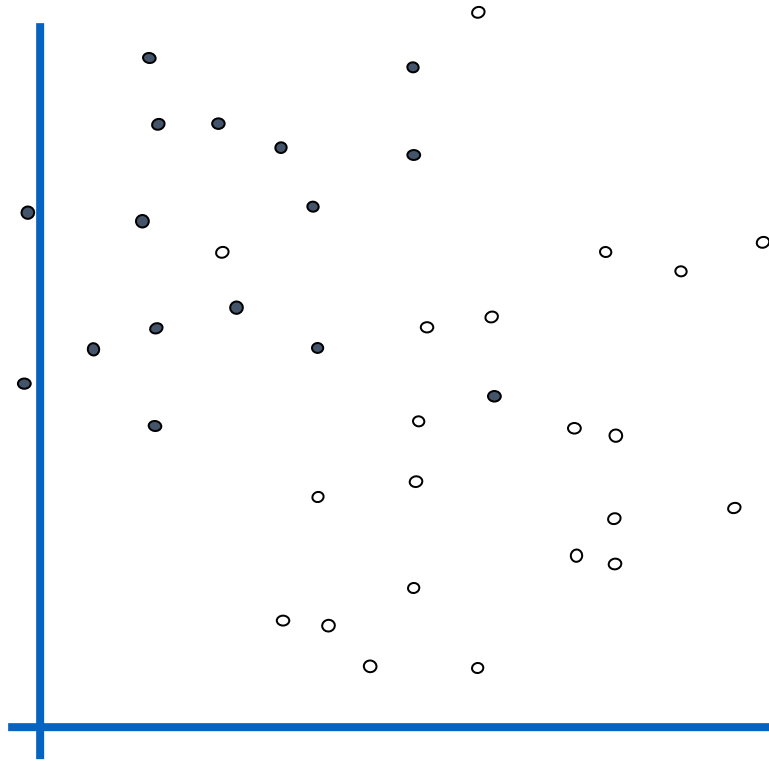
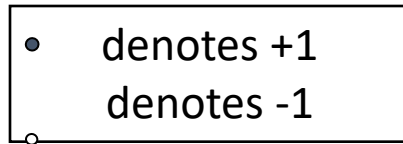


# Quadratic Programming

- Our constrained optimization problem can be solved using a method called **quadratic programming**.
- Describing **quadratic programming** in depth is outside the scope of this course.
- Our goal is simply to understand how to use quadratic programming as a black box, to solve our optimization problem.
  - This way, you can use any quadratic programming toolkit (Python has one).

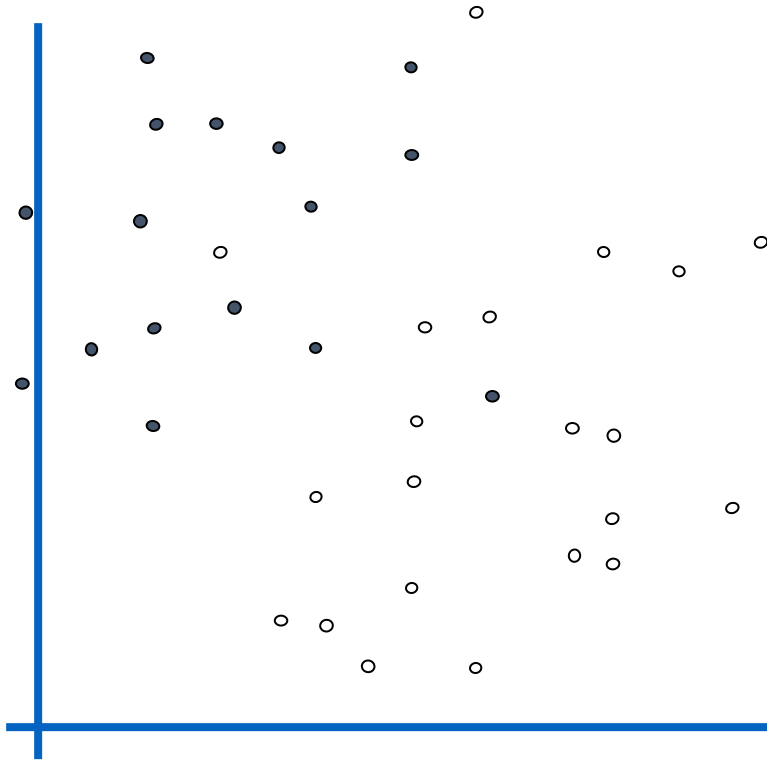
# Uh-oh!

This is going to be a problem!  
What should we do?



# Uh-oh!

• denotes +1  
○ denotes -1



This is going to be a problem!  
What should we do?

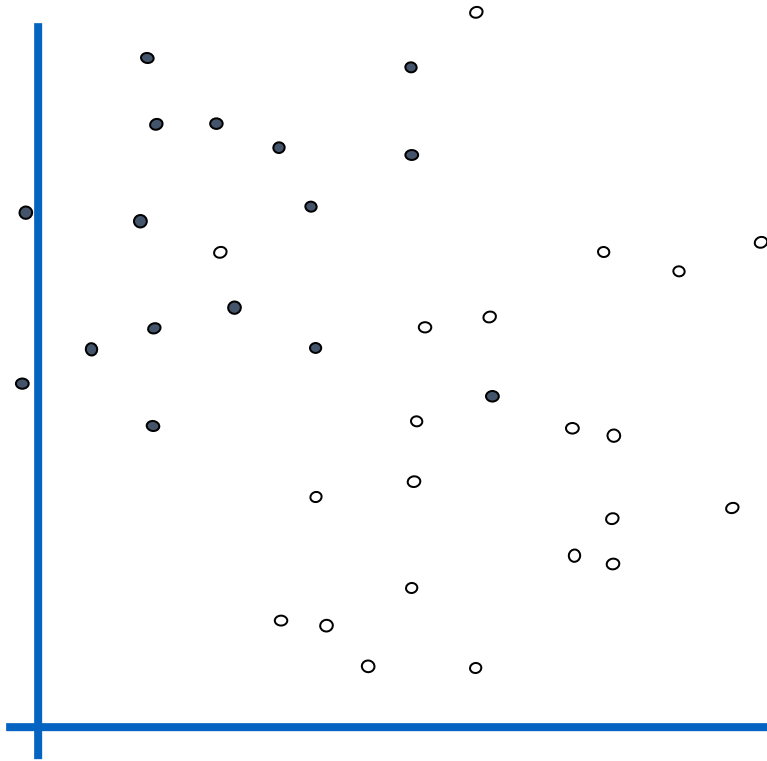
Idea 1:

Find minimum  $\mathbf{w} \cdot \mathbf{w}$ , while  
minimizing number of training  
set errors.

Problem: Two things to  
minimize makes for an ill-  
defined optimization

# Uh-oh!

- denotes +1
- denotes -1



This is going to be a problem!  
What should we do?

Idea 1.1:

Minimize

$\mathbf{w} \cdot \mathbf{w} + C (\#train\ errors)$

Tradeoff parameter

There's a serious practical problem that's about to make us reject this approach. Can you guess what it is?

# Uh-oh!

This is going to be a problem!  
What should we do?

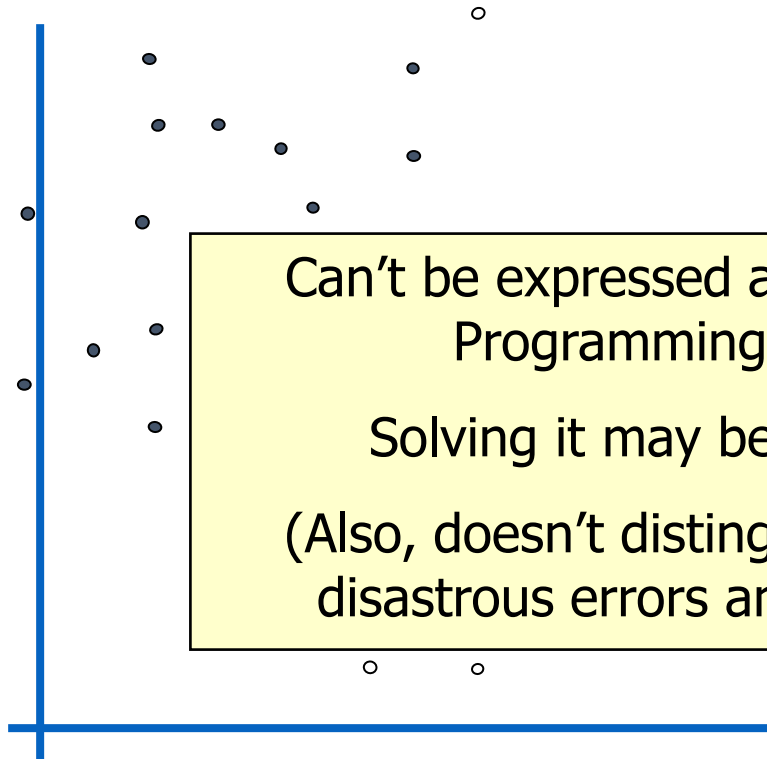
Idea 1.1:

Minimize

$$\mathbf{w} \cdot \mathbf{w} + C (\# \text{train errors})$$

Tradeoff parameter

- denotes +1
- denotes -1

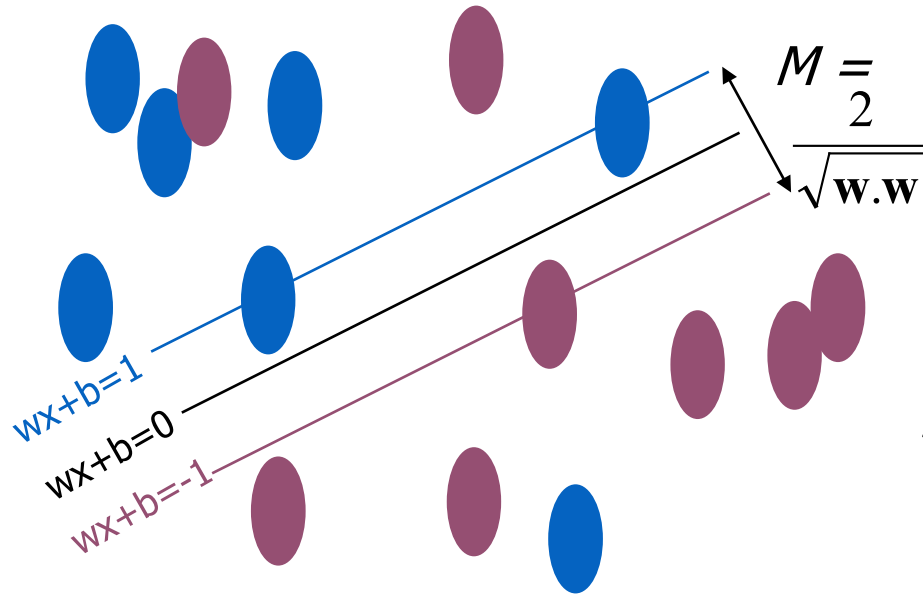


Can't be expressed as a Quadratic Programming problem.  
Solving it may be too slow.  
(Also, doesn't distinguish between disastrous errors and near misses)

There's a serious practical problem with this approach. Can you see what it is?

So... any other ideas?

# Learning Maximum Margin with Noise



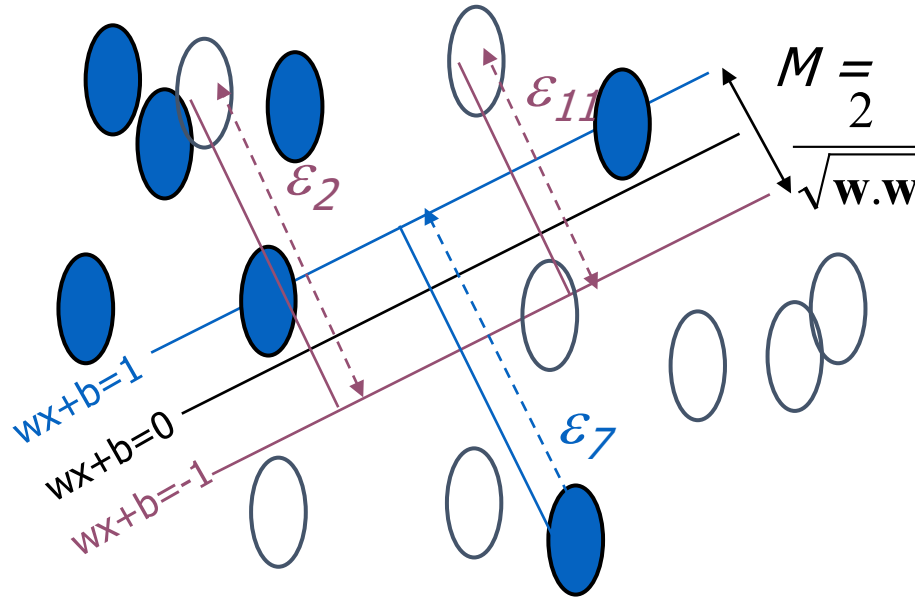
Given guess of  $\mathbf{w}$ ,  $b$  we can

- Compute sum of distances of points to their correct zones
- Compute the margin width

Assume  $R$  datapoints, each  $(\mathbf{x}_k, y_k)$  where  $y_k = +/- 1$

What should our quadratic optimization criterion be?

# Learning Maximum Margin with Noise



Given guess of  $\mathbf{w}$ ,  $b$  we can

- Compute sum of distances of points to their correct zones
- Compute the margin width

Assume  $R$  datapoints, each  $(\mathbf{x}_k, y_k)$  where  $y_k = \pm 1$

What should our quadratic optimization criterion be?

Minimize

$$\frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_{k=1}^R \epsilon_k$$

How many constraints will we have?  $2R$

What should they be?

$$\mathbf{w} \cdot \mathbf{x}_k + b \geq 1 - \epsilon_k \text{ if } y_k = 1$$

$$\mathbf{w} \cdot \mathbf{x}_k + b \leq -1 + \epsilon_k \text{ if } y_k = -1$$

# SVM Objective Function

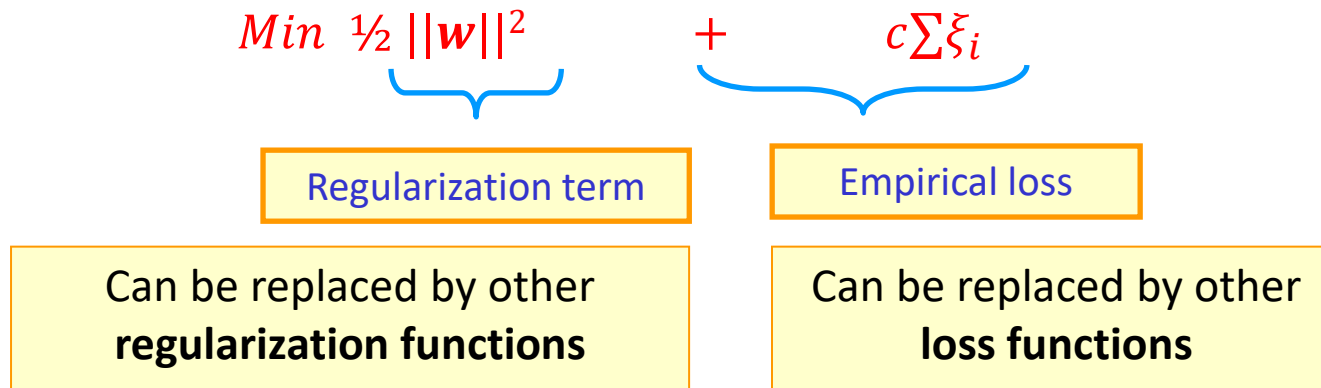
- The problem we solved is:

$$\text{Min } \frac{1}{2} ||\mathbf{w}'||^2 + c \sum \xi_i$$

- Where  $\xi_i > 0$  is called a **slack variable**, and is defined by:

- $\xi_i = \max(0, 1 - y_i \mathbf{w}'^T \mathbf{x}_i)$
- Equivalently, we can say that:  $y_i \mathbf{w}'^T \mathbf{x}_i \geq 1 - \xi_i; \xi_i \geq 0$

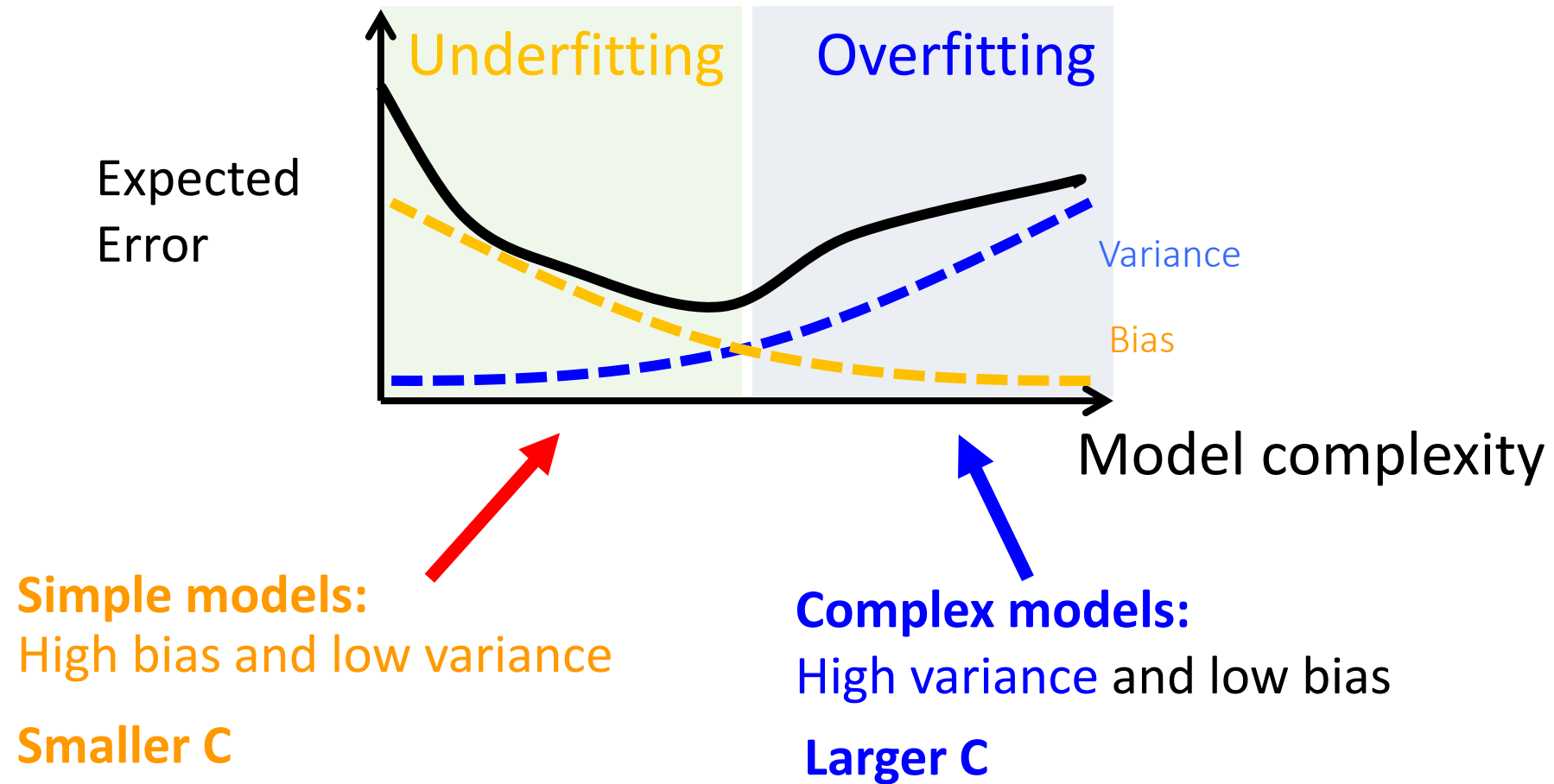
- And this can be written as:



- General Form of a learning algorithm:
  - Minimize empirical loss, and Regularize (to avoid over fitting)

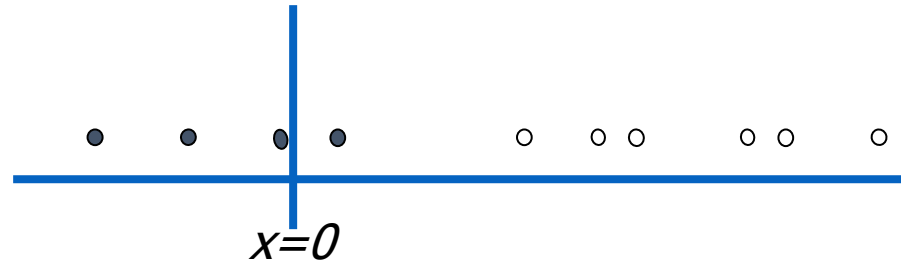


# Underfitting and Overfitting



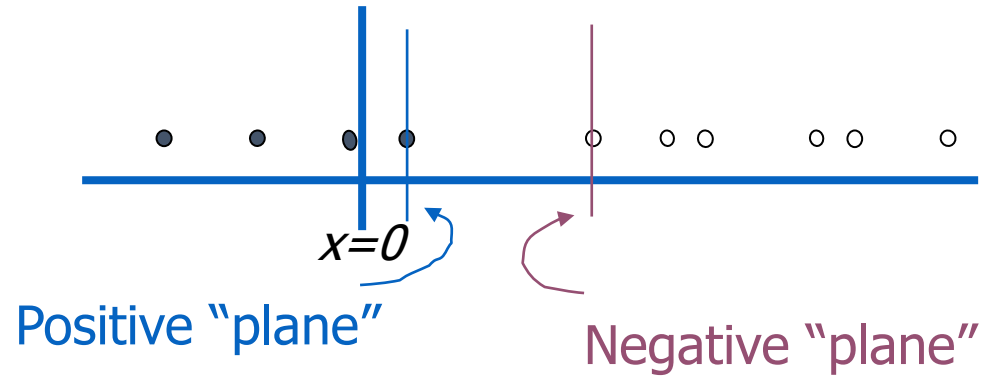
Suppose we're in 1-dimension

What would  
SVMs do with  
this data?



# Suppose we're in 1-dimension

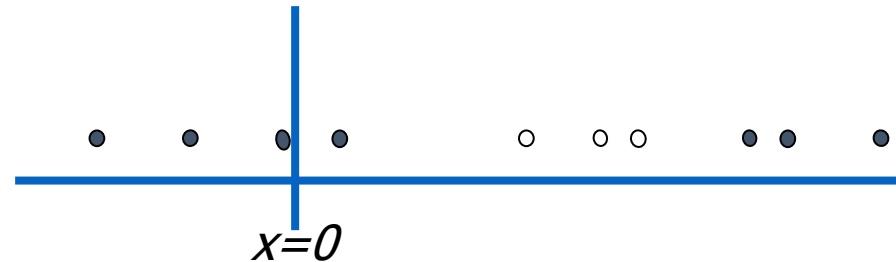
Not a big surprise



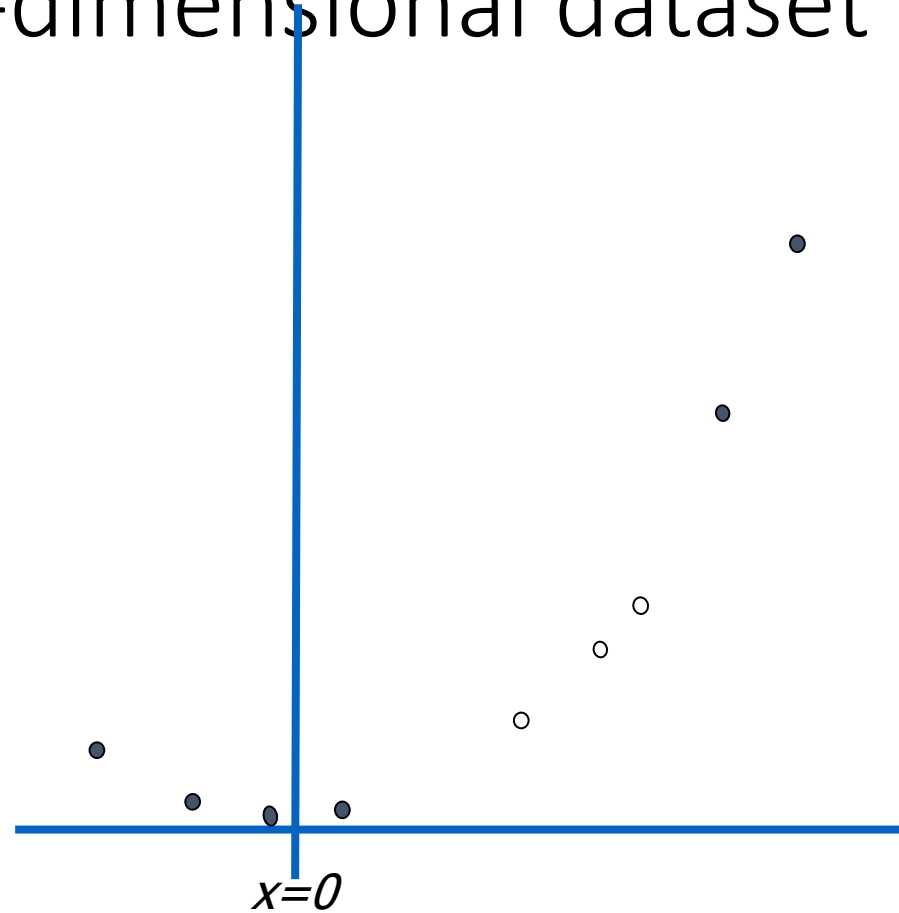
# Harder 1-dimensional dataset

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

What can be done about this?



# Harder 1-dimensional dataset

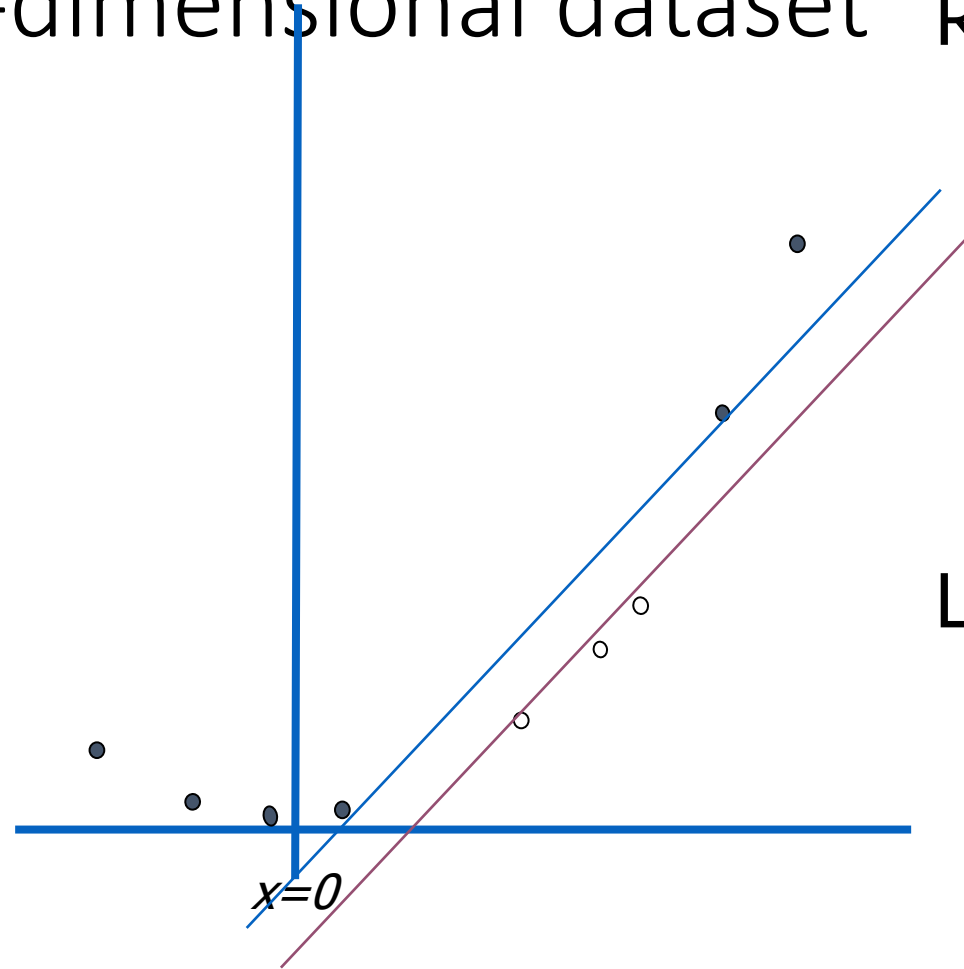


Remember how  
permitting non-  
linear basis  
functions made  
linear regression  
so much nicer?

Let's permit them  
here too

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

# Harder 1-dimensional dataset



Remember how  
permitting non-  
linear basis  
functions made  
linear regression  
so much nicer?

Let's permit them  
here too

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

# Kernels and Basis Functions

- In general, kernels make it easy to incorporate basis functions into SVMs:
  - Define  $\varphi(\mathbf{x})$  any way you like.
  - Define  $k(\mathbf{x}, \mathbf{z}) = \varphi(\mathbf{x})^T \varphi(\mathbf{z})$ .
- The kernel function represents a dot product, but in a (typically) higher-dimensional feature space compared to the original space of  $\mathbf{x}$  and  $\mathbf{z}$ .

# Common SVM basis functions

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

$\mathbf{z}_k = ( \text{radial basis functions of } \mathbf{x}_k )$

$$\mathbf{z}_k[j] = \varphi_j(\mathbf{x}_k) = \text{KernelFn}\left(\frac{|\mathbf{x}_k - \mathbf{c}_j|}{KW}\right)$$

$\mathbf{z}_k = ( \text{sigmoid functions of } \mathbf{x}_k )$



# Polynomial Kernels

- Let  $\mathbf{x}$  and  $\mathbf{z}$  be  $D$ -dimensional vectors.
- A polynomial kernel of degree  $d$  is defined as:

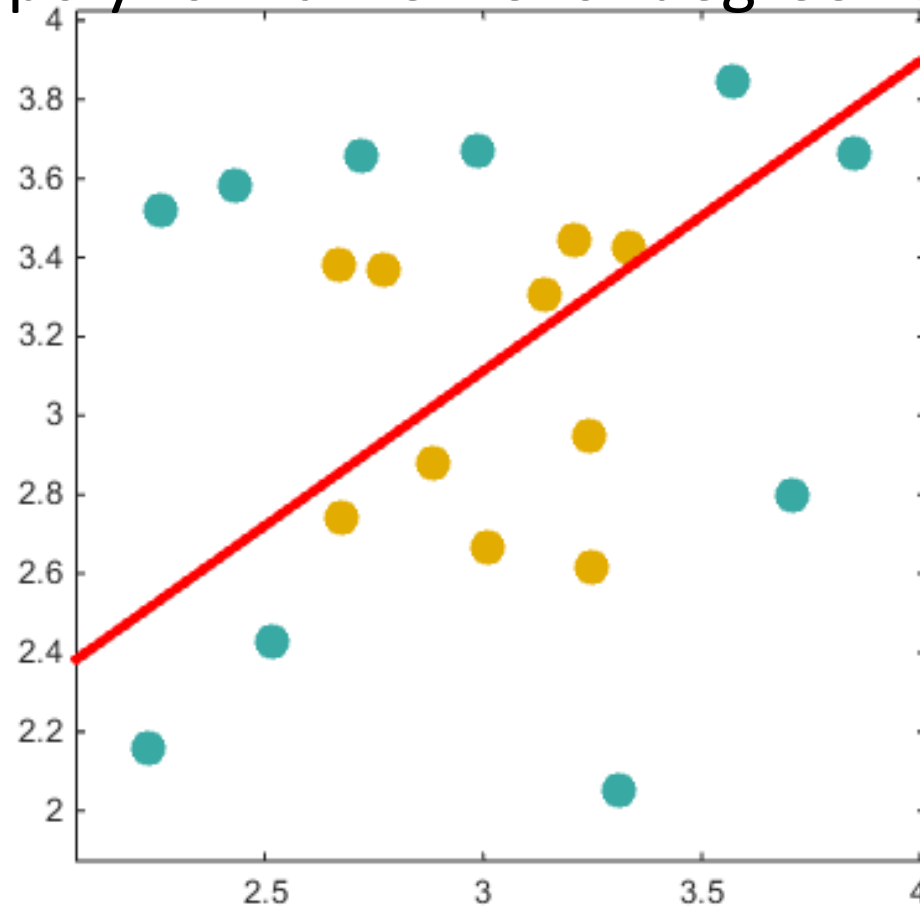
$$k(\mathbf{x}, \mathbf{z}) = (c + \mathbf{x}^T \mathbf{z})^d$$

- The kernel  $k(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^2$  that we saw a couple of slides back was a quadratic kernel.
- Parameter  $c$  controls the trade-off between influence higher-order and lower-order terms.
  - Increasing values of  $c$  give increasing influence to lower-order terms.

# Polynomial Kernels – An Easy Case

Decision boundary with polynomial kernel of degree 1.

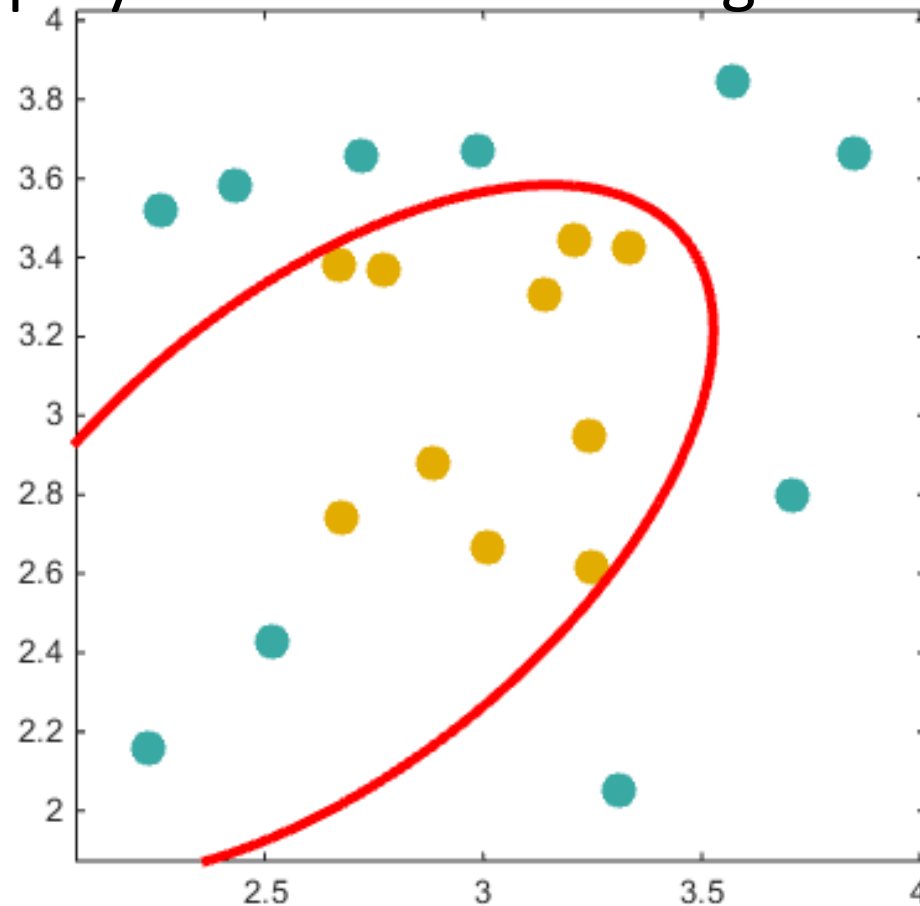
This is identical to the result using the standard dot product as kernel.



# Polynomial Kernels – An Easy Case

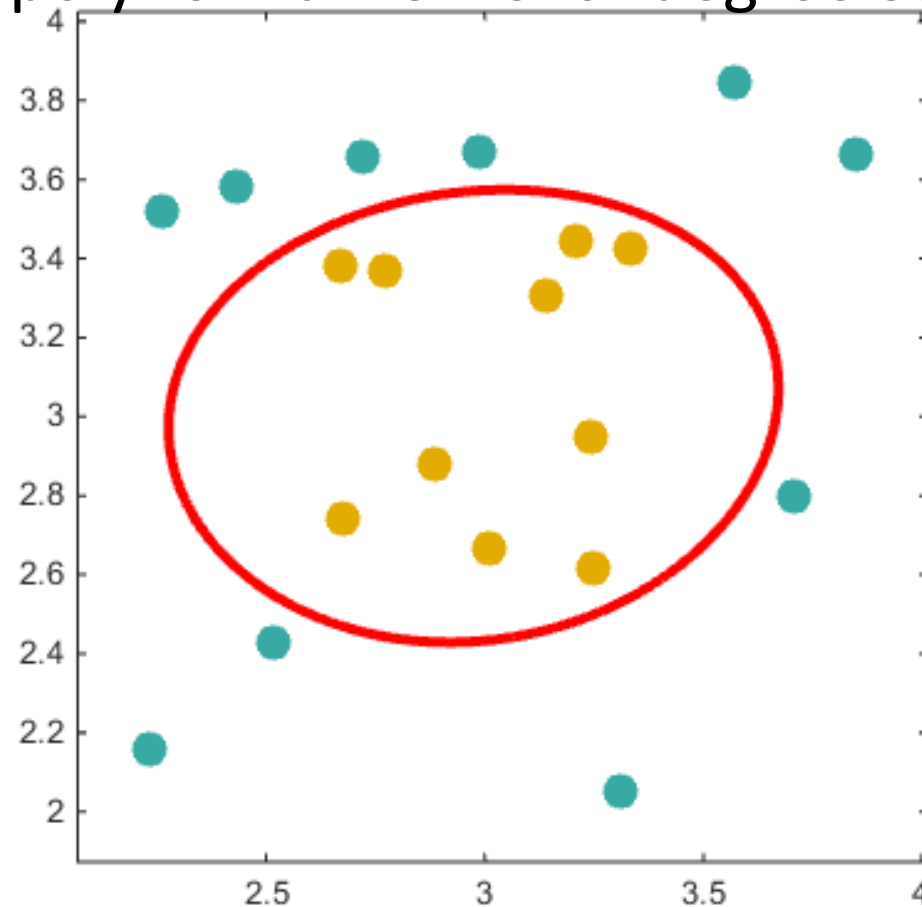
Decision boundary with polynomial kernel of degree 2.

The decision boundary is not linear anymore.



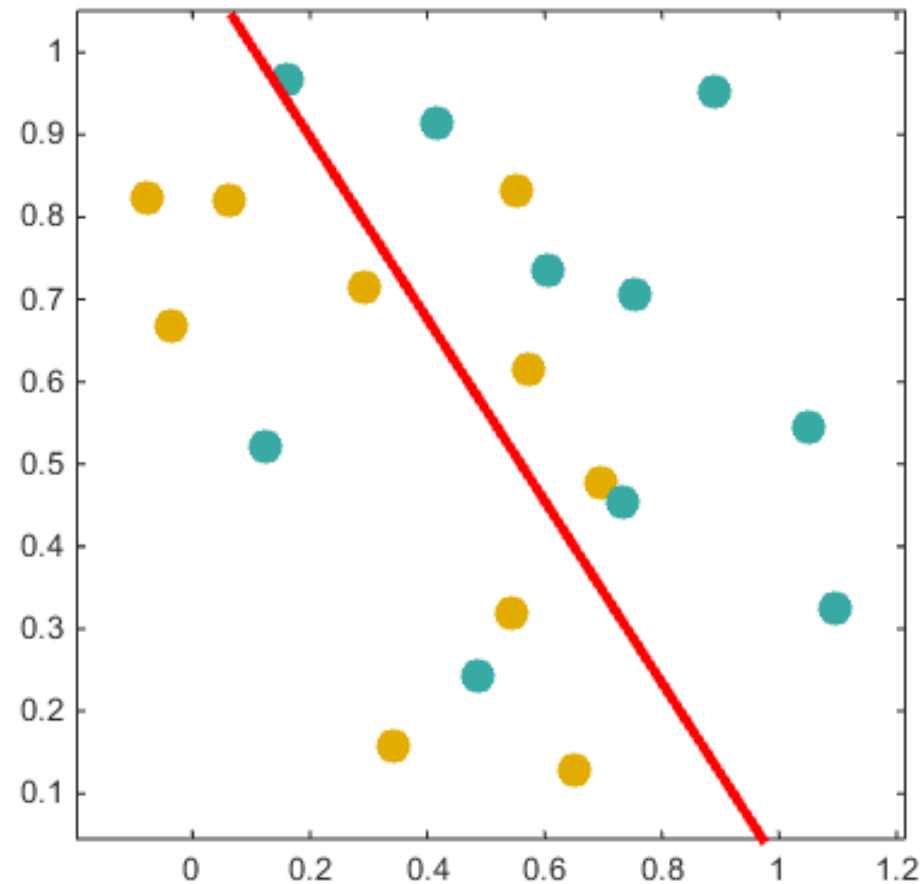
# Polynomial Kernels – An Easy Case

Decision boundary with polynomial kernel of degree 3.



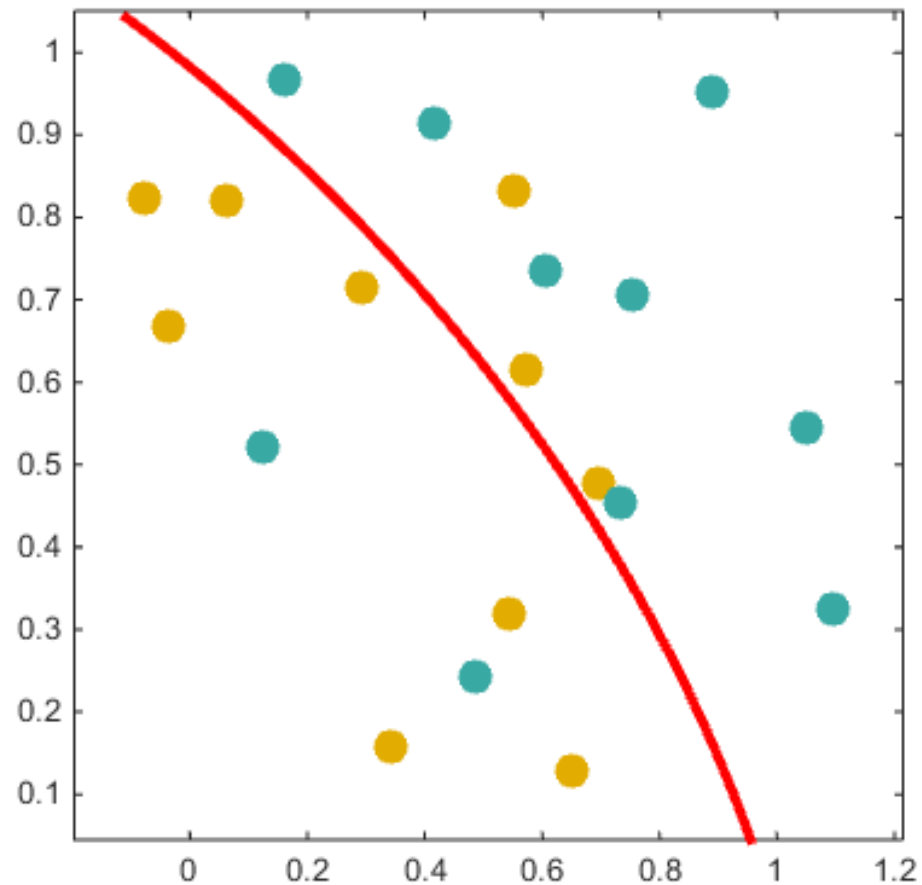
# Polynomial Kernels – A Harder Case

Decision boundary with polynomial kernel of degree 1.



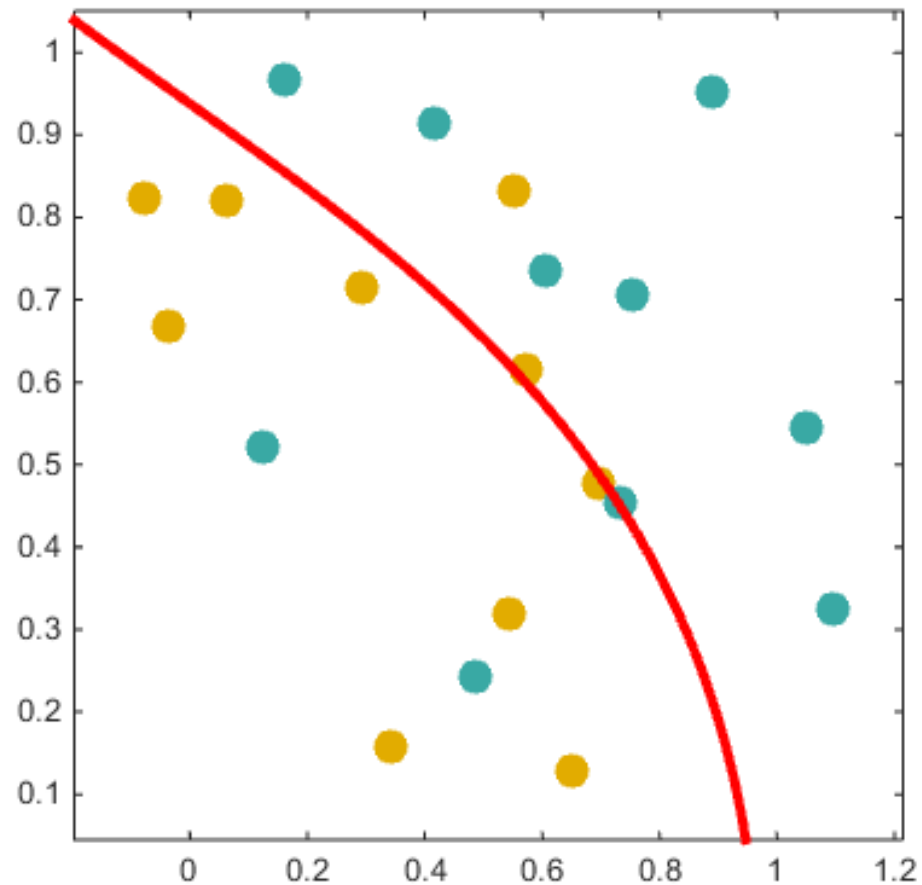
# Polynomial Kernels – A Harder Case

Decision boundary with polynomial kernel of degree 2.



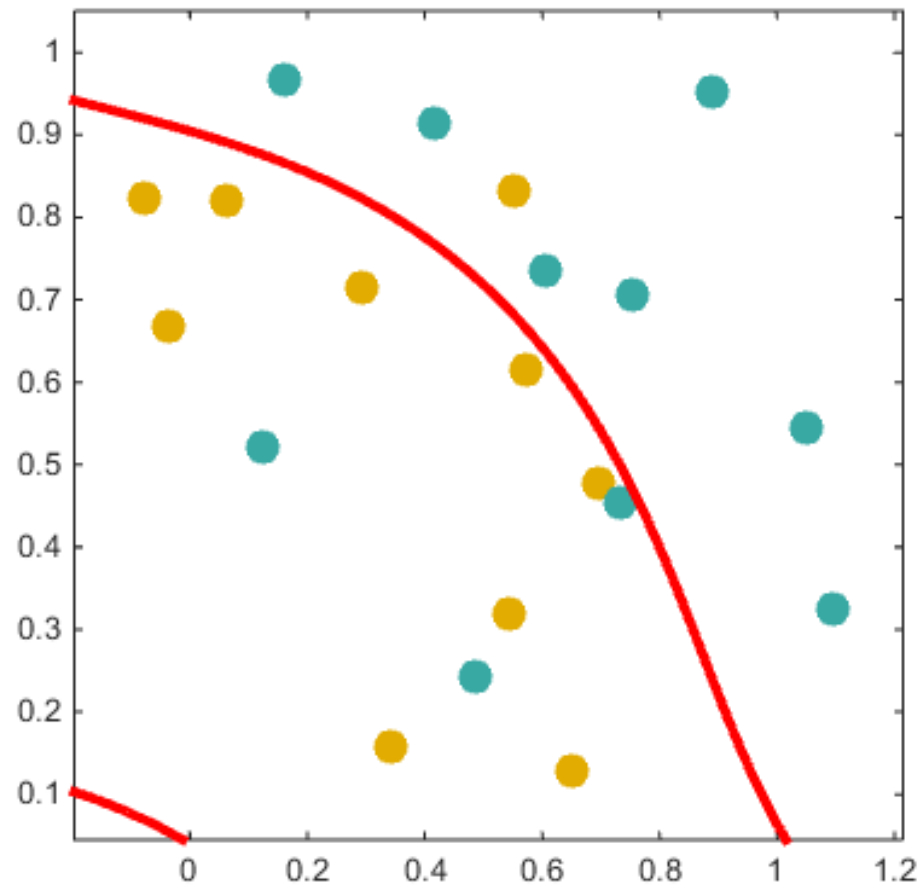
# Polynomial Kernels – A Harder Case

Decision boundary with polynomial kernel of degree 3.



# Polynomial Kernels – A Harder Case

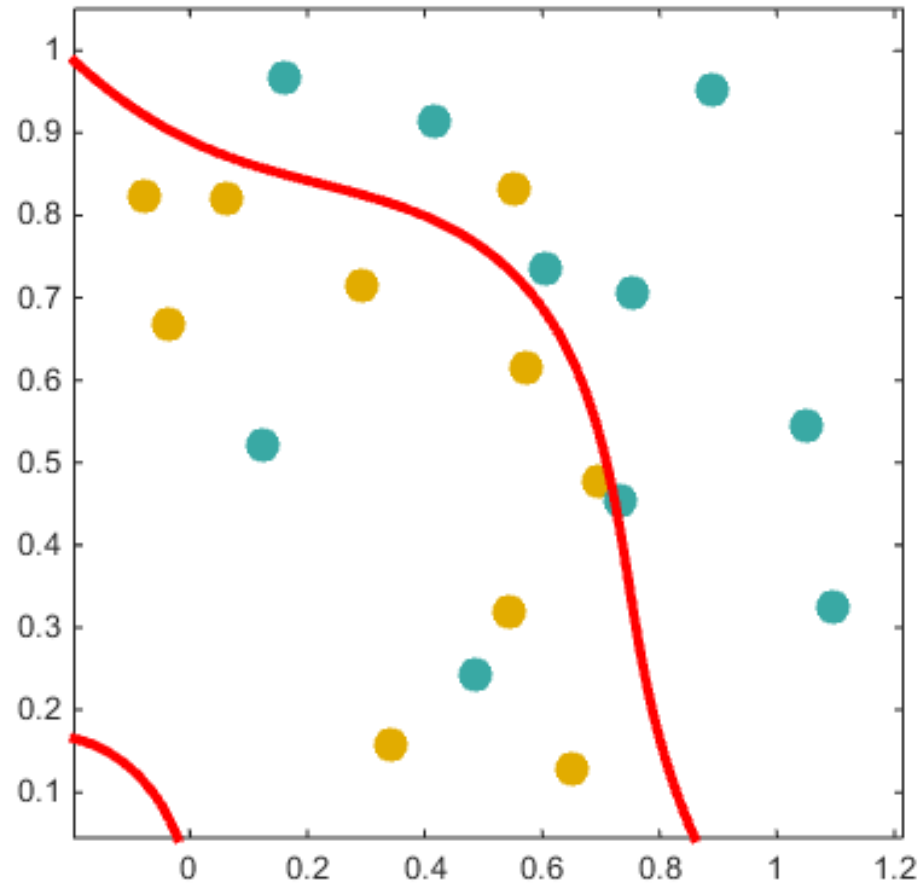
Decision boundary with polynomial kernel of degree 4.





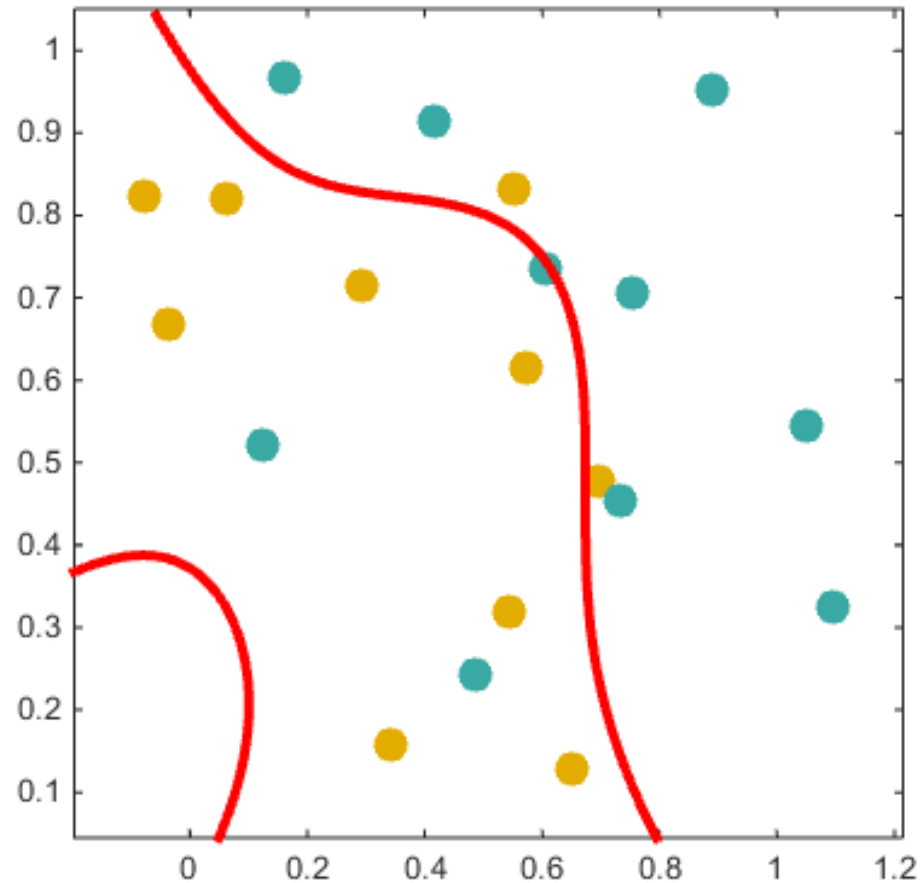
# Polynomial Kernels – A Harder Case

Decision boundary with polynomial kernel of degree 5.



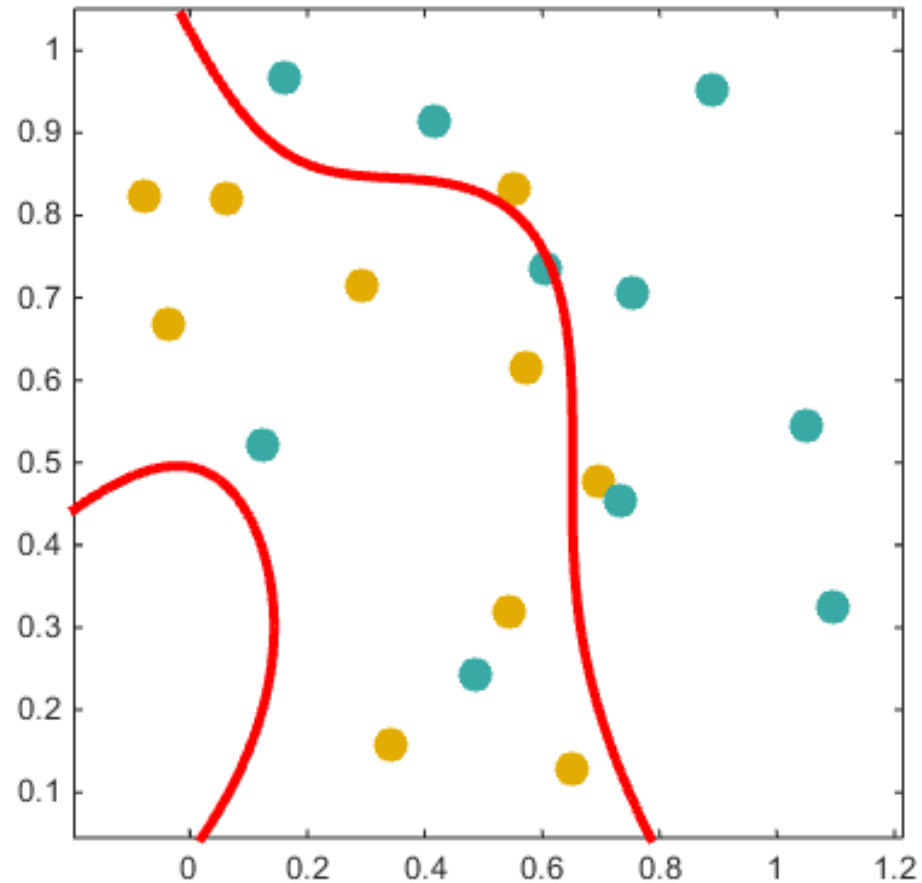
# Polynomial Kernels – A Harder Case

Decision boundary with polynomial kernel of degree 6.



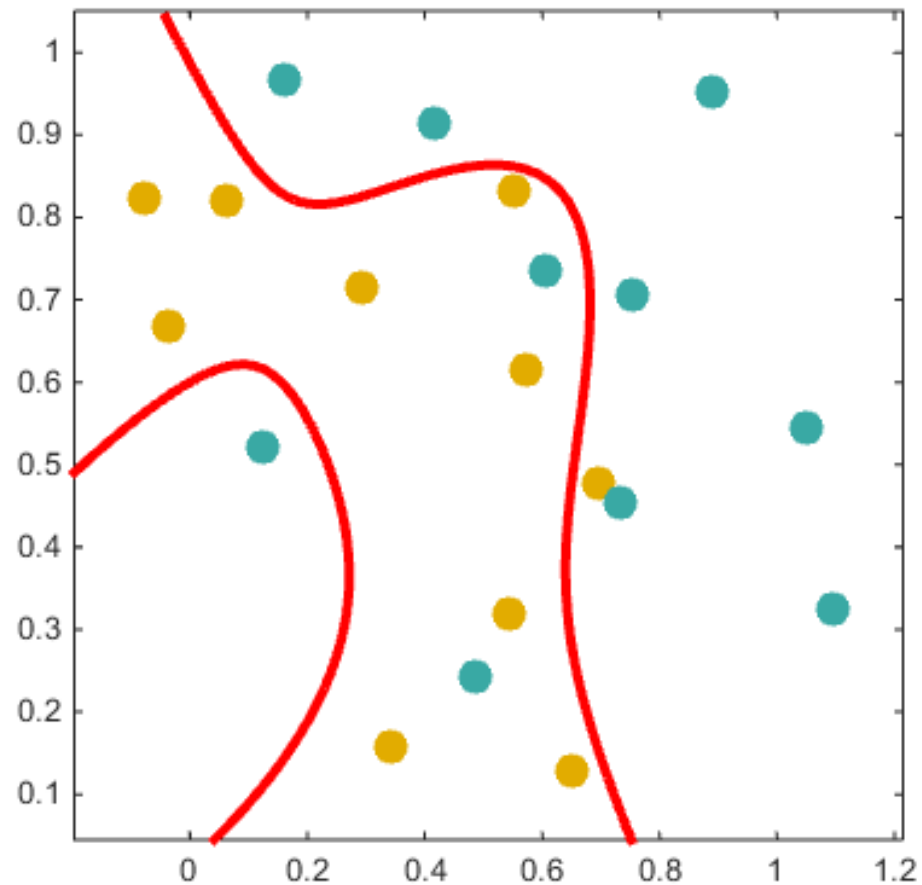
# Polynomial Kernels – A Harder Case

Decision boundary with polynomial kernel of degree 7.



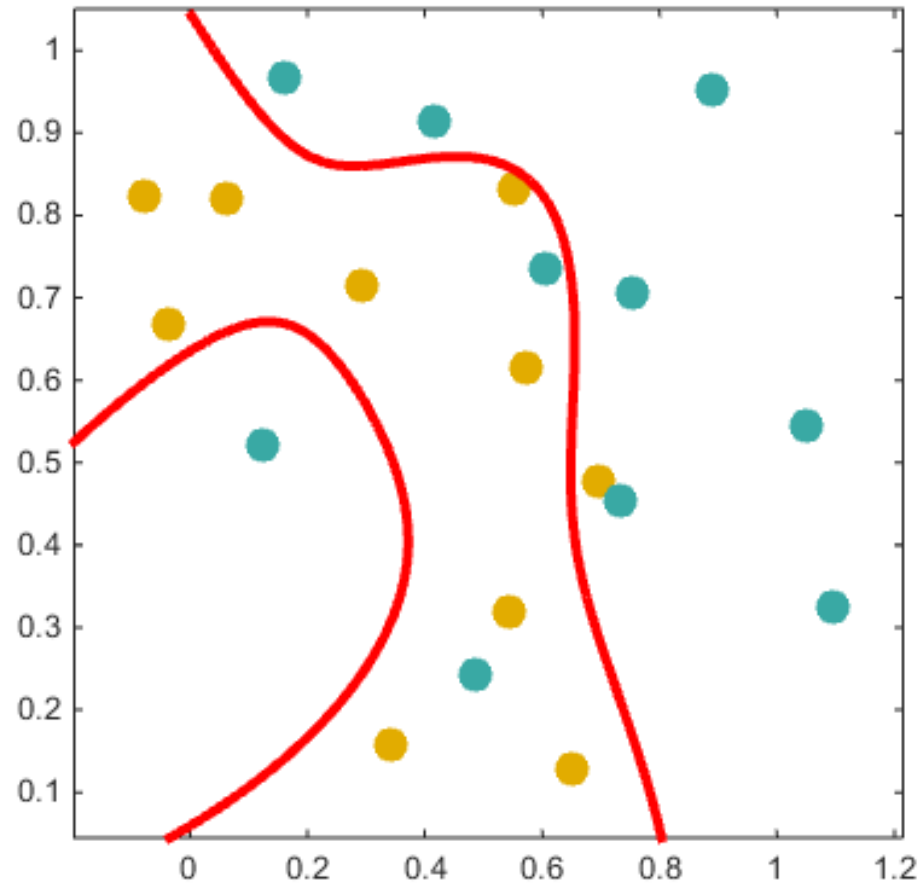
# Polynomial Kernels – A Harder Case

Decision boundary with polynomial kernel of degree 8.



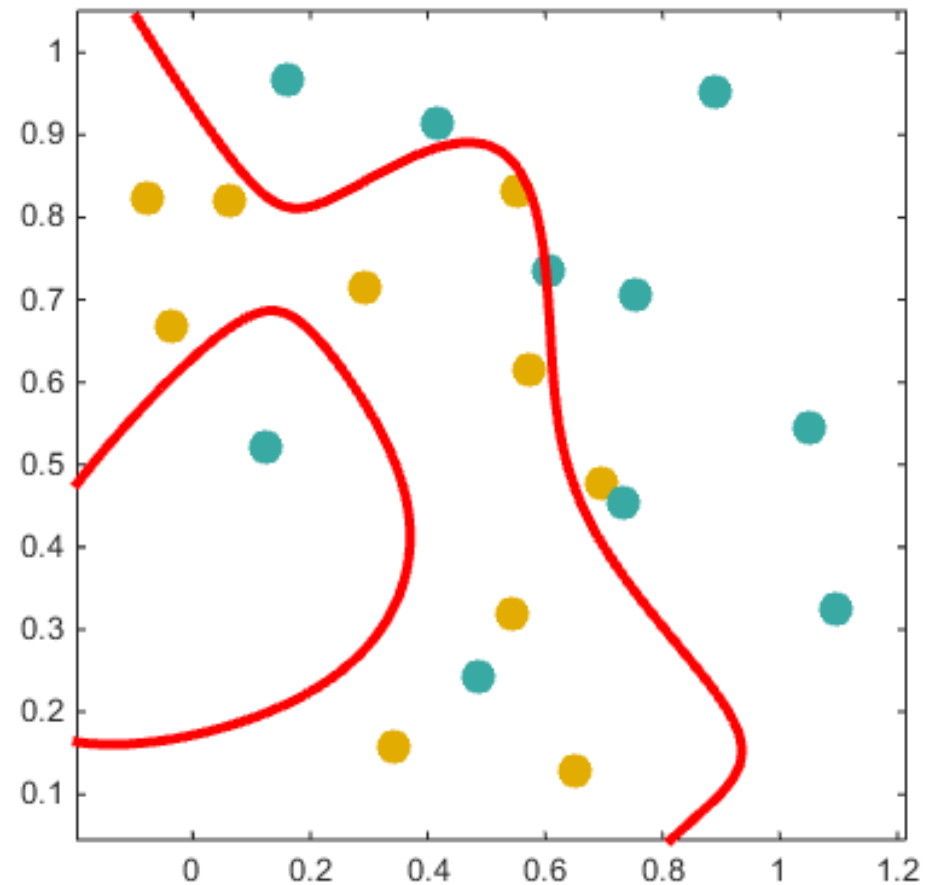
# Polynomial Kernels – A Harder Case

Decision boundary with polynomial kernel of degree 9.



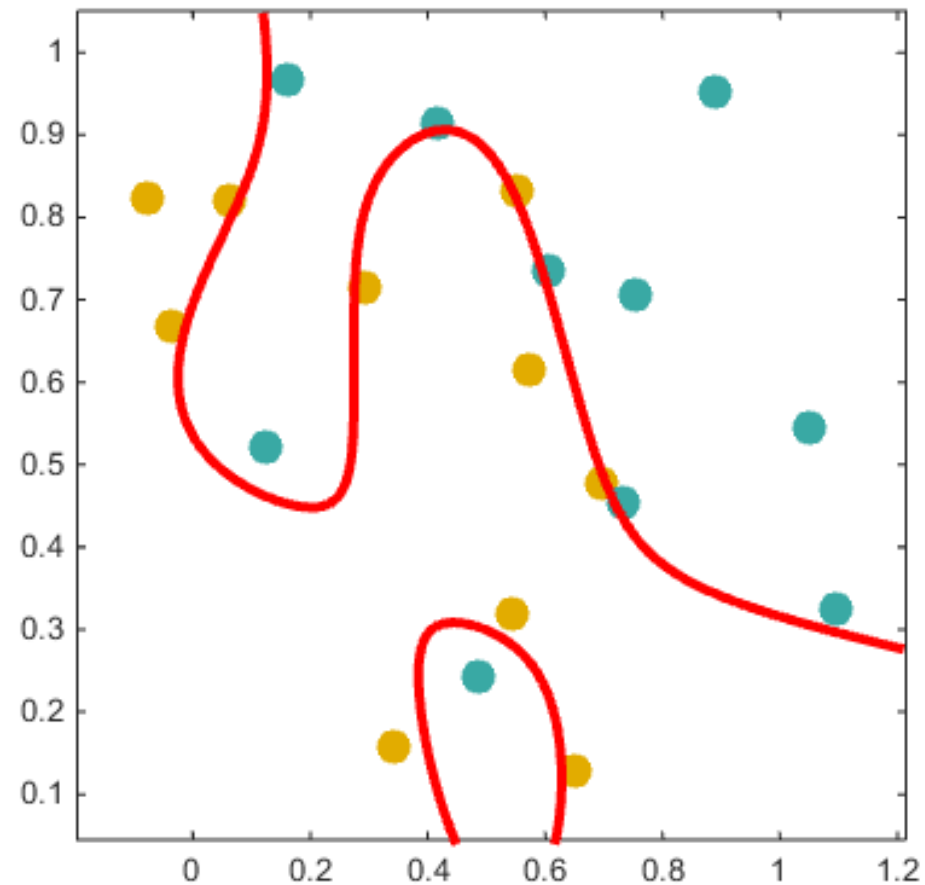
# Polynomial Kernels – A Harder Case

Decision boundary with polynomial kernel of degree 10.



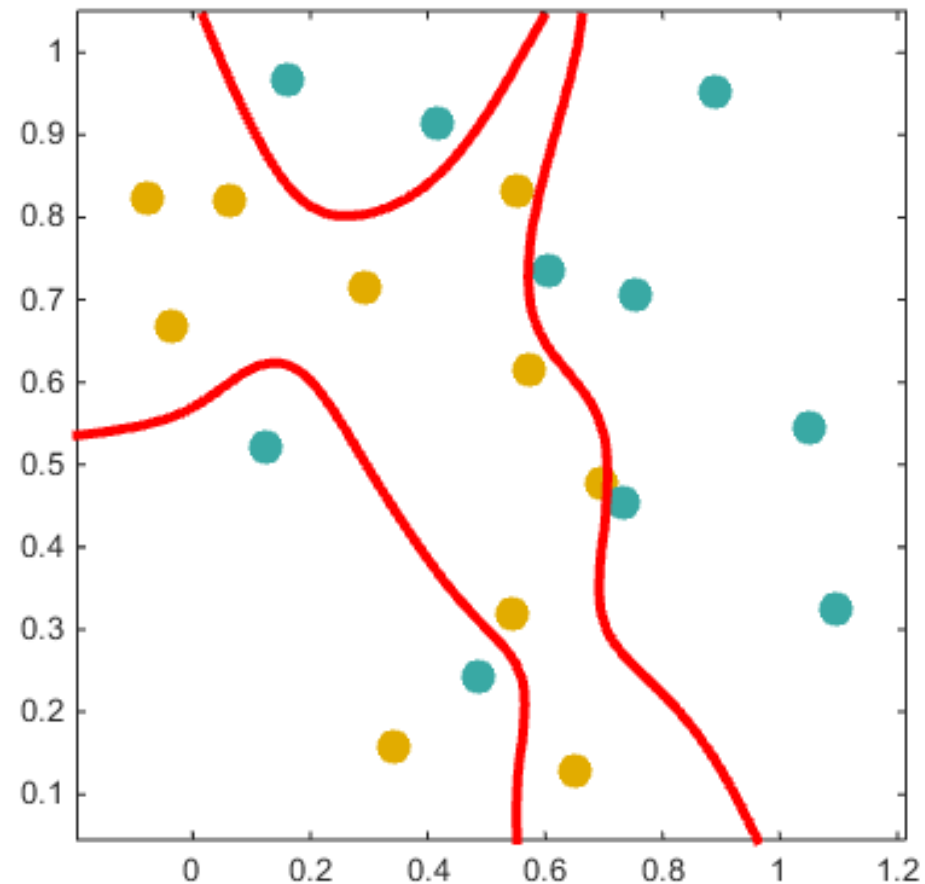
# Polynomial Kernels – A Harder Case

Decision boundary with polynomial kernel of degree 20.



# Polynomial Kernels – A Harder Case

Decision boundary with polynomial kernel of degree 100.





# RBF/Gaussian Kernels

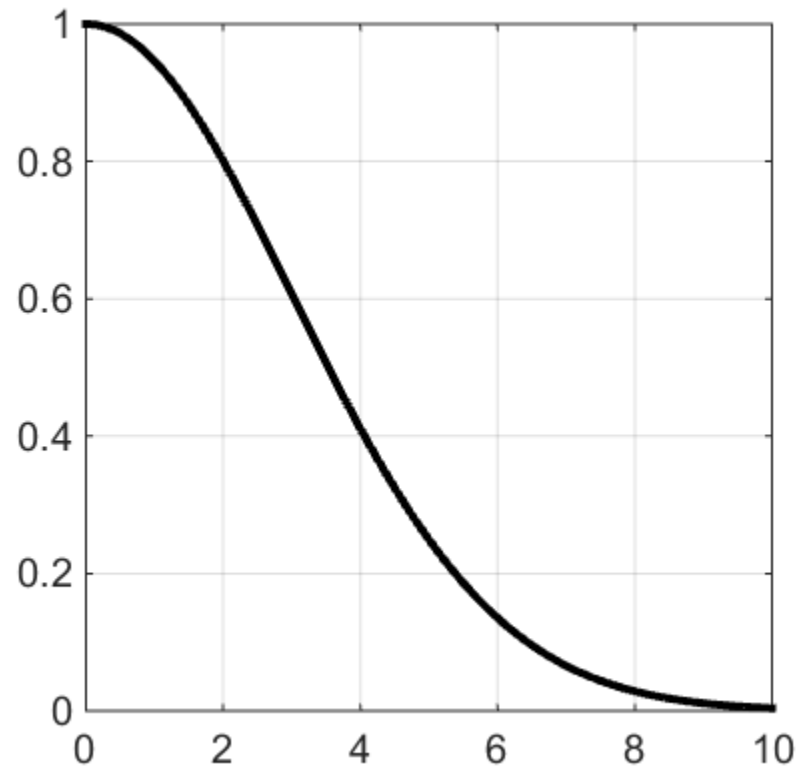
- The **Radial Basis Function (RBF) kernel**, also known as **Gaussian kernel**, is defined as:

$$k_{\sigma}(\mathbf{x}, \mathbf{z}) = e^{-\frac{\|\mathbf{x}-\mathbf{z}\|^2}{2\sigma^2}}$$

- Given  $\sigma$ , the value of  $k_{\sigma}(\mathbf{x}, \mathbf{z})$  only depends on the distance between  $\mathbf{x}$  and  $\mathbf{z}$ .
  - $k_{\sigma}(\mathbf{x}, \mathbf{z})$  decreases exponentially to the distance between  $\mathbf{x}$  and  $\mathbf{z}$ .
- Parameter  $\sigma$  is chosen manually.
  - Parameter  $\sigma$  specifies how fast  $k_{\sigma}(\mathbf{x}, \mathbf{z})$  decreases as  $\mathbf{x}$  moves away from  $\mathbf{z}$ .

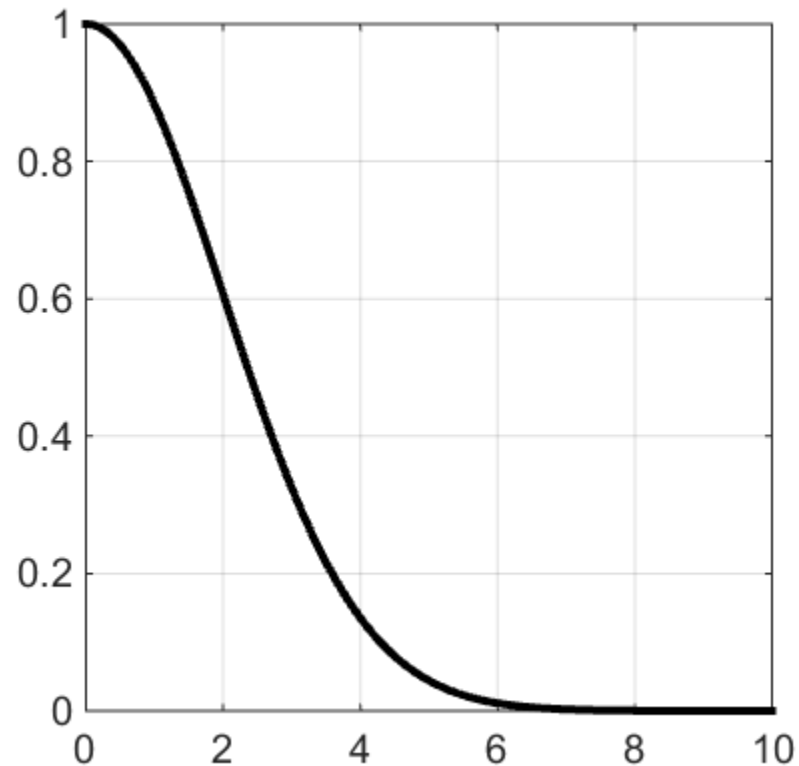
# RBF Output Vs. Distance

- X axis: distance between  $\mathbf{x}$  and  $\mathbf{z}$ .
- Y axis:  $k_{\sigma}(\mathbf{x}, \mathbf{z})$ , with  $\sigma = 3$ .



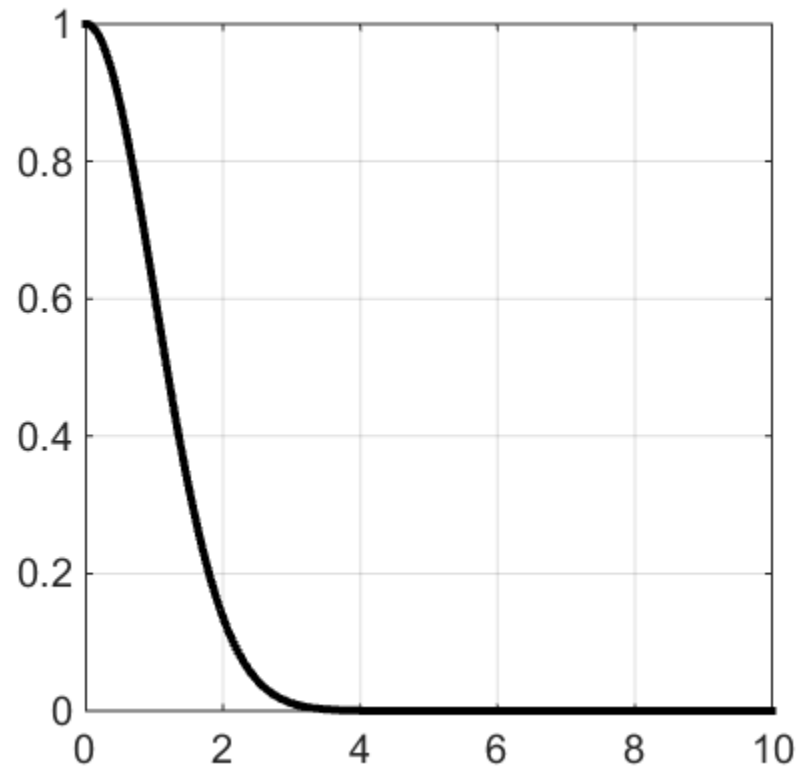
# RBF Output Vs. Distance

- X axis: distance between  $\mathbf{x}$  and  $\mathbf{z}$ .
- Y axis:  $k_{\sigma}(\mathbf{x}, \mathbf{z})$ , with  $\sigma = 2$ .



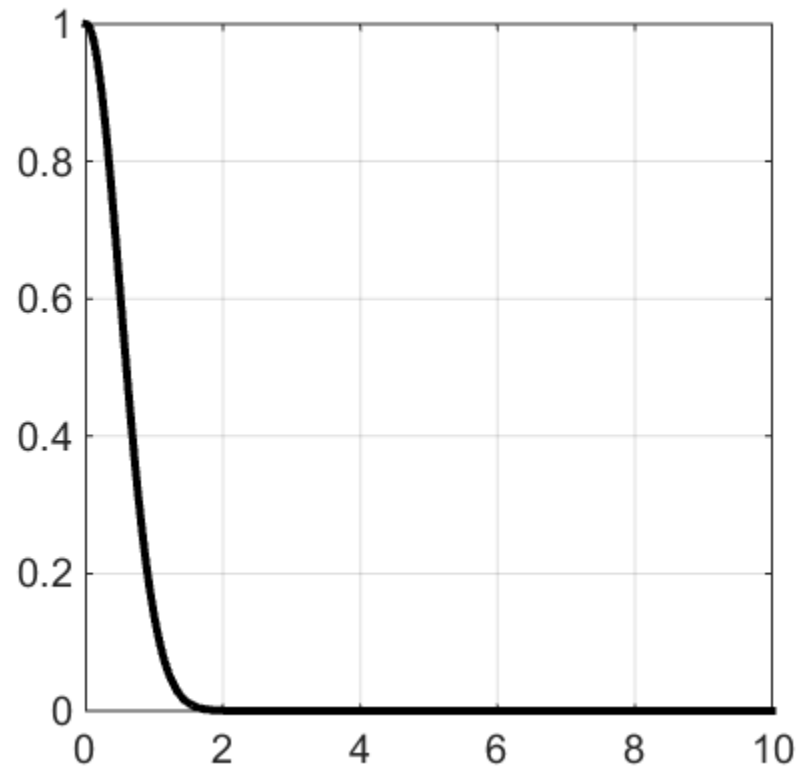
# RBF Output Vs. Distance

- X axis: distance between  $\mathbf{x}$  and  $\mathbf{z}$ .
- Y axis:  $k_{\sigma}(\mathbf{x}, \mathbf{z})$ , with  $\sigma = 1$ .



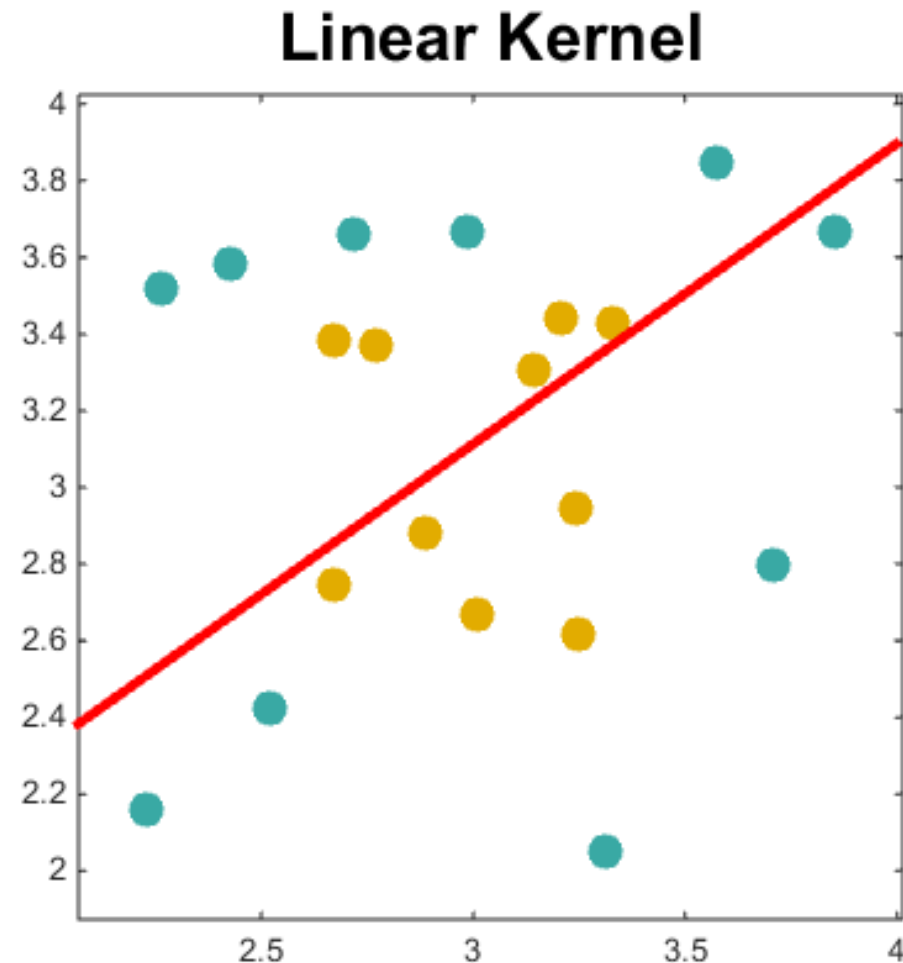
# RBF Output Vs. Distance

- X axis: distance between  $\mathbf{x}$  and  $\mathbf{z}$ .
- Y axis:  $k_{\sigma}(\mathbf{x}, \mathbf{z})$ , with  $\sigma = 0.5$ .



# RBF Kernels – An Easier Dataset

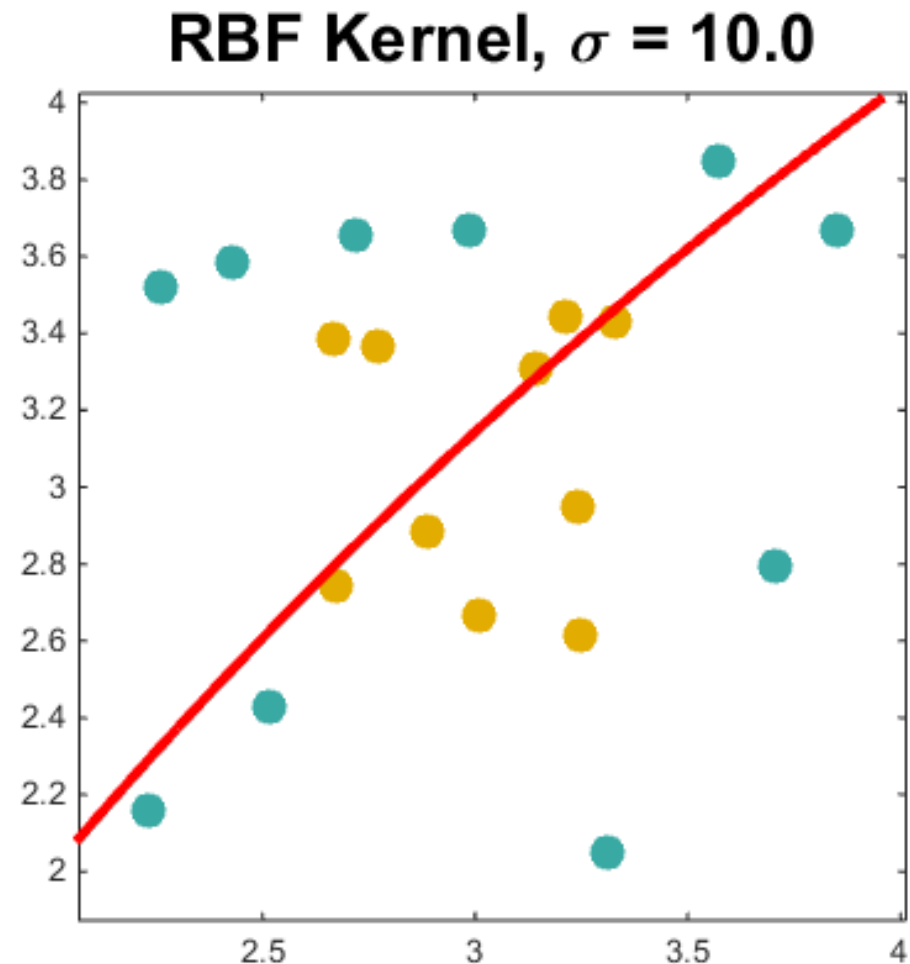
Decision boundary  
with a linear kernel.



# RBF Kernels – An Easier Dataset

Decision boundary  
with an RBF kernel.

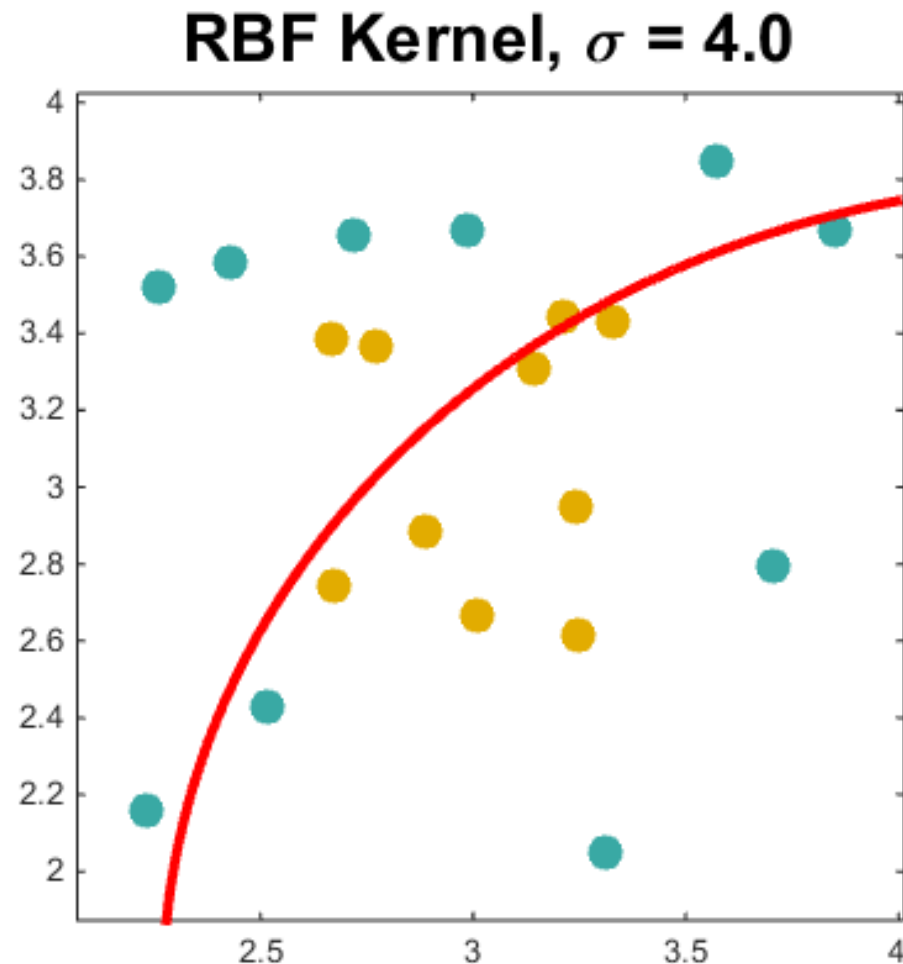
For this dataset, this is a relatively large  
value for  $\sigma$ , and it produces a boundary  
that is almost linear.



# RBF Kernels – An Easier Dataset

Decision boundary  
with an RBF kernel.

Decreasing the value of  $\sigma$  leads  
to less linear boundaries.

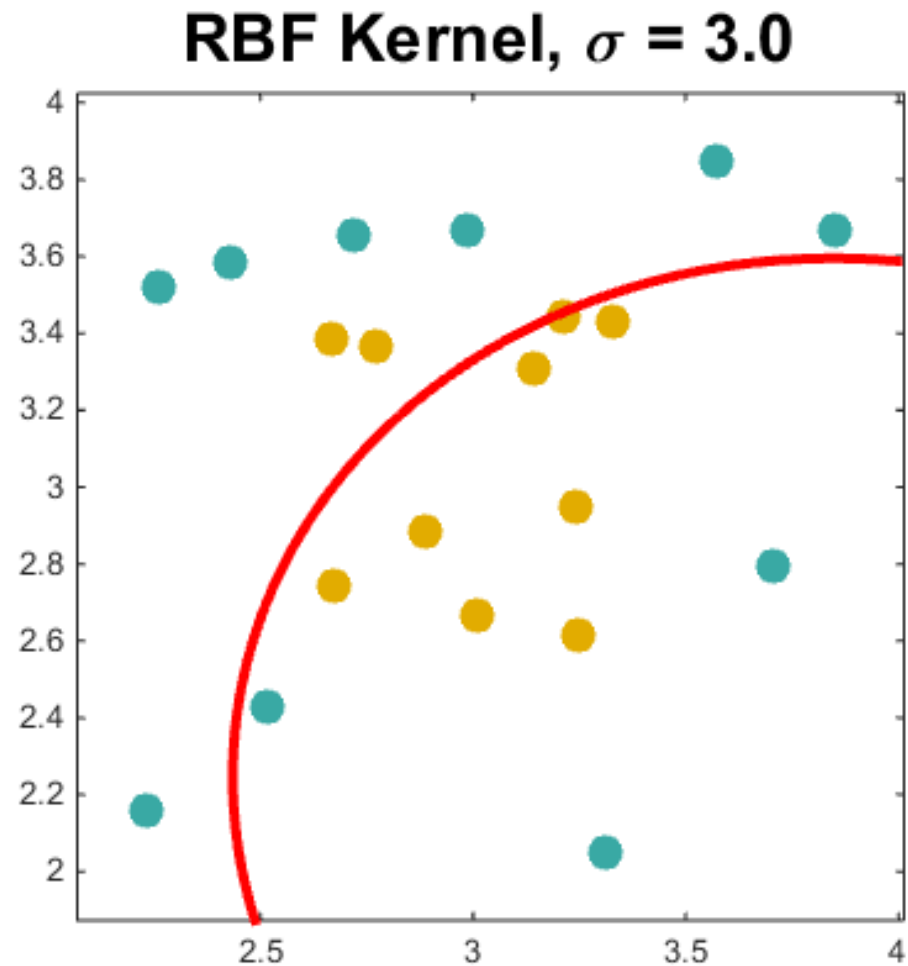




# RBF Kernels – An Easier Dataset

Decision boundary  
with an RBF kernel.

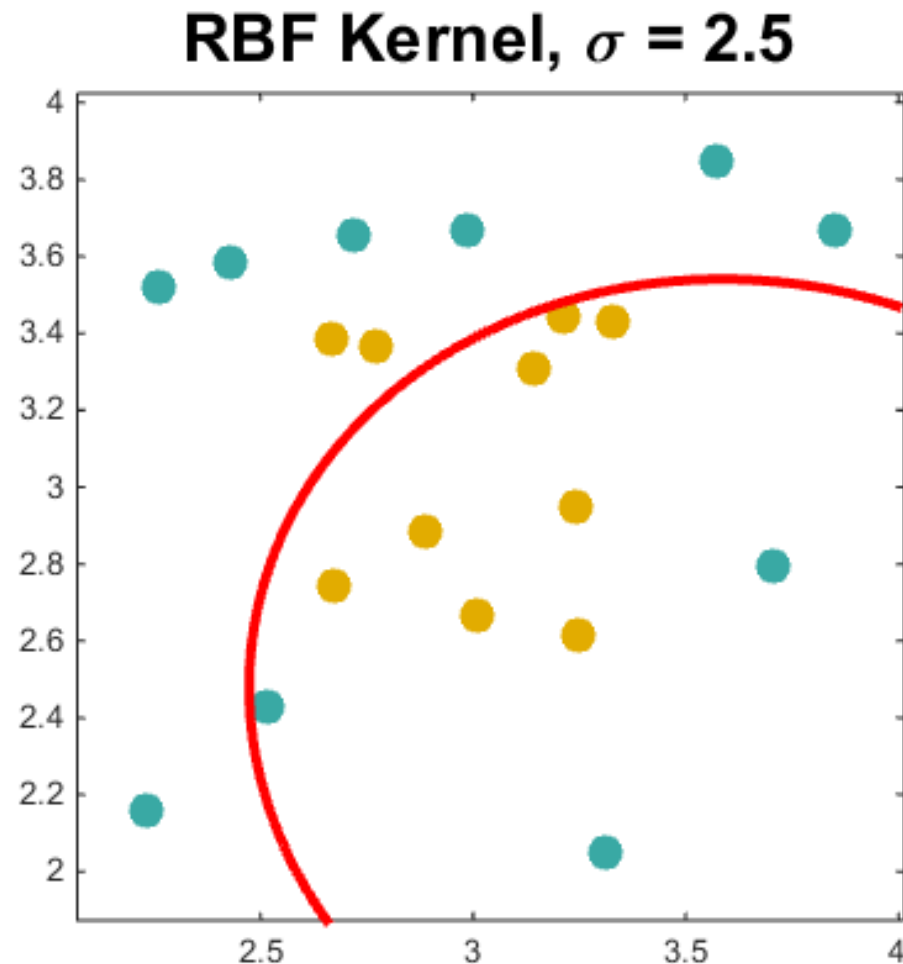
Decreasing the value of  $\sigma$  leads  
to less linear boundaries.



# RBF Kernels – An Easier Dataset

Decision boundary  
with an RBF kernel.

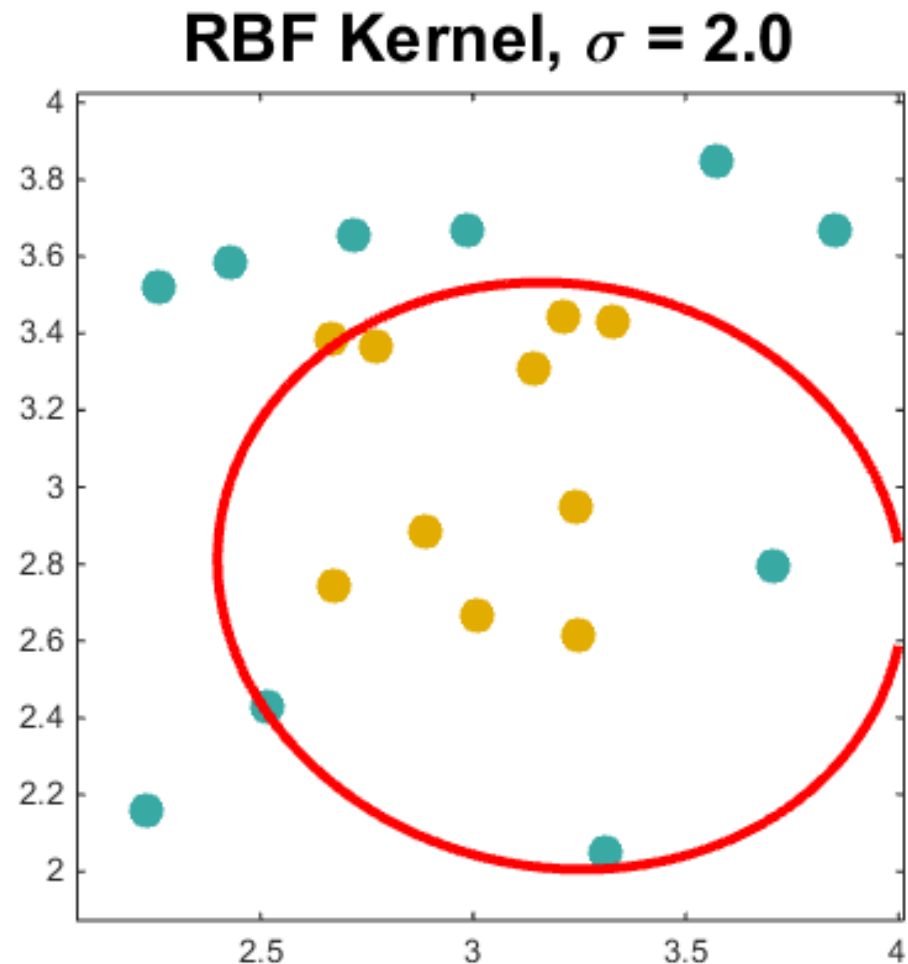
Decreasing the value of  $\sigma$  leads  
to less linear boundaries.



# RBF Kernels – An Easier Dataset

Decision boundary  
with an RBF kernel.

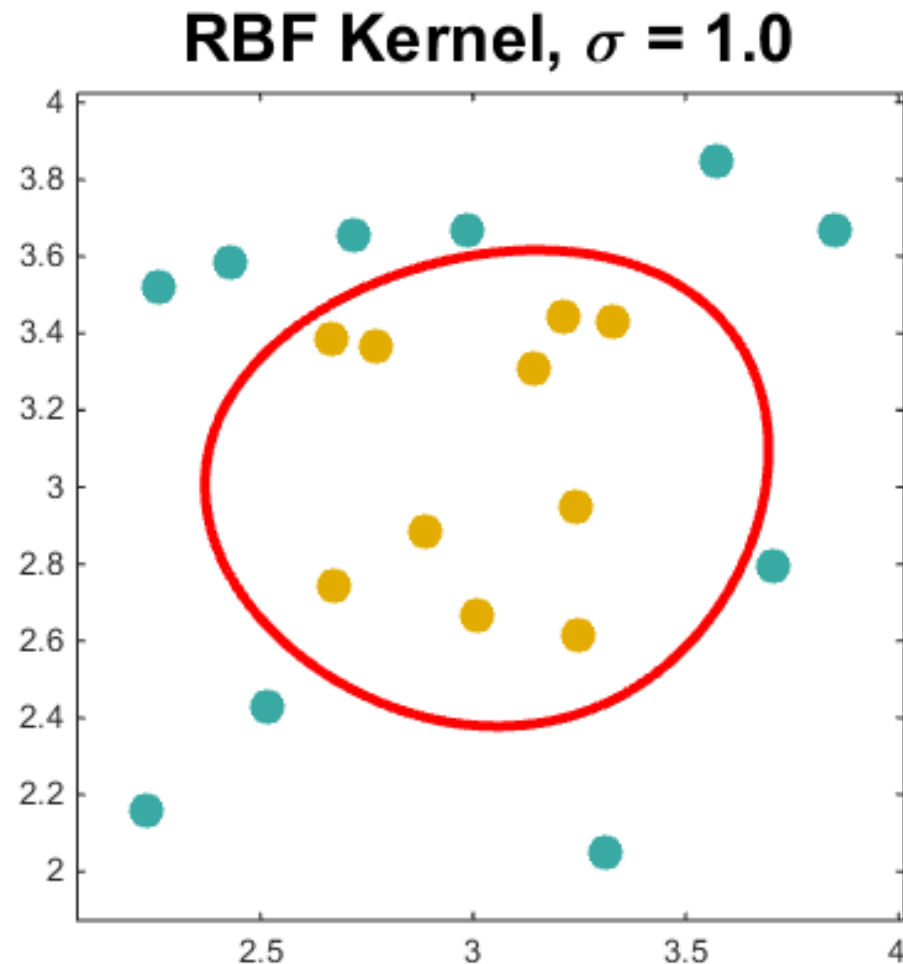
Decreasing the value of  $\sigma$  leads  
to less linear boundaries.



# RBF Kernels – An Easier Dataset

Decision boundary  
with an RBF kernel.

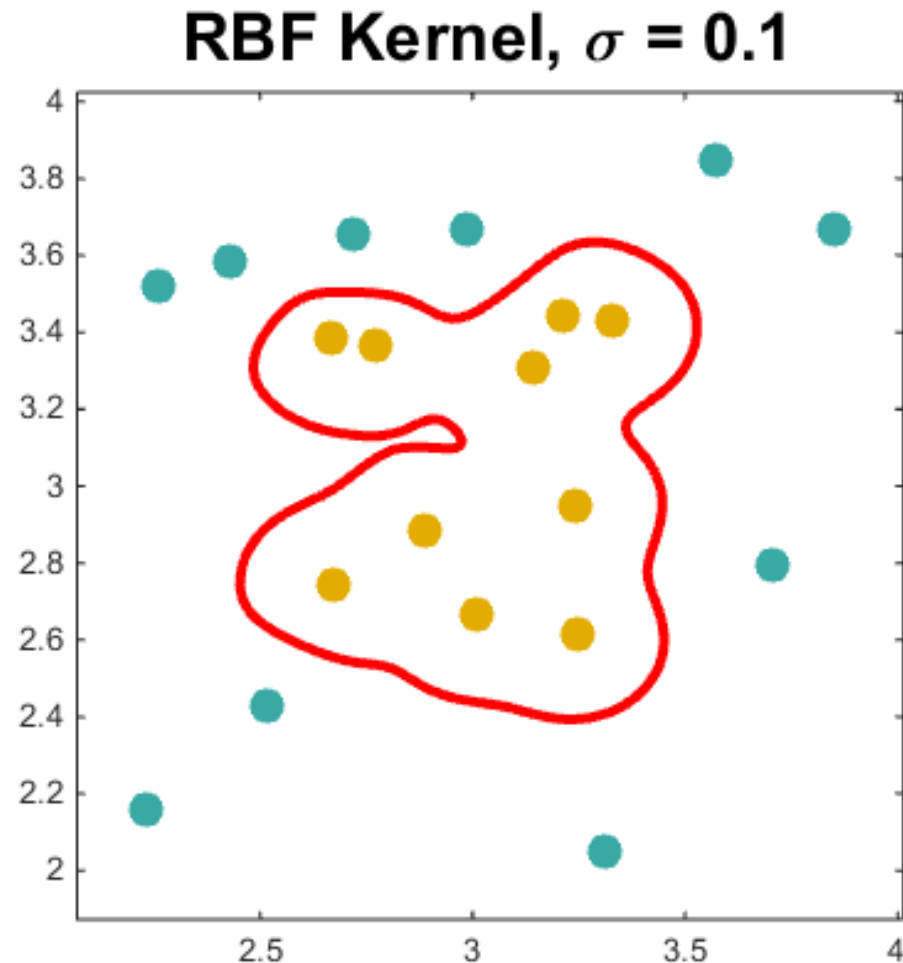
Decreasing the value of  $\sigma$  leads  
to less linear boundaries.



# RBF Kernels – An Easier Dataset

Decision boundary  
with an RBF kernel.

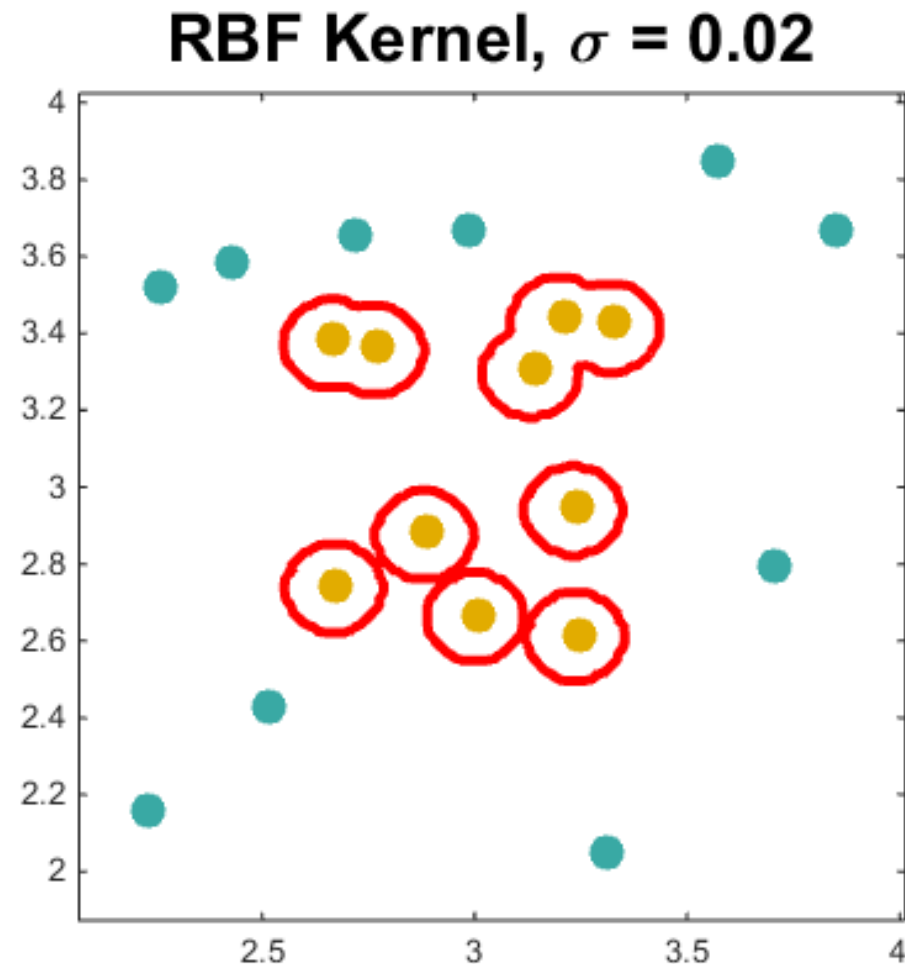
Note that smaller  
values of  $\sigma$  increase dangers of  
overfitting.



# RBF Kernels – An Easier Dataset

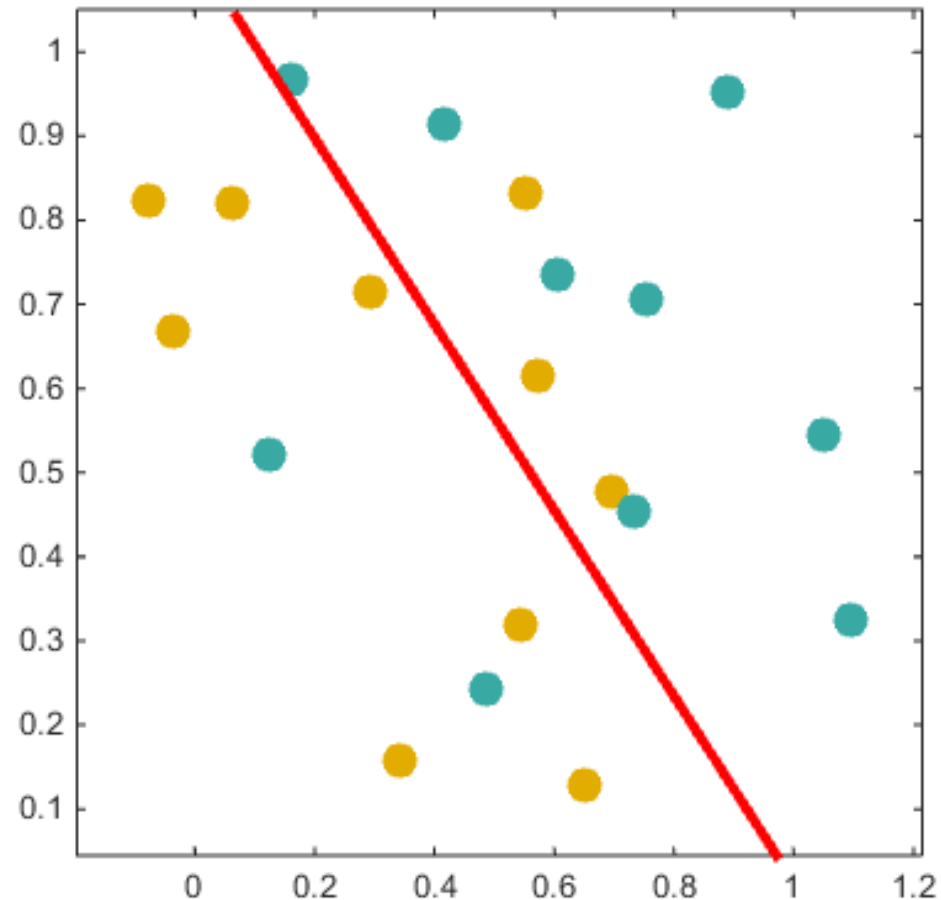
Decision boundary  
with an RBF kernel.

Note that smaller  
values of  $\sigma$  increase dangers of  
overfitting.



# RBF Kernels – A Harder Dataset

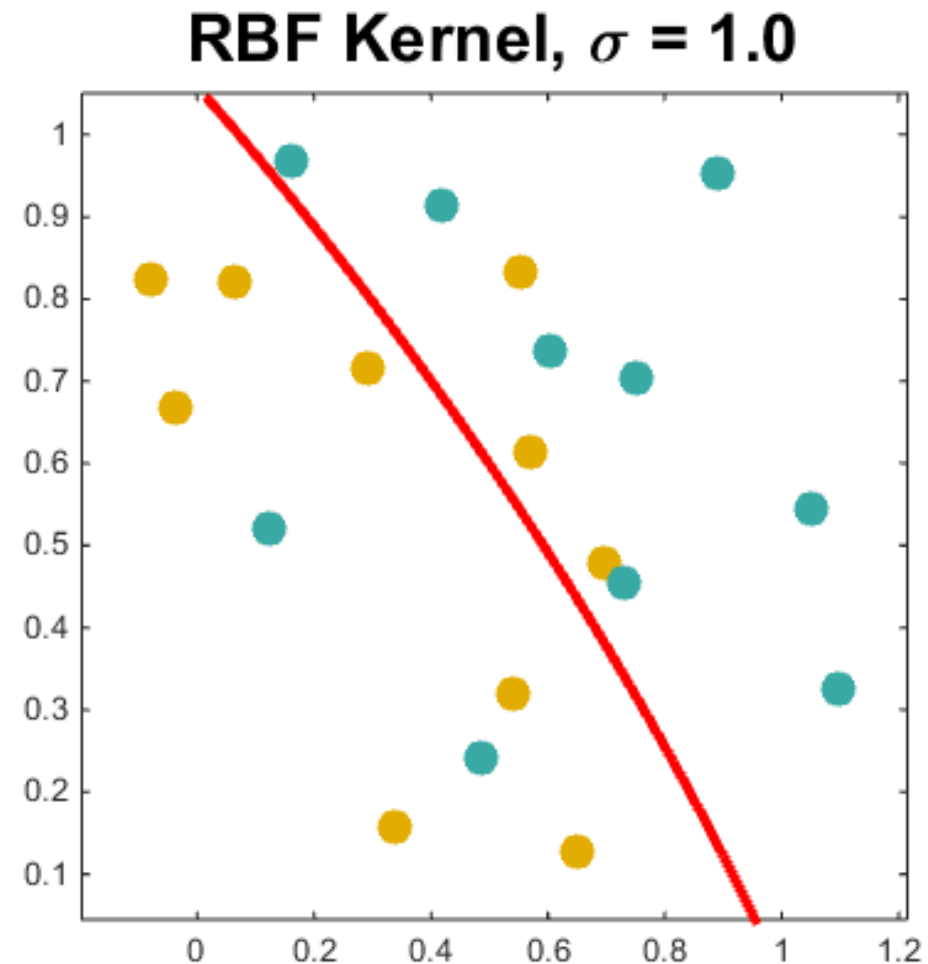
Decision boundary  
with a linear kernel.



# RBF Kernels – A Harder Dataset

Decision boundary  
with an RBF kernel.

The boundary is  
almost linear.

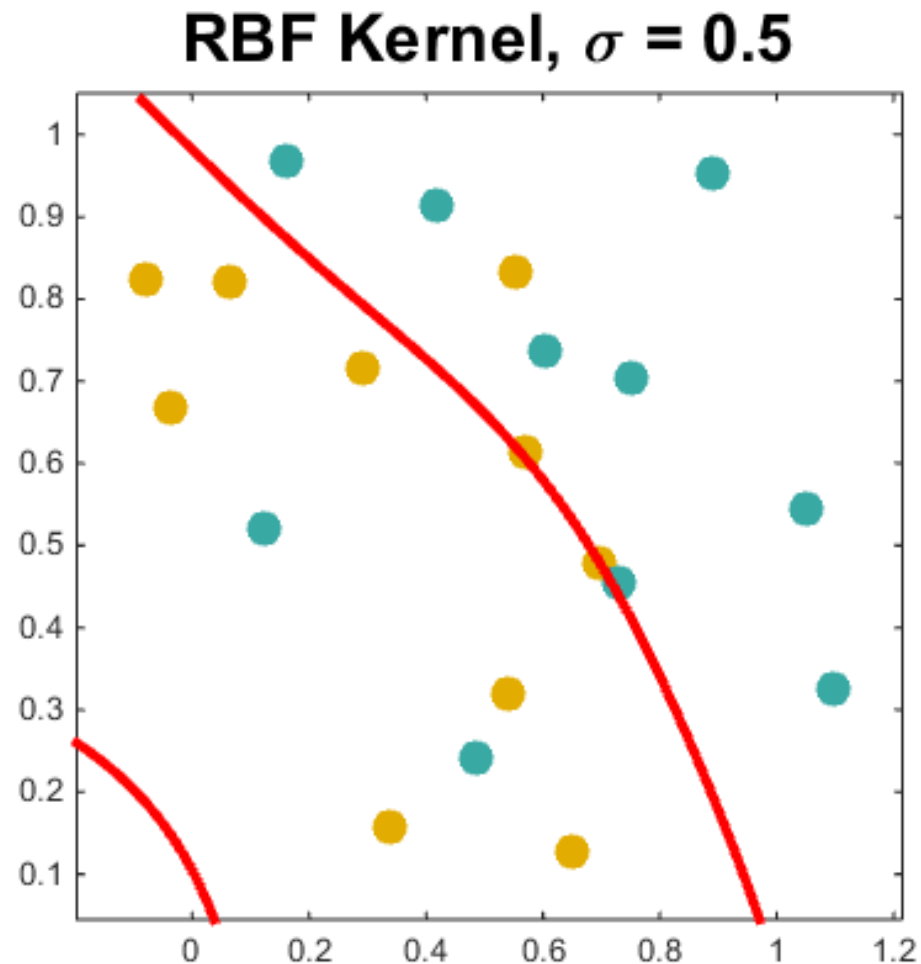




# RBF Kernels – A Harder Dataset

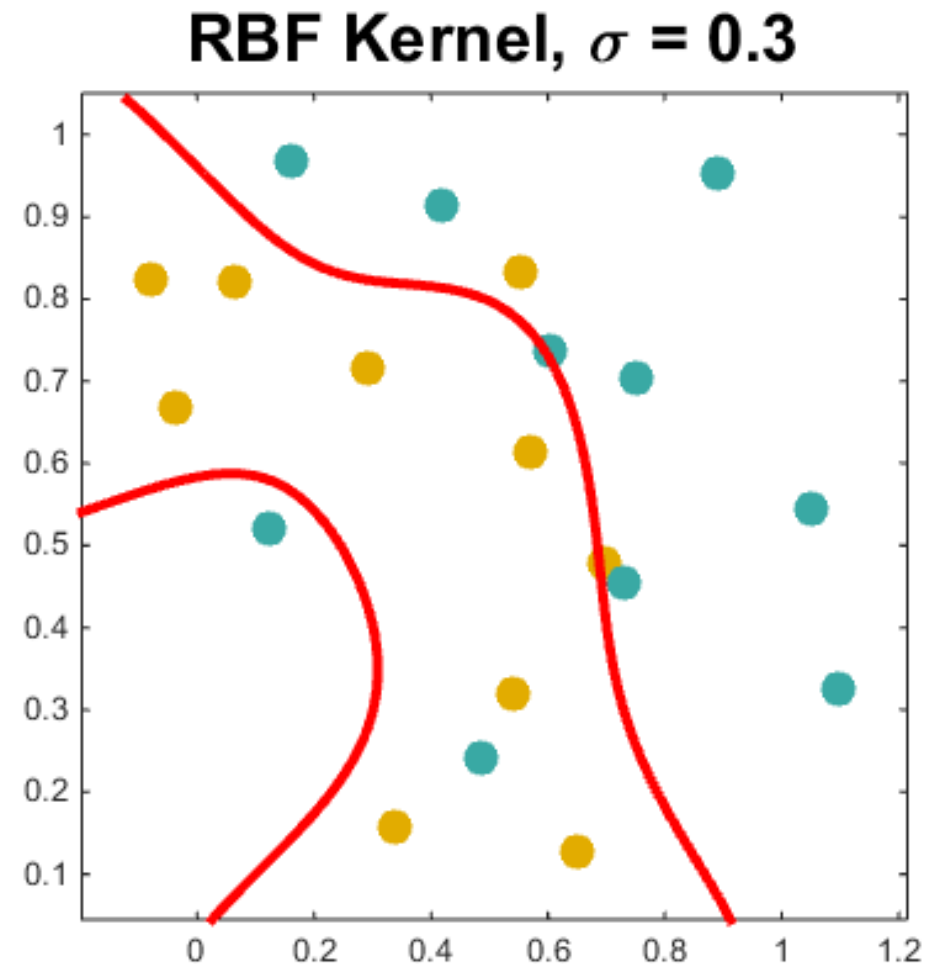
Decision boundary  
with an RBF kernel.

The boundary now  
is clearly nonlinear.



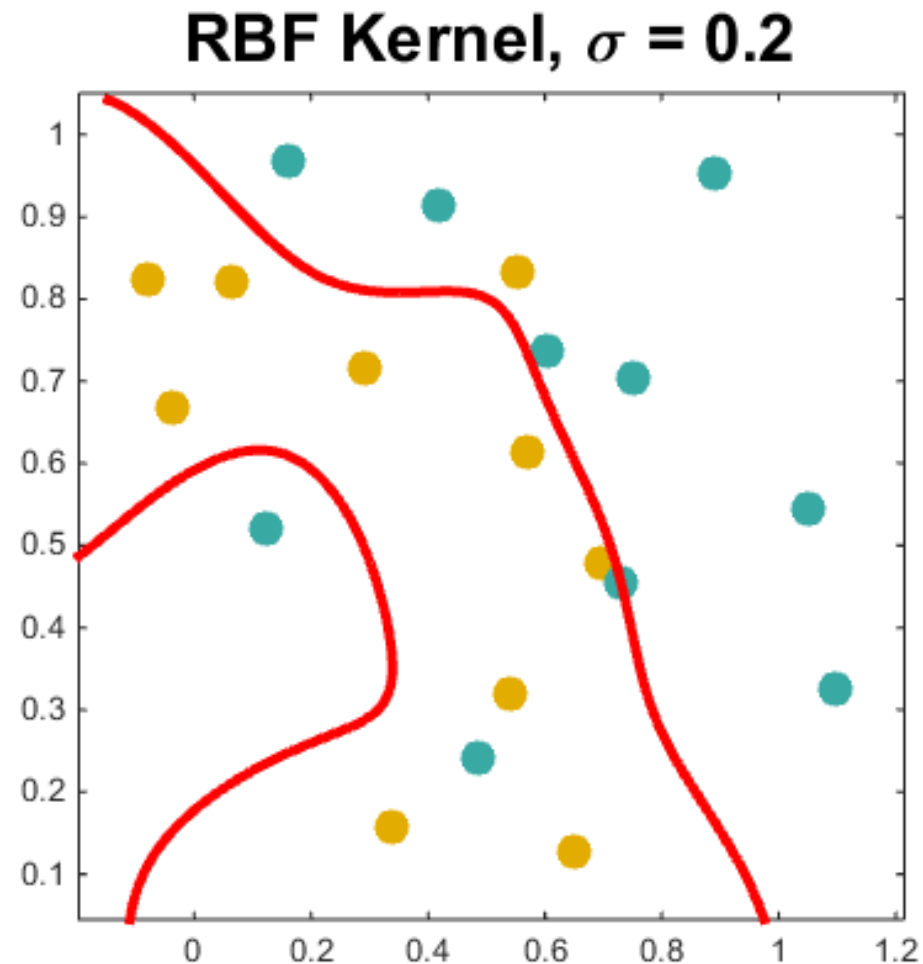
# RBF Kernels – A Harder Dataset

Decision boundary  
with an RBF kernel.



# RBF Kernels – A Harder Dataset

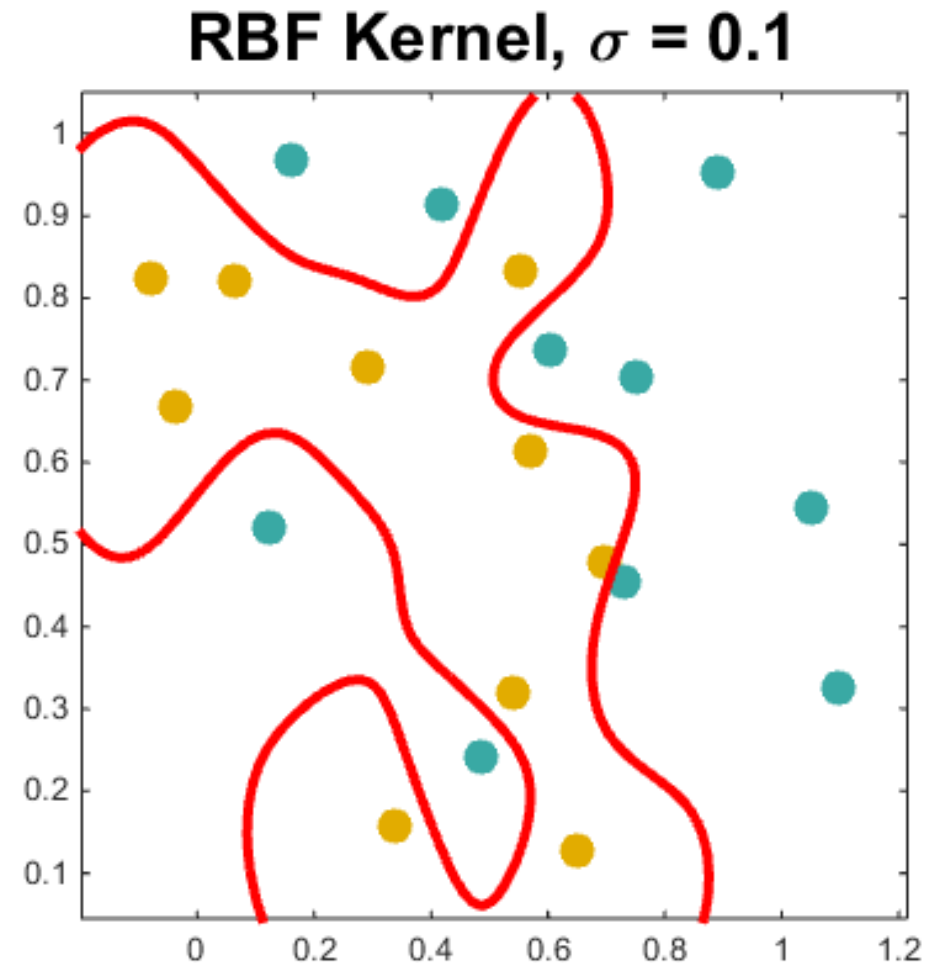
Decision boundary  
with an RBF kernel.



# RBF Kernels – A Harder Dataset

Decision boundary  
with an RBF kernel.

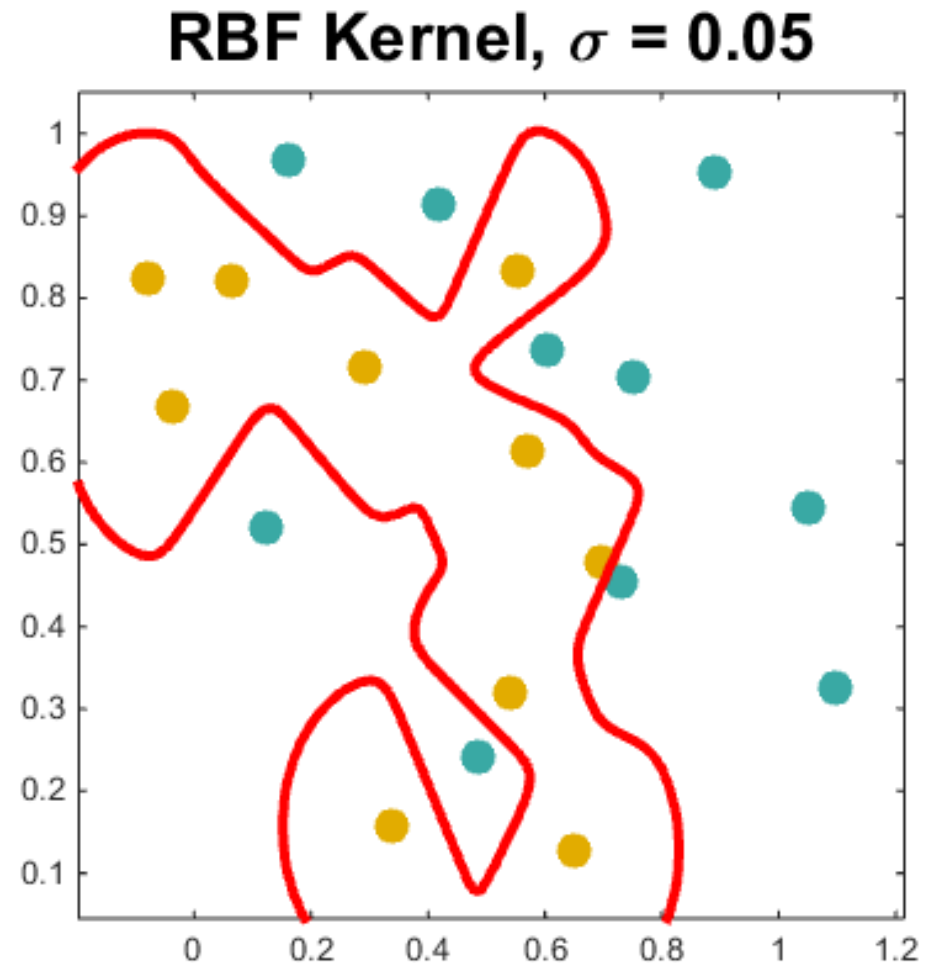
Again, smaller values  
of  $\sigma$  increase dangers of overfitting.



# RBF Kernels – A Harder Dataset

Decision boundary  
with an RBF kernel.

Again, smaller values  
of  $\sigma$  increase dangers of overfitting.



# SVM parameters

SVM has another set of parameters called hyperparameters.

- The soft margin constant  $C$ .
- Any parameters the kernel function depends on
  - linear kernel – no hyperparameter (except for  $C$ )
  - polynomial – degree
  - Gaussian – width of Gaussian

- So which kernel and which parameters should I use?
- The answer is data-dependent.
- Several kernels should be tried.
- Try linear kernel first and then see, if the classification can be improved with nonlinear kernels (tradeoff between quality of the kernel and the number of dimensions).
- Select kernel + parameters +  $C$  by crossvalidation.

# Practical Aspects

Many beginners use the following procedure:

- Transform data to the format of an SVM software (often obsolete)
- Randomly try a few kernels and parameters
- Test

Instead try

- Conduct simple data scaling/normalization
- Consider the RBF kernel first
- Use cross-validation to find the best parameter  $C$  and  $\gamma$
- Use the best parameter  $C$  and  $\gamma$  to test model



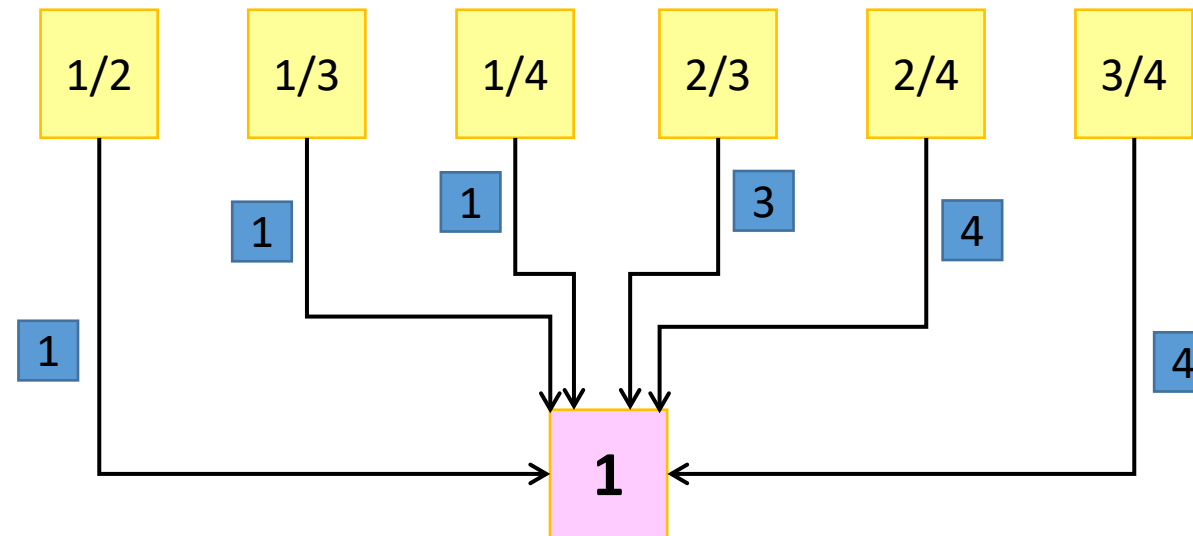
# Computational aspects

- Classification of new samples is very quick, training is longer (reasonably fast for thousands of samples).
- Linear kernel – scales linearly.
- Nonlinear kernels – scale quadratically.

# Multiclass SVM

- SVM is defined for binary classification.
- How to predict more than two classes (multiclass)?
- Simplest approach: decompose the multiclass problem into several binary problems and train several binary SVM's.

- one-versus-one approach
  - Train a binary SVM for any two classes from the training set
  - For  $k$ -class problem create  $\frac{k(k-1)}{2}$  SVM models
  - Prediction: voting procedure assigns the class to be the class with the maximum votes



- one-versus-all approach
  - For  $k$ -class problem train only  $k$  SVM models.
  - Each will be trained to predict one class (+1) vs. the rest of classes (-1)
  - Prediction:
    - Winner takes all strategy
    - Assign new example to the class with the largest output value  $f(\mathbf{x})$ .

1/rest

2/rest

3/rest

4/rest

# SVMs: Recap

- Advantages:
  - Training finds globally best solution.
    - No need to worry about local optima, or iterations.
  - SVMs can define complex decision boundaries.
- Disadvantages:
  - Training time is cubic to the number of training data. This makes it hard to apply SVMs to large datasets.
  - High-dimensional kernels increase the risk of overfitting.
    - Usually larger training sets help reduce overfitting, but SVMs cannot be applied to large training sets due to cubic time complexity.
  - Some choices must still be made manually.
    - Choice of kernel function.
    - Choice of parameter  $C$  in formula  $C\left(\sum_{n=1}^N \xi_n\right) + \frac{1}{2} \|\mathbf{w}\|^2$ .

# 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.  
<http://citeseer.nj.nec.com/burges98tutorial.html>
- The VC/SRM/SVM Bible:  
Statistical Learning Theory by Vladimir Vapnik, Wiley-Interscience; 1998