

# Support Vector Machines (SVM)

# Objectives

- ① Gain an intuition about the *purpose* and *power* of SVMs.
- ② Explore (some) of the mathematics behind SVMs.
- ③ Supercharge SVMs with kernels and soft margins.
- ④ Gain an intuition about the Bias-Variance tradeoff while using SVMs.

# Support Vector Machines

A rough history

**Maximum Margin Classifier:** (morning lecture)

1963: Vapnik, Chervonenkis

or SVC

**Soft Margins and the “Kernel Trick”:** (afternoon lecture)

1992-1995: Vapnik, Boser, Guyon, Cortes

This is the modern Support Vector Machine (SVM).

loosely referred as SVM

The real SVM

# Outline

## 1 Review

- Supervised Learning
- Notation
- Hyperplanes

## 2 Motivation

- Binary Classification
- Margin
- Maximum Margin Classifier

## 3 SVMs

- Soft Margin
- Kernels
- Misc Topics

# Supervised Learning

**High level:** What is supervised learning?

Learn an unknown function from a set of **labeled** training data.

- Our training data is limited and finite. A useful algorithm must generalize well to “unseen” data.
- Example: Children learning colors.
- Support Vector Machines (SVMs) are a *supervised* learning algorithm.

# Outline

## 1 Review

- Supervised Learning
- Notation
- Hyperplanes

## 2 Motivation

- Binary Classification
- Margin
- Maximum Margin Classifier

## 3 SVMs

- Soft Margin
- Kernels
- Misc Topics

# Notation

## Goal:

Learn a model of a function  $F : X \rightarrow Y$  from a training set  $D$ .

$D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ , where

- $x^{(j)} \in X$  is often called the “input”.
- $y^{(j)} \in Y$  is often called the “label” or “target”.

## Notation (cont.)

$$F : X \rightarrow Y$$

- Often,  $X = \mathbb{R}^n$
- Often,  $Y$  is a finite set (i.e. a classification task)

We want our learned model to *generalize* well.

Generalization error is a measure of the model's performance on all possible “unseen” data.

# Outline

## 1 Review

- Supervised Learning
- Notation
- Hyperplanes

## 2 Motivation

- Binary Classification
- Margin
- Maximum Margin Classifier

## 3 SVMs

- Soft Margin
- Kernels
- Misc Topics

# Dimensions

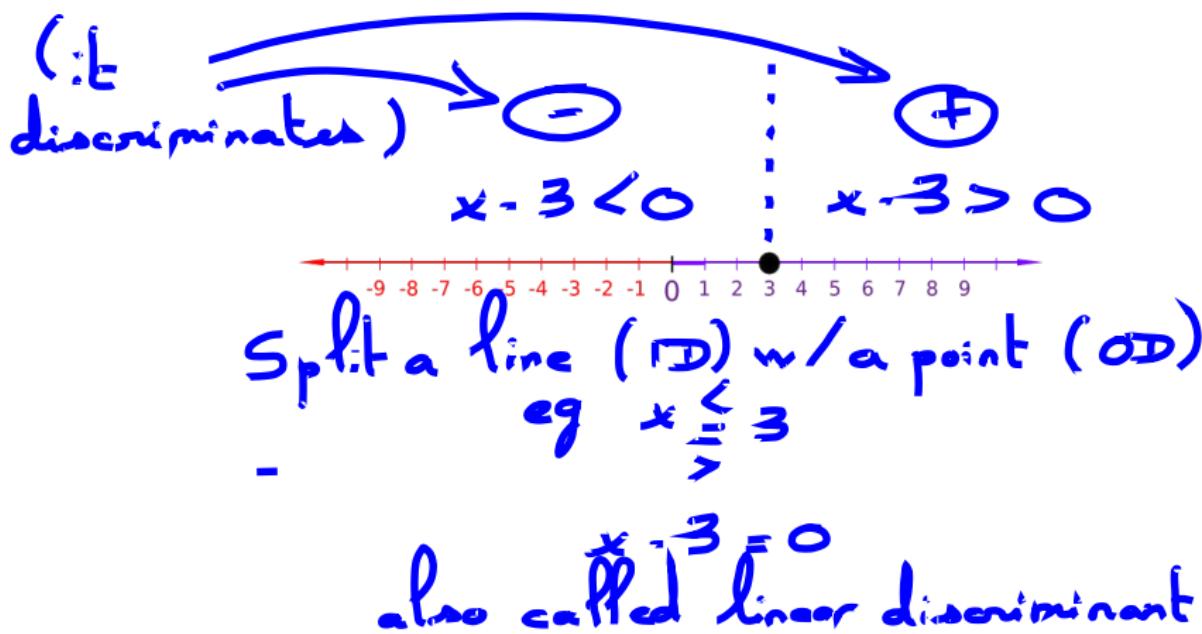
Basic stuff, I know...

						X Y Z W ↙
0	1	2	3	4	#Dim	

<sup>1</sup>By NerdBoy1392 (Own work) [CC BY-SA 3.0 (<http://creativecommons.org/licenses/by-sa/3.0>)], via Wikimedia Commons

# 1D

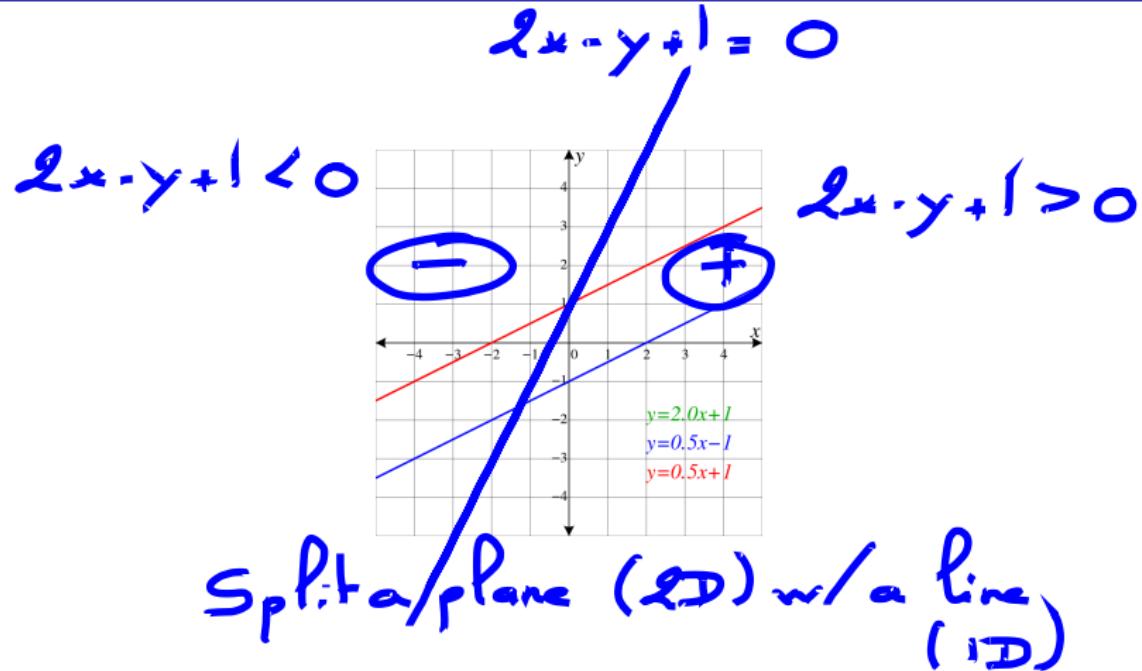
How do you split this space?



<sup>1</sup>By HakunamentaMathsIsFun at en.wikipedia [CC0], from Wikimedia Commons, Public Domain

2D

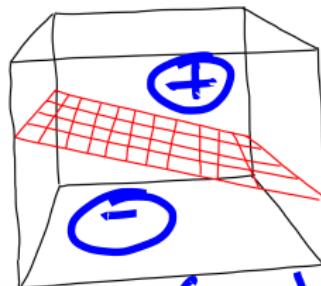
How do you split this space?



<sup>1</sup>By ElectroKid (talk • contribs). Original: HiTe. (Modification from the original work.) [CC BY-SA 1.0 (<http://creativecommons.org/licenses/by-sa/1.0/>)], via Wikimedia Commons

# 3D

How do you split this space?



Split a space (3D) w/ a plane (2D)  
Hyperplane in the form  
 $\alpha x + \beta y + \gamma z = 0$

4D, 5D, etc...

Hard to visualize... :/

In general, an  $n$ -dimensional space can be separated by an  $(n - 1)$ -dimensional *hyperplane*.

In an  $n$ -dimensional space any hyperplane can be defined by  $w \in \mathbb{R}^n$  and  $b \in \mathbb{R}$ . The hyperplane includes all  $x \in \mathbb{R}^n$  where:

$$w_0x_0 + w_1x_1 + \dots + w_{n-1}x_{n-1} - b = 0$$

usually written: *changed notation here*

*in matrix form:*

$$\boxed{w \cdot x - b = 0}$$

## How to interpret $w$ and $b$

So,  $w$  and  $b$  define a hyperplane. Is there an interpretation of  $w$  and  $b$  that can help us visualize this hyperplane?

## How to interpret $w$ and $b$

So,  $w$  and  $b$  define a hyperplane. Is there an interpretation of  $w$  and  $b$  that can help us visualize this hyperplane?

- $\frac{w}{\|w\|}$  is the hyperplane's normal vector.
- $\frac{b}{\|w\|}$  is the hyperplane's distance from the origin.

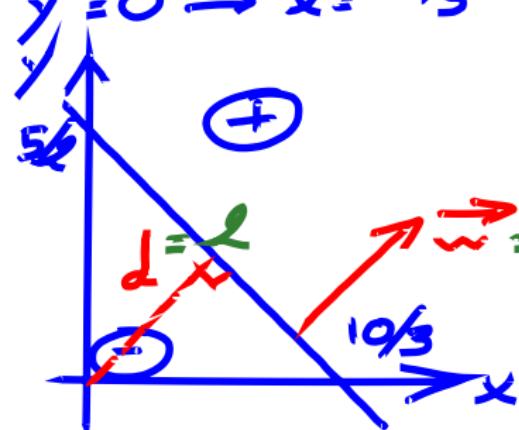
## Example in 2D

$$3x + 4y - 10 = 0 \quad \vec{w} = \begin{pmatrix} 3 \\ 4 \end{pmatrix} \quad b = 10$$

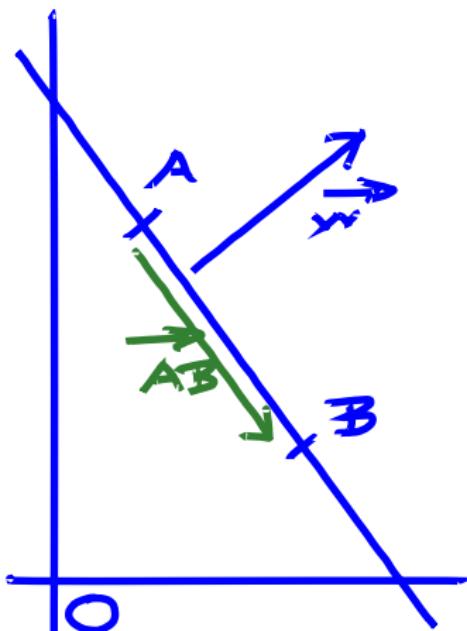
$$x = 0 \rightarrow y = \frac{5}{2}$$

$$y = 0 \rightarrow x = \frac{10}{3}$$

$$d = \frac{|b|}{\|\vec{w}\|} = \frac{10}{\sqrt{3^2 + 4^2}} = \frac{10}{\sqrt{25}} = \frac{10}{5} = 2$$



$\frac{w}{\|w\|}$  is the hyperplane's normal vector



$$\vec{w} \cdot \vec{x}_A - b = 0 \quad (1) \quad (\vec{A} \in \text{hyperplane})$$

$$\vec{w} \cdot \vec{x}_B - b = 0 \quad (2) \quad (\vec{B} \in \text{hyperplane})$$

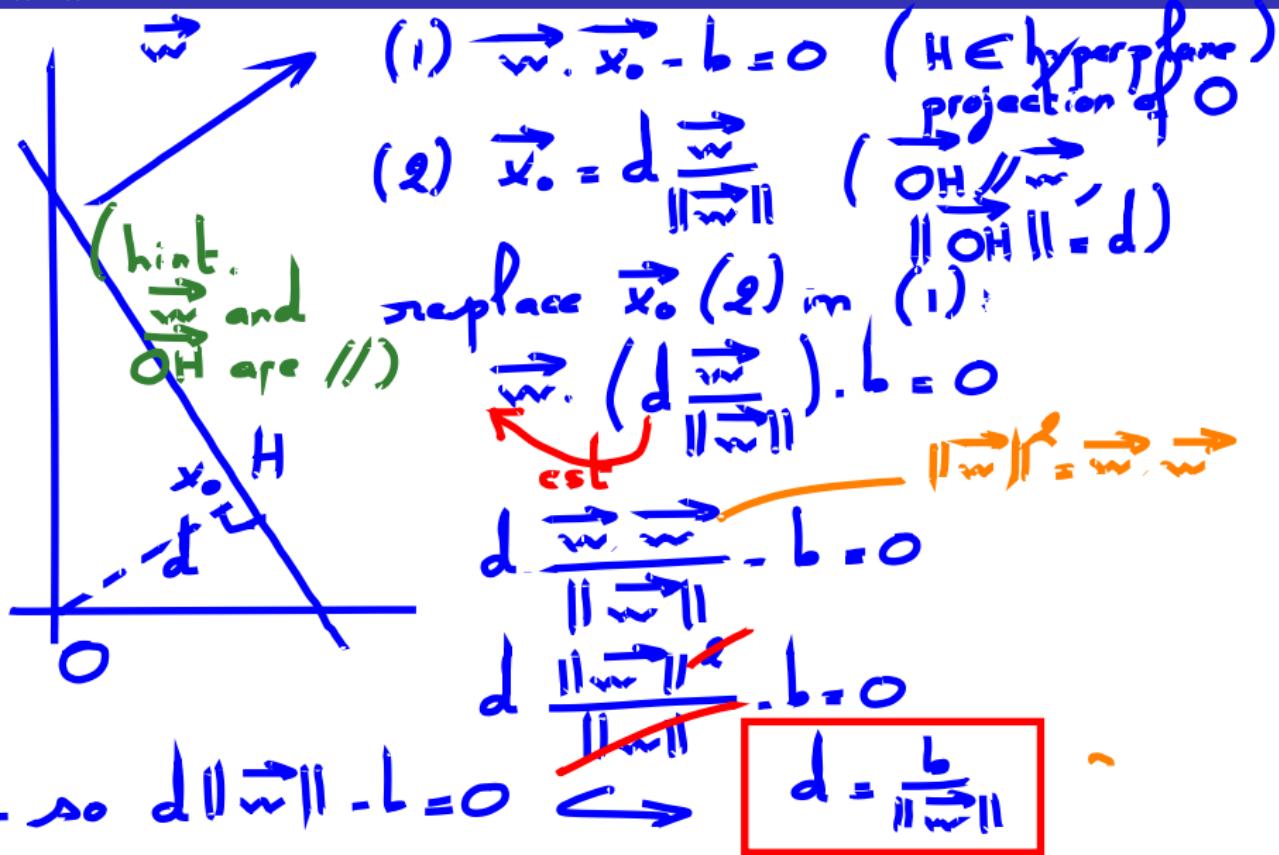
$$(2) - (1): \vec{w} \cdot \vec{OB} - b - (\vec{w} \cdot \vec{OA} + b) = 0$$

$$\vec{w} \cdot (\vec{OB} - \vec{OA}) = 0$$

$$\vec{w} \cdot \vec{AO} = 0$$
$$\vec{w} \cdot \vec{AO} = \vec{w} \cdot \vec{AB}$$
$$\vec{w} \cdot \vec{AB} = 0 \quad \forall A, B \in \text{hyperplane}$$

$\vec{w}$  is normal to the hyperplane

$\frac{b}{\|w\|}$  is the hyperplane's distance from the origin



# Outline

## 1 Review

- Supervised Learning
- Notation
- Hyperplanes

## 2 Motivation

- Binary Classification
- Margin
- Maximum Margin Classifier

## 3 SVMs

- Soft Margin
- Kernels
- Misc Topics

# Binary Classification

A supervised learning problem

Recall, we're trying to learn  $F : X \rightarrow Y$ .

- Let,  $X = \mathbb{R}^n$
- For binary classification,  $Y = \{-1, 1\}$

-1

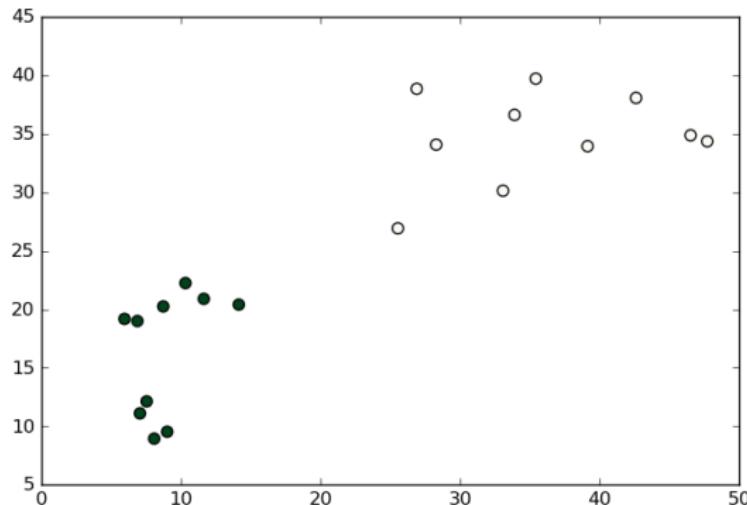
X

We'll be using -1  
instead of 0 for  
(future) convenience

Big idea: Let's have our model find a hyperplane that splits our  $n$ -dimensional data  $X$  into the set where  $y = -1$  and the set where  $y = 1$ .

# Binary Classification: Example

How many ways can we use a hyperplane to classify this dataset correctly?

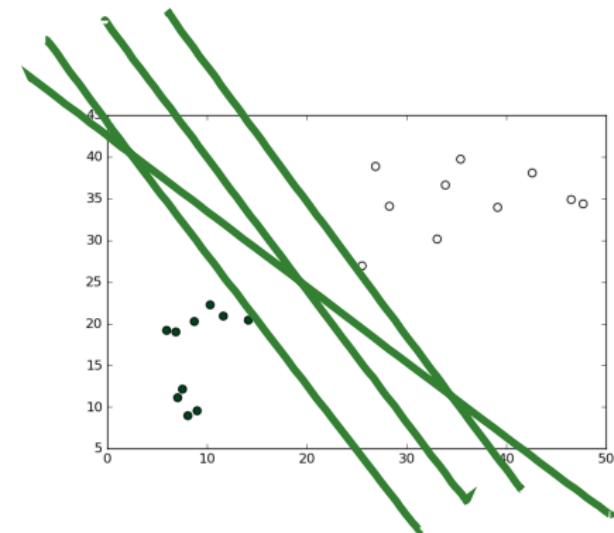


$$X = \mathbb{R}^2$$

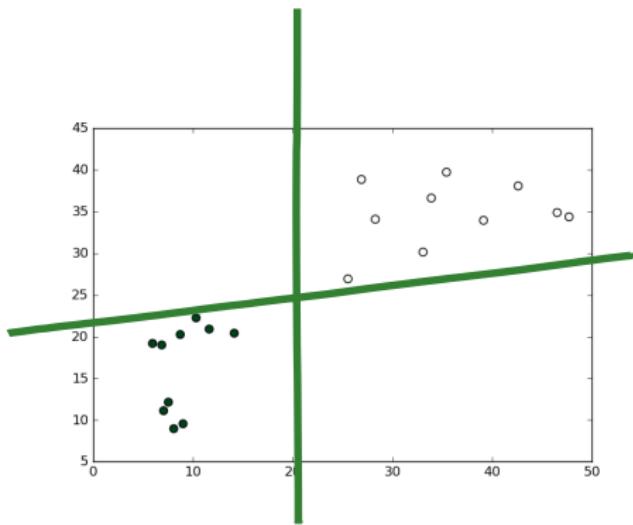
$$Y = \{-1, 1\}$$

# Binary Classification: Example

## Two Example Solutions



Which is better?



# Outline

## 1 Review

- Supervised Learning
- Notation
- Hyperplanes

## 2 Motivation

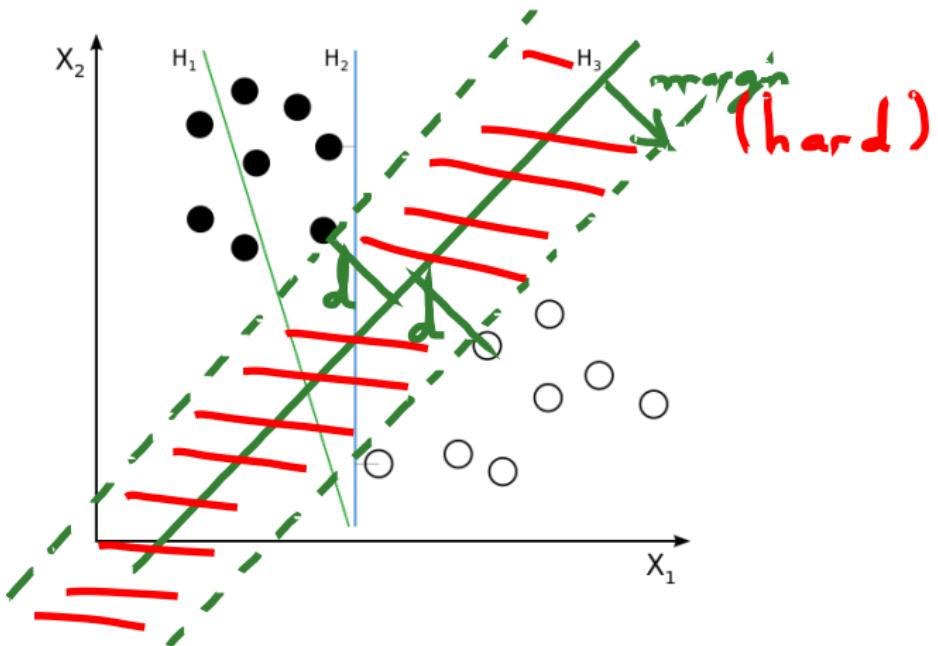
- Binary Classification
- Margin
- Maximum Margin Classifier

## 3 SVMs

- Soft Margin
- Kernels
- Misc Topics

# Defining Margin

The distance from the hyperplane to the nearest training-data point.



<sup>1</sup>By User:ZackWeinberg, based on PNG version by User:Cyc [CC BY-SA 3.0 (<http://creativecommons.org/licenses/by-sa/3.0>)], via Wikimedia Commons

# Why Maximize the Margin?

Our goal is to train a model that generalizes to “unseen” data.

Large margin means better generalization.

Intuitively, this makes sense (see previous slide)

# Outline

## 1 Review

- Supervised Learning
- Notation
- Hyperplanes

## 2 Motivation

- Binary Classification
- Margin
- Maximum Margin Classifier

## 3 SVMs

- Soft Margin
- Kernels
- Misc Topics

# Maximum Margin Classifier

## Goal

Goal: Calculate  $w$  and  $b$  of the hyperplane:

$$w \cdot x - b = 0$$

... such that the classes are split correctly and the margin is maximized.

# Maximum Margin Classifier

Goal (cont.)

First, some house cleaning: What happens to the hyperplane when we scale  $w$  and  $b$  by some factor  $c$ ?

$$\begin{aligned}\vec{w}' &= c\vec{w} \\ b' &= cb \\ \vec{w}' \cdot \vec{x} + b' &= (\vec{c}\vec{w}) \cdot \vec{x} + cb \\ &= c(\vec{w} \cdot \vec{x} + b) \\ &= c \underbrace{(\vec{w} \cdot \vec{x} + b)}_0 = 0\end{aligned}$$

$\vec{f}(\vec{w}, b)$  defines an  
hyperplane, then  
 $(c\vec{w}, cb)$  ( $c \neq 0$ )  
define the same hyperplane  
(that's a too many!)

# Maximum Margin Classifier

## Setup

*so basically just one ( $w, b$ ) pair*

We need to define a “canonical”  $w$  and  $b$ . This will help later.

Let

$$|w \cdot x^{(i)} - b| = 1$$

where  $x^{(i)}$  is the closest point to the hyperplane.

There will be a unique scaled  $w$  and  $b$  to achieve this.

# Maximum Margin Classifier

## Margin

Now, if  $x^{(i)}$  is the closest point to the hyperplane, then the distance from  $x^{(i)}$  to the hyperplane (at point  $x_0$ ) is our margin. What is that distance?

Our equations are:

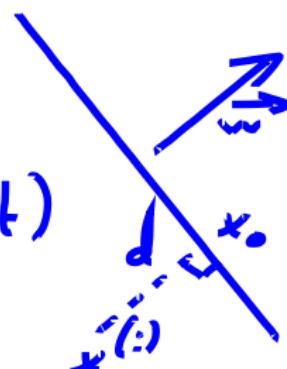
$$\|\vec{x}^{(i)} - \vec{x}_0\| = d \quad (1) \text{ (by definition)}$$

$$\vec{w} \cdot \vec{x}_0 + b = 0 \quad (2) \text{ ( } x_0 \in \text{hyperplane})$$

$$|\vec{w} \cdot \vec{x}^{(i)} + b| = 1 \quad (3) \text{ (canonical constraint)}$$

$$\|\vec{w}\| \cdot (\vec{x}^{(i)} - \vec{x}_0) = ? \text{ ways}$$

One way:  $\frac{\|\vec{w}\| \cdot (\vec{x}^{(i)} - \vec{x}_0)}{\|\vec{w}\|} = \frac{1}{\|\vec{w}\|} \|\vec{x}^{(i)} - \vec{x}_0\|$



# Maximum Margin Classifier

Margin (cont.)

Distance:

$$\|\vec{w}\| \cdot \|\vec{x} - \vec{x}_0\| =$$

both vectors are  $\parallel$

$$\|\vec{w}\| \cdot \|\vec{x} - \vec{x}_0\| = d \quad (\text{from (1)})$$

$$\frac{1}{\|\vec{w}\|} \cdot \|\vec{w}\| = d$$

Putting both ways  
together:

$$d = \frac{1}{\|\vec{w}\|}$$

$$= \frac{1}{\|\vec{w}\|} \cdot \|\vec{w}\| \cdot \left( \frac{\vec{w} \cdot \vec{x}_0 + b}{\|\vec{w}\|} \right) = \left( \frac{\vec{w} \cdot \vec{x}_0 + b}{\|\vec{w}\|} \right)$$

$$\begin{aligned} & \text{if } \vec{a} \text{ and } \vec{b} \text{ are } \parallel: \\ & \exists (\alpha, \beta, \vec{u}) \quad \vec{a} = \alpha \vec{u} \\ & \vec{b} = \beta \vec{u} \\ & \|\vec{a} \cdot \vec{b}\| = \left\| \left( \alpha \vec{u} \right) \left( \beta \vec{u} \right) \right\| \\ & = \alpha \beta \|\vec{u}\|^2 = \alpha \beta \\ & \|\vec{a}\| \cdot \|\vec{b}\| = \frac{1}{\|\vec{u}\|} \\ & \|\alpha \vec{u}\| \cdot \|\beta \vec{u}\| = \alpha \beta \|\vec{u}\|^2 \\ & \boxed{\|\vec{a} \cdot \vec{b}\| = \|\vec{a}\| \cdot \|\vec{b}\|} \end{aligned}$$

# Maximum Margin Classifier

First Attempt

$$\begin{aligned} & \text{Maximize}_{\mathbf{w}} \frac{1}{\|\mathbf{w}\|} \quad \text{subject to:} \\ & \text{minimize}_{\mathbf{w}} \|\mathbf{w}\|^2 \\ & |w \cdot x^{(i)} - b| \geq 1, \quad \text{for all } x^{(i)} \in D \end{aligned}$$

*replace w/  $y^{(i)}$*

*minimize  $\|\mathbf{w}\|^2$*

*don't like absolute values*

*replace w/  $y^{(i)}$*

*$y^{(i)} - (w \cdot x^{(i)} - b) \geq 1$*

... but we don't know how to solve this optimization problem. Let's reformulate.

# Maximum Margin Classifier

Reformulated

$$\text{Minimize } ||w||^2$$

details of optimization  
outside of class ...

subject to:

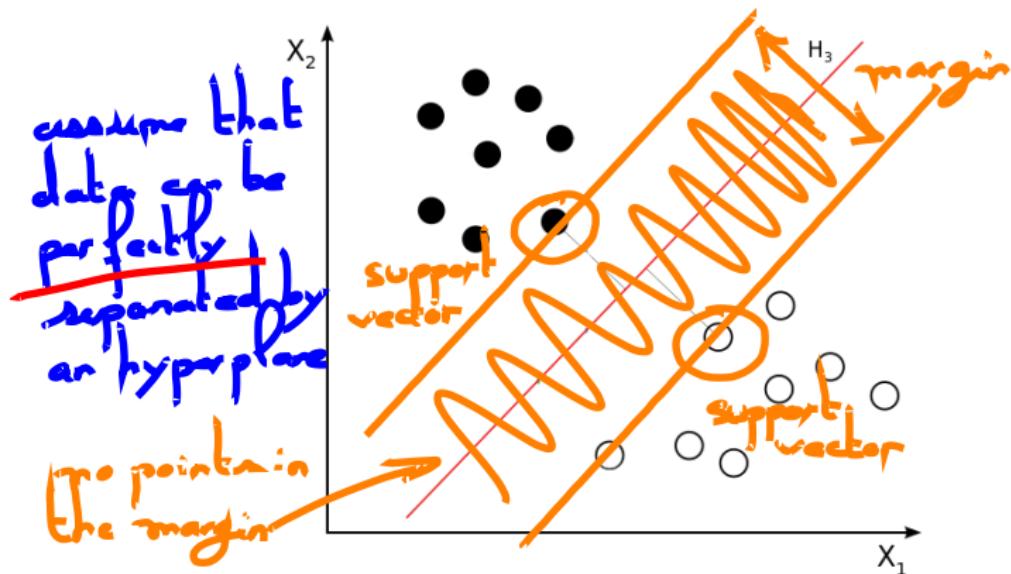
$$y^{(i)}(w \cdot x^{(i)} - b) \geq 1,$$

for all  $(y^{(i)}, x^{(i)}) \in D$

... plus more steps... and we eventually get a quadratic programming formulation.

# Support Vectors

The maximum margin hyperplane is defined only by the points that touch the margin. These are called the “support vectors”.



# sklearn's interface

## LogisticRegression vs SVC

### **LogisticRegression:**

▶ Link

### **SVC:**

▶ Link

(end of morning lecture)

# Outline

## 1 Review

- Supervised Learning
- Notation
- Hyperplanes

## 2 Motivation

- Binary Classification
- Margin
- Maximum Margin Classifier

## 3 SVMs

- Soft Margin
- Kernels
- Misc Topics

# Soft Margin Motivation

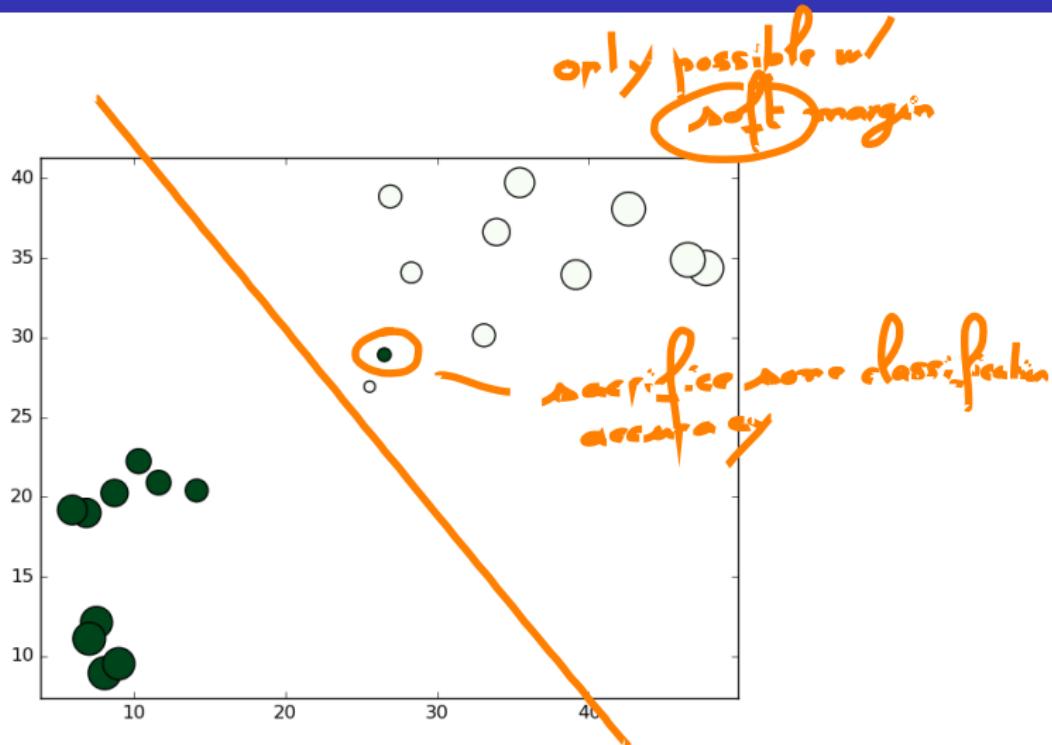
What if:

- ① Your data isn't linearly separable?
- ② Your data is noisy/has outliers?

Soft Margins address these problems.

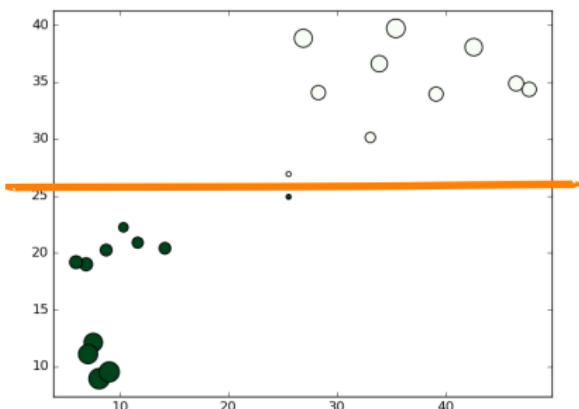
# Soft Margins

## Inseparable Data

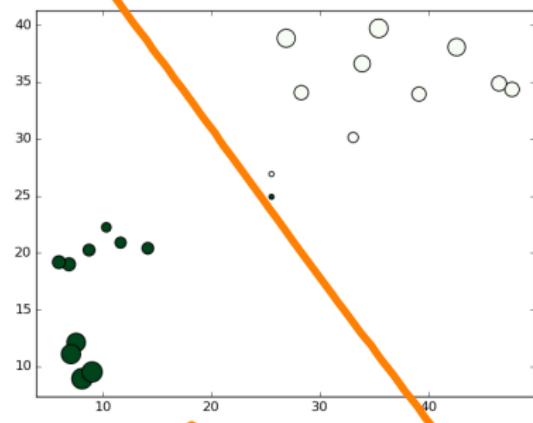


# Soft Margins

## Outliers in Data



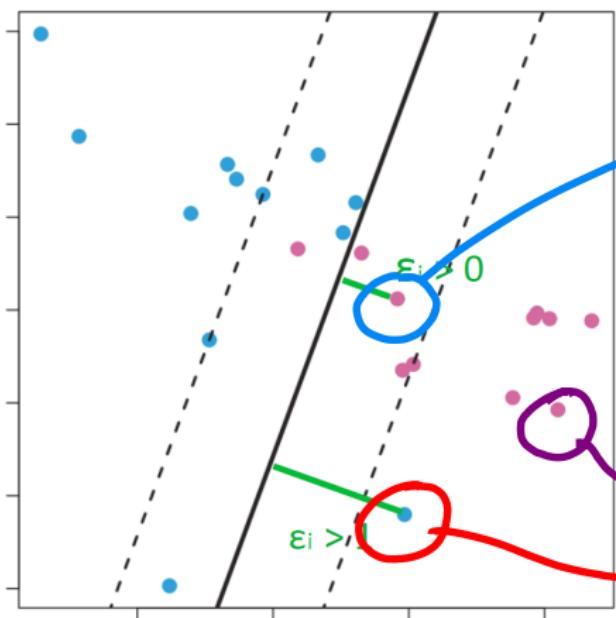
hard margin



soft margin

# Soft Margins

We need some sort of slack/budget



right side of hyperplane but violate margin  $\epsilon_i > 0$

$\epsilon_i = 0$  for being on correct side of margin  
 $\epsilon_i > 0$  for violating the margin

$\epsilon_i > 1$  for being on wrong side of hyperplane

right side of margin  $\epsilon_i = 0$   
wrong side of hyperplane  $\epsilon_i > 1$

# Soft Margin

The  $C$  hyperparameter

Support Vector Classifier  
SVC

An extension to Maximum Margin Classifiers adds a hyperparameter  $C$  to control the misclassification error penalty/margin tradeoff.

SVC = Maximum Margin Classifier  
+ soft-margin

# Support Vector Classifier

## The C hyperparameter

Minimize  $\|w\|^2$

subject to:

- $y^{(i)}(w \cdot x^{(i)} - b) \geq 1 - \varepsilon^{(i)}$

slack from each point



- $\varepsilon^{(i)} \geq 0$  ✓

- $\sum_i \varepsilon^{(i)} \leq 1/C$

Budget we can type  
ok-linear convention

# Support Vector Classifier

## The $C$ hyperparameter

Support Vector Classifier (SVC) extends Maximum Margin Classifiers by adding a hyperparameter  $C$  to control the misclassification error penalty/margin tradeoff.

**Large  $C$ :** Harder margins: value classification accuracy over a large margin

*small  $\gamma_C \rightarrow$  low margin*

**Small  $C$ :** Softer margins: value a large margin over classification accuracy

*large  $\gamma_C \rightarrow$  large margin*

# Soft Margins

scikit-learn code

```
from sklearn.svm import SVC  
...  
svc = SVC(C=1.0, kernel='linear')  
svc.fit(x, y)
```

SVC supports the  $C$  parameter as the soft-margin hyperparameter.

# Outline

## 1 Review

- Supervised Learning
- Notation
- Hyperplanes

## 2 Motivation

- Binary Classification
- Margin
- Maximum Margin Classifier

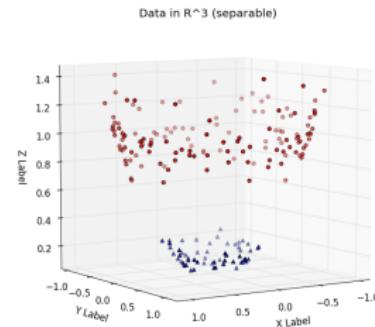
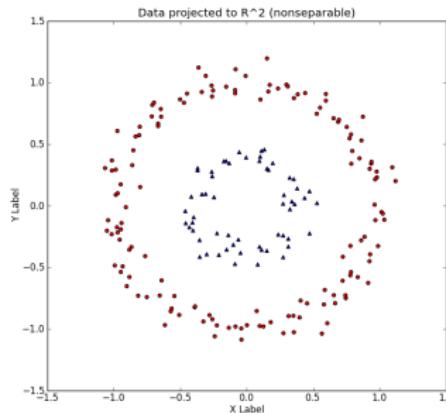
## 3 SVMs

- Soft Margin
- Kernels
- Misc Topics

# The “Kernel Trick”

The idea...

Idea: If data is inseparable in its input space, maybe it will be separable in a higher-dimensional space.



# The “Kernel Trick”

Back to some math...

We have not discussed exactly how the SVC is computed but it turns out that the optimization problem to maximize the margin involves only the *dot products* of the observations (as opposed to the observations themselves):

$$\mathcal{L}(\alpha) = \sum_{i=1}^m \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \alpha^{(i)} \alpha^{(j)} (\mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)})$$

*Dot / inner product*

# The “Kernel Trick”

Creating a kernel...

What if we never took the dot product, but we instead replaced it with a generalization of the dot product of the form (a “kernel function”)?

$$x^{(i)} \cdot x^{(j)} \rightarrow K(x^{(i)}, x^{(j)})$$

# The “Kernel Trick”

Why is this so cool?

SVM = SVC + “kernels”

$$\mathcal{L}(\alpha) = \sum_{i=1}^m \alpha^{(i)} - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \alpha^{(i)} \alpha^{(j)} (\mathcal{K}(x^{(i)}, x^{(j)}))$$

- Saves some computation.
- Opens new possibilities. A kernel can operate in infinite dimensions!

# The Linear Kernel

$$K(x^{(i)}, x^{(j)}) = x^{(i)} \cdot x^{(j)}$$

(Basically SVC)

## The Polynomial Kernel

expand feature space by adding new features

$$K(x^{(i)}, x^{(j)}) = (1 + x^{(i)} \cdot x^{(j)})^d$$

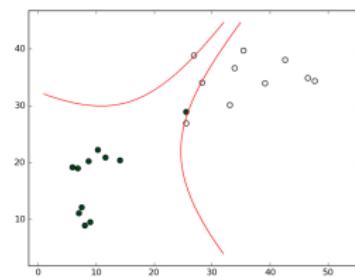
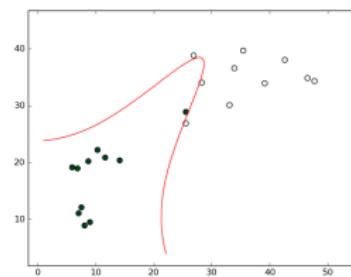
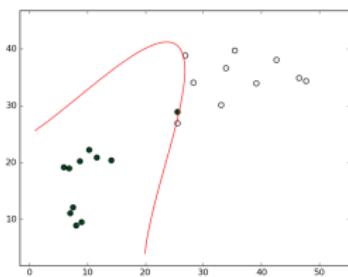
kernel is a  
powerful way  
to accommodate  
non-linear  
class  
boundaries

- Equivalent to the dot product in the  $d$ -order  $\phi$  space
- Requires an extra hyper-parameter,  $d$ , for “degree”

$$K(x^{(i)}, x^{(j)}) = 1 + 2x^{(i)} \cdot x^{(j)} + (x^{(i)})^2 (x^{(j)})^2$$

# The Polynomial Kernel

## Example



```
svc = SVC(C=10000.0,  
kernel='poly',  
degree=3)
```

```
svc = SVC(C=10000.0,  
kernel='poly',  
degree=5)
```

```
svc = SVC(C=10000.0,  
kernel='poly',  
degree=10)
```

# The RBF Kernel ("Gaussian")

(Radial Basis Function)

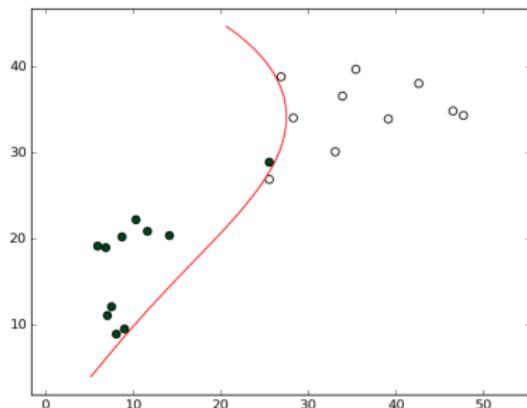
$$K(x^{(i)}, x^{(j)}) = \exp(-\gamma \|x^{(i)} - x^{(j)}\|^2)$$

geometric  
euclidean  
space

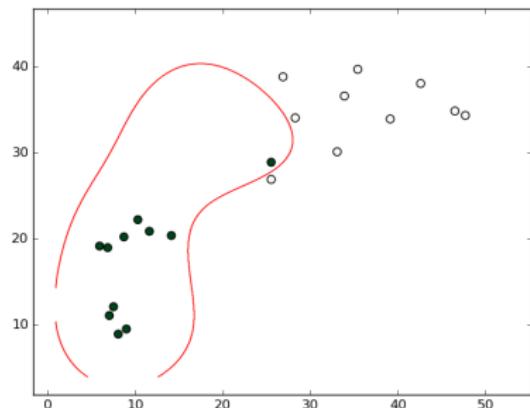
- Equivalent to the dot product in the Hilbert space of infinite dimensions
- Requires an extra hyper-parameter,  $\gamma$ , "gamma"

# The RBF Kernel

## Examples



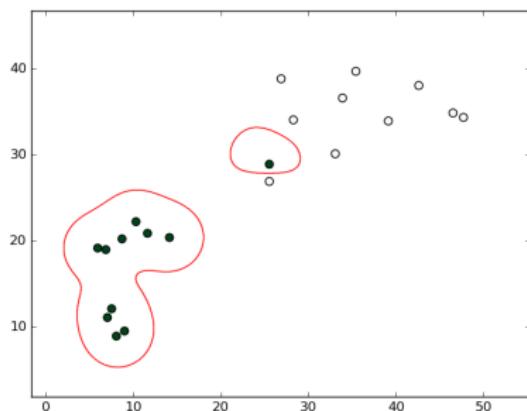
```
svc = SVC(C=10000.0, kernel='rbf',  
gamma=0.001)
```



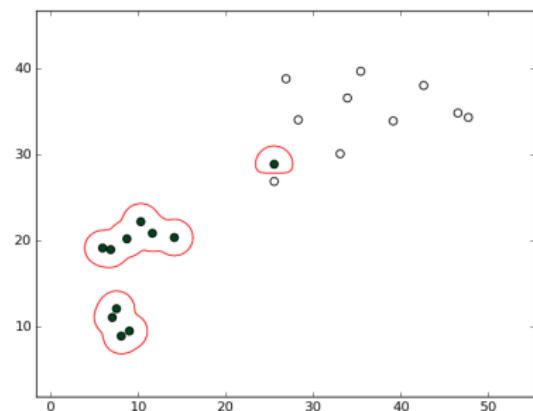
```
svc = SVC(C=10000.0, kernel='rbf',  
gamma=0.01)
```

# The RBF Kernel

## More Examples



```
svc = SVC(C=10000.0, kernel='rbf',  
gamma=0.1)
```



```
svc = SVC(C=10000.0, kernel='rbf',  
gamma=1.0)
```

# Revisiting the Bias-Variance Tradeoff with SVM...

## High-Bias Models

Makes many assumptions and prefers to solve problems a certain way.

E.g. A linear SVM looks for dividing hyperplanes in the input space *only*.

For complex data, high-bias models often *underfit* the data.

## High-Variance Models

Makes fewer assumptions and has more representational power.

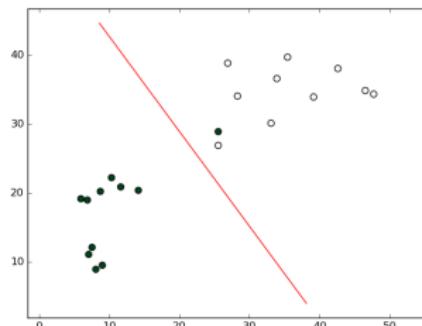
E.g. An RBF SVM looks for dividing hyperplanes in an infinite-dimensional space.

For simple data, high-variance models often *overfit* the data.

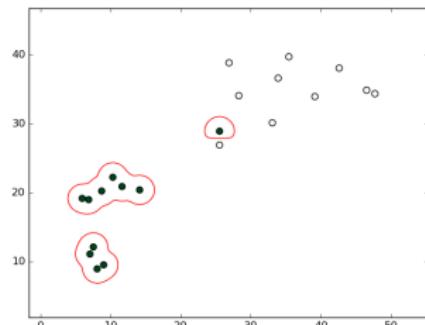
# Bias-Variance Tradeoff

Example

low variance /  
higher bias



high bias /  
low variance



```
svc = SVC(C=0.01, kernel="linear")
```

```
svc = SVC(C=10000.0, kernel='rbf',  
          gamma=1.0)
```

# Outline

## 1 Review

- Supervised Learning
- Notation
- Hyperplanes

## 2 Motivation

- Binary Classification
- Margin
- Maximum Margin Classifier

## 3 SVMs

- Soft Margin
- Kernels
- Misc Topics

# SVMs vs Logistic Regression

(some rules of thumb)

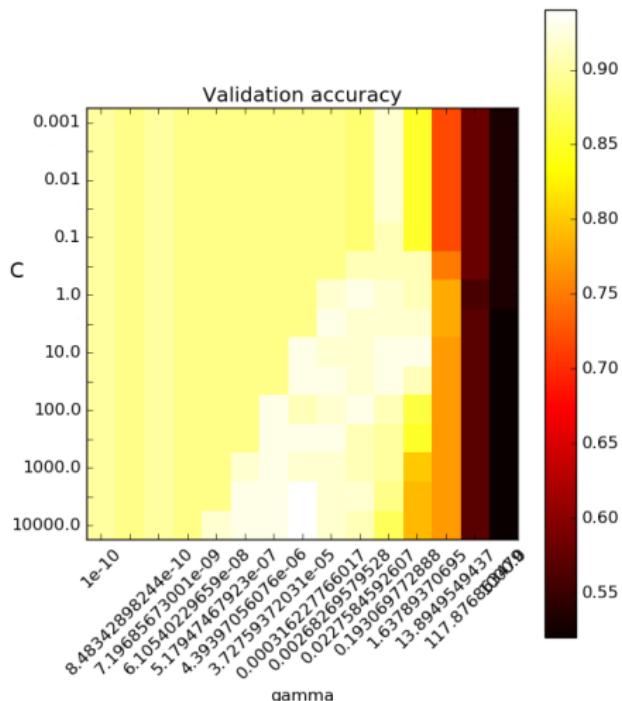
- ① Logistic Regression maximizes the *Binomial Log Likelihood* function. SVMs maximize the *margin*.
- ② When classes are nearly separable, SVMs tends to do better than Logistic Regression. Otherwise, Logistic Regression (with Ridge) and SVMs are similar.
- ③ If you want to estimate probabilities however, Logistic Regression is the better choice.
- ④ With kernels, SVMs work well. Logistic Regression works fine with kernels but can get computationally too expensive.

# Grid Search

## Hyperparameter Tuning

Let's find  $C$  and  $\gamma$  by searching through values we expect might work well.

Use cross-validation accuracy to determine which values are best.



# Grid Search

code

```
svc_rbf = SVC(kernel='rbf')

param_space = {'C': np.logspace(-3, 4, 15),
               'gamma': np.logspace(-10, 3, 15)}

grid_search = GridSearchCV(svc_rbf, param_space,
                           scoring='accuracy', cv=10)
grid_search.fit(x, y)

print grid_search.grid_scores_
print grid_search.best_params_
print grid_search.best_score_
print grid_search.best_estimator_
```