

# Do not distribute course material

You may not and may not allow others to reproduce or distribute lecture notes and course materials publicly whether or not a fee is charged.

# Lecture Support Vector Machines continued

---

PROF. LINDA SELLIE

SOME SLIDES FROM PROF. RANGAN

You are not responsible for the  
material covered in class on March  
29th

# Outline

---

- ❑ Notation change, intuition, and finding how to compare hyperplanes - mathematically how do compare hyperplanes to find the one with the maximum margin. Can we turn this way of comparing hyperplanes into an objective function
- ❑ Support vector machines
  - ★ hard margin - find the objective function when the data is linearly separable
  - ★ Dealing with non-linear data - “Soft” margins for SVM - New objective function for the case where the data is not linearly separable
  - ★ dual formula - a clever trick  $g(\mathbf{x}) = \text{sign} \left( w_0 + \sum_{i \in I} \alpha^{(i)} y^{(i)} \underline{\phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x})} \right) = \text{sign} \left( w_0 + \sum_{i \in I} \alpha^{(i)} y^{(i)} \underline{K(\mathbf{x}^{(i)}, \mathbf{x})} = \Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x}) \right)$
  - ★ Dealing with non-linear data - feature transformation with the kernel trick - Show two popular feature maps
- ❑ Multiclass

# SVMs

- Maximizes distance of training data to boundary (built in regularization)
- Generalizes to nonlinear decision boundaries (I.e. feature transformation)
- Performs well with high dimensional data

## Duel formula motivation: Polynomial Kernel

$$g(\mathbf{x}) = \text{sign} \left( \sum_{i \in I} \alpha^{(i)} y^{(i)} (\mathbf{x}^{(i)T} \mathbf{x}) + w_0 \right)$$

$$\mathbf{x} = [x_1, x_2]^T$$

$$\Phi_2 : R^2 \rightarrow R^6$$

$$\Phi_2(\mathbf{x}) = [1, x_1, x_2, x_1 x_2, x_1^2, x_2^2]^T$$

$$\Phi_2(\mathbf{x})^T \Phi_2(\mathbf{x}') = [1, x_1, x_2, x_1 x_2, x_1^2, x_2^2] \cdot [1, x'_1, x'_2, x'_1 x'_2, x'_1^2, x'_2^2]$$

$$\Phi_2(\mathbf{x})^T \Phi_2(\mathbf{x}') = 1 + x_1 x'_1 + x_2 x'_2 + x_1 x_2 x'_1 x'_2 + x_1^2 x'_1^2 + x_2^2 x'_2^2$$

$$\Phi_2(-0.75, -0.25) = (1, -0.75, -0.25, (-0.75) \cdot (-0.25), (-0.75)^2, (-0.25)^2)$$

$$\Phi_2(-1, 1) = (1, -1, 1, -1 \cdot 1, (-1)^2, 1^2)$$

$$\Phi_2(-0.75, -0.25) \cdot \Phi_2(-1, 1) = 1 + 0.75 - 0.25 + 0.75 \cdot (-0.25) + (-0.75)^2 + (-0.25)^2$$

But if we use a different (but similar) transformation:

$$\Phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2, x_1^2, x_2^2]^T$$

$$\begin{aligned} K(\mathbf{x}, \mathbf{x}') &= \Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}') = 1 + 2x_1 x'_1 + 2x_2 x'_2 + 2x_1 x_2 x'_1 x'_2 + x_1^2 x'_1^2 + x_2^2 x'_2^2 \\ &= (1 + \mathbf{x} \cdot \mathbf{x}')^2 \end{aligned}$$

$$\Phi(-0.75, -0.25) = (1, \sqrt{2} \cdot (-0.75), \sqrt{2} \cdot (-0.25), \sqrt{2}(-0.75) \cdot (-0.25), (-0.75)^2, (-0.25)^2)$$

$$\Phi(-1, 1) = (1, -\sqrt{2}, \sqrt{2}, -\sqrt{2}, (-1)^2, 1^2)$$

$$K((-1, 1), (-0.75, -0.25)) = \frac{2.25}{2.25}$$

$$= 1 + 2(0.75) + 2(-0.25) - 2(0.75)(0.25) + (-0.75)^2 + (-0.25)^2$$

$$= (1 + (0.75 - 0.25))^2$$

Class exercise

So, instead of computing  $\Phi(\mathbf{x})^T \Phi(\mathbf{x}')$ ,

we compute  $\mathbf{x} \cdot \mathbf{x}'$  in the original feature space, and then compute  $(1 + \mathbf{x} \cdot \mathbf{x}')^2$

The polynomial kernel can be generalized to higher dimensions. A degree k polynomial is  $K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x} \cdot \mathbf{x}')^k$

Learn more at: [https://en.wikipedia.org/wiki/Polynomial\\_kernel](https://en.wikipedia.org/wiki/Polynomial_kernel)

## Duel formula motivation: Polynomial Kernel

$$\mathbf{x} = [x_1, x_2]^T$$

$$\Phi_2 : R^2 \rightarrow R^6$$

$$\Phi_2(\mathbf{x}) = [1, x_1, x_2, x_1x_2, x_1^2, x_2^2]^T$$

$$\Phi_2(\mathbf{x})^T \Phi_2(\mathbf{x}') = [1, x_1, x_2, x_1x_2, x_1^2, x_2^2] \cdot [1, x'_1, x'_2, x'_1x'_2, x'^2_1, x'^2_2]$$

$$\Phi_2(\mathbf{x})^T \Phi_2(\mathbf{x}') = 1 + x_1x'_1 + x_2x'_2 + x_1x_2x'_1x'_2 + x_1^2x'^2_1 + x_2^2x'^2_2$$

$$g(\mathbf{x}) = \text{sign}\left( \sum_{i \in I} \alpha^{(i)} y^{(i)} (\mathbf{x}^{(i)T} \mathbf{x}) + w_0 \right)$$

If d=200 then in z-space we have a feature vector of size approx 20,000

$$\Phi_2(-0.75, -0.25) = (1, -0.75, -0.25, (-0.75) \cdot (-0.25), (-0.75)^2, (-0.25)^2)$$

$$\Phi_2(-1, 1) = (1, -1, 1, -1 \cdot 1, (-1)^2, 1^2)$$

$$\Phi_2(-0.75, -0.25) \cdot \Phi_2(-1, 1) = 1 + 0.75 - 0.25 + 0.75 \cdot (-0.25) + (-0.75)^2 + (-0.25)^2$$

But if we use a different (but similar) transformation:

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

$$\begin{aligned} K(\mathbf{x}, \mathbf{x}') &= \Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}') = 1 + 2x_1x'_1 + 2x_2x'_2 + 2x_1x_2x'_1x'_2 + x_1^2x'^2_1 + x_2^2x'^2_2 \\ &= (1 + \mathbf{x} \cdot \mathbf{x}')^2 \end{aligned}$$

So, instead of computing  $\Phi(\mathbf{x})^T \Phi(\mathbf{x}')$ ,

we compute  $\mathbf{x} \cdot \mathbf{x}'$  in the original feature space, and then compute  $(1 + \mathbf{x} \cdot \mathbf{x}')^2$

The polynomial kernel can be generalized to higher dimensions. A degree k polynomial is  $K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x} \cdot \mathbf{x}')^k$

If d=200 then in z-space we have a feature vector of size approx 20,000, but we do roughly the same amount of work to compute  $K(\mathbf{x}, \mathbf{x}')$  as if we hadn't transformed date features

$$\Phi(-0.75, -0.25) = (1, \sqrt{-0.75}, \sqrt{-0.25}, -0.75 \cdot -0.25, -0.75^2, -0.25^2)$$

$$\Phi(-1, 1) = (1, -\sqrt{2}, \sqrt{2}, -1 \cdot 1, (-1)^2, 1^2)$$

$$K((-1, 1), (-0.75, -0.25)) =$$

Class exercise

$$= 1 + 2(0.75) + 2(-0.25) - 2(0.75)(0.25) + (-0.75)^2 + (-0.25)^2$$

$$= (1 + (0.75 - 0.25))^2$$

Learn more at: [https://en.wikipedia.org/wiki/Polynomial\\_kernel](https://en.wikipedia.org/wiki/Polynomial_kernel)

# Duality

Transforming our constrained objective function into a new objective function. (Thus we don't solve the original objective function)

- New objective function scales well with high-dimensional data
- Our hypothesis will be a function of the training examples (thus we can use the kernel function)

For SVM's the original constrained objective function (primal problem) and the dual formulation are equivalent problems. (This is not true in general when transforming a problem into its dual.)

# Learning in a transformed feature space:

Pair share: if we use  $\Phi(\mathbf{x})$ , do we need to actually perform the transformation?

## Primal

$$\min_{\mathbf{w}, w_0} \|\mathbf{w}\|_2^2$$

subject to:  $y^{(i)}(\mathbf{w}^T \Phi(\mathbf{x}^{(i)}) + w_0) \geq 1$  for all  $i = 1 \dots N$

## Dual

$$\min_{\alpha} \sum_{i,j} \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) - \sum_i \alpha^{(i)}$$

$$\sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0 \quad \alpha^{(i)} \geq 0$$

We only need to know the result of  $\Phi(\mathbf{x}^{(j)})^T \Phi(\mathbf{x}^{(i)})$

## Prediction

$$g(\mathbf{x}) = \text{sign}(\mathbf{w}^T \Phi(\mathbf{x}) + w_0)$$

$$g(\mathbf{x}) = \text{sign} \left( \sum_{i \in I} \alpha^{(i)} y^{(i)} K(\mathbf{x}^{(i)}, \mathbf{x}) + w_0 \right)$$

Support Vector  $\mathbf{x}^{(i)}$

We only need to know the result of  $\Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x})$

By computing the value of  $K$  without writing the transformation, we can have a computational advantage

# Learning in a transformed feature space:

Pair share: if we use  $\Phi(\mathbf{x})$ , do we need to actually perform the transformation?

## Primal

$$\min_{\mathbf{w}, w_0} \|\mathbf{w}\|_2^2$$

subject to:  $y^{(i)}(\mathbf{w}^T \Phi(\mathbf{x}^{(i)}) + w_0) \geq 1$  for all  $i = 1 \dots N$

## Dual

$$\min_{\alpha} \sum_{i,j} \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$$

$$\sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0$$

$$\alpha^{(i)} \geq 0$$

Kernel trick: evaluate  $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$  without competing  $\Phi(\mathbf{x}^{(j)})$  or  $\Phi(\mathbf{x}^{(i)})$

We only need to know the result of  $\Phi(\mathbf{x}^{(j)})^T \Phi(\mathbf{x}^{(i)})$

## Prediction

$$g(\mathbf{x}) = \text{sign}(\mathbf{w}^T \Phi(\mathbf{x}) + w_0)$$

$$g(\mathbf{x}) = \text{sign} \left( \sum_{i \in I} \alpha^{(i)} y^{(i)} K(\mathbf{x}^{(i)}, \mathbf{x}) + w_0 \right)$$

Support Vector  $\mathbf{x}^{(i)}$

We only need to know the result of  $\Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x})$

By computing the value of  $K$  without writing the transformation, we can have a computational advantage

# Lagrangian (hard margin case)

Primal SVM

- $$\min_{\mathbf{w}, w_0, \xi} \frac{1}{2} \|\mathbf{w}\|_2^2$$

Subject to  $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1$  for all  $i = 1 \dots N$

We will transform the hard margin case, the soft margin case is similar.  
We will assume the data is linearly separable

Lagrangian

- $$L(\mathbf{w}, w_0, \alpha^{(i)}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^N \alpha^{(i)} (1 - y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + w_0))$$

where  $\alpha^{(i)} \geq 0$

Primal variables:  $\mathbf{w}, w_0$   
Dual variables:  $\alpha^{(i)}, i \in \{1, \dots, N\}$   
(Lagrangian multipliers)

Dual formulation of SVM

$$\max_{\alpha^{(i)}, i \in \{1, \dots, N\}} \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i,j} y^{(i)} y^{(j)} \alpha^{(i)} \alpha^{(j)} \mathbf{x}^{(i)} T \mathbf{x}^{(j)}$$

subject to  $\sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0$  and  $\alpha^{(i)} \geq 0$

$$\mathbf{w} = \sum_{i=1}^N \alpha^{(i)} y^{(i)} \mathbf{x}^{(i)}$$

Find  $\alpha^{(j)} \neq 0$  solve  
 $w_0 = y^{(j)} - \mathbf{w}^T \mathbf{x}^{(j)}$

Steps:

- \* Minimize the Lagrangian function over the primal variables  $\mathbf{w}, w_0$  by setting  $\nabla_{\mathbf{w}} L = 0, \frac{\partial L}{\partial w_0} = 0$ .
- \* Solve for primal variables wrt dual variables.
- \* Remove primal variables in the Lagrangian by substituting their equivalent dual variables

# Lagrangian (hard margin case)

Primal SVM

- $$\min_{\mathbf{w}, w_0, \xi} \frac{1}{2} \|\mathbf{w}\|_2^2$$

Subject to  $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + w_0) \geq 1$  for all  $i = 1 \dots N$

We will transform the hard margin case, the soft margin case is similar.  
We will assume the data is linearly separable

Lagrangian

- $$L(\mathbf{w}, w_0, \alpha^{(i)}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^N \alpha^{(i)} (1 - y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + w_0))$$

where  $\alpha^{(i)} \geq 0$

Primal variables:  $\mathbf{w}, w_0$   
Dual variables:  $\alpha^{(i)}, i \in \{1, \dots, N\}$   
(Lagrangian multipliers)

Steps:

\* Minimize the Lagrangian function

variables  $\mathbf{w}, w_0$  by  
 $\frac{\partial L}{\partial w_0} = 0$ .

variables wrt dual

variables in the  
substituting their  
variables

Dual formulation of SVM

$$\max_{\alpha^{(i)}, i \in \{1, \dots, N\}} \sum_{i=1}^N \alpha^{(i)} - \frac{1}{2} \sum_{i,j} y^{(i)} y^{(j)} \alpha^{(i)} \alpha^{(j)} \mathbf{x}^{(i)} T \mathbf{x}^{(j)}$$

subject to  $\sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0$  and  $\alpha^{(i)} \geq 0$

Three key ideas from transformation you need to know (you don't need to know how the transformation works):

- $\alpha^{(i)}(1 - y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + w_0)) = 0$
- $\mathbf{w} = \sum_{i=1}^N \alpha^{(i)} y^{(i)} \mathbf{x}^{(i)}$
- Find  $\alpha^{(j)} \neq 0$  solve  $w_0 = y^{(j)} - \mathbf{w}^T \mathbf{x}^{(j)}$

# Recovering $\mathbf{w}, w_0$

Dual formulation returns  $\alpha = [\alpha^{(1)}, \alpha^{(2)}, \dots, \alpha^{(N)}]^T$

Recovering  $\mathbf{w}$ :

$$\mathbf{w} = \sum_{i=1}^N \alpha^{(i)} y^{(i)} \mathbf{x}^{(i)} = \sum_{i \in I} \alpha^{(i)} y^{(i)} \mathbf{x}^{(i)}$$

Only depends on support vectors! ( $\alpha^{(i)} > 0$ )

$$I = \{i \mid \alpha^{(i)} \neq 0\}$$

If  $\alpha^{(i)} \neq 0$  then  $\mathbf{x}^{(i)}$  is a support vector,  
The number of nonzero  $\alpha^{(i)} \neq 0$  could be less than d

Recovering  $w_0$ : There exists some  $\alpha^{(j)} \neq 0$

$$\alpha^{(j)} \left( \underbrace{y^{(j)}(w_0 + \mathbf{w}^T \mathbf{x}^{(j)}) - 1}_0 \right) = 0$$

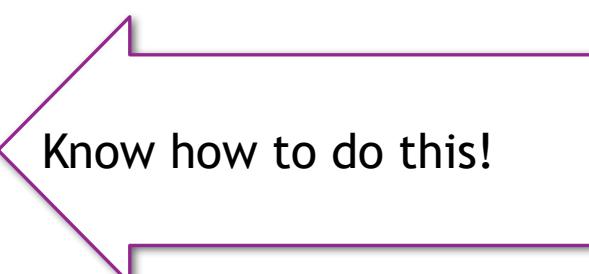
Solve for  $w_0$ :  $y^{(j)}(w_0 + \mathbf{w}^T \mathbf{x}^{(j)}) = 1$

$$w_0 = y^{(j)} - \mathbf{w}^T \mathbf{x}^{(j)} = y^{(j)} - \sum_{i \in I} \alpha^{(i)} y^{(i)} \mathbf{x}^{(i)} \mathbf{x}^{(j)}$$

# Creating the Dual SVM

Another objective function for finding the maximum margin hyperplane.

- Steps to create the dual SVM
  - Create the Lagrangian
  - Minimize the Lagrangian over the primal variables by taking the gradient and setting the gradient to zero (thus we can write the primal variables in terms of the dual variables)
  - Substitute the primal variables with their equivalent dual formulas
  - Maximize the Lagrangian
  - Find  $\mathbf{w}, w_0$  from dual variables



Know how to do this!

# Lagrangian

Example: Given training examples  $(\mathbf{x}^T, y)$ :

$$((1, 2.5), 1), ((2, 2), 1), ((1, 1), -1)$$

$$\min_{w_0, w_1, w_2} \frac{1}{2} [w_1, w_2] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

Subject to:

$$y^{(1)} (w_0 + w_1 x_1^{(1)} + w_2 x_2^{(1)}) \geq 1,$$

$$y^{(2)} (w_0 + w_1 x_1^{(2)} + w_2 x_2^{(2)}) \geq 1,$$

$$y^{(3)} (w_0 + w_1 x_1^{(3)} + w_2 x_2^{(3)}) \geq 1$$

$$\min_{w_0, w_1, w_2} \frac{1}{2} [w_1, w_2] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

Subject to:

$$(1 - (w_0 + 1w_1 + 2.5w_2)) \leq 0$$

$$(1 - (w_0 + 2w_1 + 2w_2)) \leq 0$$

$$(1 - (-w_0 + -1w_1 + -1w_2)) \leq 0$$

The constrained quadratic optimization function can be rewritten as:

$$\begin{aligned} \min_{w_0, w_1, w_2} \max_{\alpha^{(1)} \geq 0} \max_{\alpha^{(2)} \geq 0} \max_{\alpha^{(3)} \geq 0} \frac{1}{2} [w_1, w_2] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} &+ \alpha^{(1)} (1 - (w_0 + 1w_1 + 2.5w_2)) + \alpha^{(2)} (1 - (w_0 + 2w_1 + 2w_2)) \\ &+ \alpha^{(3)} (1 - (-w_0 + -1w_1 + -1w_2)) \end{aligned}$$

$$L(\mathbf{w}, w_0, \alpha) = \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^N \alpha^{(i)} (1 - y^{(i)}(w_0 + \mathbf{w}^T \mathbf{x}^{(i)}))$$

$$L(\mathbf{w}, w_0, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha^{(i)} (y^{(i)}(w_0 + \mathbf{w}^T \mathbf{x}^{(i)}) - 1)$$

# Background on solving constrained optimization

---

LAGRANGE DUALITY

# Background: The Lagrangian

[https://commons.wikimedia.org/wiki/File:Joseph-Louis\\_Lagrange.jpeg](https://commons.wikimedia.org/wiki/File:Joseph-Louis_Lagrange.jpeg)



Simplified problem

minimize  $\mathbf{u}$ :

$$\frac{1}{2} \mathbf{u}^T \mathbf{u}$$

primal problem

subject to:

$$\begin{aligned} \mathbf{a}^T \mathbf{u} &\geq c & \text{where } \mathbf{a}^T \mathbf{u} \geq c \text{ is } a_1 u_1 + a_2 u_2 + \dots + a_r u_r \geq c \\ \mathbf{a}'^T \mathbf{u} &\geq c' \end{aligned}$$

If there is a *valid* solution, this is equal to:

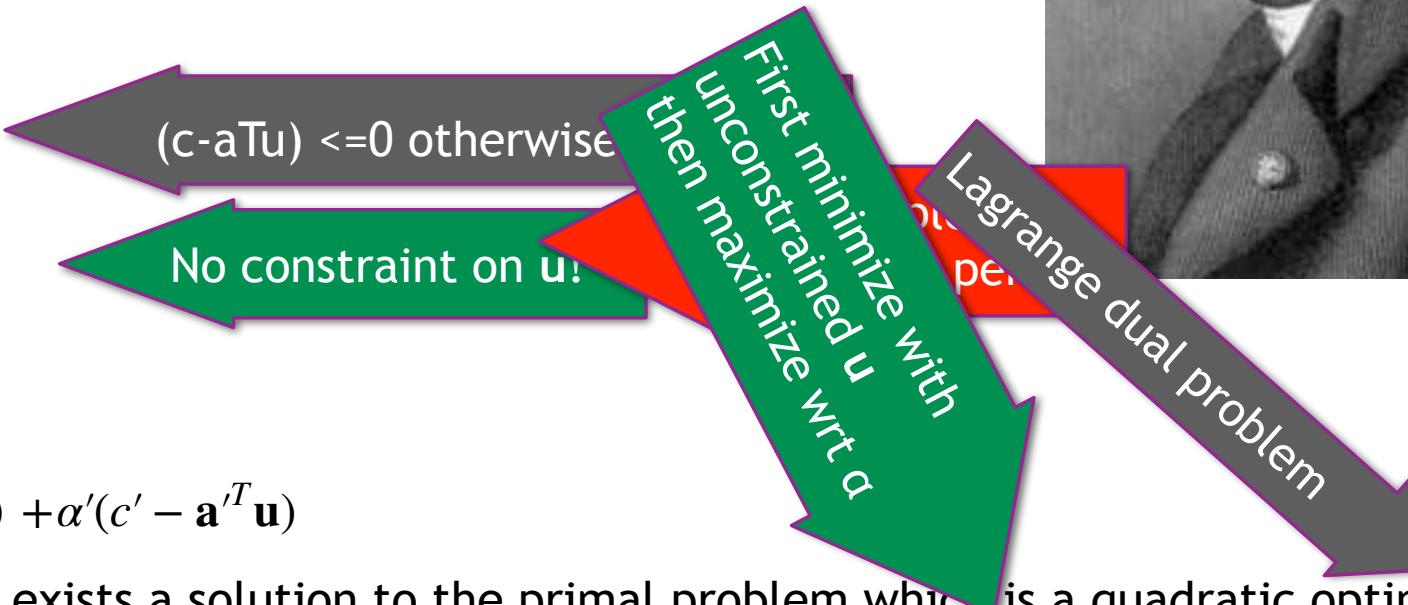
$$\begin{aligned} \text{minimize } \mathbf{u}: \quad & \frac{1}{2} \mathbf{u}^T \mathbf{u} + \max_{\alpha \geq 0} \alpha(c - \mathbf{a}^T \mathbf{u}) \\ & + \max_{\alpha' \geq 0} \alpha'(c' - \mathbf{a}'^T \mathbf{u}) \end{aligned}$$

The Lagrangian function

$$L(\mathbf{u}, \alpha) = \frac{1}{2} \mathbf{u}^T \mathbf{u} + \alpha(c - \mathbf{a}^T \mathbf{u}) + \alpha'(c' - \mathbf{a}'^T \mathbf{u})$$

The optimization function

$$\min_{\mathbf{u}} \max_{\alpha \geq 0} L(\mathbf{u}, \alpha)$$



If there exists a solution to the primal problem which is a quadratic optimization with linear constraints there is **strong duality**:  $\min_{\mathbf{u}} \max_{\alpha \geq 0} L(\mathbf{u}, \alpha) = \max_{\alpha \geq 0} \min_{\mathbf{u}} L(\mathbf{u}, \alpha)$

SVM's can appear magical because we can transform the feature space  $\Phi(\mathbf{x}^{(i)})$  (e.g. polynomial, RBF) and never have to compute  $\Phi(\mathbf{x}^{(i)})$ , instead we only need to compute the kernel  $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x}^{(j)})$  (In the description below we will not transform the feature space to keep the explanation clearer - but you can replace each  $\mathbf{x}^{(i)}$  by  $\Phi(\mathbf{x}^{(i)})$  yourself)

To show how this works, we transform our original problem into its dual formalization.

Given our constrained quadratic optimization problem:  $\min \frac{1}{2} \|\mathbf{w}\|_2^2$  subject to  $y^{(i)}(w_0 + \mathbf{w}^T \mathbf{x}^{(i)}) \geq 1$  for all  $i=1,\dots,N$

We create the Lagrangian (each constraint is multiplied by a Lagrangian multiplier  $\alpha^{(i)}$ , for  $i = 1,\dots,N$ )

$$L(\mathbf{w}, w_0, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha^{(i)} (y^{(i)}(w_0 + \mathbf{w}^T \mathbf{x}^{(i)}) - 1) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \mathbf{w}^T \sum_{i=1}^N \alpha^{(i)} y^{(i)} \mathbf{x}^{(i)} - w_0 \sum_{i=1}^N \alpha^{(i)} y^{(i)} + \sum_{i=1}^N \alpha^{(i)}$$

$\alpha^{(i)} \geq 0$

The original problem is now the same as minimization of  $L(\mathbf{w}, w_0, \alpha)$ . We first minimize with respect to  $\mathbf{w}$ ,  $w_0$ ,  $(\alpha = [\alpha^{(1)}, \dots, \alpha^{(N)}])$ . We take the gradient and find the optimal by setting the gradient to zero

$$\nabla_{\mathbf{w}} L(\mathbf{w}, w_0, \alpha) = \mathbf{w} = \sum_{i=1}^N \alpha^{(i)} y^{(i)} \mathbf{x}^{(i)}$$

**w is a linear sum of the training examples**

$$\nabla_{w_0} L(\mathbf{w}, w_0, \alpha) = - \sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0$$

**convex optimization problem**

Plugging these values back into  $L(\mathbf{w}, w_0, \alpha)$  we get the dual normalization

$$L(\alpha) = \frac{1}{2} \sum_{j=1}^N \alpha^{(j)} y^{(j)} \mathbf{x}^{(j)T} \sum_{i=1}^N \alpha^{(i)} y^{(i)} \mathbf{x}^{(i)} - \sum_{j=1}^N \alpha^{(j)} y^{(j)} \mathbf{x}^{(j)T} \sum_{i=1}^N \alpha^{(i)} y^{(i)} \mathbf{x}^{(i)} + \sum_{i=1}^N \alpha^{(i)}$$

subject to  $\sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0$  and  $\alpha^{(i)} \geq 0$

# Now we can send this to a quadratic programming solver!

There are many specialized QP solvers designed specifically for SVM's. We will not cover them :)

Maximize:

$$L(\alpha) = -\frac{1}{2} \sum_{i,j=1}^N \alpha^{(j)} \alpha^{(i)} y^{(j)} y^{(i)} \mathbf{x}^{(j)T} \mathbf{x}^{(i)} + \sum_{i=1}^N \alpha^{(i)}$$

subject to  $\sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0$  and  $\alpha^{(i)} \geq 0$

Quadratic programming solvers usually are for minimization

$$\min_{\alpha} \frac{1}{2} \sum_{i,j=1}^N \alpha^{(j)} \alpha^{(i)} y^{(j)} y^{(i)} \mathbf{x}^{(j)T} \mathbf{x}^{(i)} - \sum_{i=1}^N \alpha^{(i)}$$

subject to  $\sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0$  and  $\alpha^{(i)} \geq 0$

$$\min_{\alpha} \frac{1}{2} \boldsymbol{\alpha}^T \begin{bmatrix} y^{(1)} y^{(1)T} \mathbf{x}^{(1)T} \mathbf{x}^{(1)} & y^{(1)} y^{(2)T} \mathbf{x}^{(1)T} \mathbf{x}^{(2)} & \dots & y^{(1)} y^{(N)T} \mathbf{x}^{(1)T} \mathbf{x}^{(N)} \\ y^{(2)} y^{(1)T} \mathbf{x}^{(2)T} \mathbf{x}^{(1)} & y^{(2)} y^{(2)T} \mathbf{x}^{(2)T} \mathbf{x}^{(2)} & \dots & y^{(2)} y^{(N)T} \mathbf{x}^{(2)T} \mathbf{x}^{(N)} \\ \vdots & \vdots & & \vdots \\ y^{(N)} y^{(1)T} \mathbf{x}^{(N)T} \mathbf{x}^{(1)} & y^{(N)} y^{(2)T} \mathbf{x}^{(N)T} \mathbf{x}^{(2)} & \dots & y^{(N)} y^{(N)T} \mathbf{x}^{(N)T} \mathbf{x}^{(N)} \end{bmatrix} \boldsymbol{\alpha} + (-\mathbf{1}^T) \boldsymbol{\alpha}$$

quadratic term

linear term

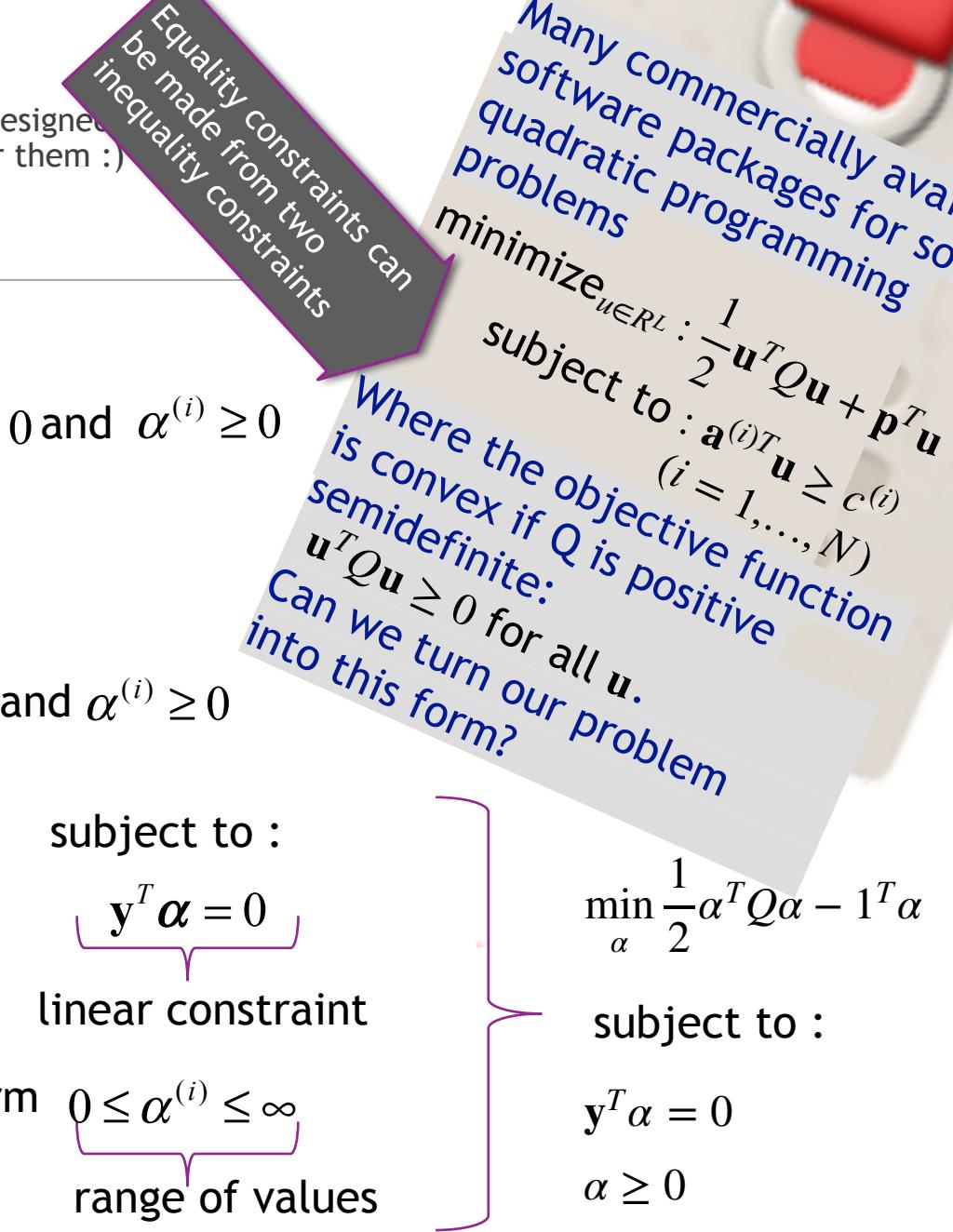
subject to :

$\mathbf{y}^T \boldsymbol{\alpha} = 0$

linear constraint

$0 \leq \alpha^{(i)} \leq \infty$

range of values



Instead of computing  $\mathbf{x}^{(i)T} \mathbf{x}^{(j)}$  we could have computed  $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$  as long as  $K$  is symmetric and  $Q$  is pos. Semidefinite

# We can perform the same calculations in Z-Space

---

Maximize:

$$L(\alpha) = -\frac{1}{2} \sum_{i,j=1}^N \alpha^{(j)} \alpha^{(i)} y^{(j)} y^{(i)} K(\mathbf{x}^{(j)}, \mathbf{x}^{(i)}) + \sum_{i=1}^N \alpha^{(i)} \quad \text{subject to } \sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0 \text{ and } \alpha^{(i)} \geq 0$$

Quadratic programming solvers usually are for minimization

$$\min_{\alpha} \frac{1}{2} \sum_{i,j=1}^N \alpha^{(j)} \alpha^{(i)} y^{(j)} y^{(i)} K(\mathbf{x}^{(j)}, \mathbf{x}^{(i)}) - \sum_{i=1}^N \alpha^{(i)} \quad \text{subject to } \sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0 \text{ and } \alpha^{(i)} \geq 0$$

$$\min_{\alpha} \frac{1}{2} \alpha^T \begin{bmatrix} y^{(1)} y^{(1)} K(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) & y^{(1)} y^{(2)} K(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) & \cdots & y^{(1)} y^{(N)} K(\mathbf{x}^{(1)}, \mathbf{x}^{(N)}) \\ y^{(2)} y^{(1)} K(\mathbf{x}^{(2)}, \mathbf{x}^{(1)}) & y^{(2)} y^{(2)} K(\mathbf{x}^{(2)}, \mathbf{x}^{(2)}) & \cdots & y^{(2)} y^{(N)} K(\mathbf{x}^{(2)}, \mathbf{x}^{(N)}) \\ \vdots & \vdots & & \vdots \\ y^{(N)} y^{(1)} K(\mathbf{x}^{(N)}, \mathbf{x}^{(1)}) & y^{(N)} y^{(2)} K(\mathbf{x}^{(N)}, \mathbf{x}^{(2)}) & \cdots & y^{(N)} y^{(N)} K(\mathbf{x}^{(N)}, \mathbf{x}^{(N)}) \end{bmatrix} \alpha - \sum_{i=1}^N \alpha^{(i)}$$

quadratic term

subject to :  $\mathbf{y}^T \alpha = 0$   
 linear constraint  $\alpha + (-1^T) \alpha = 0$   
 linear term  $\sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0$   
 range of values  $0 \leq \alpha^{(i)} \leq \infty$

min  $\frac{1}{2} \alpha^T Q \alpha - \mathbf{1}^T \alpha$   
 subject to :  $\mathbf{y}^T \alpha = 0$   
 $\alpha \geq 0$

For symmetric  $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x}^{(j)})$  and positive semidefinite Q



# “Kernel Trick”

Write the algorithm so the data only appeared when it was part of the dot product  $\Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x}^{(j)})$

- For some special types of transformations, we can efficiently compute  $\Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x}^{(j)})$  without transforming  $\Phi(\mathbf{x}^{(i)})$  or  $\Phi(\mathbf{x}^{(j)})$ 
  - Remember the polynomial transformation?
  - Define  $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x}^{(j)})$
  - We never need to create  $\mathbf{z}^{(j)} = \Phi(\mathbf{x}^{(j)})$

The kernel must satisfy Mercer's conditions (beyond the scope of this course)

This means...if I can compute,  $\Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x}^{(j)})$  efficiently I never need to worry about the length of  $\Phi(\mathbf{x}^{(j)})$

# Learning in a transformed feature space:

Pair share: if we use  $\Phi(\mathbf{x})$ , do we need to actually perform the transformation?

## Primal

$$\min_{\mathbf{w}, w_0} \|\mathbf{w}\|_2^2$$

subject to:  $y^{(i)}(\mathbf{w}^T \Phi(\mathbf{x}^{(i)}) + w_0) \geq 1$  for all  $i = 1 \dots N$

## Dual

$$\min_{\alpha} \sum_{i,j} \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) - \sum_i \alpha^{(i)}$$

$$\sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0$$

$$\alpha^{(i)} \geq 0$$

We only need to know the result of  $\Phi(\mathbf{x}^{(j)})^T \Phi(\mathbf{x}^{(i)})$

## Prediction

$$g(\mathbf{x}) = \text{sign}(\mathbf{w}^T \Phi(\mathbf{x}) + w_0)$$

$$g(\mathbf{x}) = \text{sign} \left( \sum_{i \in I} \alpha^{(i)} y^{(i)} K(\mathbf{x}^{(i)}, \mathbf{x}) + w_0 \right)$$

Support Vector  $\mathbf{x}^{(i)}$

We only need to know the result of  $\Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x})$

By computing the value of  $K$  without writing the transformation, we can have a computational advantage

# Learning in a transformed feature space:

Pair share: if we use  $\Phi(\mathbf{x})$ , do we need to actually perform the transformation?

## Primal

$$\min_{\mathbf{w}, w_0} \|\mathbf{w}\|_2^2$$

subject to:  $y^{(i)}(\mathbf{w}^T \Phi(\mathbf{x}^{(i)}) + w_0) \geq 1$  for all  $i = 1 \dots N$

## Dual

$$\min_{\alpha} \sum_{i,j} \alpha^{(i)} \alpha^{(j)} y^{(i)} y^{(j)} K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$$

$$\sum_{i=1}^N \alpha^{(i)} y^{(i)} = 0$$

$$\alpha^{(i)} \geq 0$$

Kernel trick: evaluate  $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$  without competing  $\Phi(\mathbf{x}^{(j)})$  or  $\Phi(\mathbf{x}^{(i)})$

We only need to know the result of  $\Phi(\mathbf{x}^{(j)})^T \Phi(\mathbf{x}^{(i)})$

## Prediction

$$g(\mathbf{x}) = \text{sign}(\mathbf{w}^T \Phi(\mathbf{x}) + w_0)$$

$$g(\mathbf{x}) = \text{sign} \left( \sum_{i \in I} \alpha^{(i)} y^{(i)} K(\mathbf{x}^{(i)}, \mathbf{x}) + w_0 \right)$$

Support Vector  $\mathbf{x}^{(i)}$

We only need to know the result of  $\Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x})$

By computing the value of  $K$  without writing the transformation, we can have a computational advantage

# Soft margin support vectors

MARGIN SUPPORT VECTORS

$$y^{(i)}(w_0 + \mathbf{w}^T \mathbf{x}^{(i)}) = 1$$

$$\xi^{(i)} = 0$$

$$\alpha^{(i)} > 0$$

To decide on the value of C, use cross validation!

A good first guess is:  $C = \frac{1}{N}$

NON-MARGIN SUPPORT VECTORS

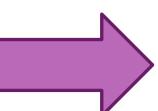
$$y^{(i)}(w_0 + \mathbf{w}^T \mathbf{x}^{(i)}) < 1 \quad \xi^{(i)} > 0$$

as long as  
you are  
violating the  
margin you  
are a  
support  
vector

# Outline

---

- ❑ Notation change, intuition, and finding how to compare hyperplanes - **mathematically how do compare hyperplanes to find the one with the maximum margin.** Can we turn this way of comparing hyperplanes into an objective function
- ❑ Support vector machines
  - ★ hard margin - **find the objective function when the data is linearly separable**
  - ★ Dealing with non-linear data - “Soft” margins for SVM - **New objective function for the case where the data is not linearly separable**
  - ★ dual formula - **a clever trick**
  - ★ Dealing with non-linear data - feature transformation with the kernel trick - **Show two popular feature maps**



# Kernel SVM

# What if the data isn't linearly separable (part 2)?

---

SOME OF THE MATERIAL IN THE NEXT FEW SLIDES ARE FROM SLIDES  
CREATED BY PROF LISA HELLERSTEIN AND PROF LINDA SELLIE

SOME SLIDES WERE INSPIRED BY SLIDES FROM CMU 18-661

# How to deal with non-linearly separable data

- Transform features?

- By adding enough new features we can make any data linearly separable...

$$\Phi(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2 \\ x_1x_2 \\ x_1^2 \\ x_2^2 \end{bmatrix}$$

We can transform  $\mathbf{x}$  into  $\Phi(\mathbf{x})$  in both the primal and dual SVM formulations

- What could be the problem?

- high computation costs
  - overfitting

- Dual formulation of support vectors are designed to deal (somewhat) with these issues

# Standard Choices of Kernel functions

- ➊ Dot product (no transformation):

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

- ➋ Polynomial kernels:

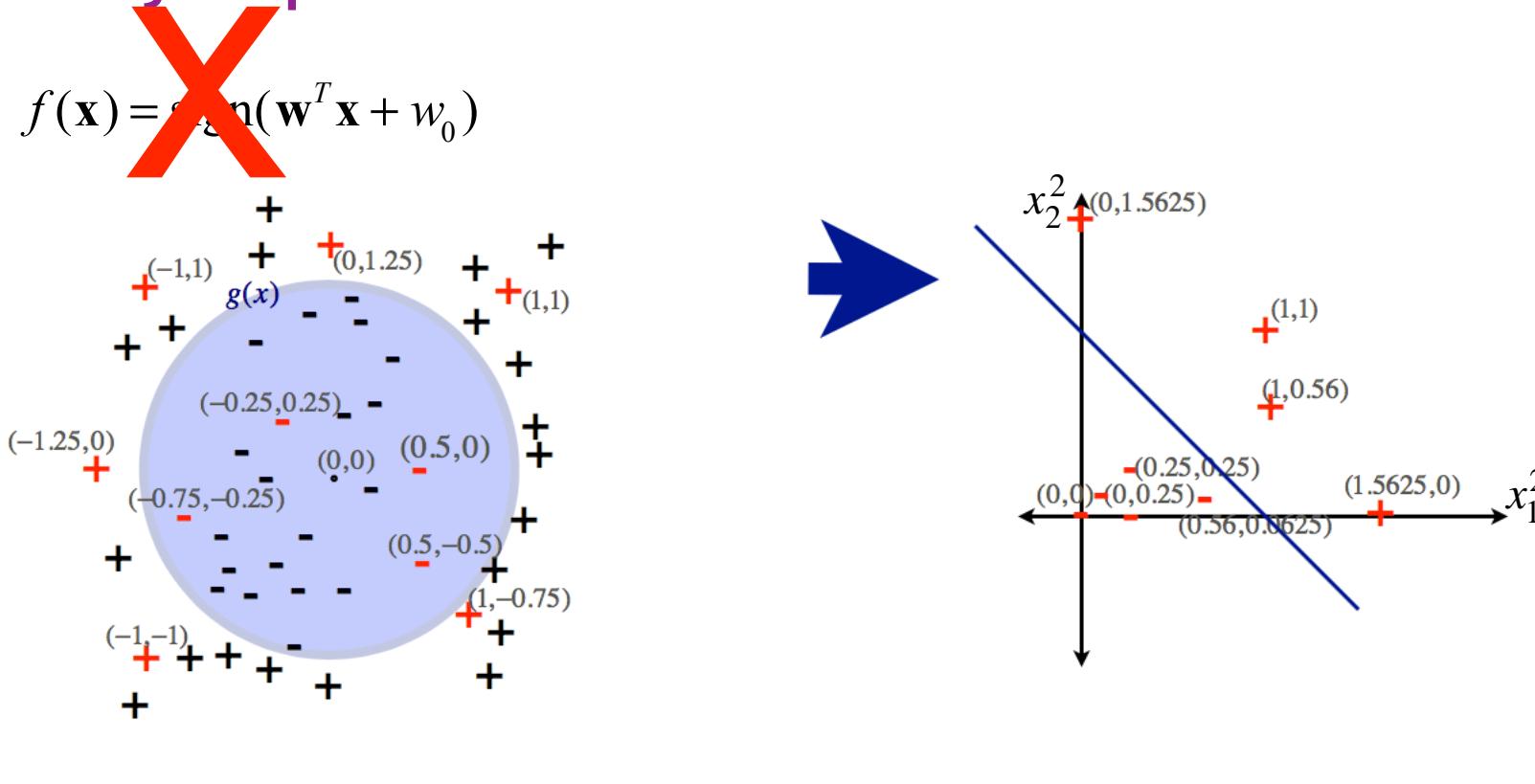
$$K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = (1 + \mathbf{x}^{(i)T} \mathbf{x}^{(j)})^d$$

- ➌ Radial basis function:

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

# Map data into a different space

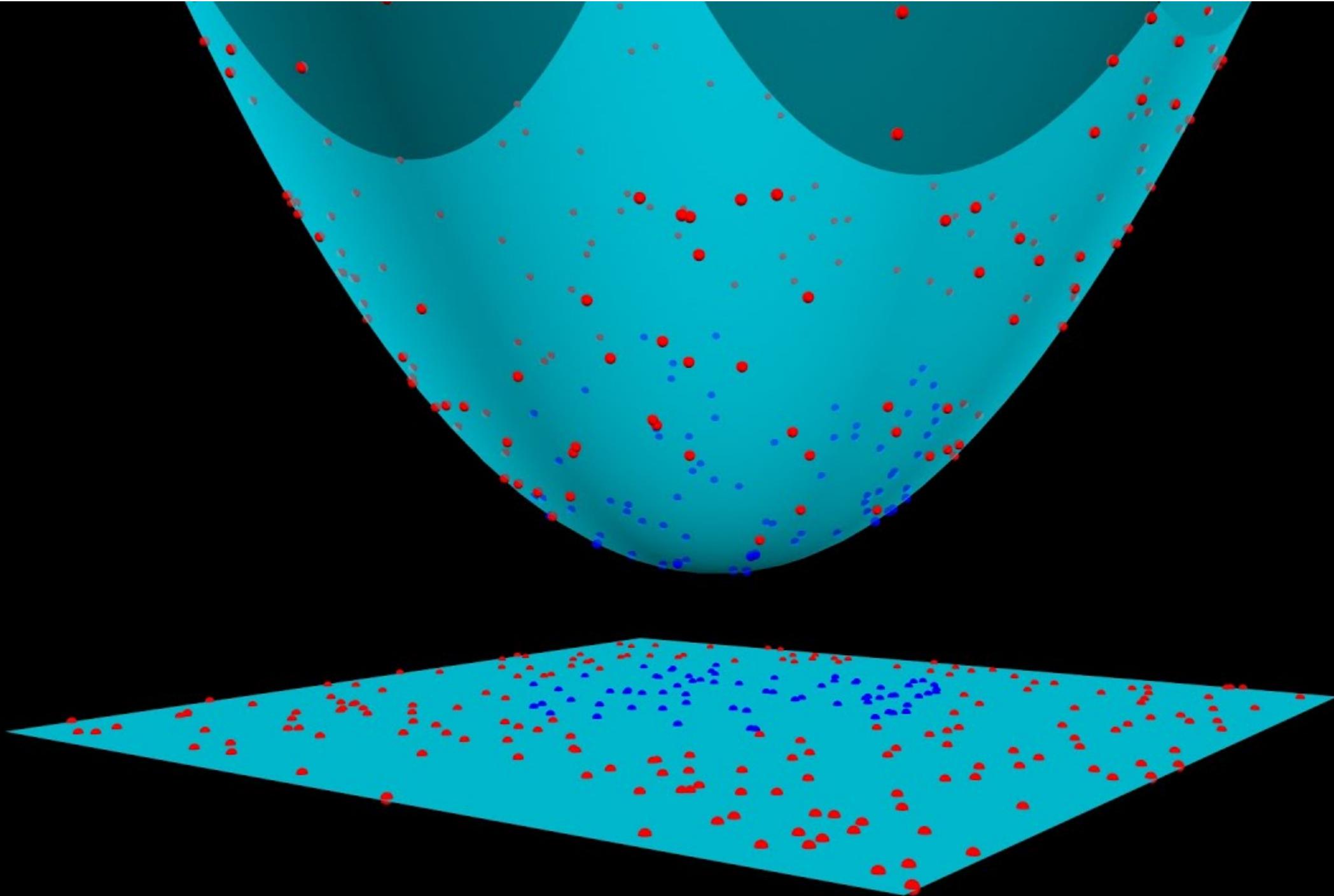
linearly separable?



If we transform the features using  $\Phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2, x_1^2, x_2^2]^T$

$$\text{sign}\left( \left( \sum_{i \in I} \alpha^{(i)} y^{(i)} \Phi(\mathbf{x}^{(i)}) \right)^T \Phi(\mathbf{x}) + w_0 \right) \rightarrow g(\mathbf{x}) = \text{sign}\left( \sum_{i \in I} \alpha^{(i)} y^{(i)} K(\mathbf{x}^{(i)}, \mathbf{x}) + w_0 \right)$$

$K(\mathbf{x}^{(i)}, \mathbf{x}) = (1 + \mathbf{x}^{(i)T} \mathbf{x})^2$



# Another feature transformation idea

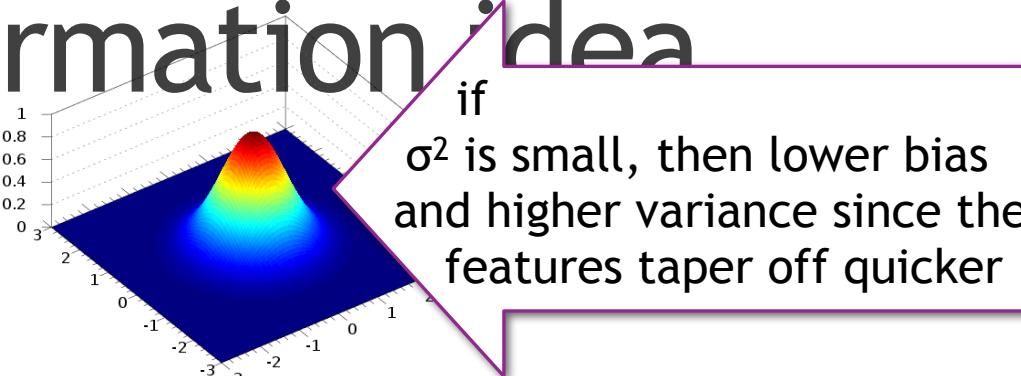
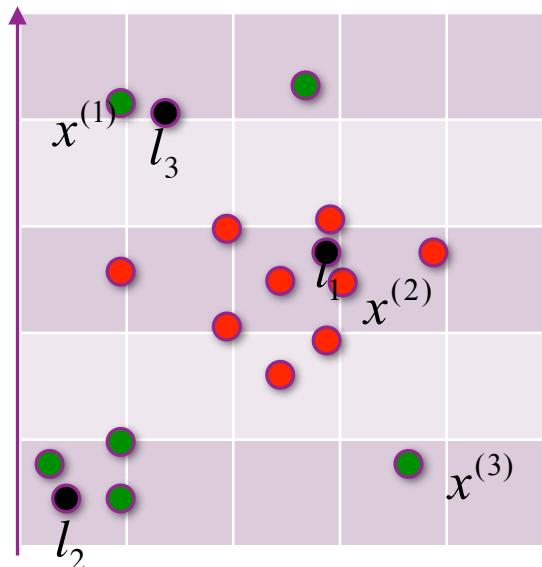
□ Define new features,  $f_1$ ,  $f_2$ , and  $f_3$  for every example on how close it is to chosen “landmarks”  $l^{(1)}$ ,  $l^{(2)}$ , and  $l^{(3)}$

□ For an example  $\mathbf{x}$  we create new features  $z_1$ ,  $z_2$ ,  $z_3$ :

$$z_1 = \exp\left\{-\frac{\|\mathbf{x} - l^{(1)}\|_2^2}{2\sigma^2}\right\}$$

$$z_2 = \exp\left\{-\frac{\|\mathbf{x} - l^{(2)}\|_2^2}{2\sigma^2}\right\}$$

$$z_3 = \exp\left\{-\frac{\|\mathbf{x} - l^{(3)}\|_2^2}{2\sigma^2}\right\}$$



if  $\sigma^2$  is small, then lower bias and higher variance since the features taper off quicker

□  $\Phi(\mathbf{x}^{(1)}) = (z_1, z_2, z_3)$  approx  $(0, 0, 1)$

□  $\Phi(\mathbf{x}^{(2)}) = (z_1, z_2, z_3)$  approx  $(1, 0, 0)$

□ The larger the sigma - the more weight is given to features far away from the landmark

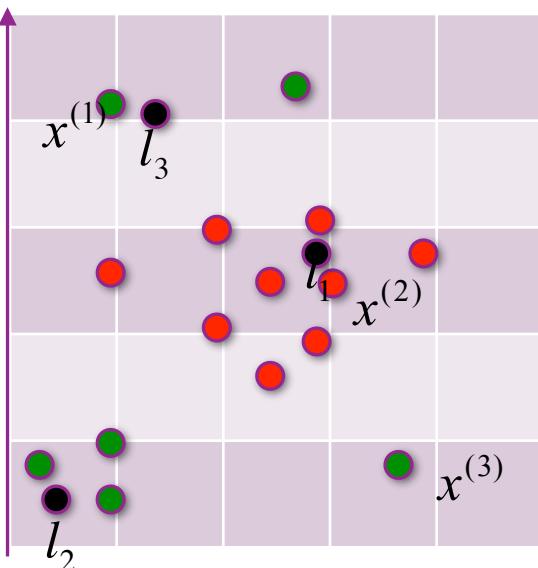
# Another feature transformation

- Define new features,  $f_1$ ,  $f_2$ , and  $f_3$  for every example  $x^{(i)}$  based on how close it is to chosen “landmarks”  $l^{(1)}$ ,  $l^{(2)}$ , and  $l^{(3)}$
- For an example  $\mathbf{x}$  we create new features  $z_1$ ,  $z_2$ ,  $z_3$ :

$$z_1 = \exp\left\{-\frac{\|\mathbf{x} - l^{(1)}\|_2^2}{2\sigma^2}\right\}$$

$$z_2 = \exp\left\{-\frac{\|\mathbf{x} - l^{(2)}\|_2^2}{2\sigma^2}\right\}$$

$$z_3 = \exp\left\{-\frac{\|\mathbf{x} - l^{(3)}\|_2^2}{2\sigma^2}\right\}$$



$\|\mathbf{x} - l^{(i)}\|_2^2$   
square of Euclidean  
distance of  $\mathbf{x}$  and  $l^{(i)}$

if  $\sigma^2$  is large, then  
higher bias and lower  
variance since the  
features taper off  
slower.  $\sigma^2$  is the  
radius of influence

□  $\Phi(\mathbf{x}^{(1)}) = (z_1, z_2, z_3) \approx (0, 0, 1)$

□  $\Phi(\mathbf{x}^{(2)}) = (z_1, z_2, z_3) \approx (1, 0, 0)$

□ The larger the sigma - the more weight is given to features far away from the landmark

# Gaussian-RBF Kernel

$$\Phi(\mathbf{x}) = \exp\left(-\frac{(\mathbf{x})^2}{2\sigma^2}\right) \left(1, \sqrt{\frac{2^1}{1!}}(\mathbf{x})^1, \sqrt{\frac{2^2}{2!}}(\mathbf{x})^2, \sqrt{\frac{2^3}{3!}}(\mathbf{x})^3, \dots\right)$$

*infinite-dimensions transformation*

$$K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp\left(-\frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|_2^2}{2\sigma^2}\right) = \exp\left(-\frac{(\mathbf{x}^{(i)})^2}{2\sigma^2}\right) \exp\left(\frac{2\mathbf{x}^{(i)}\mathbf{x}^{(j)}}{2\sigma^2}\right) \exp\left(-\frac{(\mathbf{x}^{(j)})^2}{2\sigma^2}\right)$$

*Taylor series*

$$= \exp\left(-\frac{(\mathbf{x}^{(i)})^2}{2\sigma^2}\right) \left(\sum_{k=0}^{\infty} \frac{2^k (\mathbf{x}^{(i)})^k (\mathbf{x}^{(j)})^k}{k!}\right) \exp\left(-\frac{(\mathbf{x}^{(j)})^2}{2\sigma^2}\right)$$
$$= \Phi(\mathbf{x}^{(i)})^T \Phi(\mathbf{x}^{(j)})$$

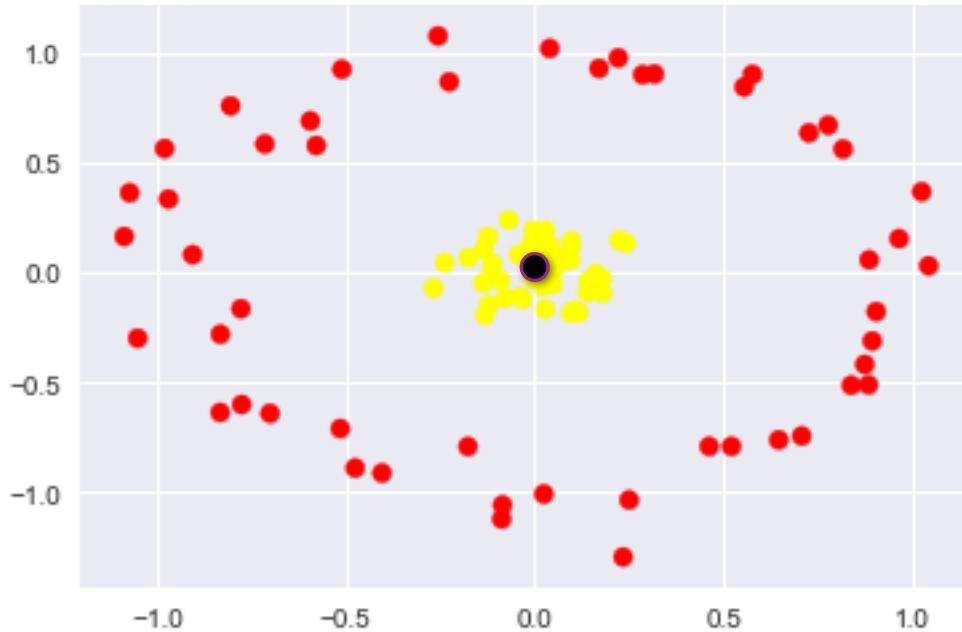
The Taylor series expansion of  $e^x$

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!}$$

see: <http://people.math.sc.edu/girardi/m142/handouts/10sTaylorPolySeries.pdf>

## Example

We choose to make a landmark of a center point,  $l$ , and compute a new feature for example  $x^{(i)}$



## Lets transform the feature space by using a RBF (radial basis function) Kernel SVM

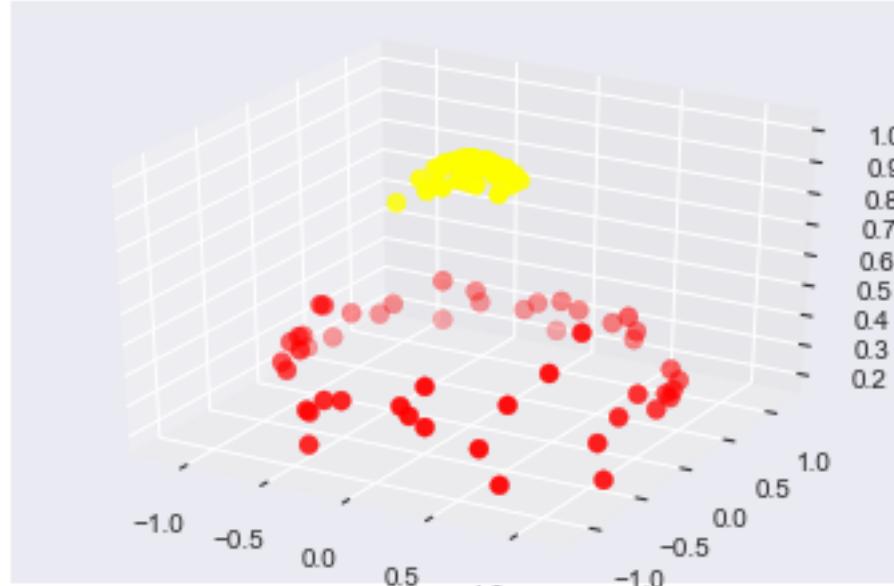
We choose to make a landmark of a center point,  $l$ , and compute a new feature for example  $x^{(i)}$

$$z_1^{(i)} = \exp\left(-\frac{\|\mathbf{x}^{(i)} - \ell\|_2^2}{(2\sigma^2)}\right)$$

For every example  $x^{(i)}$ , where  $z^{(i)}$  will have only one feature.

Now all the examples have only one feature (instead of 2 features).

We will show the original features and the new feature  $z$



In the new space, it is easy to see we can separate the points

# Gaussian Radial Basis Function

- For every example make a landmark:

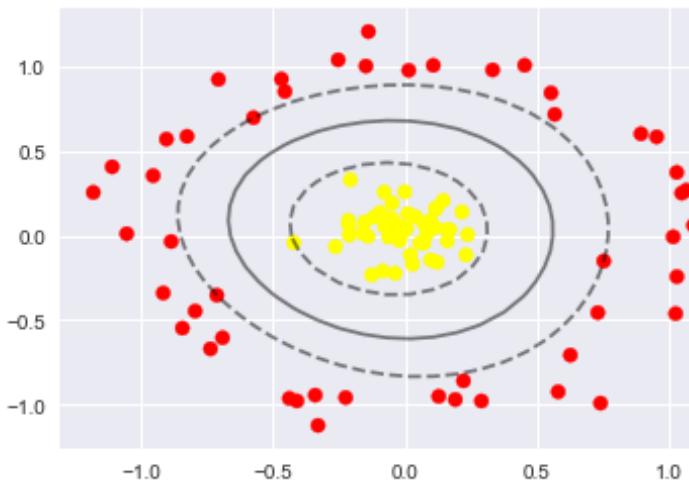
$$l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(N)} = x^{(N)}$$

- Then for every  $x$ , we have created  $N$  new features  $z_1, z_2, \dots, z_n$ :

$$K(x^{(i)}, l^{(j)}) = z_j^{(i)} = \exp\left(-\frac{\|x^{(i)} - l^{(j)}\|_2^2}{2\sigma^2}\right)$$

- The number of features is now  $N$

$$z = \begin{bmatrix} \exp(\|\mathbf{x} - l^{(1)}\|_2^2 / 2\sigma^2) \\ \exp(\|\mathbf{x} - l^{(2)}\|_2^2 / 2\sigma^2) \\ \vdots \\ \exp(\|\mathbf{x} - l^{(N-1)}\|_2^2 / 2\sigma^2) \\ \exp(\|\mathbf{x} - l^{(N)}\|_2^2 / 2\sigma^2) \end{bmatrix}$$



support vectors  
[-0.7116806, -0.35318537]  
[-0.57297015, 0.69623782]  
[ 0.7511019, -0.1515738 ]  
[-0.42044048, -0.04330787]

# Gaussian Radial Basis Function

- For every example make a landmark  $l^{(1)} = \mathbf{x}^{(1)}, l^{(2)} = \mathbf{x}^{(2)}, \dots, l^{(N)} = \mathbf{x}^{(N)}$
- Then for every  $\mathbf{x}$ , we have created  $N$  new features  $z_1, z_2, \dots, z_n$ :

$$K(\mathbf{x}^{(i)}, \mathbf{l}^{(j)}) = z_j^{(i)} = \exp\left(-\frac{\|\mathbf{x}^{(i)} - \mathbf{l}^{(j)}\|_2^2}{2\sigma^2}\right)$$

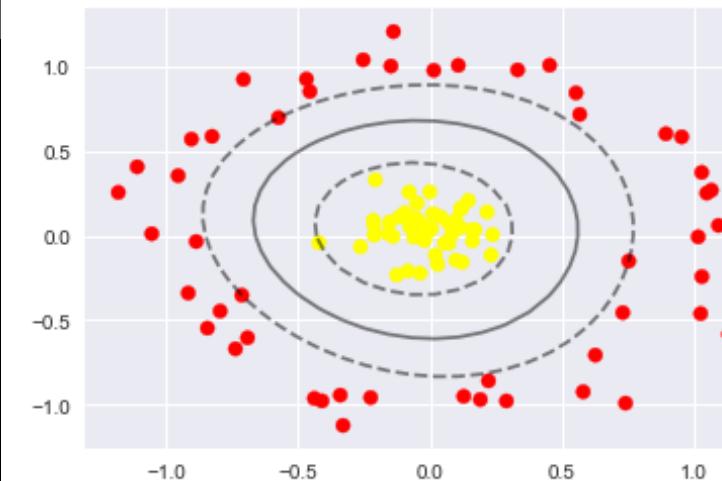
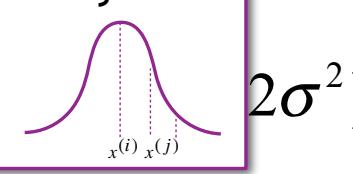
- The number of features is now  $N$

Perform feature scaling!

$$\exp\left(\|\mathbf{x} - l^{(1)}\|_2^2 / 2\sigma^2\right)$$

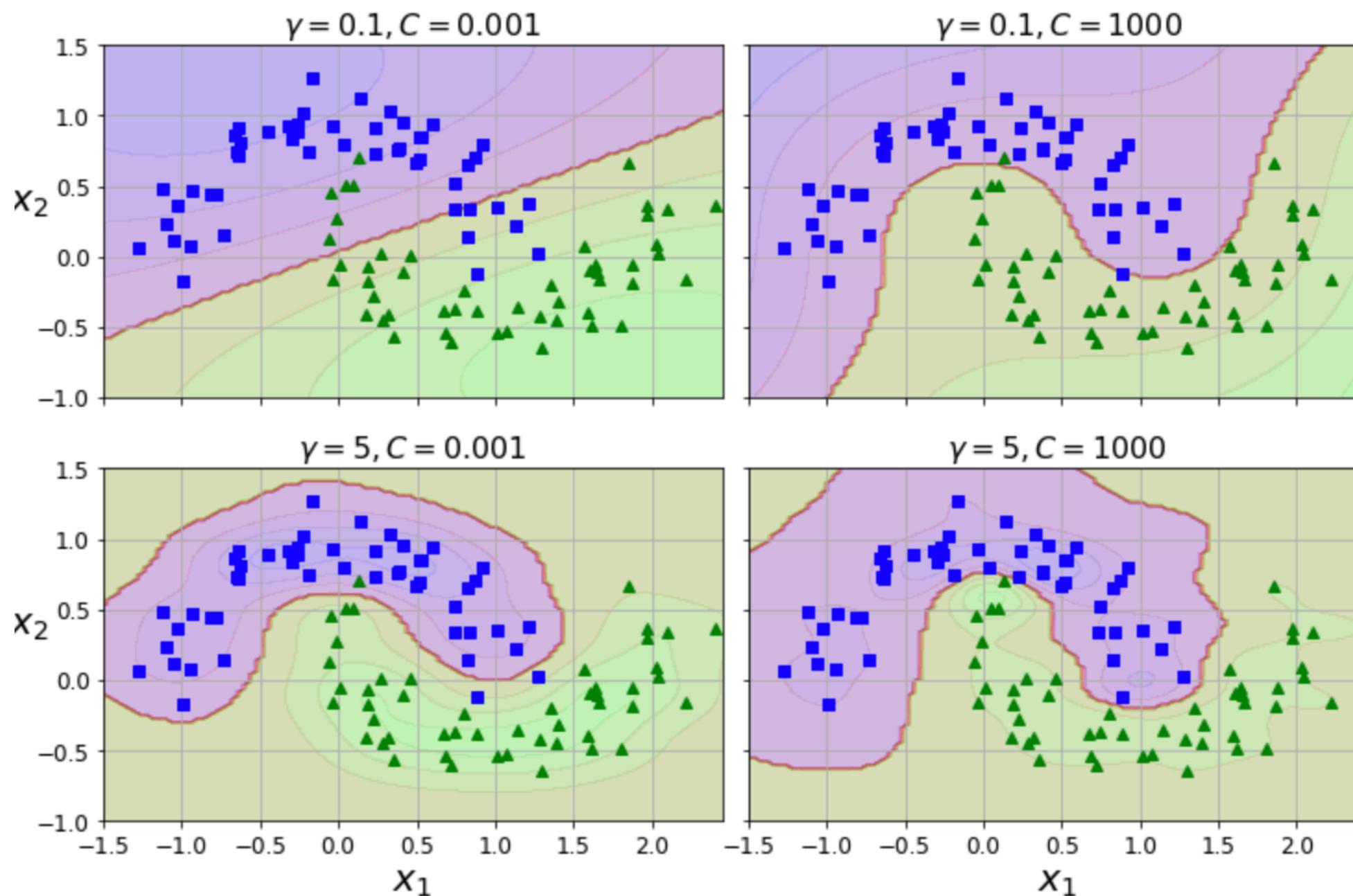
$$\dots \exp\left(\|\mathbf{x} - l^{(N)}\|_2^2 / 2\sigma^2\right)$$

Similarity measure.  $\mathbf{x}^{(i)}, \mathbf{x}^{(j)}$  are “very similar” if they are at most  $\sigma$  distance



support vectors  
[-0.7116806, -0.35318537]  
[-0.57297015, 0.69623782]  
[ 0.7511019, -0.1515738 ]  
[-0.42044048, -0.04330787]

$$\gamma = \frac{1}{2\sigma^2}$$



Read the next slides on how to use the sklearn implementation of SVM on your own

---

# Sklearn

---

LinearSVC: optimizer for the primal problem

SVC: optimizer for the dual problem



## 1.4.6. Kernel functions

The *kernel function* can be any of the following:

- linear:  $\langle x, x' \rangle$ .
- polynomial:  $(\gamma \langle x, x' \rangle + r)^d$ , where  $d$  is specified by parameter `degree`,  $r$  by `coef0`.
- rbf:  $\exp(-\gamma \|x - x'\|^2)$ , where  $\gamma$  is specified by parameter `gamma`, must be greater than 0.
- sigmoid  $\tanh(\gamma \langle x, x' \rangle + r)$ , where  $r$  is specified by `coef0`.

Different kernels are specified by the `kernel` parameter:

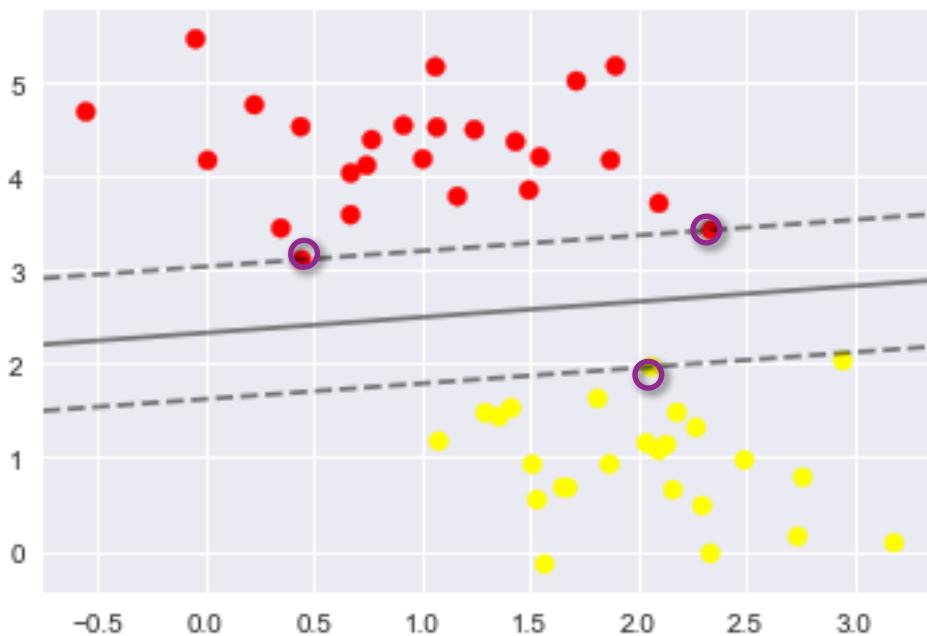
```
>>> linear_svc = svm.SVC(kernel='linear')
>>> linear_svc.kernel
'linear'
>>> rbf_svc = svm.SVC(kernel='rbf')
>>> rbf_svc.kernel
'rbf'
```

## Using sklearn's support vector classifier to find the maximum margin on the training examples X

```
: from sklearn.svm import SVC # "Support vector classifier"
model = SVC(kernel='linear', C=1E10)
model.fit(x, y)
```

By choosing the parameter C to be very large we assure no misclassification points

The classifier finds the optimal model for the training examples



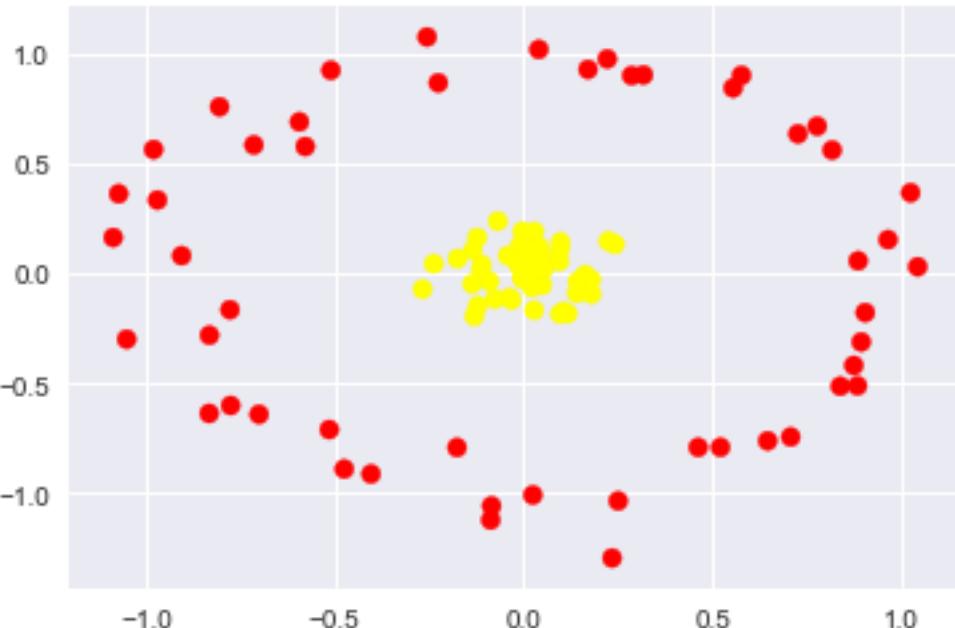
We can observe what the support vectors are:

```
model.support_vectors_
```

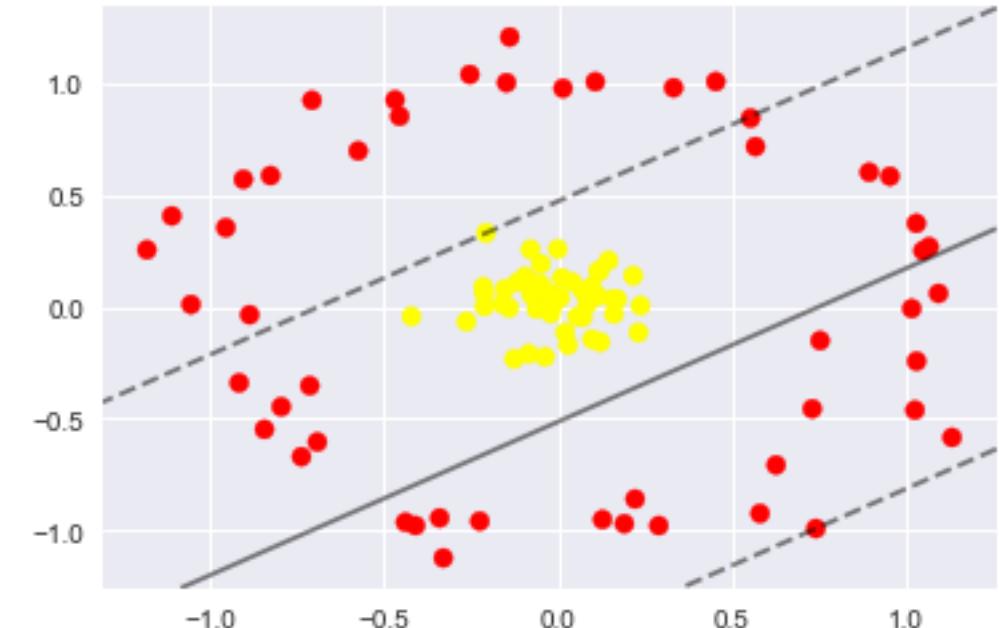
```
array([[ 0.44359863,  3.11530945],
       [ 2.33812285,  3.43116792],
       [ 2.06156753,  1.96918596]])
```

## Suppose the data is not linearly separable

Lets try to fit it with a linear support vector machine (we will allow for some misclassification)...



```
clf = SVC(kernel='linear').fit(X, y)
```



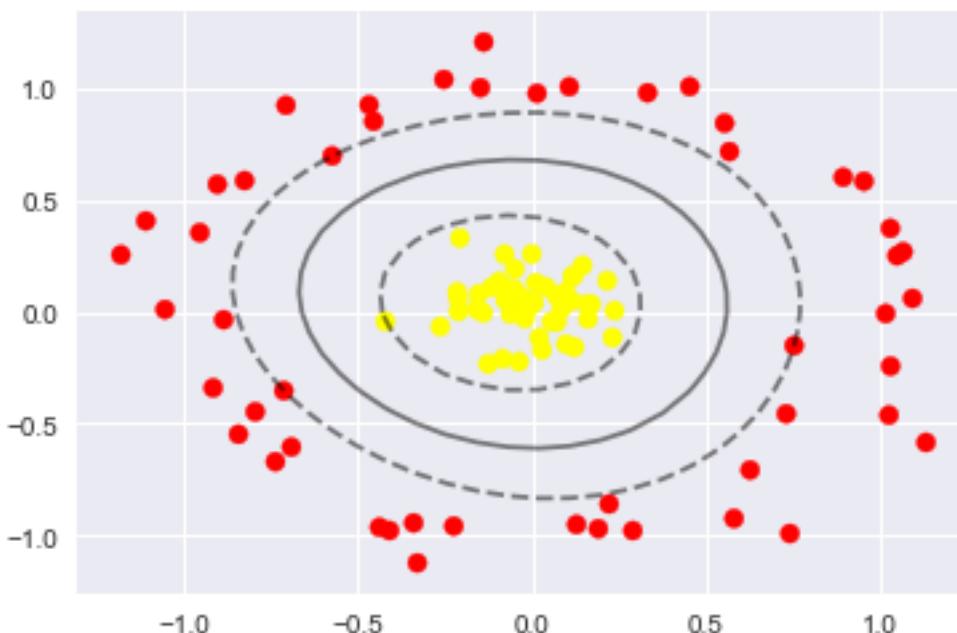
## How can we find the right point to transform the points to become linearly separable?

Easiest answer is to compute the distance to every point in the data set!

Problem? The number of features becomes N, the number of training examples

In sklearn, we can apply the radial basis function kernel (RBF kernel)

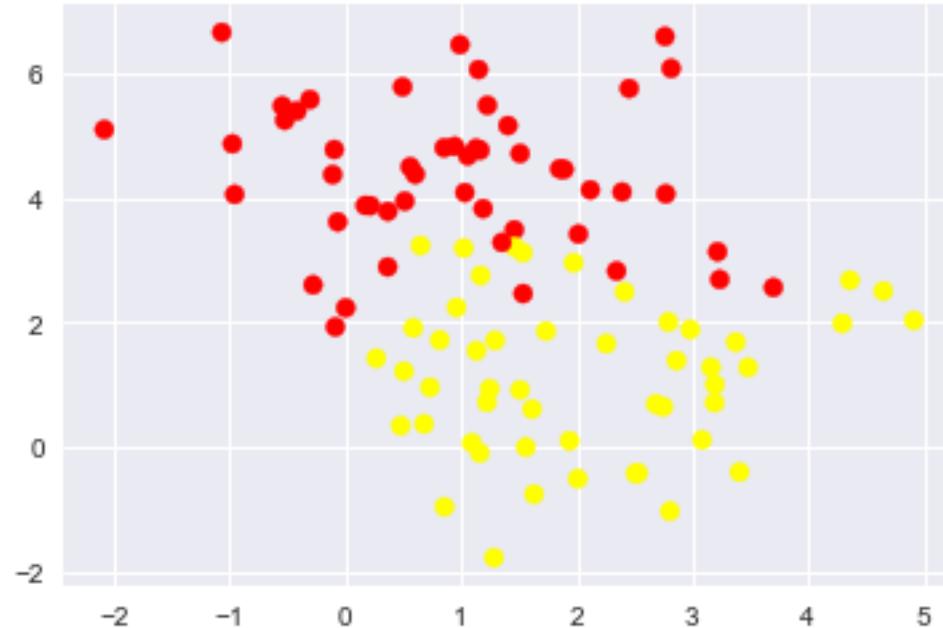
```
clf = SVC(kernel='rbf', C=1E6)
clf.fit(X, y)
```



```
clf.support_vectors_
```

```
array([[-0.7116806 , -0.35318537],
       [-0.57297015,  0.69623782],
       [ 0.7511019 , -0.1515738 ],
       [-0.42044048, -0.04330787]])
```

**Another way to work with non separable data is to allowing for errors**  
- Soft Margin SVM



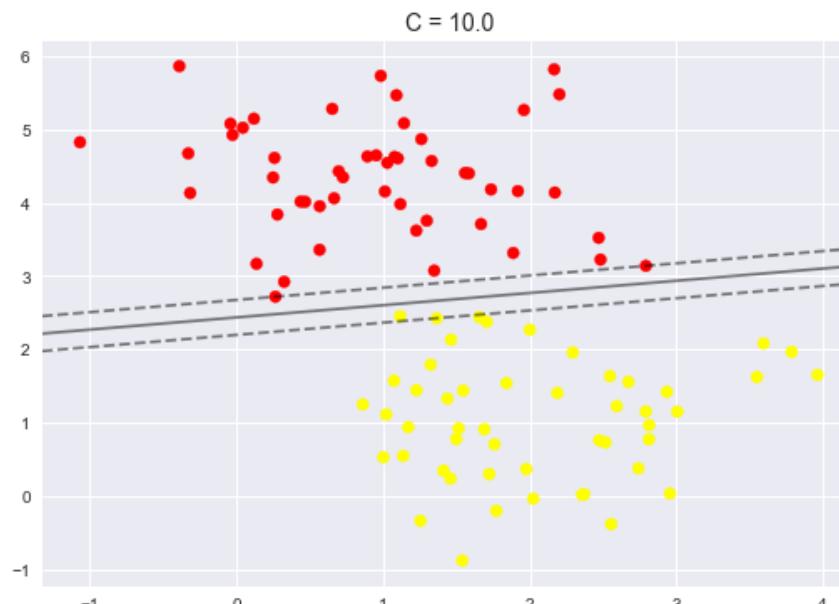
We create a soft- margin classifier by adding a cost to points that are not correctly classified and outside the margin

The primal objective function:

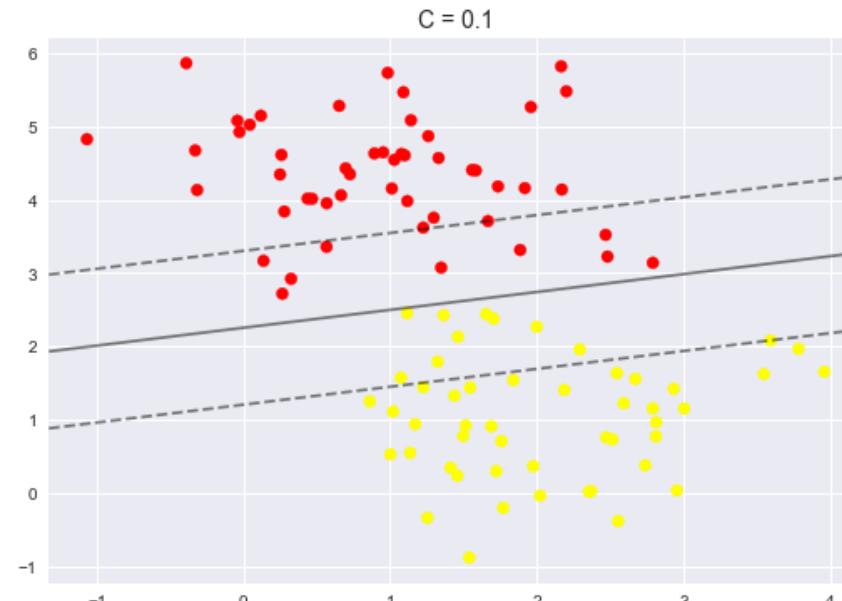
$$\min_{w_0, \mathbf{w}, \{\xi^{(i)}\}_{i=1}^N} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi^{(i)}$$

subject to  $y^{(i)}(w_0 + \mathbf{w}^T \mathbf{x}^{(i)}) \geq 1 - \xi^{(i)}$  for all  $i = 1, \dots, N$   $\xi^{(i)} \geq 0$

The tuning parameter  $C$  softens ( $C$  smaller) or hardens ( $C$  larger)

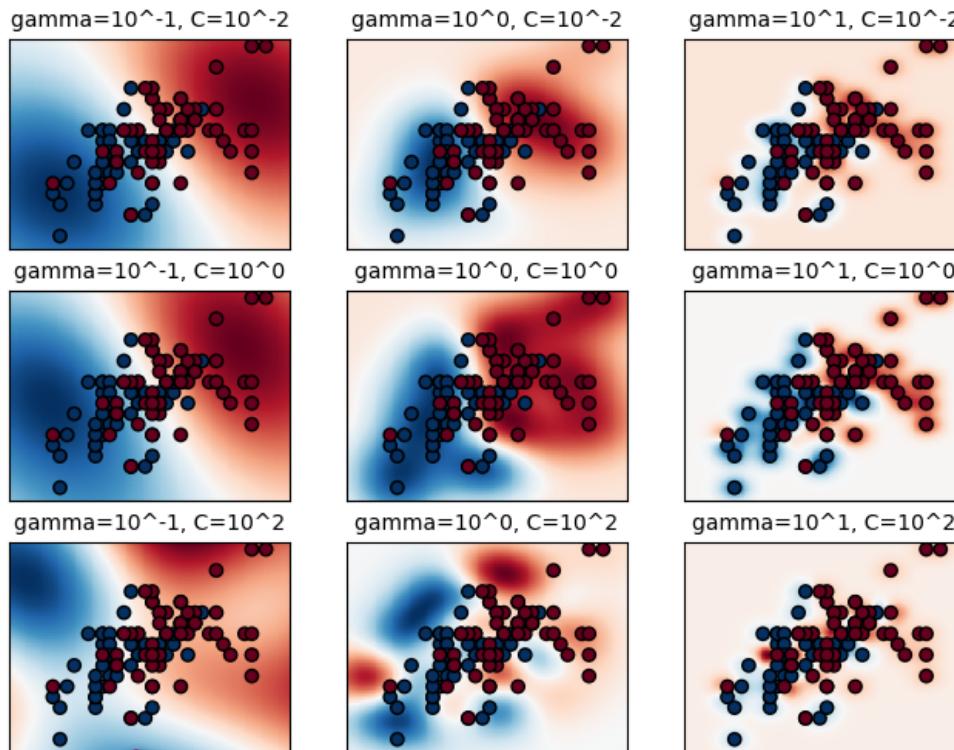


`model = SVC(kernel='linear', C=10).fit(X, y)`



`model = SVC(kernel='linear', C=0.1).fit(X, y)`

# SVC



Under-fitting for  
Small gamma

Overfitting for  
large gamma

```
from sklearn.svm import SVC
```

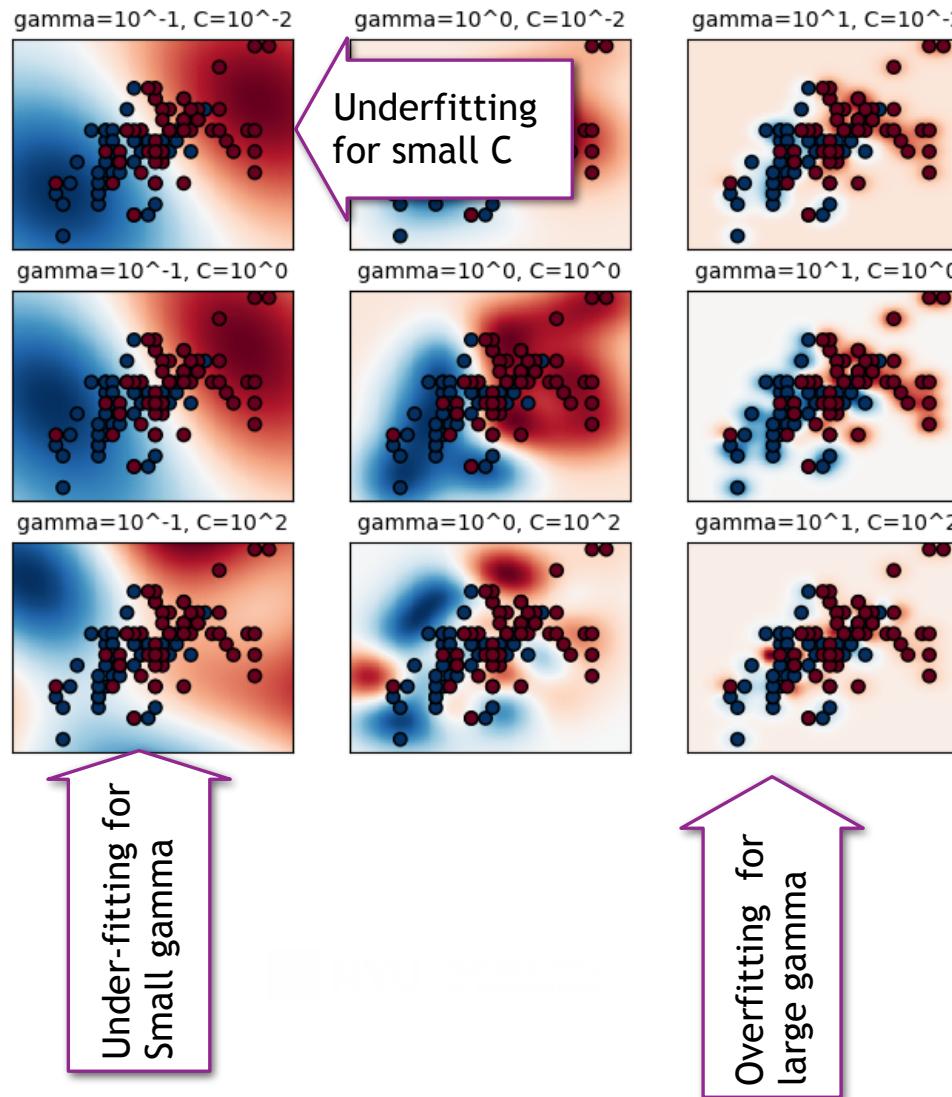
```
C_2d_range = [1e-2, 1, 1e2]
gamma_2d_range = [1e-1, 1, 1e1]
classifiers = []
for C in C_2d_range:
    for gamma in gamma_2d_range:
        clf = SVC(C=C, gamma=gamma)
        clf.fit(X_2d, y_2d)
        classifiers.append((C, gamma, clf))
```

“Intuitively, the `gamma` parameter defines how far the influence of a single training example reaches, with low values meaning ‘far’ and high values meaning ‘close’.”

“The `C` parameter trades off correct classification of training examples against maximization of the decision function’s margin”

The radius of the RBF kernel alone acts as a good structural regularizer.

# SVC



```
from sklearn.svm import SVC
```

```
C_2d_range = [1e-2, 1, 1e2]
gamma_2d_range = [1e-1, 1, 1e1]
classifiers = []
for C in C_2d_range:
    for gamma in gamma_2d_range:
        clf = SVC(C=C, gamma=gamma)
        clf.fit(X_2d, y_2d)
        classifiers.append((C, gamma, clf))
```

“Intuitively, the `gamma` parameter defines how far the influence of a single training example reaches, with low values meaning ‘far’ and high values meaning ‘close’.”

“The `C` parameter trades off correct classification of training examples against maximization of the decision function’s margin”

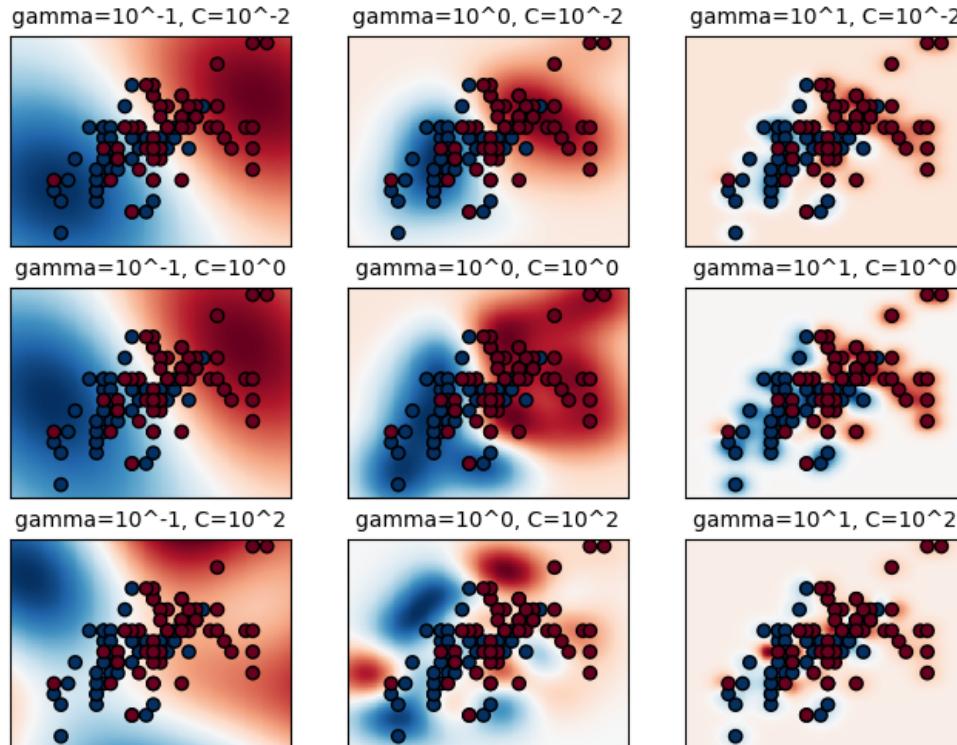
The radius of the RBF kernel alone acts as a good structural regularizer.

Ack!!! There are too  
many parameters to  
tune

---

sklearn's [RandomizedSearchCV](#) and [GridSearchCV](#) are classes for parameter tuning that methodically builds and evaluates different combinations of parameters as a grid

# GridSearchCV



```
from sklearn.svm import SVC
from sklearn.model_selection import StratifiedShuffleSplit
from sklearn.model_selection import GridSearchCV

# Train classifiers
#
# For an initial search, a logarithmic grid with basis
# 10 is often helpful. Using a basis of 2, a finer
# tuning can be achieved but at a much higher cost.

C_range = np.logspace(-2, 10, 13)
gamma_range = np.logspace(-9, 3, 13)
param_grid = dict(gamma=gamma_range, C=C_range)
cv = StratifiedShuffleSplit(n_splits=5, test_size=0.2, random_state=42)
grid = GridSearchCV(SVC(), param_grid=param_grid, cv=cv)
grid.fit(X, y)

print("The best parameters are %s with a score of %0.2f"
      % (grid.best_params_, grid.best_score_))
```

These also  
need to be  
imported