



Department of Computer Science  
UNIVERSITY OF COLORADO **BOULDER**



## Machine Learning: Chenhao Tan

University of Colorado Boulder

LECTURE 18

Slides adapted from Jordan Boyd-Graber, Chris Ketelsen

## Logistics

---

- Project proposal is due on Friday!

## Overview

---

Hinge-loss view of soft-margin SVM

Kernels

Examples

## Outline

---

Hinge-loss view of soft-margin SVM

Kernels

Examples

## Recap: Karush-Kuhn-Tucker (KKT) conditions

### Primal and dual feasibility

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, \xi_i \geq 0, C \geq \alpha_i \geq 0, \beta_i \geq 0$$

### Stationarity

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i, \sum_{i=1}^m \alpha_i y_i = 0, \alpha_i + \beta_i = C$$

### Complementary slackness

$$\alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i] = 0, \beta_i \xi_i = 0$$

## Soft-margin SVM

---

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

subject to

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, i \in [1, m]$$

$$\xi_i \geq 0, i \in [1, m]$$

What is  $\xi_i$ ?

## Soft-margin SVM

---

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

subject to

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, i \in [1, m]$$

$$\xi_i \geq 0, i \in [1, m]$$

What is  $\xi_i$ ?

$$\xi_i = \begin{cases} 0, & \text{if } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 \\ 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + b), & \text{otherwise} \end{cases}$$

## Soft-margin SVM

---

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \ell^{(\text{hin})}(y_i, \mathbf{w} \cdot \mathbf{x}_i + b)$$

You can solve this with gradient descent.



## Outline

---

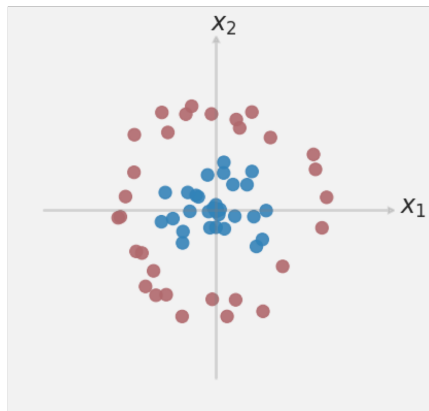
Hinge-loss view of soft-margin SVM

Kernels

Examples

Can you solve this with linear separator?

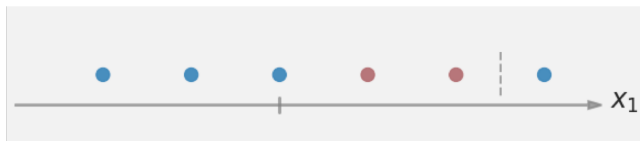
---



## Can you solve this with linear separator?

---

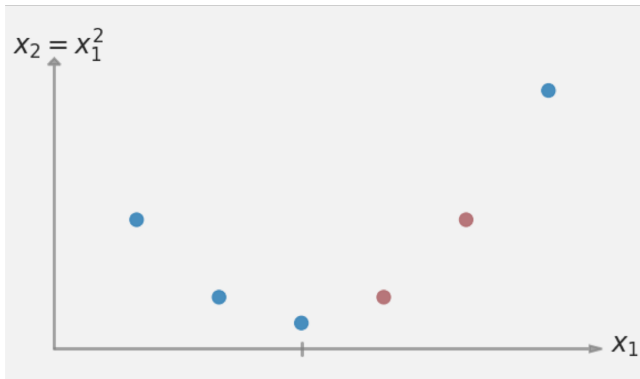
What can we do if the data is clearly not linearly separable?



## Can you solve this with linear separator?

---

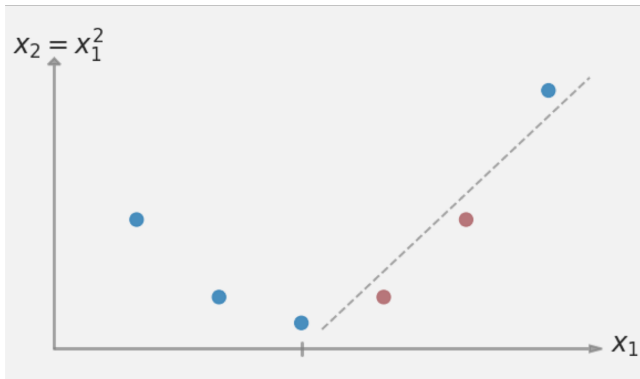
Add a dimension.



## Can you solve this with linear separator?

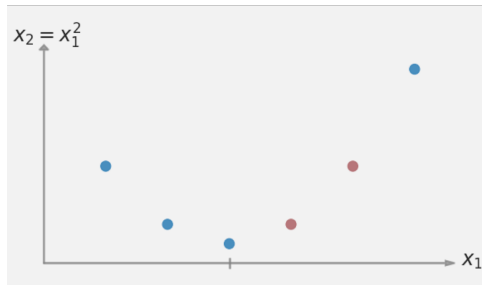
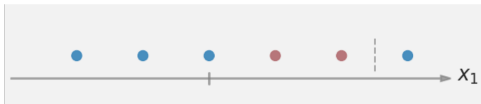
---

Add a dimension.



## Derived features

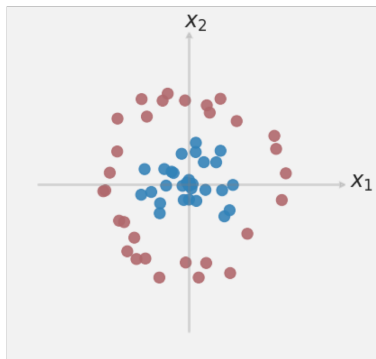
We started with the original feature vector,  $\mathbf{x} = (x_1)$ ,  
and we created a new derived feature vector,  $\phi(\mathbf{x}) = (x_1, x_1^2)$ .



## What about the previous problem?

---

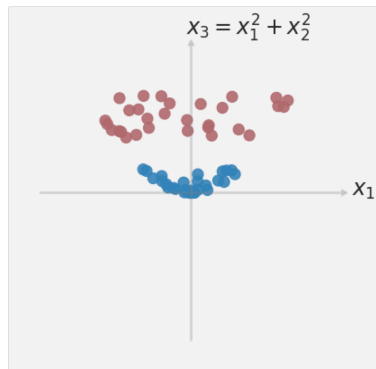
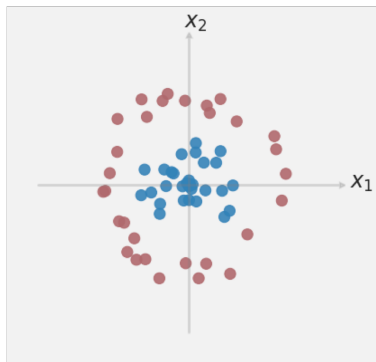
Definitely not separable in two dimensions.



## What about the previous problem?

---

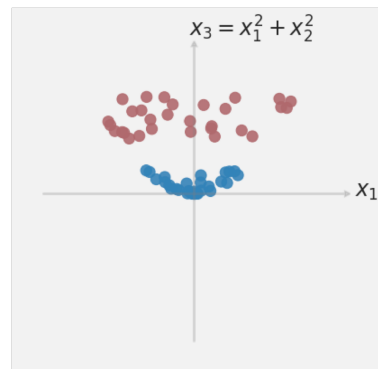
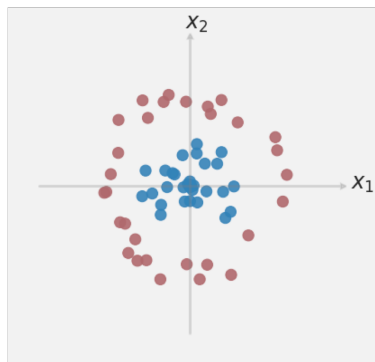
Definitely not separable in two dimensions.  
But in three dimensions, it becomes easily separable.





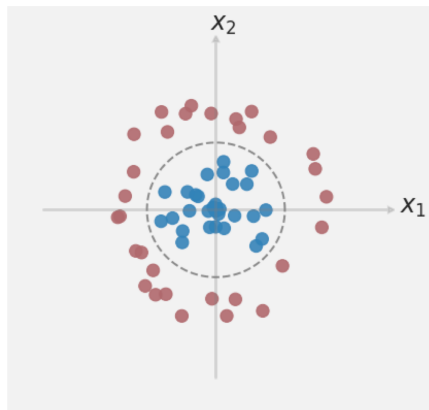
## Derived features

We started with the original feature vector,  $\mathbf{x} = (x_1, x_2)$ ,  
and we created a new derived feature vector,  $\phi(\mathbf{x}) = (x_1, x_2, x_1^2 + x_2^2)$ .



## Derived features

We started with the original feature vector,  $\mathbf{x} = (x_1, x_2)$ ,  
and we created a new derived feature vector,  $\phi(\mathbf{x}) = (x_1, x_2, x_1^2 + x_2^2)$ .



## What's special about SVMs?

---

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

## What's special about SVMs?

---

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

- This dot product is basically just how much  $x_i$  looks like  $x_j$ . Can we generalize that?

## What's special about SVMs?

---

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

- This dot product is basically just how much  $x_i$  looks like  $x_j$ . Can we generalize that?
- Kernels!

## What's special about SVMs?

---

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j))$$

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i \cdot \mathbf{x}_j)$$

## What does the kernel trick buy us?

---

Polynomial kernel:

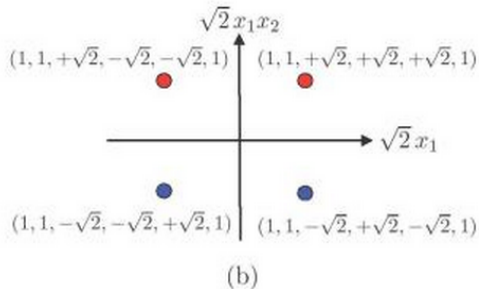
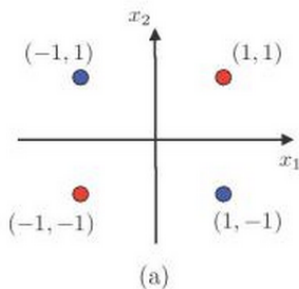
$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}' + c)^d$$

## What does the kernel trick buy us?

Polynomial kernel:

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}' + c)^d$$

When  $d = 2, c = 1$ :





## What does the kernel trick buy us?

---

Polynomial kernel:

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}' + 1)^2$$

What is the corresponding  $\phi(\mathbf{x})$ , where  $\mathbf{x} \in \mathbb{R}^k$ ?

What is the complexity of storing  $\phi(\mathbf{x})$  and computing  $\phi(\mathbf{x}) \cdot \phi(\mathbf{x}')$ ?

What about using the kernel function?

## What's a kernel?

---

- A function  $K : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$  is a kernel over  $\mathcal{X}$ .
- This is equivalent to taking the dot product  $\langle \phi(x_1), \phi(x_2) \rangle$  for some mapping
- **Mercer's Theorem:** So long as the function is continuous and symmetric, then  $K$  admits an expansion of the form

$$K(x, x') = \sum_{n=0}^{\infty} a_n \phi_n(x) \phi_n(x')$$

## What's a kernel?

---

- A function  $K : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$  is a kernel over  $\mathcal{X}$ .
- This is equivalent to taking the dot product  $\langle \phi(x_1), \phi(x_2) \rangle$  for some mapping
- **Mercer's Theorem:** So long as the function is continuous and symmetric, then  $K$  admits an expansion of the form

$$K(x, x') = \sum_{n=0}^{\infty} a_n \phi_n(x) \phi_n(x')$$

- The computational cost is just in computing the kernel

## Kernel Matrix

---

The important property of the kernel matrix  $\mathbf{K} = [K(x_i, x_j)]_{ij} \in \mathbb{R}^{m \times m}$  is symmetric positive semidefinite.

## Kernel Matrix

---

The important property of the kernel matrix  $\mathbf{K} = [K(x_i, x_j)]_{ij} \in \mathbb{R}^{m \times m}$  is symmetric positive semidefinite.

$$\mathbf{K}^T = \mathbf{K}$$

## Kernel Matrix

---

The important property of the kernel matrix  $\mathbf{K} = [K(x_i, x_j)]_{ij} \in \mathbb{R}^{m \times m}$  is symmetric positive semidefinite.

$$\mathbf{K}^T = \mathbf{K}$$

$$\forall \mathbf{x}, \mathbf{x}^T \mathbf{K} \mathbf{x} \geq 0$$

## Kernel Matrix

---

The important property of the kernel matrix  $\mathbf{K} = [K(x_i, x_j)]_{ij} \in \mathbb{R}^{m \times m}$  is symmetric positive semidefinite.

$$\mathbf{K}^T = \mathbf{K}$$

$$\forall \mathbf{x}, \mathbf{x}^T \mathbf{K} \mathbf{x} \geq 0$$

Also known as Gram matrix.

## Gaussian Kernel

---

$$K(x, x') = \exp\left(-\frac{\|x' - x\|^2}{2\sigma^2}\right)$$



## Gaussian Kernel

---

$$K(x, x') = \exp \left( -\frac{\|x' - x\|^2}{2\sigma^2} \right)$$

which can be rewritten as

$$K(x, x') = \sum_n \frac{(x \cdot x')^n}{\sigma^n n!}$$

(All polynomials!)

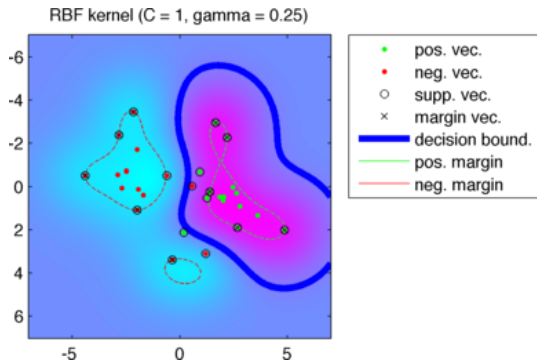
## Gaussian Kernel

$$K(x, x') = \exp\left(-\frac{\|x' - x\|^2}{2\sigma^2}\right)$$

which can be rewritten as

$$K(x, x') = \sum_n \frac{(x \cdot x')^n}{\sigma^n n!}$$

(All polynomials!)



## How does it affect optimization

---

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

- Replace all dot product with kernel evaluations  $K(\mathbf{x}_i, \mathbf{x}_j)$
- Makes computation more expensive, overall structure is the same

## Outline

---

Hinge-loss view of soft-margin SVM

Kernels

Examples

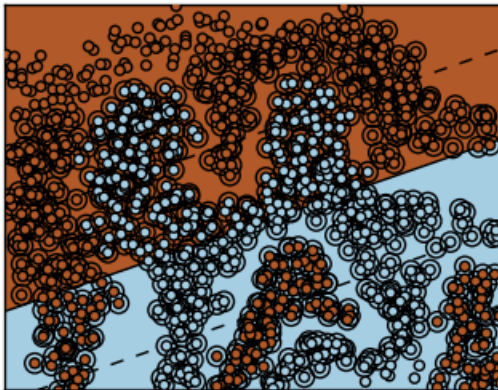
## Kernelized SVM

---

```
X, Y = read_data("ex8a.txt")  
clf = svm.SVC(kernel=kk, degree=dd, gamma=gg)  
clf.fit(X, Y)
```

## Linear Kernel Doesn't Work

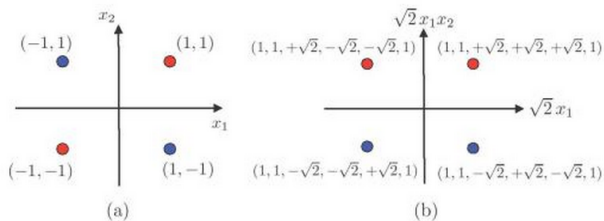
---



## Polynomial Kernel

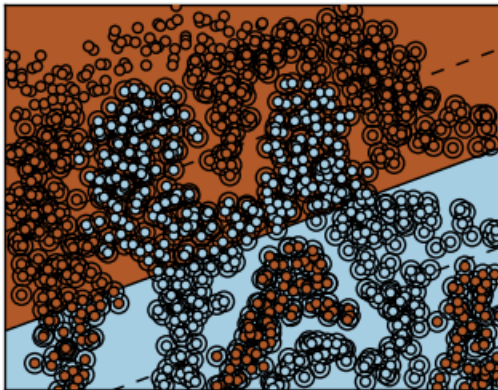
$$K(x, x') = (x \cdot x' + c)^d$$

When  $d = 2$ :



## Polynomial Kernel $d = 1, c = 5$

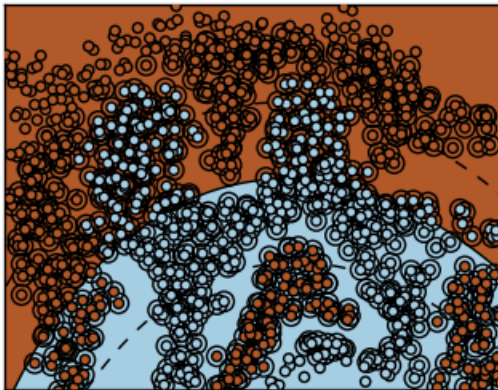
---





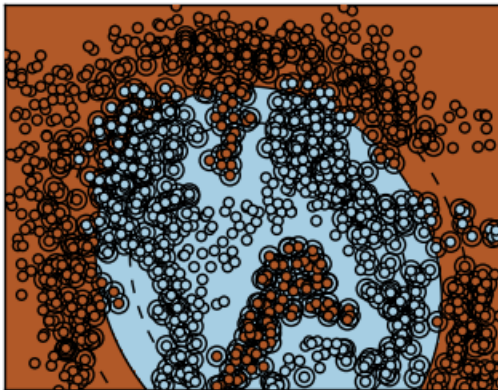
## Polynomial Kernel $d = 2, c = 5$

---



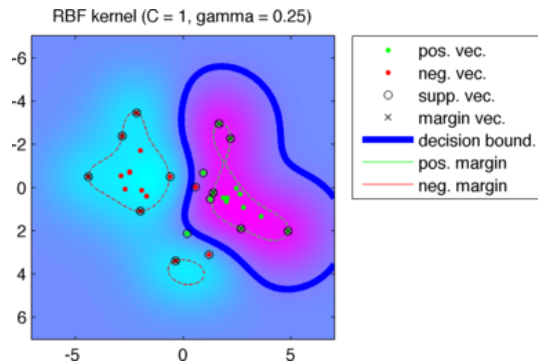
## Polynomial Kernel $d = 3, c = 5$

---



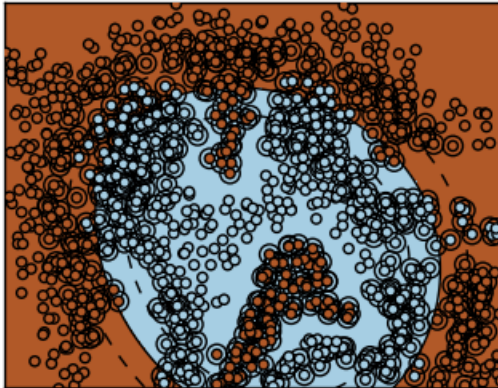
## Gaussian Kernel

$$K(x, x') = \exp \left( -\gamma \|x' - x\|^2 \right)$$



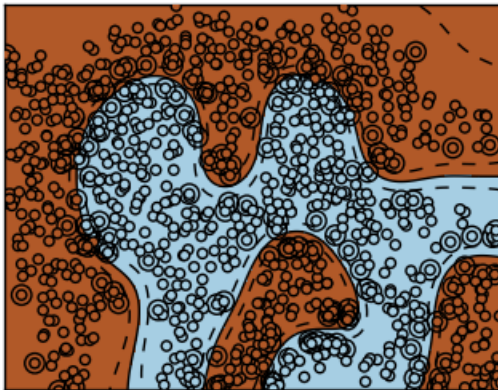
## RBF Kernel $\gamma = 2$

---



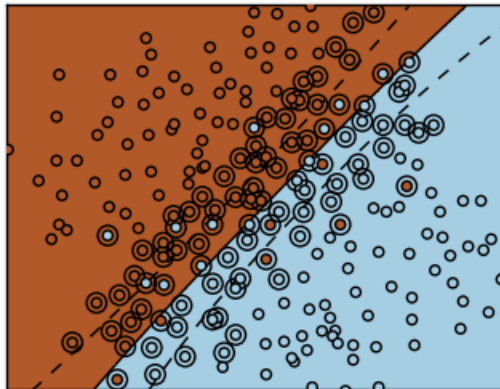
## RBF Kernel $\gamma = 100$

---



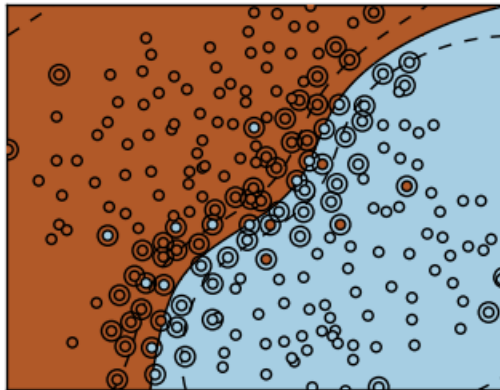
## RBF Kernel $\gamma = 1$

---



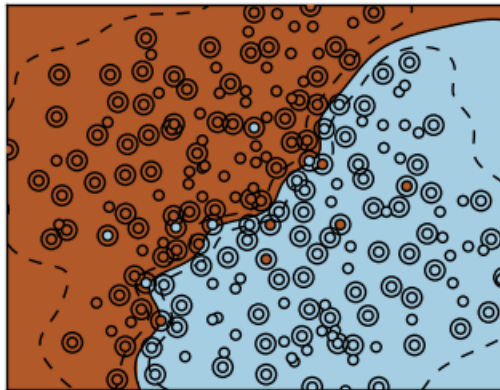
## RBF Kernel $\gamma = 10$

---



## RBF Kernel $\gamma = 100$

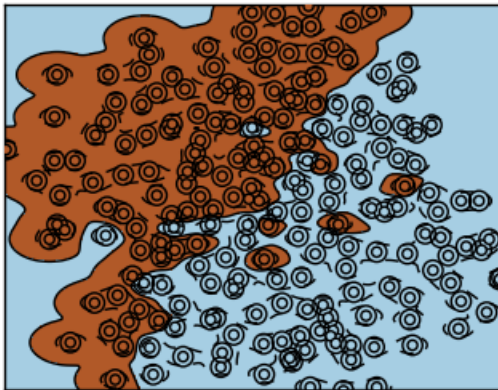
---





## RBF Kernel $\gamma = 1000$

---



## **Be careful!**

---

- Which has the lowest training error?
- Which one would generalize best?

## Recap

---

- This completes our discussion of SVMs
- Workhorse method of machine learning
- Flexible, fast, effective

## Recap

---

- This completes our discussion of SVMs
- Workhorse method of machine learning
- Flexible, fast, effective
- Kernels: applicable to wide range of data, inner product trick keeps method simple