# Support Vector Machines

# +

# Kernels

Matt Gormley
Lecture 27
Apr. 22, 2020

# Reminders

- **Homework 8: Reinforcement Learning**
  - **Out: Fri, Apr 10**
  - **Due: Wed, Apr 22 at 11:59pm**
- **Homework 9: Learning Paradigms**
  - **Out: Wed, Apr. 22**
  - **Due: Wed, Apr. 29 at 11:59pm**
  - **Can only be submitted up to 3 days late, so we can return grades before final exam**

- **Today's In-Class Poll**
  - **http://poll.mlcourse.org**

# CONSTRAINED OPTIMIZATION

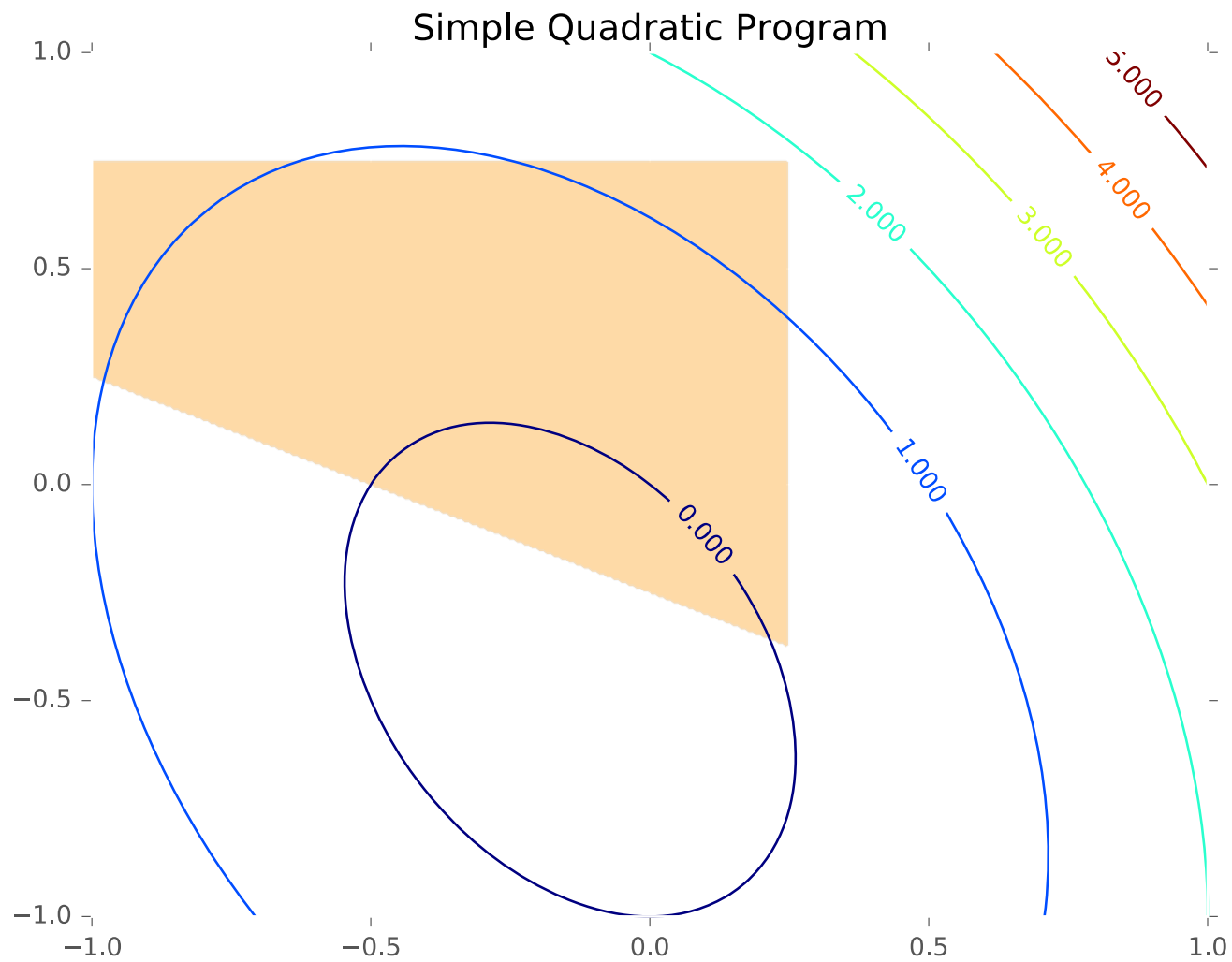# Constrained Optimization

**Unconstrained**

$$\min_{\vec{\theta}} \; J(\vec{\theta})$$

**Constrained**

$$\min_{\vec{\theta}} \; J(\vec{\theta})$$

$$\text{s.t.} \; g(\vec{\theta}) \leq \vec{b}$$

# Quadratic Program

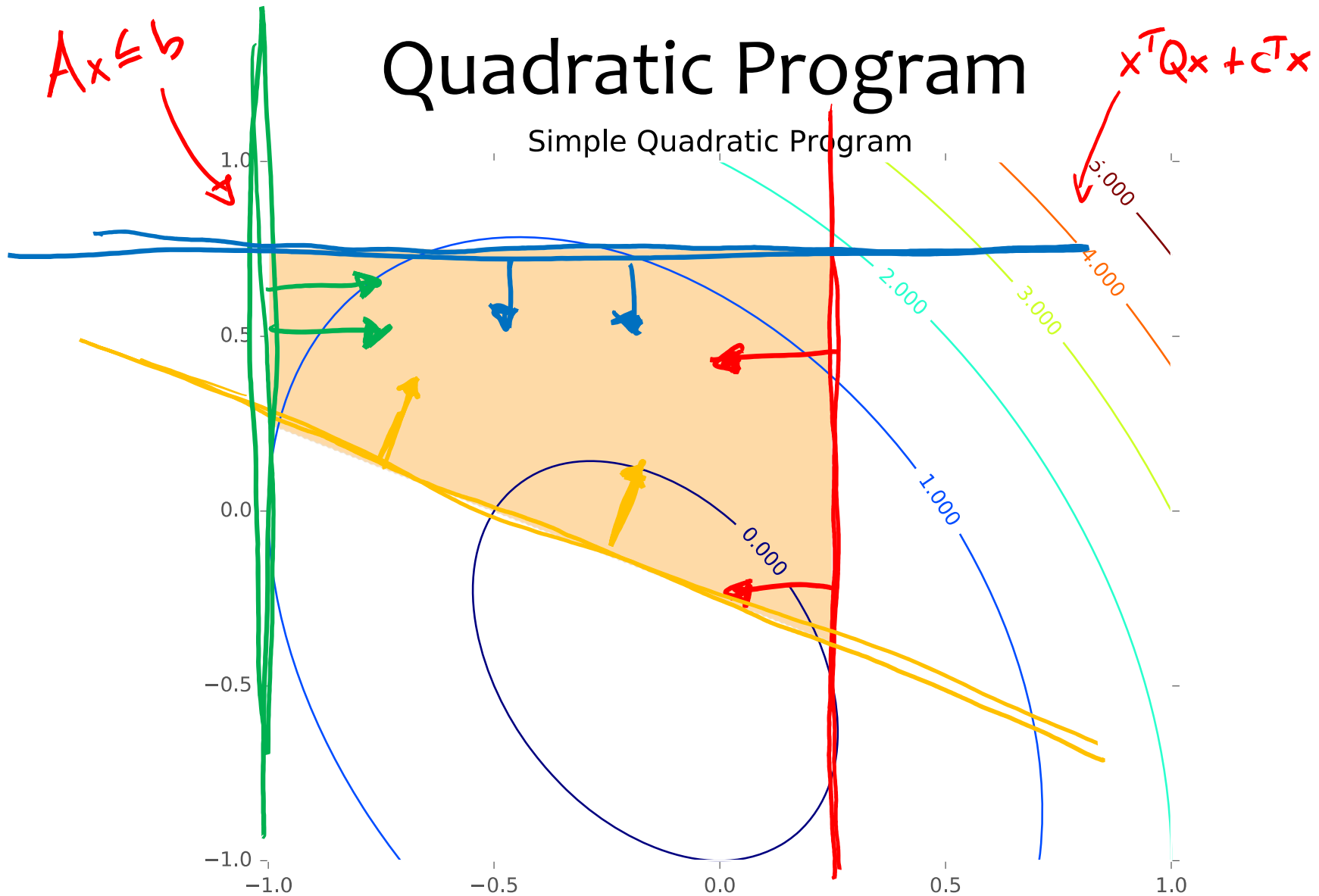

Simple Quadratic Program
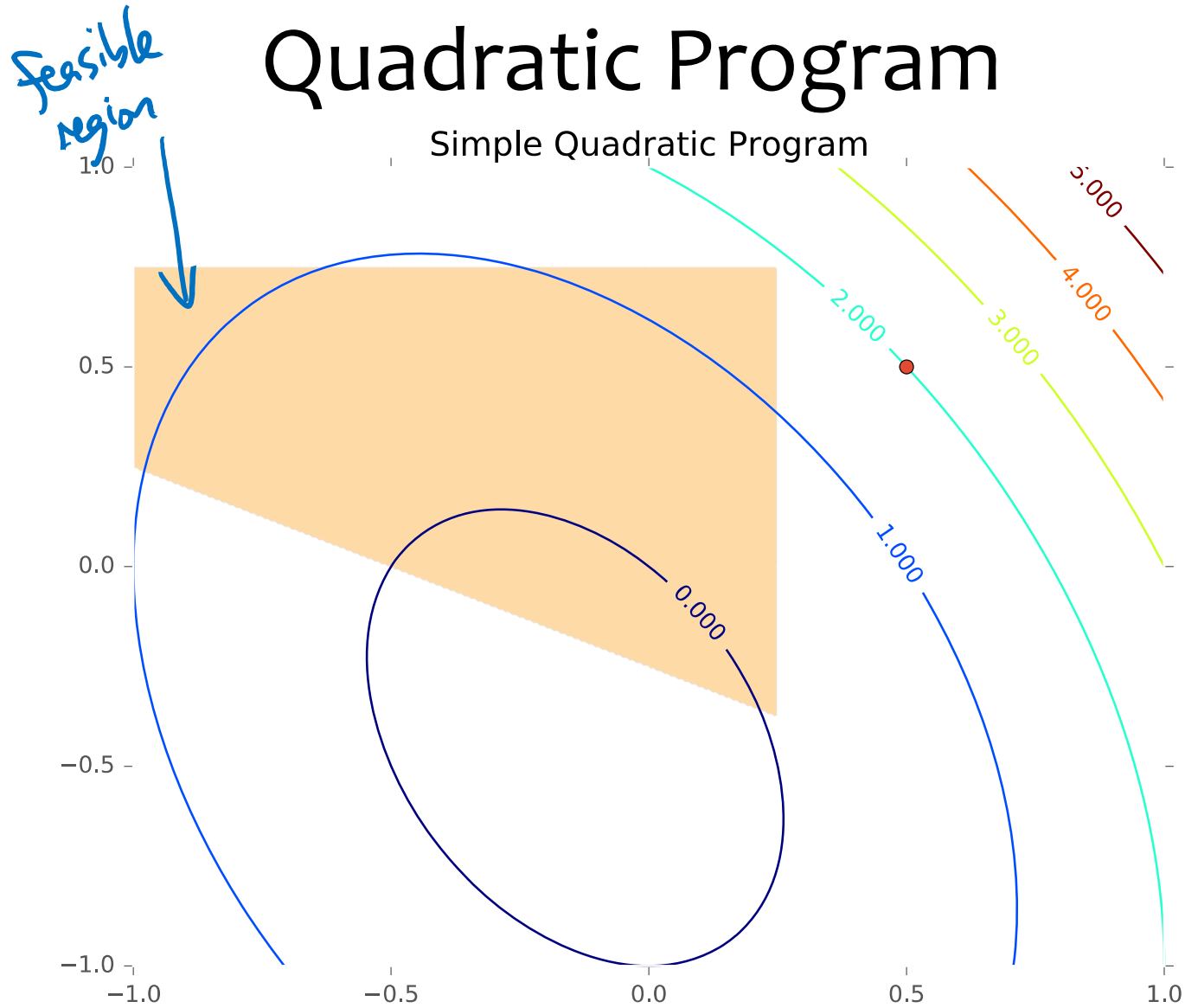
# SVM: Optimization Background

*Whiteboard*

– Constrained Optimization

– Linear programming

– Quadratic programming

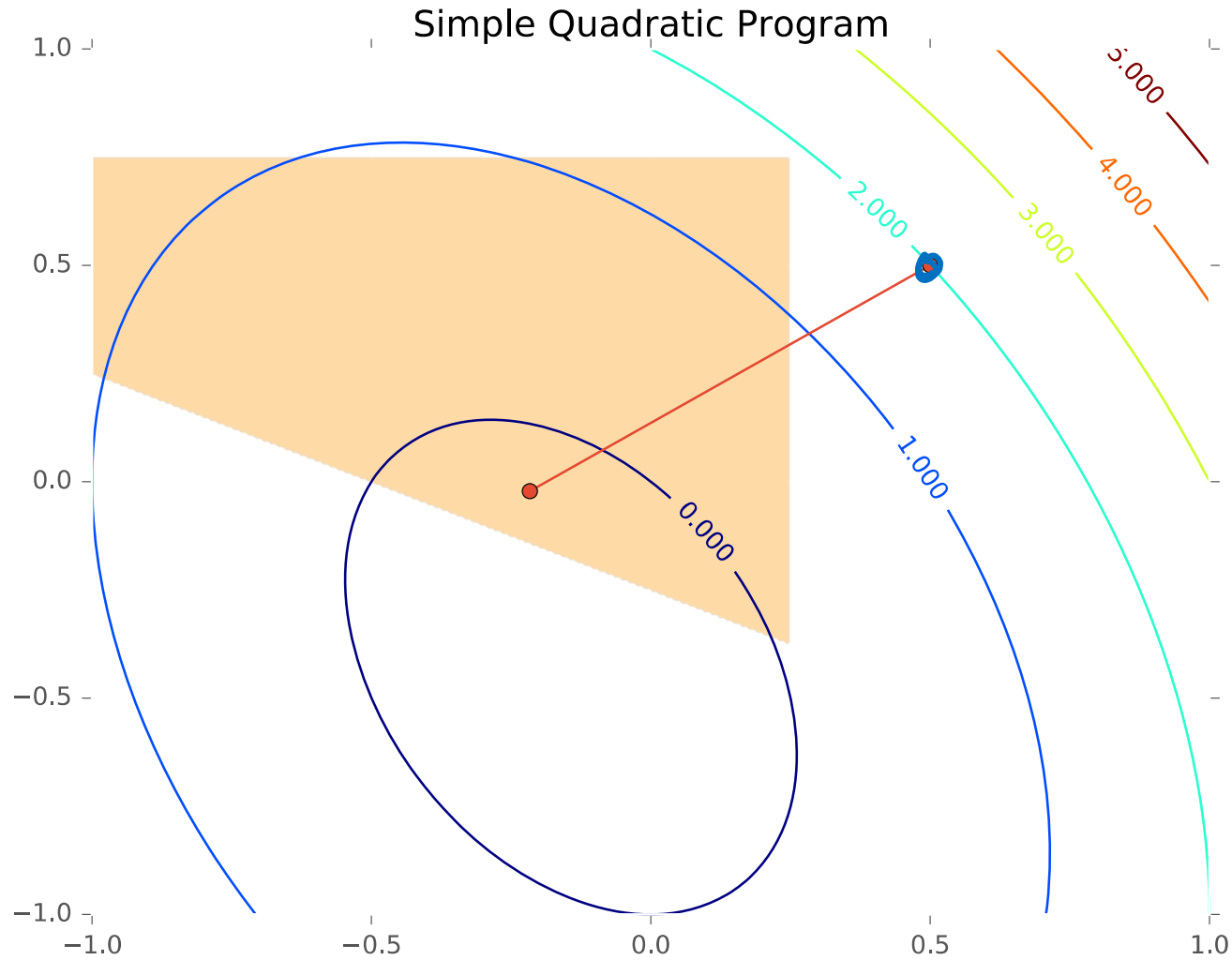– Example: 2D quadratic function with linear constraints

# Quadratic Program

## Simple Quadratic Program



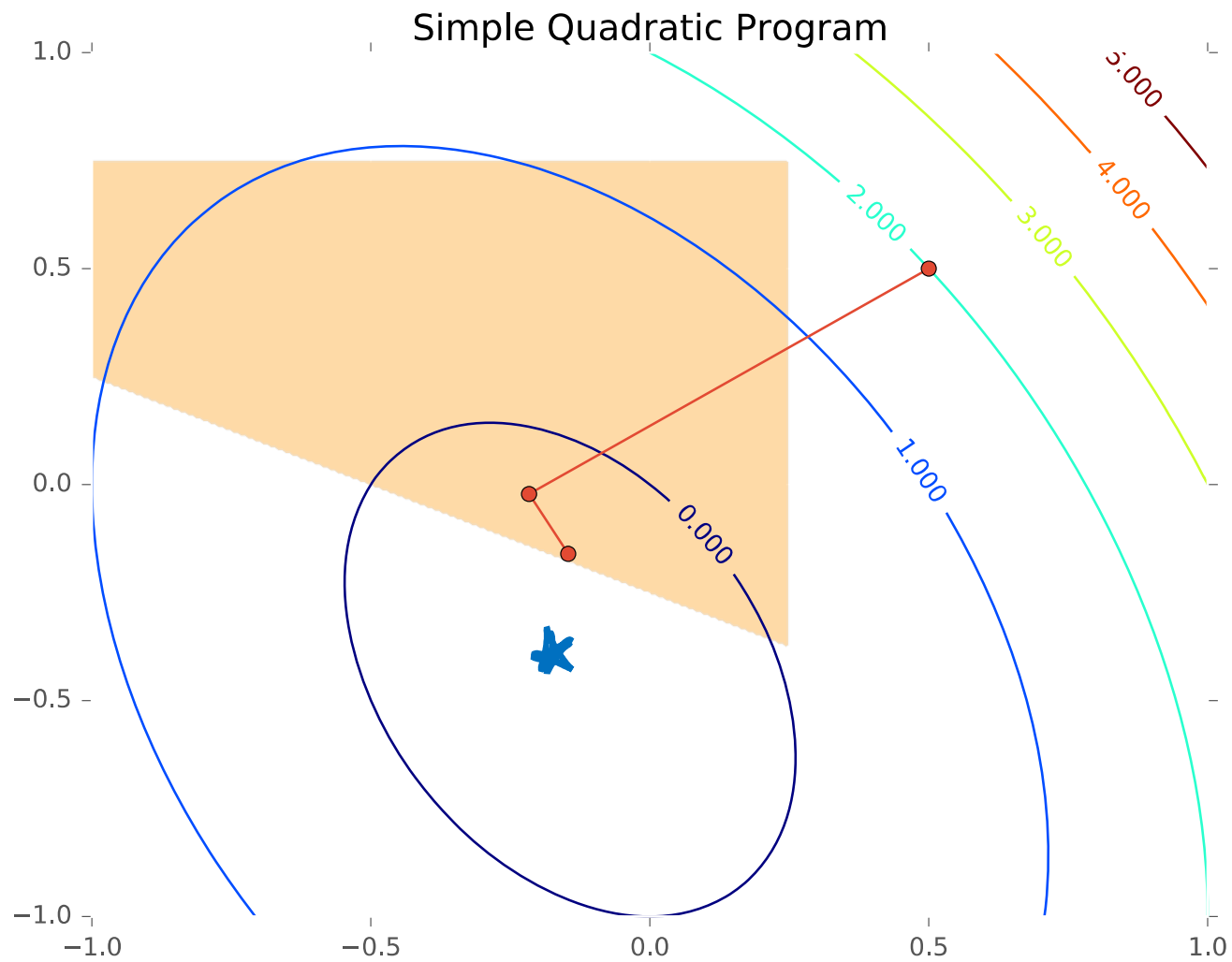$Ax \leq b$

$x^TQx + c^Tx$

11

# Quadratic Program



Simple Quadratic Program

feasible region

12

# Quadratic Program

## Simple Quadratic Program

# Quadratic Program



Simple Quadratic Program

# Quadratic Program

## Simple Quadratic Program

# SUPPORT VECTOR MACHINE (SVM)

# Example: Building Walls
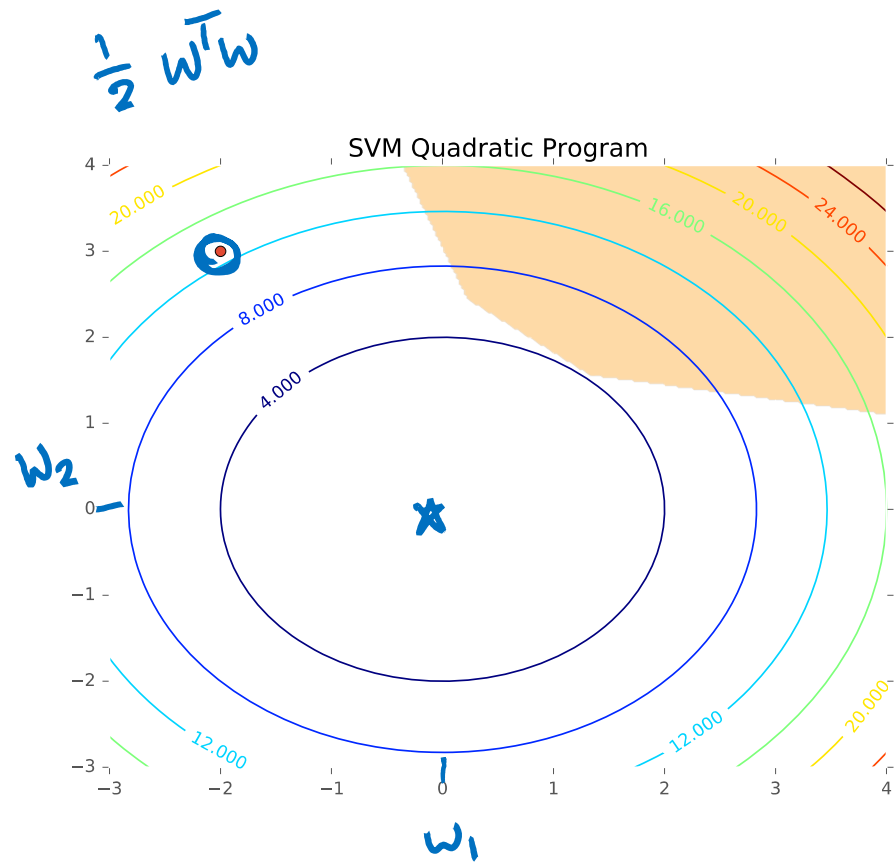
https://www.facebook.com/Mondobloxx/
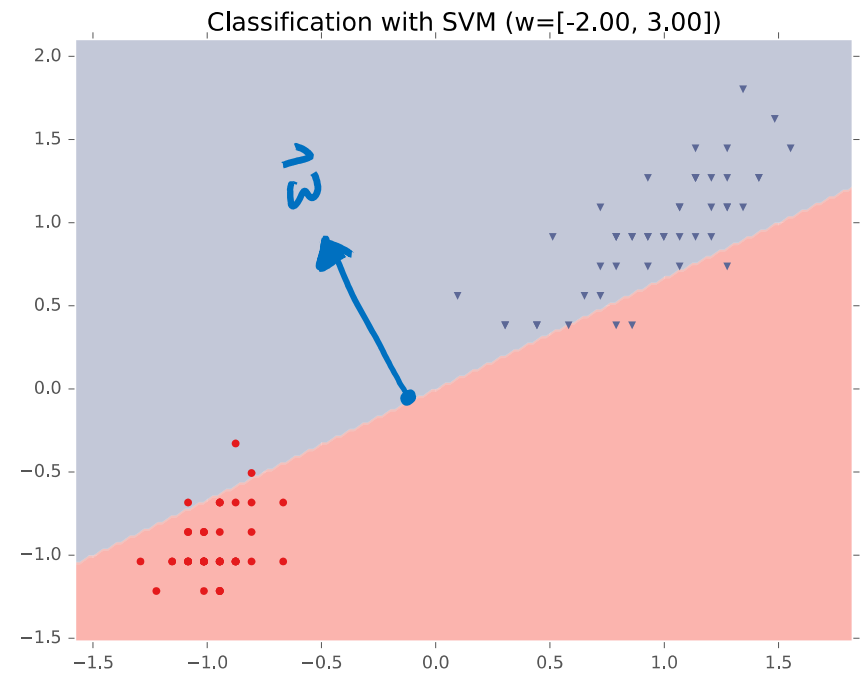
# SVM

*Whiteboard*

– SVM Primal (Linearly Separable Case)

This section borrows ideas from Nina Balcan's SVM lectures at CMU and Patrick Winston's "widest street" SVM lecture at MIT (https://www.youtube.com/watch?v=_PwhiWxHK8o).

# SVM QP

$\frac{1}{2} w^T w$

### SVM Quadratic Program



$w_2$

$w_1$

b is not shown

### Classification with SVM (w=[-2.00, 3.00])



$\vec{w}$

# SVM QP

$$y^{(i)}\left(w^\top x^{(i)} + b\right) \geq 1$$



SVM Quadratic Program



Classification with SVM (w=[0.37, 1.51])

# SVM QP



SVM Quadratic Program

Classification with SVM (w=[0.62, 1.58])

# SVM QP



SVM Quadratic Program

Classification with SVM (w=[1.04, 1.77])

# SVM QP

>1



SVM Quadratic Program

Classification with SVM (w=[1.28, 1.62])

# SVM QP



SVM Quadratic Program

Classification with SVM (w=[1.28, 1.60])

Black Box QP Solver

# Support Vector Machines (SVMs)

**Hard-margin SVM (Primal)**

$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|_2^2$$

$$\text{s.t. } y^{(i)}(\mathbf{w}^T\mathbf{x}^{(i)} + b) \geq 1, \quad \forall i = 1,\dots,N$$

**Hard-margin SVM (Lagrangian Dual)**

$$\max_{\boldsymbol{\alpha}} \sum_{i=1}^{N} \alpha_i - \frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N} \alpha_i\alpha_j y^{(i)}y^{(j)}\mathbf{x}^{(i)}\cdot\mathbf{x}^{(j)}$$

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

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

- Instead of minimizing the primal, we can maximize the dual problem
- For the SVM, these two problems give the same answer (i.e. the minimum of one is the maximum of the other)
- *Definition:* **support vectors** are those points x$^{(i)}$ for which $\alpha^{(i)} \neq 0$

# METHOD OF LAGRANGE MULTIPLIERS

# Method of Lagrange Multipliers

Method of Lagrange Multipliers <span style="color:green">(case w/inequalities)</span>

Goal: $\min \; f(\vec{x}) \; \text{s.t.} \; g(\vec{x}) \leq c$

① Construct Lagrangian

$$L(\vec{x}, \lambda) = f(\vec{x}) - \lambda(g(\vec{x}) - c)$$

② Solve
$$\min_{\vec{x}} \max_{\lambda} L(\vec{x}, \lambda)$$

$$\nabla_x f(x) = \lambda \nabla_g g(x)$$

$$\boxed{\nabla L(\vec{x}, \lambda) = 0} \quad \text{s.t.} \; \lambda \geq 0, \; g(\vec{x}) \leq c$$

Equivalent to solving:
$$\nabla f(\vec{x}) = \lambda \nabla_g(\vec{x}) \quad \text{s.t.} \; \lambda \geq 0, \; g(\vec{x}) \leq c$$

# Method of Lagrange Multipliers

# Method of Lagrange Multipliers

# Method of Lagrange Multipliers

# Method of Lagrange Multipliers

# Method of Lagrange Multipliers

# Method of Lagrange Multipliers



$\nabla f(x,y) = \lambda \nabla g(x,y)$

lagrange multipliers

$f(x,y) = -4$

$f(x,y) = -3$

$f(x,y) = -2$

$f(x,y) = -1$

$f(x,y) = 0$

$x^2 + y^2 = 1$

# Method of Lagrange Multipliers

# SVM DUAL

# Method of Lagrange Multipliers

*Whiteboard*

- – Lagrangian Duality
- – Example: SVM Dual

# Support Vector Machines (SVMs)

Hard-margin SVM (Primal)

$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|_2^2$$

$$\text{s.t. } y^{(i)}(\mathbf{w}^T\mathbf{x}^{(i)} + b) \geq 1, \quad \forall i = 1,\ldots,N$$

Hard-margin SVM (Lagrangian Dual)

$$\max_{\boldsymbol{\alpha}} \sum_{i=1}^{N} \alpha_i - \frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N} \alpha_i\alpha_j y^{(i)}y^{(j)}\mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$$

$$\text{s.t. } \alpha_i \geq 0, \quad \forall i = 1,\ldots,N$$

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

$$\vec{w} = \sum_{i=1}^{N} \alpha_i \, y^{(i)} x^{(i)}$$

- Instead of minimizing the primal, we can maximize the dual problem
- For the SVM, these two problems give the same answer (i.e. the minimum of one is the maximum of the other)
- *Definition*: **support vectors** are those points $x^{(i)}$ for which $\alpha^{(i)} \neq 0$ ← points on the margin

# SVM EXTENSIONS

# Soft-Margin SVM

Hard-margin SVM (Primal)

$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|_2^2$$

$$\text{s.t. } y^{(i)}(\mathbf{w}^T\mathbf{x}^{(i)} + b) \geq 1, \quad \forall i = 1, \ldots, N$$

Soft-margin SVM (Primal)

$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|_2^2 + C\left(\sum_{i=1}^{N} e_i\right)$$

$$\text{s.t. } y^{(i)}(\mathbf{w}^T\mathbf{x}^{(i)} + b) \geq 1 - e_i, \quad \forall i = 1, \ldots, N$$

$$e_i \geq 0, \quad \forall i = 1, \ldots, N$$

- **Question**: If the dataset is not linearly separable, can we still use an SVM?

- **Answer**: Not the hard-margin version. It will never find a feasible solution.

In the soft-margin version, we add "**slack variables**" that **allow some points to violate** the large-margin constraints.

The constant C dictates **how large** we should allow the slack variables to be

# Soft-Margin SVM

**Hard-margin SVM (Primal)**

$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|_2^2$$

$$\text{s.t. } y^{(i)}(\mathbf{w}^T\mathbf{x}^{(i)} + b) \geq 1, \quad \forall i = 1,\ldots,N$$

**Soft-margin SVM (Primal)**

$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|_2^2 + C\left(\sum_{i=1}^{N} e_i\right)$$

$$\text{s.t. } y^{(i)}(\mathbf{w}^T\mathbf{x}^{(i)} + b) \geq 1 - e_i, \quad \forall i = 1,\ldots,N$$

$$e_i \geq 0, \quad \forall i = 1,\ldots,N$$

# Soft-Margin SVM

**Hard-margin SVM (Primal)**

$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|_2^2$$

$$\text{s.t. } y^{(i)}(\mathbf{w}^T\mathbf{x}^{(i)} + b) \geq 1, \quad \forall i = 1,\dots,N$$

**Hard-margin SVM (Lagrangian Dual)**

$$\max_{\boldsymbol{\alpha}} \sum_{i=1}^{N}\alpha_i - \frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N}\alpha_i\alpha_j y^{(i)}y^{(j)}\mathbf{x}^{(i)}\cdot\mathbf{x}^{(j)}$$

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

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

**Soft-margin SVM (Primal)**

$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|_2^2 + C\left(\sum_{i=1}^{N}e_i\right)$$

$$\text{s.t. } y^{(i)}(\mathbf{w}^T\mathbf{x}^{(i)} + b) \geq 1 - e_i, \quad \forall i = 1,\dots,N$$

$$e_i \geq 0, \quad \forall i = 1,\dots,N$$

**Soft-margin SVM (Lagrangian Dual)**

$$\max_{\boldsymbol{\alpha}} \sum_{i=1}^{N}\alpha_i - \frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N}\alpha_i\alpha_j y^{(i)}y^{(j)}\mathbf{x}^{(i)}\cdot\mathbf{x}^{(j)}$$

$$\text{s.t. } 0 \leq \alpha_i \leq C, \quad \forall i = 1,\dots,N$$

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

**We can also work with the dual of the soft-margin SVM**

# Multiclass SVMs

The SVM is **inherently** a **binary** classification method, but can be extended to handle K-class classification in many ways.

1. ***one-vs-rest:***
   - build K binary classifiers
   - train the $k^{th}$ classifier to predict whether an instance has label k or something else
   - predict the class with largest score

2. ***one-vs-one:***
   - build (K choose 2) binary classifiers
   - train one classifier for distinguishing between each pair of labels
   - predict the class with the most "votes" from any given classifier

# Learning Objectives

**Support Vector Machines**

*You should be able to...*

1. Motivate the learning of a decision boundary with large margin
2. Compare the decision boundary learned by SVM with that of Perceptron
3. Distinguish unconstrained and constrained optimization
4. Compare linear and quadratic mathematical programs
5. Derive the hard-margin SVM primal formulation
6. Derive the Lagrangian dual for a hard-margin SVM
7. Describe the mathematical properties of support vectors and provide an intuitive explanation of their role
8. Draw a picture of the weight vector, bias, decision boundary, training examples, support vectors, and margin of an SVM
9. Employ slack variables to obtain the soft-margin SVM
10. Implement an SVM learner using a black-box quadratic programming (QP) solver

# KERNELS

# Kernels: Motivation

Most real-world problems exhibit data that is not linearly separable.

Example: pixel representation for Facial Recognition:



**Q:** When your data is **not linearly separable,** how can you still use a linear classifier?

**A:** Preprocess the data to produce **nonlinear features**

# Kernels: Motivation

- Motivation #1: Inefficient Features
  - Non-linearly separable data requires **high dimensional** representation
  - Might be **prohibitively expensive** to compute or store

- Motivation #2: Memory-based Methods
  - k-Nearest Neighbors (KNN) for facial recognition allows a **distance metric** between images -- no need to worry about linearity restriction at all

# Kernel Methods

Φ

- **Key idea:**
  1. Rewrite the algorithm so that we only work with **dot products** $x^T z$ of feature vectors
  2. Replace the **dot products** $x^T z$ with a **kernel function** $k(x, z)$

- The kernel $k(x,z)$ can be **any** legal definition of a dot product:

$$k(x, z) = \varphi(x)^T \varphi(z) \text{ for any function } \varphi: X \rightarrow \mathbf{R}^D$$

  So we only compute the $\varphi$ dot product **implicitly**

- This **"kernel trick"** can be applied to many algorithms:
  - classification: perceptron, SVM, ...
  - regression: ridge regression, ...
  - clustering: k-means, ...

# SVM: Kernel Trick

## Hard-margin SVM (Primal)

$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|_2^2$$

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

- Suppose we do some feature engineering
- Our feature function is $\phi$
- We apply $\phi$ to each input vector $\mathbf{x}$

$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|_2^2$$

$$\text{s.t. } y^{(i)}(\mathbf{w}^T\phi\left(\mathbf{x}^{(i)}\right) + b) \geq 1, \quad \forall i$$

## Hard-margin SVM (Lagrangian Dual)

$$\max_{\boldsymbol{\alpha}} \sum_{i=1}^{N} \alpha_i - \frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N} \alpha_i\alpha_j y^{(i)}y^{(j)}\mathbf{x}^{(i)}\cdot\mathbf{x}^{(j)}$$

$$\text{s.t. } \alpha_i \geq 0, \quad \forall i = 1,\ldots,N$$

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

$$\max_{\boldsymbol{\alpha}} \sum_{i=1}^{N} \alpha_i - \frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N} \alpha_i\alpha_j y^{(i)}y^{(j)}\phi\left(\mathbf{x}^{(i)}\right)\cdot\phi\left(\mathbf{x}^{(j)}\right)$$

$$\text{s.t. } \alpha_i \geq 0, \quad \forall i = 1,\ldots,N$$

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

# SVM: Kernel Trick

Hard-margin SVM (Lagrangian Dual)

$$\max_{\boldsymbol{\alpha}} \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j y^{(i)} y^{(j)} \boxed{\phi\left(\mathbf{x}^{(i)}\right) \cdot \phi\left(\mathbf{x}^{(j)}\right)}$$

$$\text{s.t. } \alpha_i \geq 0, \quad \forall i = 1, \ldots, N$$

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

We could replace the dot product of the two feature vectors in the transformed space with a function k(x,z) where $k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \phi\left(\mathbf{x}^{(i)}\right) \cdot \phi\left(\mathbf{x}^{(j)}\right)$

# SVM: Kernel Trick

Hard-margin SVM (Lagrangian Dual)

$$\max_{\boldsymbol{\alpha}} \sum_{i=1}^{N} \alpha_i - \frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N} \alpha_i\alpha_j y^{(i)}y^{(j)}\boxed{k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})}$$

$$\text{s.t. } \alpha_i \geq 0, \quad \forall i = 1,\ldots,N$$

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

We could replace the dot product of the two feature vectors in the transformed space with a function k(x,z) where $k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \phi\left(\mathbf{x}^{(i)}\right) \cdot \phi\left(\mathbf{x}^{(j)}\right)$

# Kernel Methods

- **Key idea:**
  1. Rewrite the algorithm so that we only work with **dot products** $x^Tz$ of feature vectors
  2. Replace the **dot products** $x^Tz$ with a **kernel function** $k(x, z)$

- The kernel $k(x,z)$ can be **any** legal definition of a dot product:

  $$k(x, z) = \varphi(x)^T\varphi(z) \text{ for any function } \varphi: X \rightarrow \mathbf{R}^D$$

  So we only compute the $\varphi$ dot product **implicitly**

- This **"kernel trick"** can be applied to many algorithms:
  – classification: perceptron, SVM, …
  – regression: ridge regression, …
  – clustering: k-means, …

# Kernel Methods

**Q:** These are just non-linear features, right?

**A:** Yes, but…

**Q:** Can't we just compute the feature transformation $\varphi$ explicitly?

**A:** That depends…

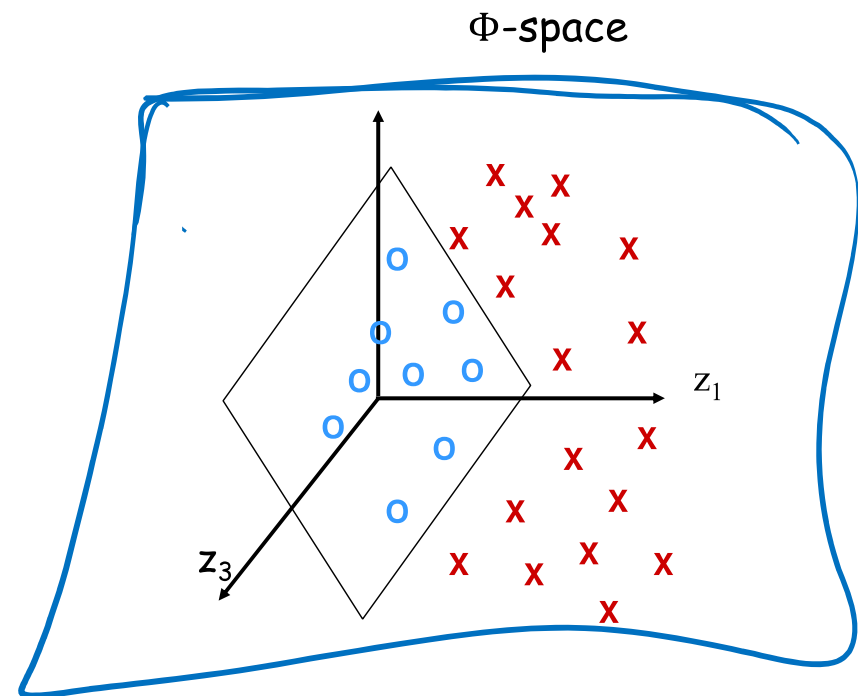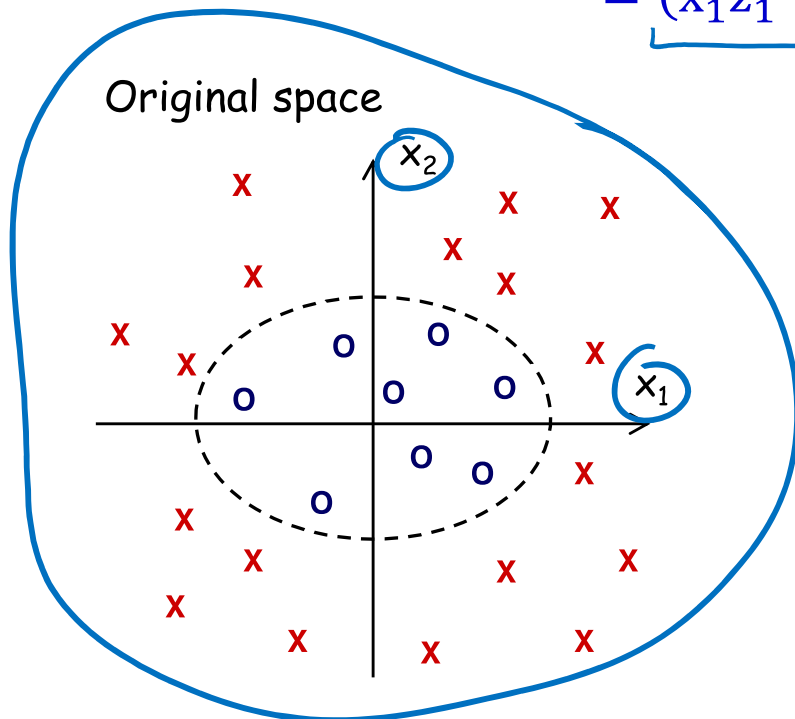**Q:** So, why all the hype about the kernel trick?

**A:** Because the **explicit features** might either be **prohibitively expensive** to compute or **infinite length** vectors

# Example: Polynomial Kernel

For $n=2$, $d=2$, the kernel $K(x, z) = (x \cdot z)^d$ corresponds to

$$\phi : R^2 \rightarrow R^3, (x_1, x_2) \rightarrow \Phi(x) = (x_1^2, x_2^2, \sqrt{2}x_1 x_2)$$

$$\phi(x) \cdot \phi(z) = (x_1^2, x_2^2, \sqrt{2}x_1 x_2) \cdot (z_1^2, z_2^2, \sqrt{2}z_1 z_2)$$
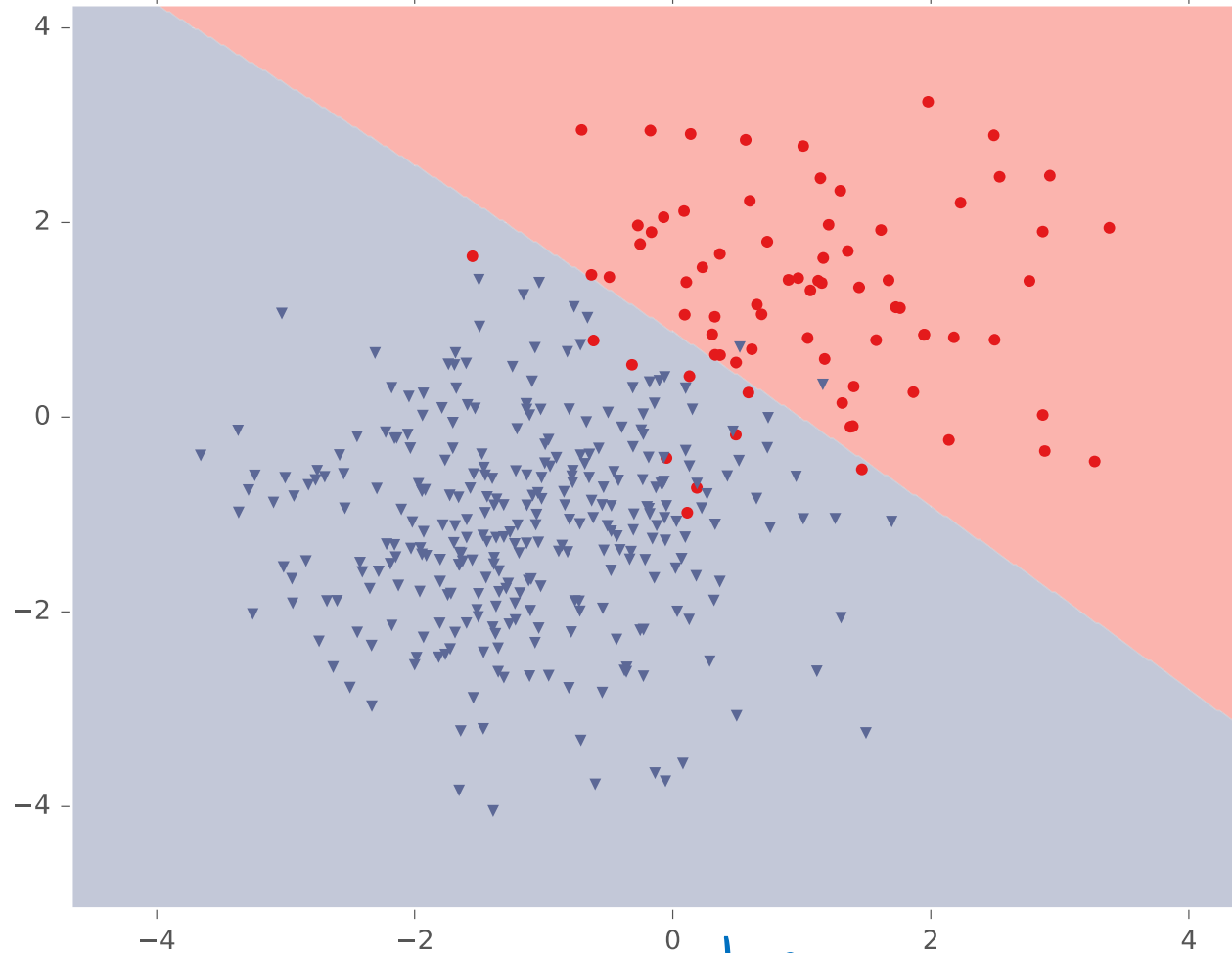
$$= (x_1 z_1 + x_2 z_2)^2 = (x \cdot z)^2 = K(x, z)$$



Original space

$\Phi$-space

# Kernel Examples

| Name | Kernel Function (implicit dot product) | Feature Space (explicit dot product) |
|---|---|---|
| Linear | $K(\mathbf{x}, \mathbf{z}) = \mathbf{x}^T \mathbf{z}$ | Same as original input space |
| Polynomial (v1) | $K(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^d$ | All polynomials **of** degree d |
| Polynomial (v2) | $K(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^d$ | All polynomials **up to** degree d |
| Gaussian | $K(\mathbf{x}, \mathbf{z}) = \exp(-\frac{\|\mathbf{x} - \mathbf{z}\|_2^2}{2\sigma^2})$ | Infinite dimensional space |
| Hyperbolic Tangent (Sigmoid) Kernel | $K(\mathbf{x}, \mathbf{z}) = \tanh(\alpha \mathbf{x}^T \mathbf{z} + c)$ | (With SVM, this is equivalent to a 2-layer neural network) |

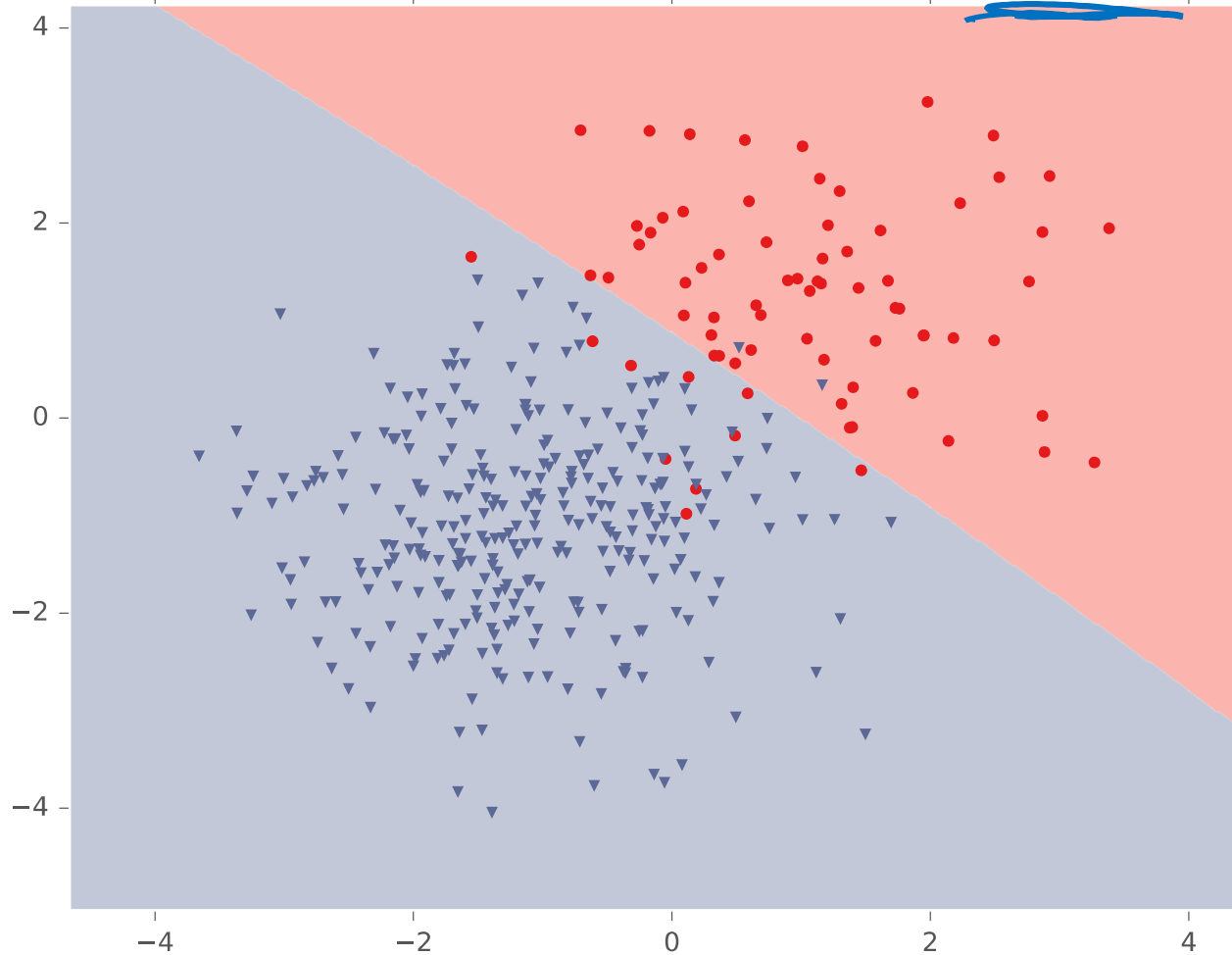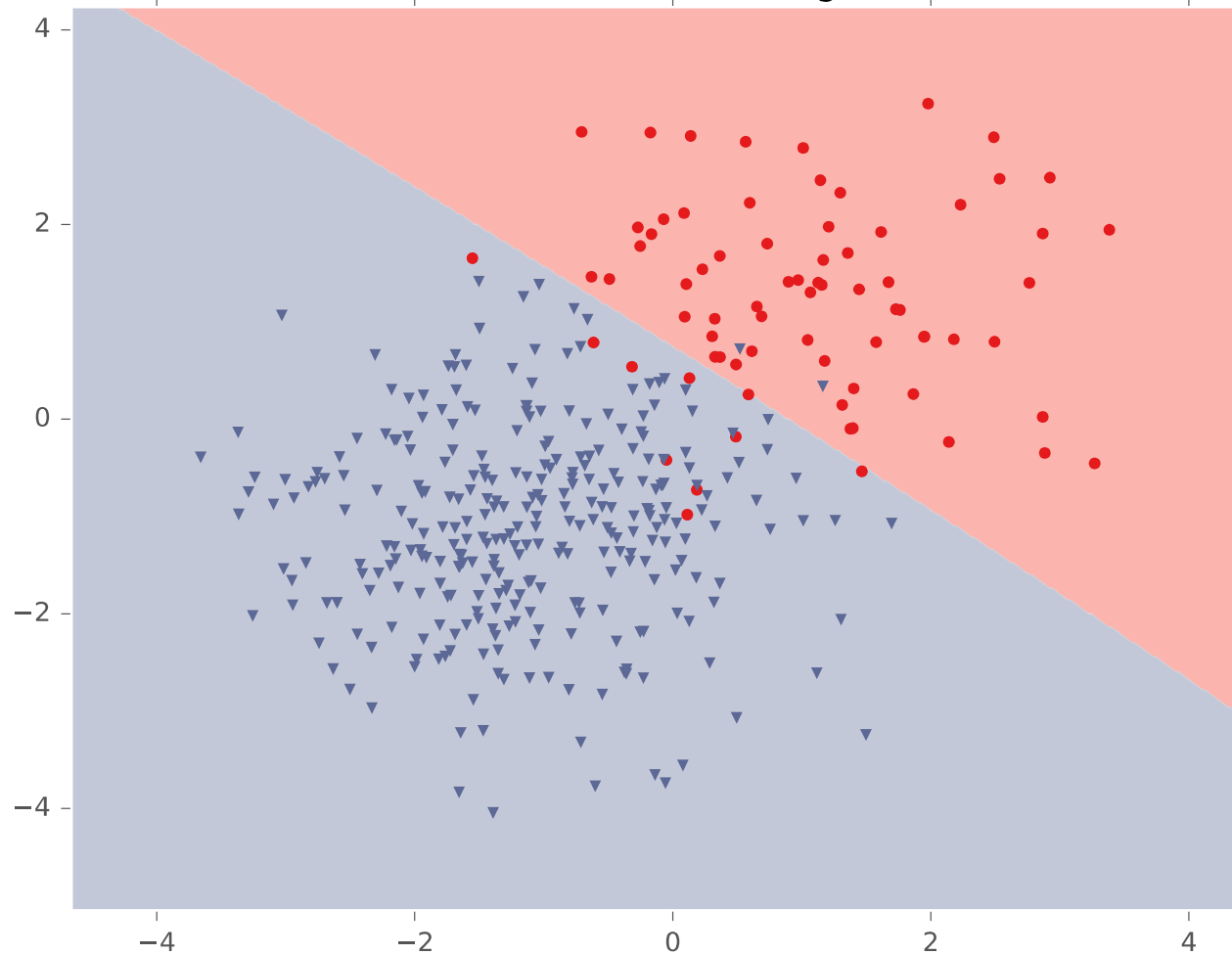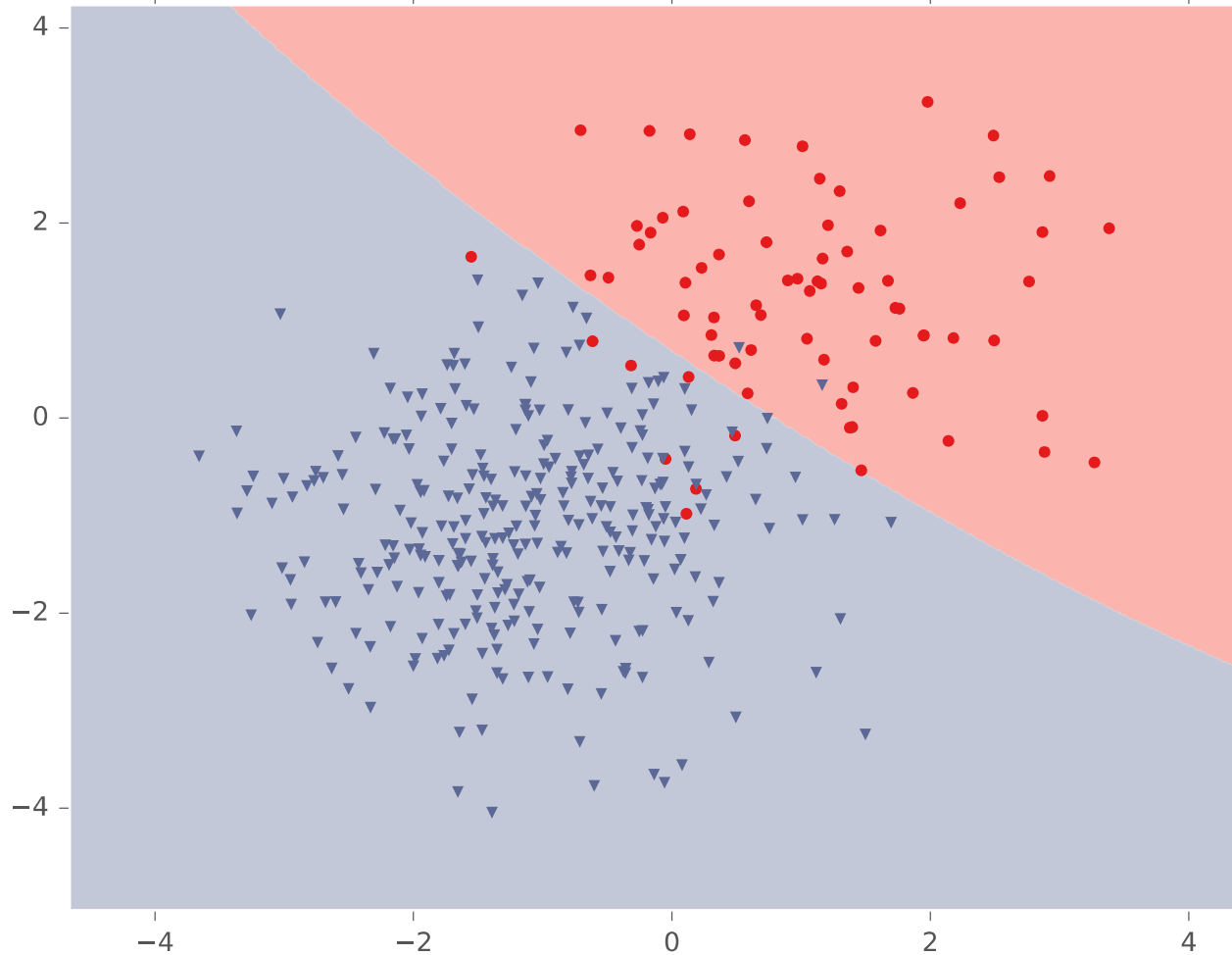# RBF Kernel Example

Classification with SVM (kernel=rbf, gamma=0.010000)



**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$

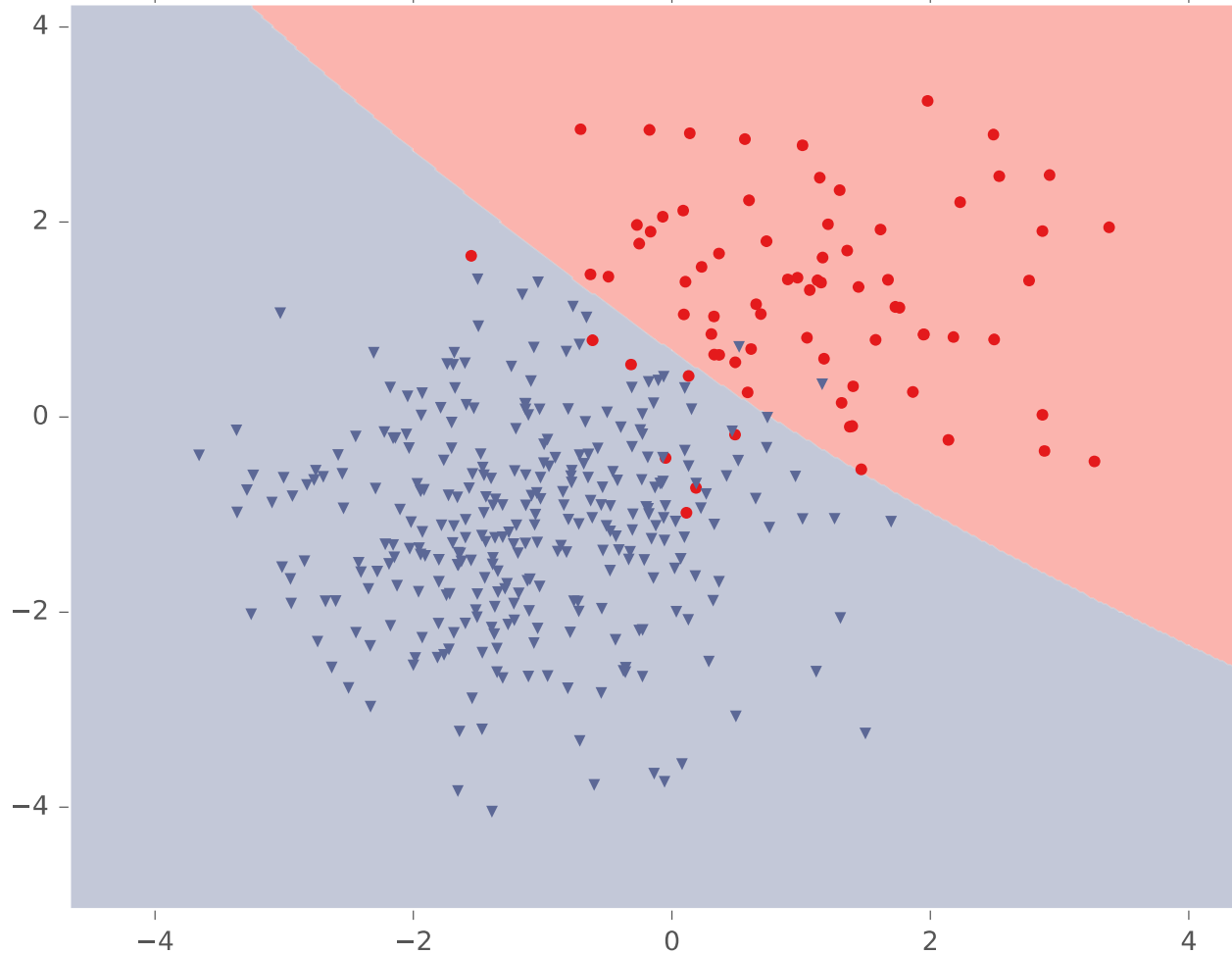hyp.

67

# RBF Kernel Example



Classification with SVM (kernel=rbf, gamma=0.010000)

**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$

68

# RBF Kernel Example
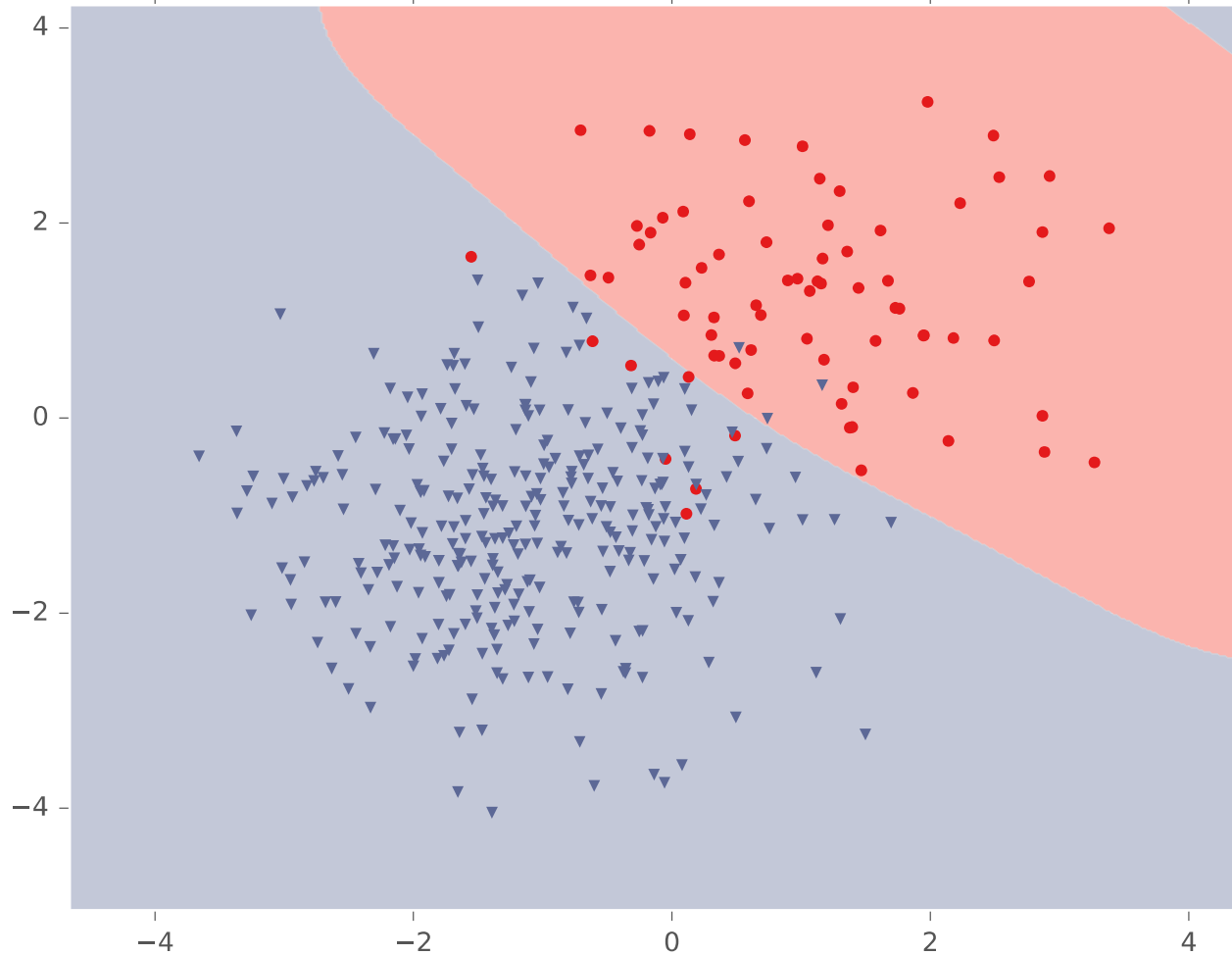
Classification with SVM (kernel=rbf, gamma=0.020000)



**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$

# RBF Kernel Example
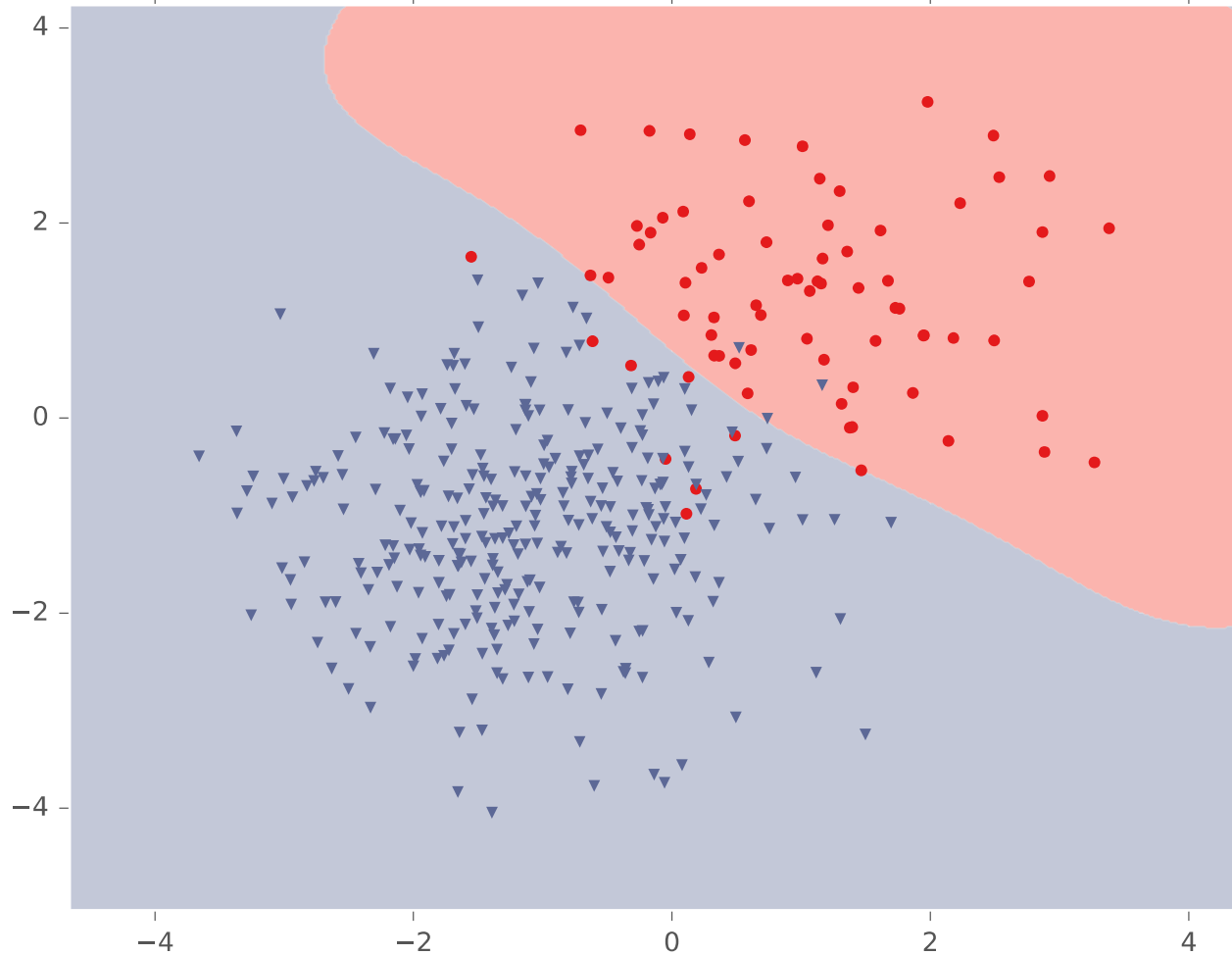
Classification with SVM (kernel=rbf, gamma=0.040000)



**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$

# RBF Kernel Example

Classification with SVM (kernel=rbf, gamma=0.080000)



**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$

# RBF Kernel Example



Classification with SVM (kernel=rbf, gamma=0.160000)

**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$
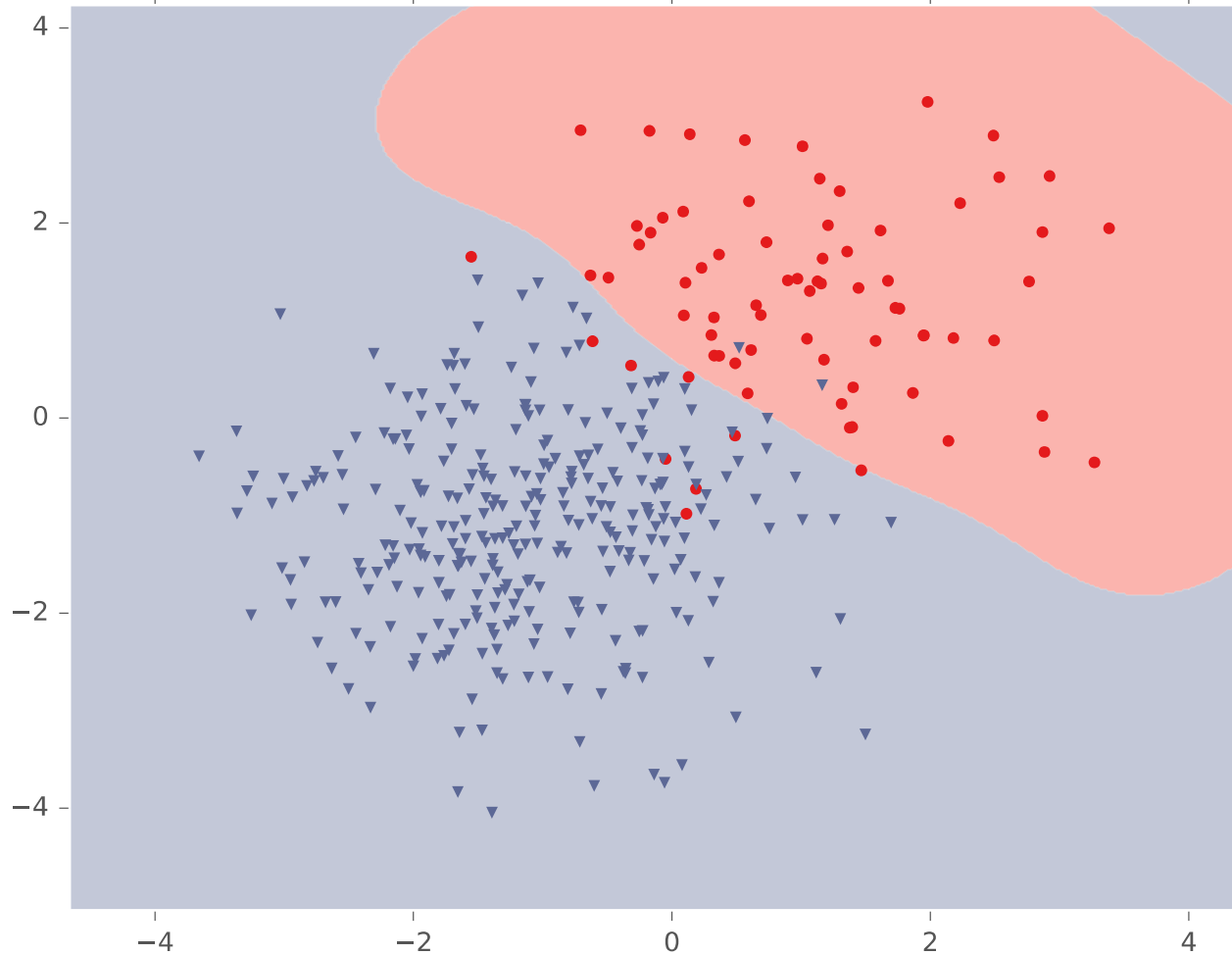
# RBF Kernel Example

## Classification with SVM (kernel=rbf, gamma=0.320000)



**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$
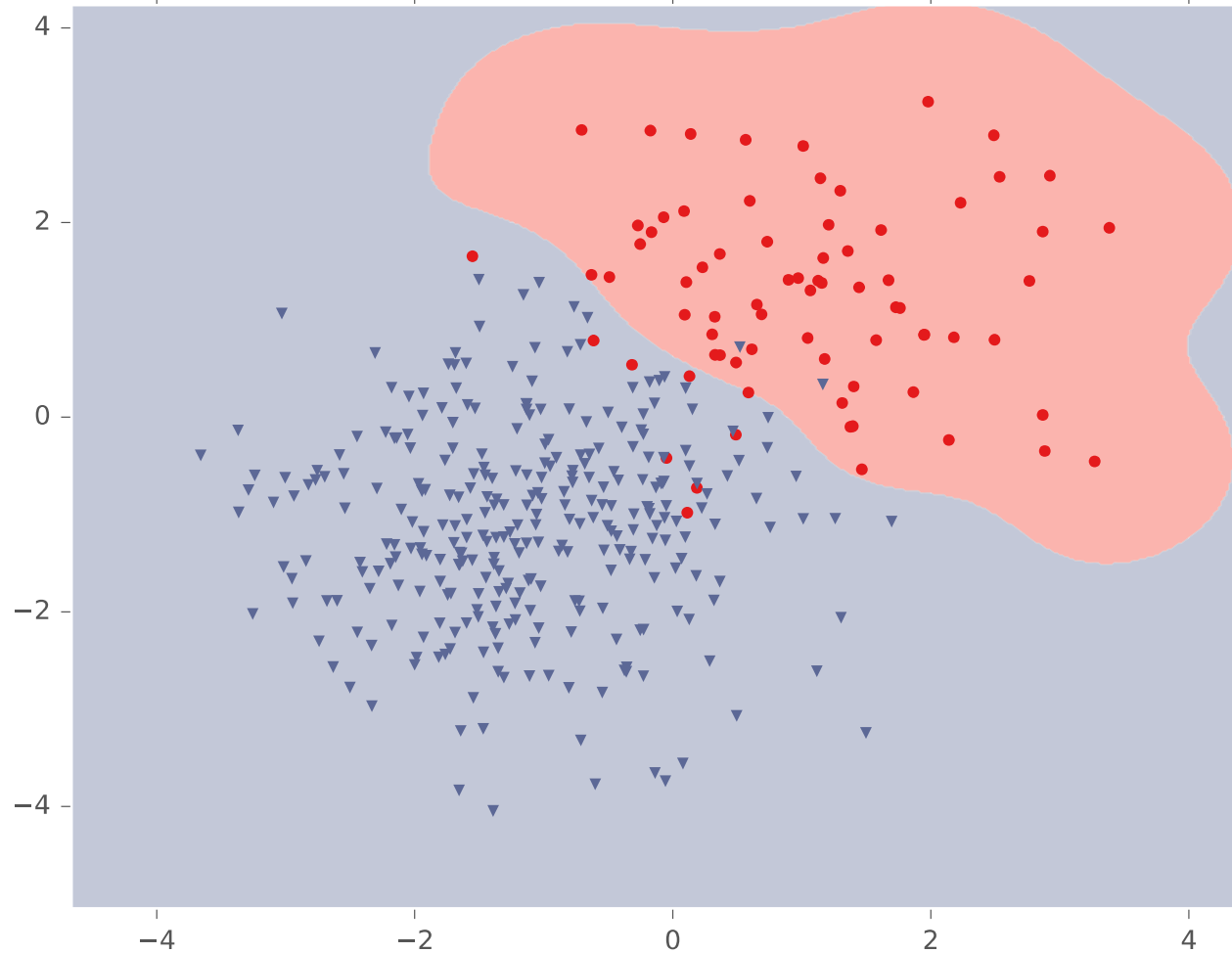
# RBF Kernel Example

Classification with SVM (kernel=rbf, gamma=0.640000)



**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$
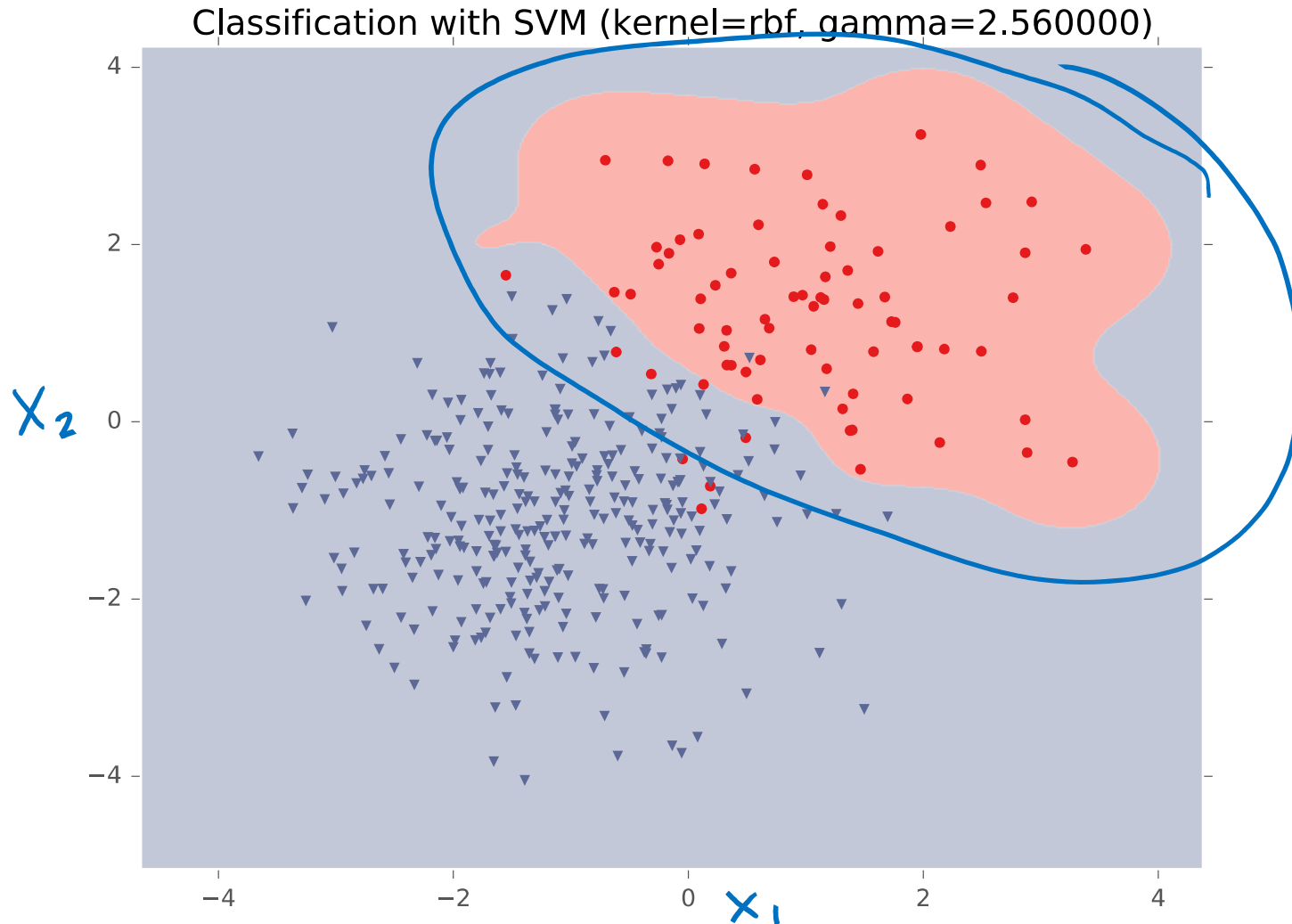
# RBF Kernel Example

Classification with SVM (kernel=rbf, gamma=1.280000)



**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$

# RBF Kernel Example

## Classification with SVM (kernel=rbf, gamma=2.560000)



$X_2$

$X_1$

**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$

# RBF Kernel Example

Classification with SVM (kernel=rbf, gamma=5.120000)



**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma ||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$
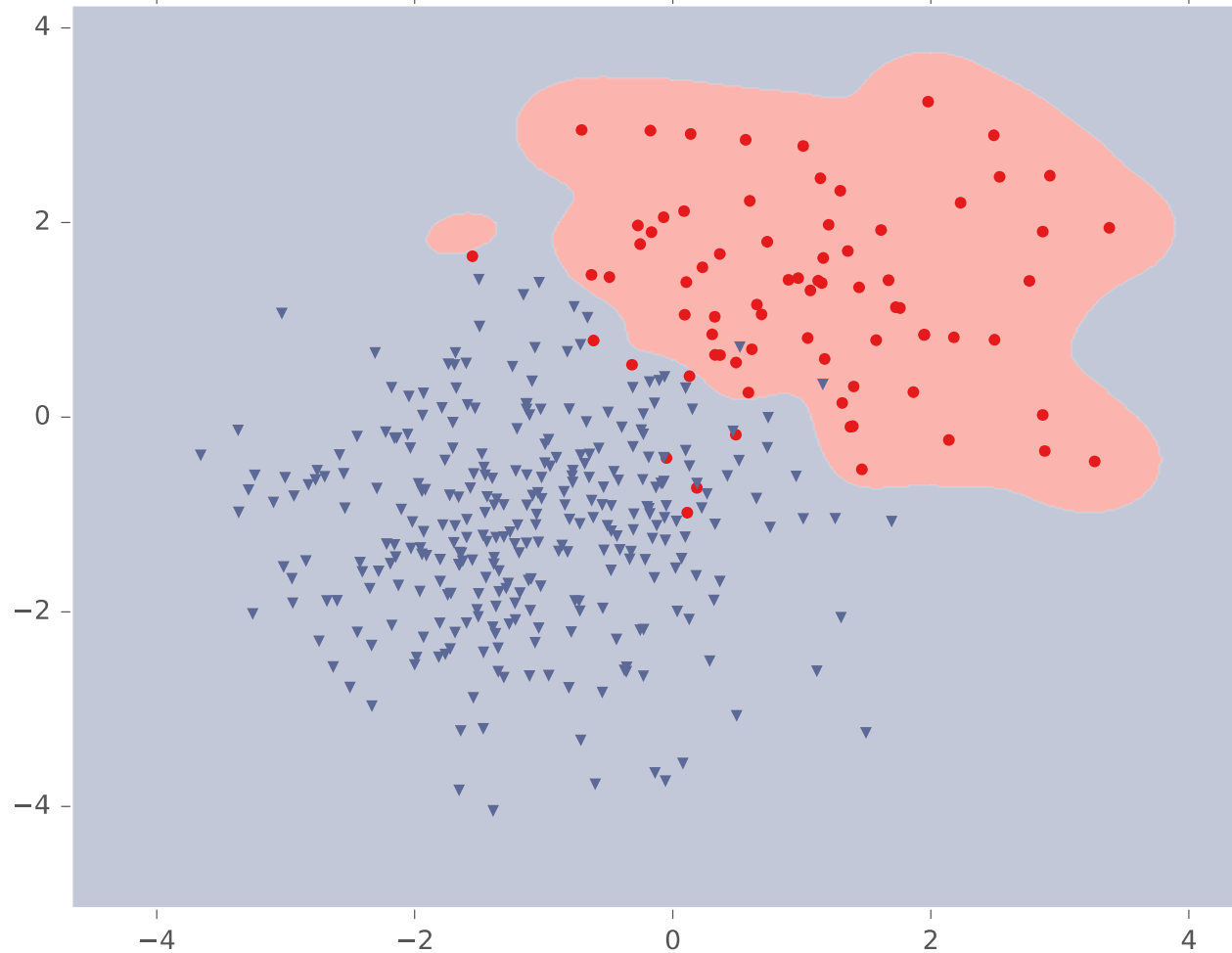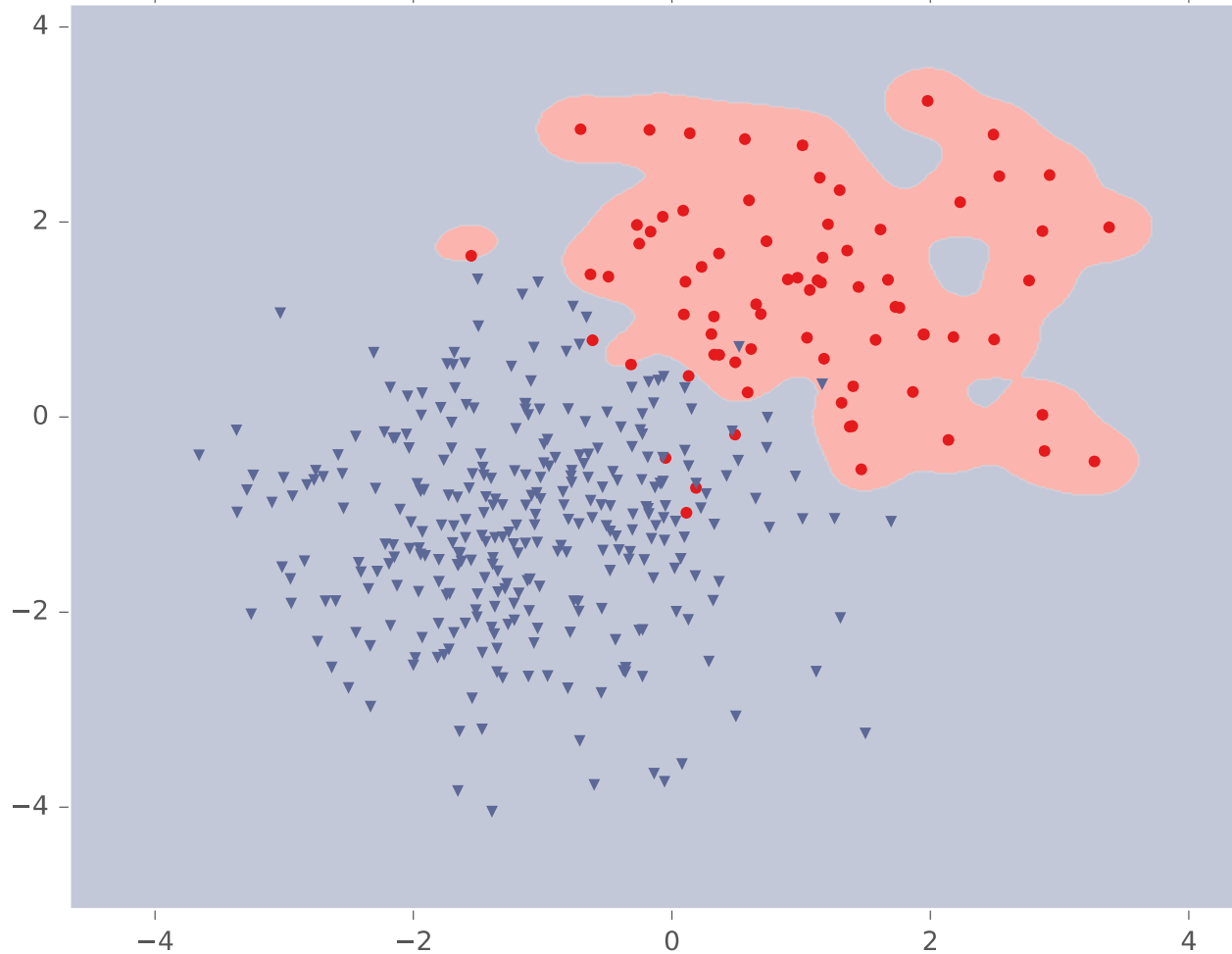
# RBF Kernel Example

Classification with SVM (kernel=rbf, gamma=10.000000)



**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$

# RBF Kernel Example

## KNN vs. SVM
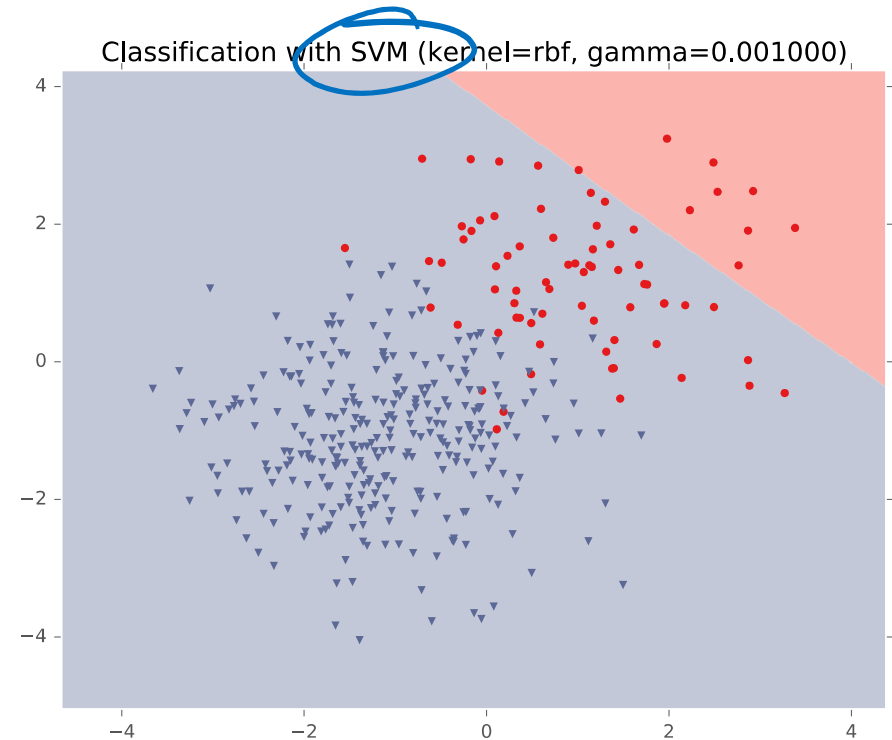
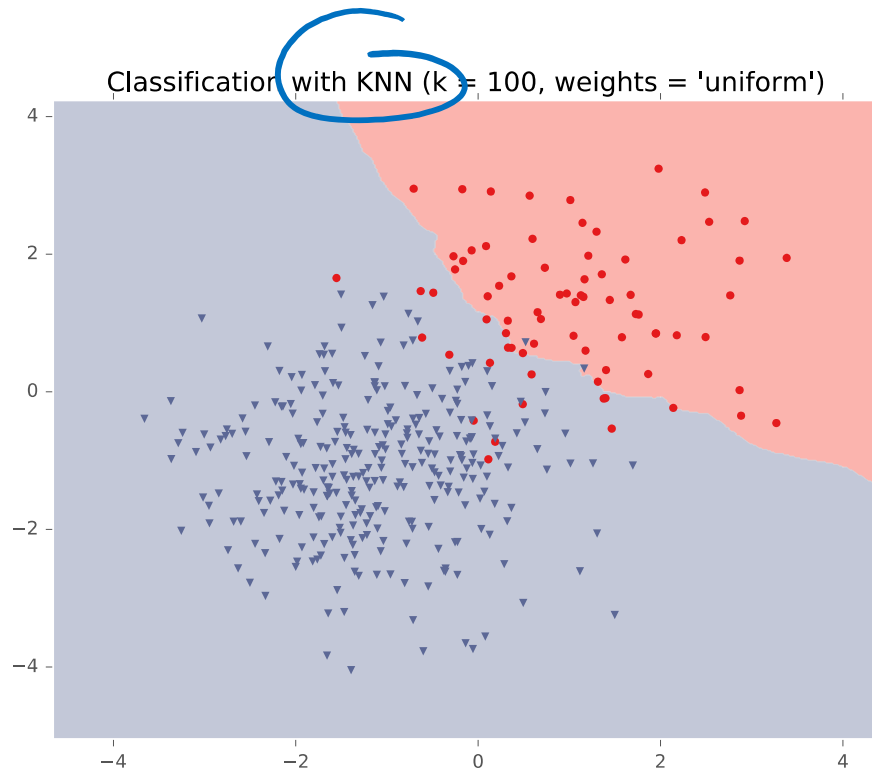Classification with KNN (k = 100, weights = 'uniform')

Classification with SVM (kernel=rbf, gamma=0.001000)

**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$

79

# RBF Kernel Example

## KNN vs. SVM



Classification with KNN (k = 16, weights = 'uniform')

Classification with SVM (kernel=rbf, gamma=0.040000)

**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$
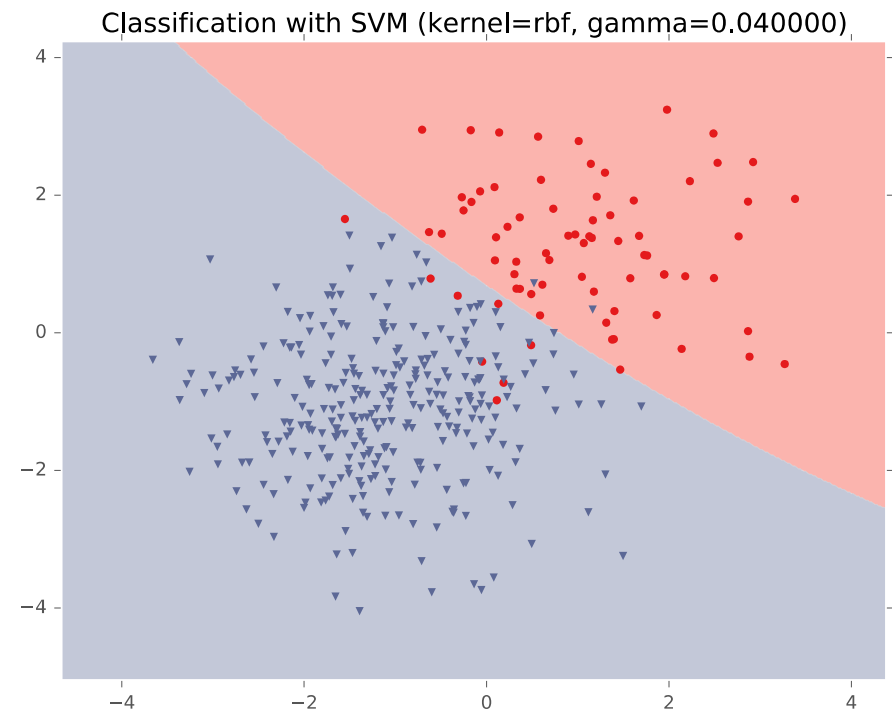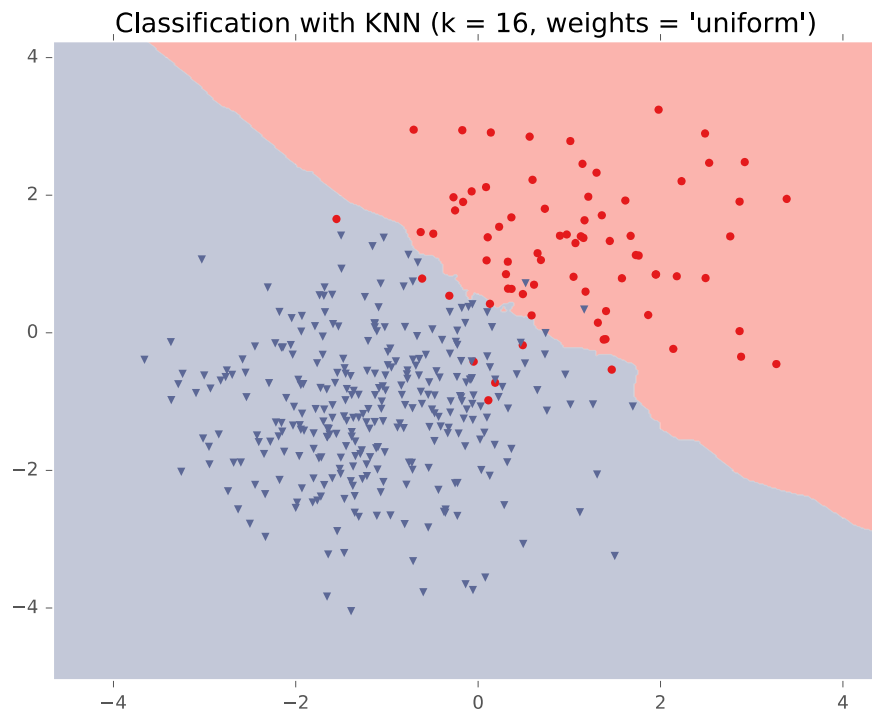
# RBF Kernel Example

## KNN vs. SVM



Classification with KNN (k = 4, weights = 'uniform')
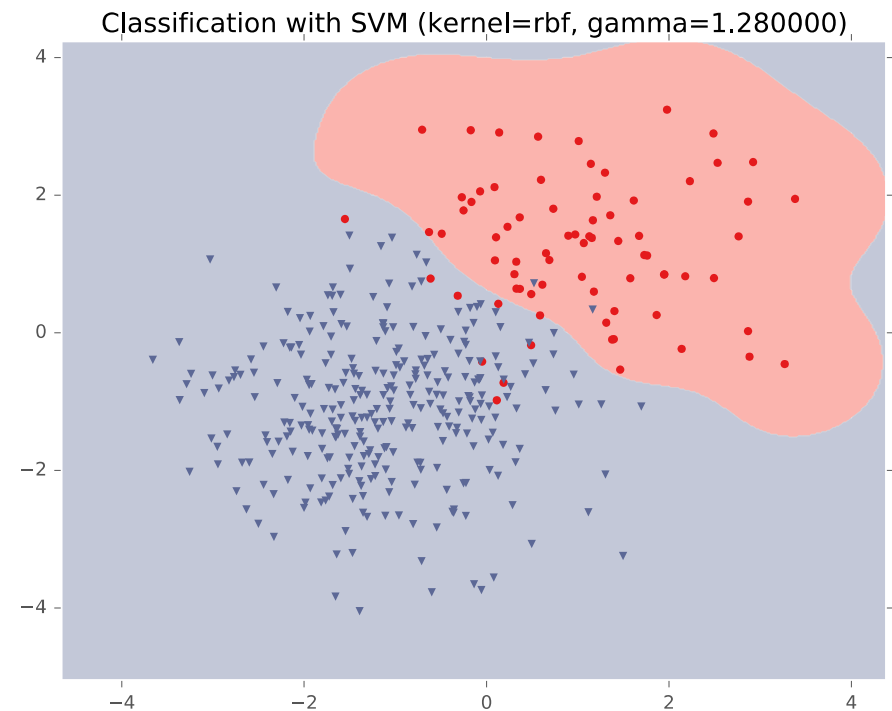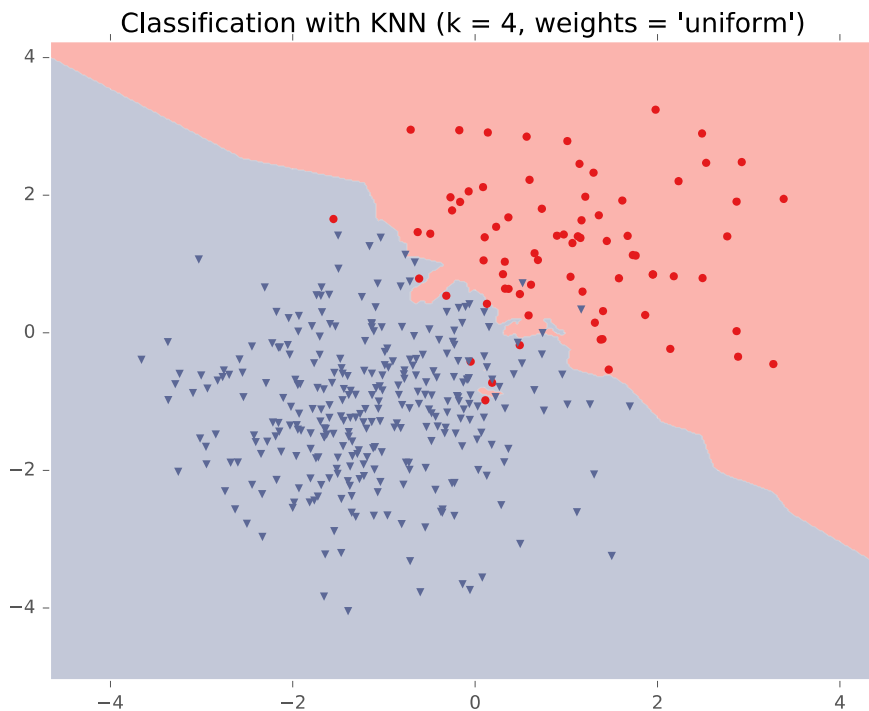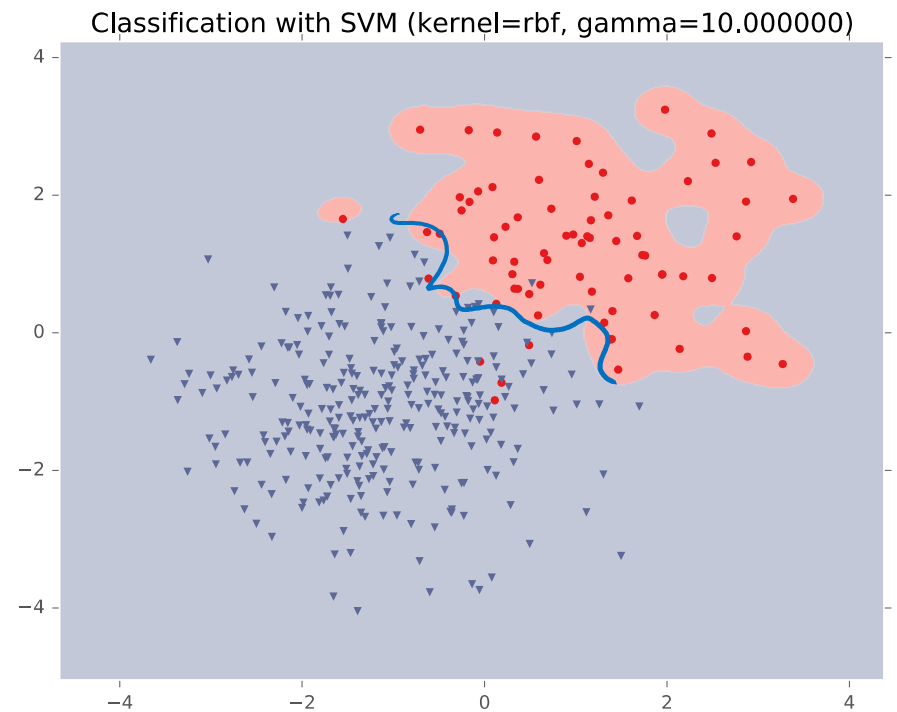
Classification with SVM (kernel=rbf, gamma=1.280000)

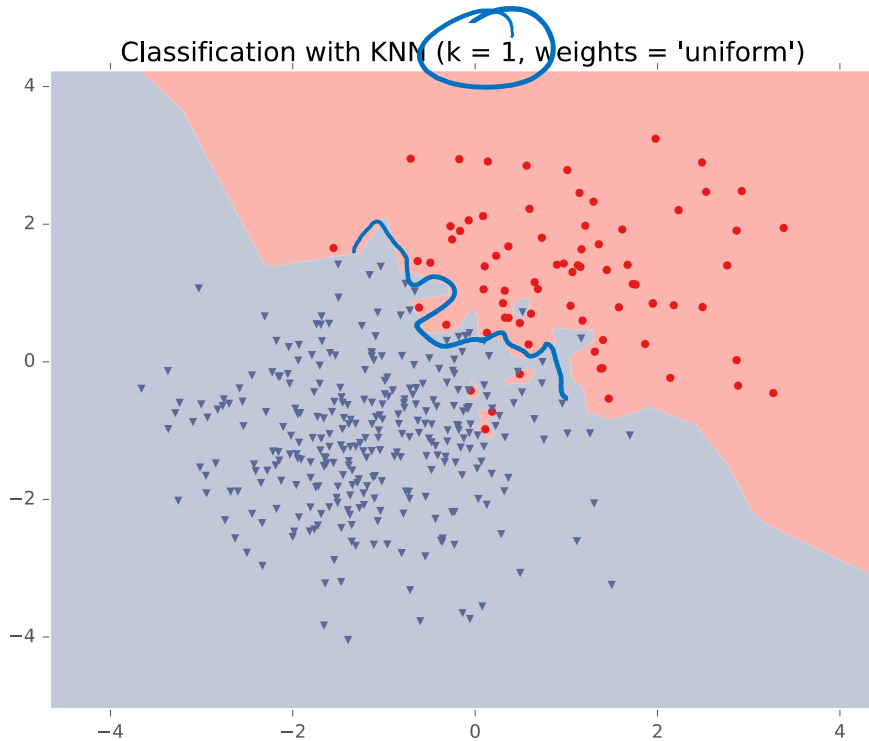**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$

# RBF Kernel Example

**KNN vs. SVM**

Classification with KNN (k = 1, weights = 'uniform')

Classification with SVM (kernel=rbf, gamma=10.000000)

**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$

82

Φ

# Kernel Methods

- **Key idea:**
  1. Rewrite the algorithm so that we only work with **dot products** $x^Tz$ of feature vectors
  2. Replace the **dot products** $x^Tz$ with a **kernel function** $k(x, z)$

- The kernel $k(x,z)$ can be **any** legal definition of a dot product:

$$k(x, z) = \varphi(x)^T\varphi(z) \text{ for any function } \varphi: X \rightarrow \mathbf{R}^D$$

So we only compute the $\varphi$ dot product **implicitly**

- This **"kernel trick"** can be applied to many algorithms:
  - classification: perceptron, SVM, ...
  - regression: ridge regression, ...
  - clustering: k-means, ...

# SVM + Kernels: Takeaways

- Maximizing the margin of a linear separator is a **good training criteria**

- Support Vector Machines (SVMs) learn a **max-margin linear classifier**

- The SVM optimization problem can be solved with **black-box Quadratic Programming (QP) solvers**

- Learned decision boundary is defined by its **support vectors**

- Kernel methods allow us to work in a transformed feature space **without explicitly representing that space**

- The **kernel-trick** can be applied to **SVMs**, as well as many other algorithms

# Learning Objectives

## Kernels

*You should be able to…*

1. Employ the kernel trick in common learning algorithms
2. Explain why the use of a kernel produces only an implicit representation of the transformed feature space
3. Use the "kernel trick" to obtain a computational complexity advantage over explicit feature transformation
4. Sketch the decision boundaries of a linear classifier with an RBF kernel