# Announcements

## Assignments

- HW8: due Thu, 12/3, 11:59 pm
- HW9
  - Out Friday
  - Due Wed, 12/9, 11:59 pm
  - The two slip days are free (last possible submission Fri, 12/11, 11:59 pm)

## Final Exam

- Mon, 12/14
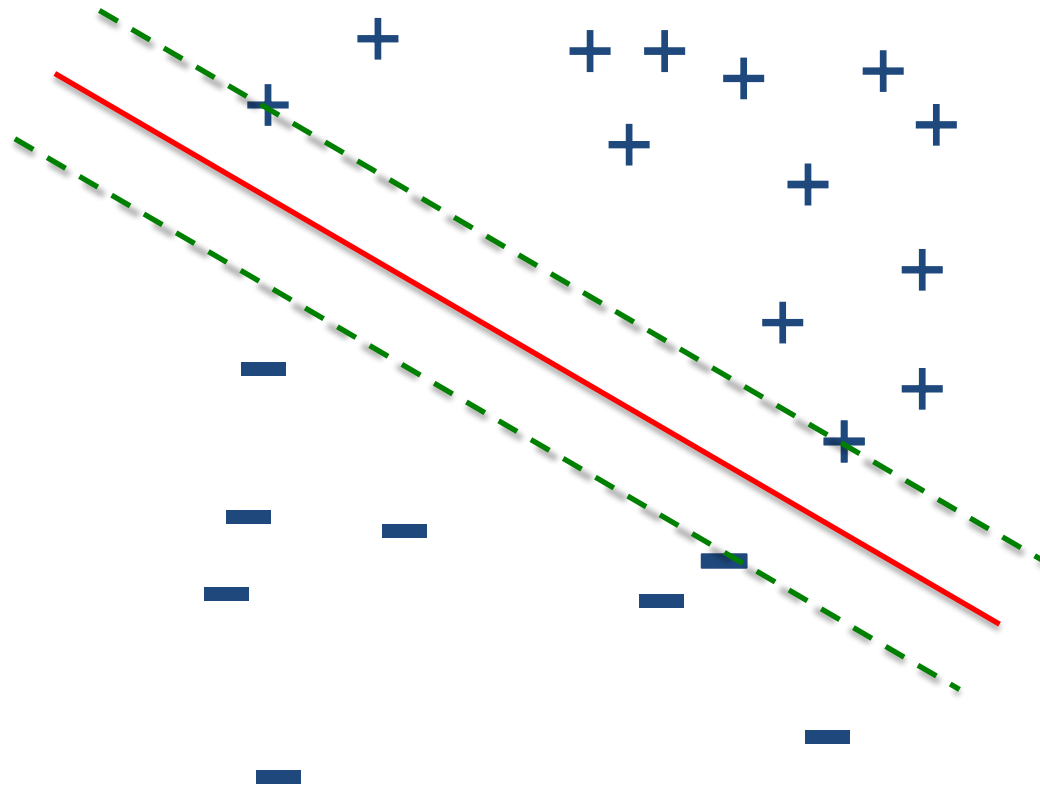- Stay tuned to Piazza for more details

# Introduction to Machine Learning

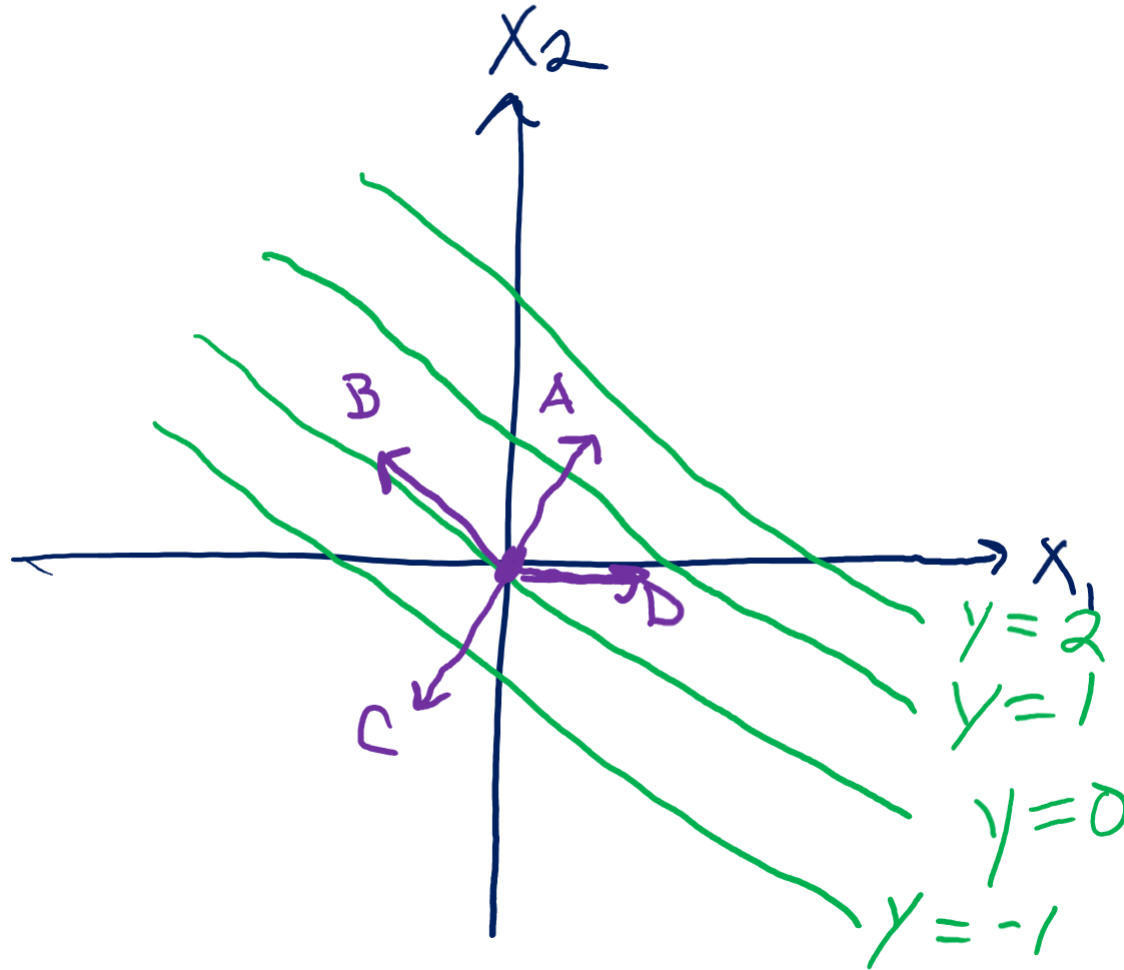# Support Vector Machines

Instructor: Pat Virtue

# Support Vector Machines

Find linear separator with maximum margin
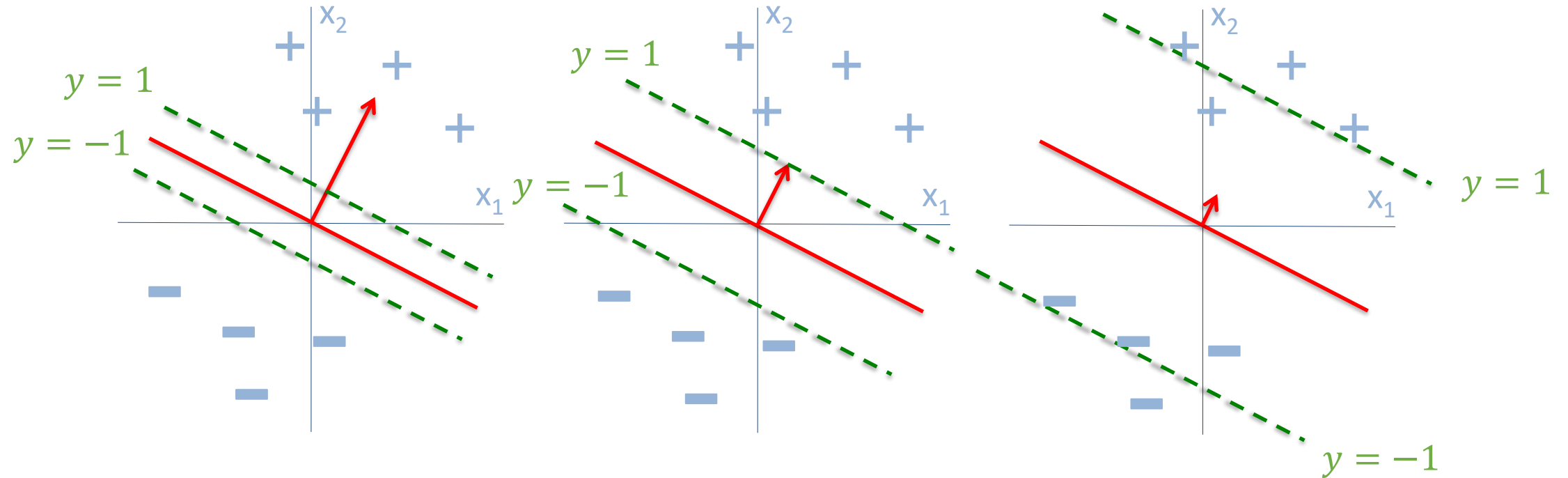
# Previous Piazza Poll

As the magnitude of w increases, will the distance between the contour lines of $y = \boldsymbol{w}^T \boldsymbol{x} + b$ increase or decrease?

# Support Vector Machines

Find linear separator with maximum margin

# Linear Separability

**Data**

$$\mathcal{D} = \left\{ \boldsymbol{x}^{(i)}, y^{(i)} \right\}_{i=1}^{N} \quad \boldsymbol{x} \in \mathbb{R}^M, \quad y \in \{-1, +1\}$$

**Linearly separable iff:**

$$\exists \, \boldsymbol{w}, b \quad s.t. \quad \boldsymbol{w}^T \boldsymbol{x}^{(i)} + b > 0 \quad \text{if} \quad y^{(i)} = +1 \quad \text{and}$$

$$\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b < 0 \quad \text{if} \quad y^{(i)} = -1$$

# Linear Separability

Data

$$\mathcal{D} = \{\boldsymbol{x}^{(i)}, y^{(i)}\}_{i=1}^{N} \quad \boldsymbol{x} \in \mathbb{R}^{M}, \quad y \in \{-1, +1\}$$

Linearly separable iff:

$$\exists\, \boldsymbol{w}, b \quad s.t. \quad \boldsymbol{w}^T \boldsymbol{x}^{(i)} + b > 0 \quad \text{if} \quad y^{(i)} = +1 \quad \text{and}$$

$$\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b < 0 \quad \text{if} \quad y^{(i)} = -1$$

$$\Leftrightarrow \exists\, \boldsymbol{w}, b \quad s.t. \quad y^{(i)}\left(\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b\right) > 0$$

$$\Leftrightarrow \exists\, \boldsymbol{w}, b, c \quad s.t. \quad y^{(i)}\left(\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b\right) \geq c \quad \text{and} \quad c > 0$$

# Piazza Poll 1

Are these two statements equivalent?

$\exists \, \boldsymbol{w}, b, c \quad s.t. \quad y^{(i)}\left(\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b\right) \geq \, c \quad$ and $\quad c > 0$

$\exists \, \boldsymbol{w}, b \qquad s.t. \quad y^{(i)}\left(\boldsymbol{w}^T \boldsymbol{x}^{(i)} + b\right) \geq 1$

# Linear Separability

Data

$$\mathcal{D} = \{x^{(i)}, y^{(i)}\}_{i=1}^{N} \quad x \in \mathbb{R}^M, \quad y \in \{-1, +1\}$$

Linearly separable iff:

$$\exists \, w, b \quad s.t. \quad w^T x^{(i)} + b > 0 \quad \text{if} \quad y^{(i)} = +1 \quad \text{and}$$

$$w^T x^{(i)} + b < 0 \quad \text{if} \quad y^{(i)} = -1$$

$$\Leftrightarrow \exists \, w, b \quad s.t. \quad y^{(i)}\left(w^T x^{(i)} + b\right) > 0$$

$$\Leftrightarrow \exists \, w, b, c \quad s.t. \quad y^{(i)}\left(w^T x^{(i)} + b\right) \geq c \quad \text{and} \quad c > 0$$

$$\Leftrightarrow \exists \, w, b \quad s.t. \quad y^{(i)}\left(w^T x^{(i)} + b\right) \geq 1$$
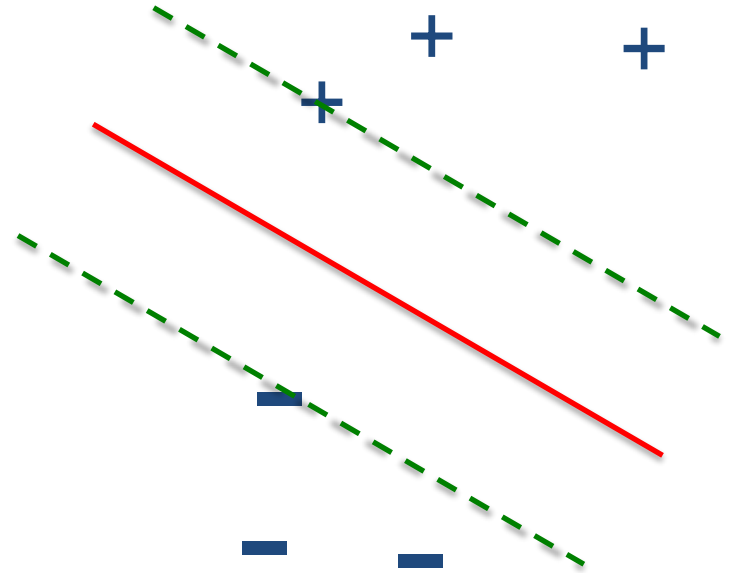
# Support Vector Machines

Find linear separator with maximum margin

Let $x_+$ and $x_-$ be hypothetical points on the +/- margin from the decision boundary

$$\exists\, w, b \qquad s.t. \quad y^{(i)}\left(w^T x^{(i)} + b\right) \geq 1$$

$$\Leftrightarrow \exists\, w, b \qquad s.t. \quad w^T x_+ + b \geq +1 \quad \text{and}$$

$$w^T x_- + b \leq -1$$

Consider the vector from $x_-$ to $x_+$ and its projection onto the vector $w$:

# Support Vector Machines

$$width = \frac{w^T}{\|w\|_2}(x_+ - x_-)$$

Find linear separator with maximum margin

$$\max_{w,b} \quad \text{"width"}$$

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

# Support Vector Machines

Find linear separator with maximum margin

$$\underset{\boldsymbol{w},b}{\text{argmax}} \quad \text{width}$$

$$\text{width} = \frac{2}{\|\boldsymbol{w}\|_2}$$

# Support Vector Machines

Find linear separator with maximum margin
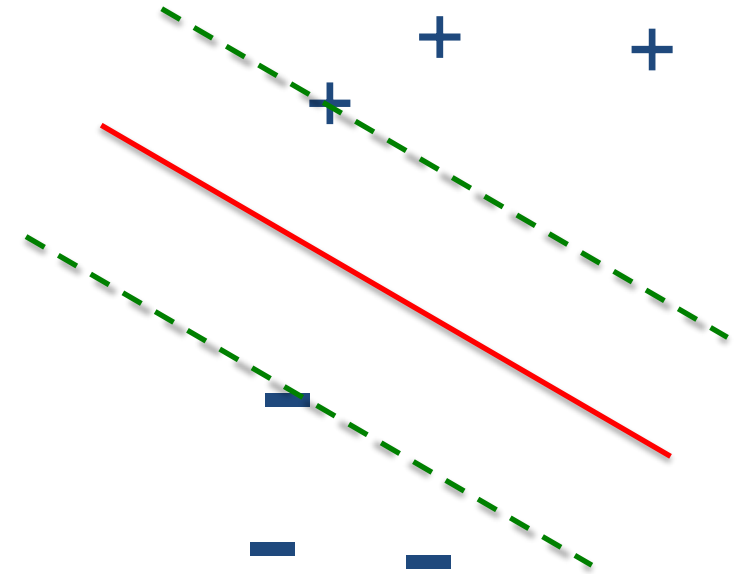
$$\text{width} = \frac{2}{\|\boldsymbol{w}\|_2}$$

$$\underset{\boldsymbol{w},b}{\text{argmax}} \quad \text{width}$$

$$\Leftrightarrow \underset{\boldsymbol{w},b}{\text{argmax}} \quad \frac{2}{\|\boldsymbol{w}\|_2}$$

$$\Leftrightarrow \underset{\boldsymbol{w},b}{\text{argmin}} \quad \frac{1}{2}\|\boldsymbol{w}\|_2$$

$$\Leftrightarrow \underset{\boldsymbol{w},b}{\text{argmin}} \quad \frac{1}{2}\|\boldsymbol{w}\|_2^2$$

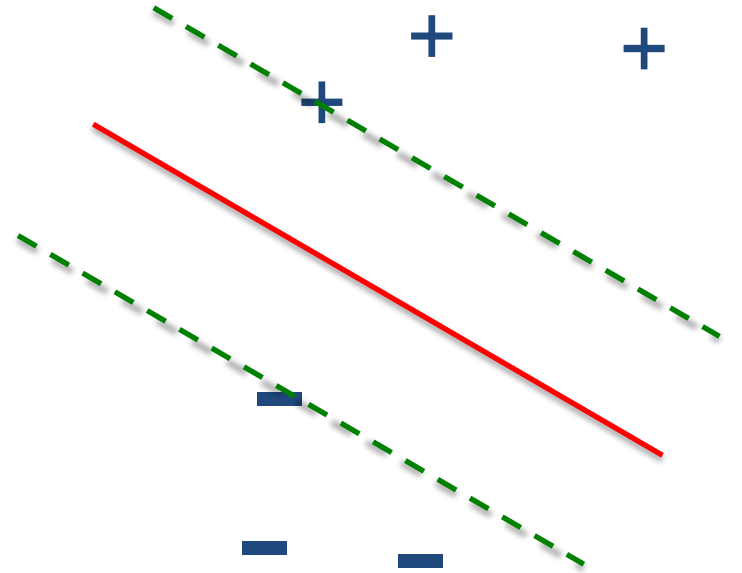$$\Leftrightarrow \underset{\boldsymbol{w},b}{\text{argmin}} \quad \frac{1}{2}\boldsymbol{w}^T\boldsymbol{w}$$

# SVM Optimization

Quadratic program!

$$\min_{\mathbf{w}, b} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

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

Quadratic Program

$$\min_{\mathbf{x}} \quad \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{c}^T \mathbf{x}$$

$$\text{s.t.} \quad \mathbf{A}\mathbf{x} \preceq \mathbf{b}$$

# SVM Optimization

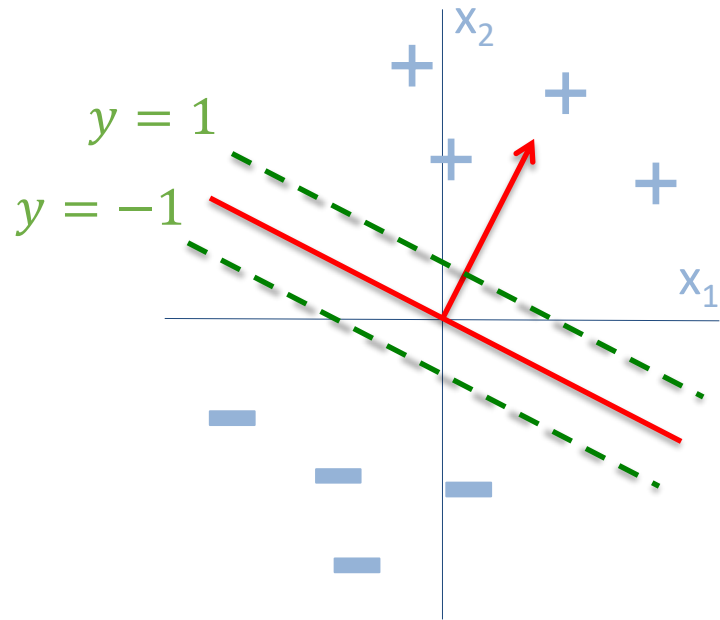How did we go from maximizing margin to minimizing $\|\boldsymbol{w}\|_2$?

# SVM Optimization

How did we go from maximizing margin to minimizing $\|\boldsymbol{w}\|_2$?
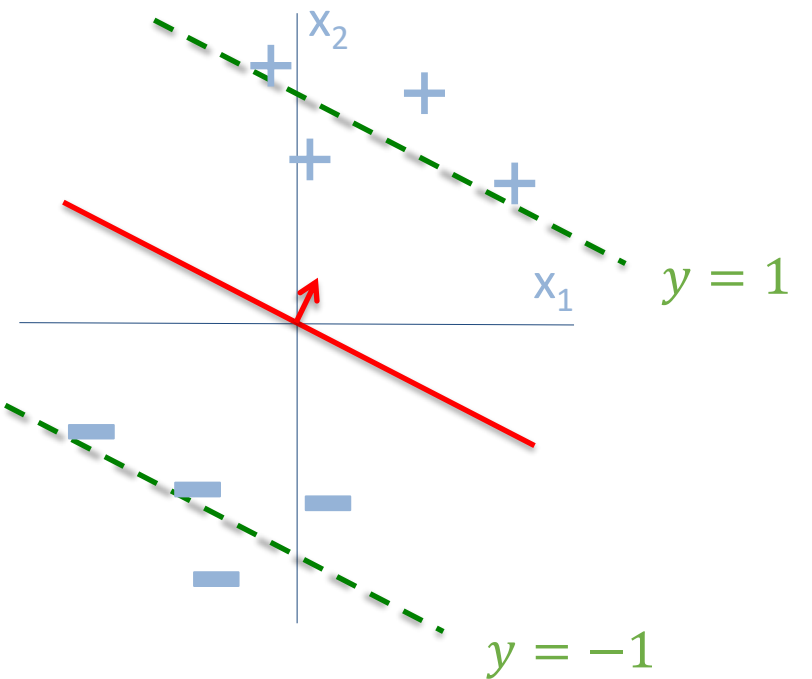


$$\|\boldsymbol{w}\|_2 = 2$$

$$\frac{1}{\|\boldsymbol{w}\|_2} = \frac{1}{2}$$

$$\|\boldsymbol{w}\|_2 = 1$$

$$\frac{1}{\|\boldsymbol{w}\|_2} = 1$$

$$\|\boldsymbol{w}\|_2 = \frac{1}{2}$$

$$\frac{1}{\|\boldsymbol{w}\|_2} = 2$$

# Linear Separability

## Data

$$\mathcal{D} = \left\{ \boldsymbol{x}^{(i)}, y^{(i)} \right\}_{i=1}^{N} \quad \boldsymbol{x} \in \mathbb{R}^{M}, \quad y \in \{-1, +1\}$$

## Linearly separable iff:

$$\exists \, \boldsymbol{w}, b \quad s.t. \quad \boldsymbol{w}^{T} \boldsymbol{x}^{(i)} + b > 0 \quad \text{if} \quad y^{(i)} = +1 \quad \text{and}$$

$$\boldsymbol{w}^{T} \boldsymbol{x}^{(i)} + b < 0 \quad \text{if} \quad y^{(i)} = -1$$

# Support Vector Machines

Find linear separator with maximum margin

$$\max_{\boldsymbol{w},b} \quad \text{"width"}$$

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

$$width = \frac{\boldsymbol{w}^T}{\|\boldsymbol{w}\|_2}(\boldsymbol{x}_+ - \boldsymbol{x}_-)$$

# SVM Optimization

Quadratic program!

$$\min_{\mathbf{w},b} \quad \frac{1}{2}\|\mathbf{w}\|_2^2$$
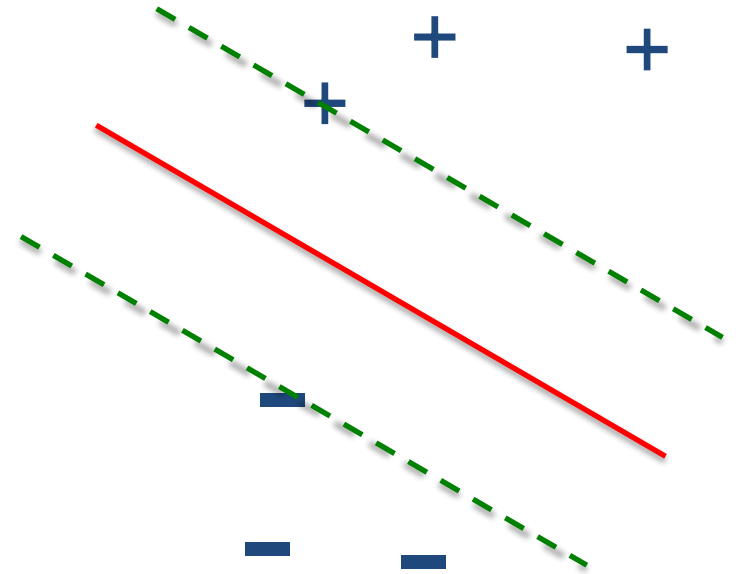
$$\text{s.t.} \quad y^{(i)}\left(\mathbf{w}^T\mathbf{x}^{(i)} + b\right) \geq 1 \quad \forall\, i$$

Quadratic Program

$$\min_{\mathbf{x}} \quad \mathbf{x}^T\mathbf{Q}\mathbf{x} + \mathbf{c}^T\mathbf{x}$$

$$\text{s.t.} \quad \mathbf{A}\mathbf{x} \preceq \mathbf{b}$$

# Constrained Optimization

## Linear Program

$$\min_{x} \quad c^T x$$

$$\text{s.t.} \quad Ax \preceq b$$

## Solvers

- Simplex
- Interior point methods

## Quadratic Program

$$\min_{x} \quad x^T Q x + c^T x$$

$$\text{s.t.} \quad Ax \preceq b$$

## Solvers

- Conjugate gradient
- Ellipsoid method
- Interior point methods

# Constrained Optimization

## Linear Program

$$\min_{x} \quad c^T x$$

$$\text{s.t.} \quad Ax \preceq b$$

## Solvers

- Simplex

- Interior point methods

## Quadratic Program

$$\min_{x} \quad x^T Q x + c^T x$$

$$\text{s.t.} \quad Ax \preceq b$$

## Special Case

- If $Q$ is positive-definite, the problem is convex

- $Q$ is positive-definite if:
$$v^T Q v > 0 \quad \forall\, v \in \mathbb{R}^M \backslash 0$$

- A symmetric $Q$ is positive-definite if all of its eigenvalues are positive

# Support Vector Machines

Next steps

- Different optimization formulation
    - Primal → dual
    - "Support vectors"

- Support non-linear classification
    - Feature maps
    - Kernel trick

- Support non-separable data
    - Hard-margin SVM → soft-margin SVM

# Method of Lagrange Multipliers

## Goal

$$\min_{\boldsymbol{x}} \quad f(\boldsymbol{x})$$

$$\text{s.t.} \quad g(\boldsymbol{x}) \leq c$$

## Step 1: Construct Lagrangian

$$\mathcal{L}(\boldsymbol{x}, \lambda) = f(\boldsymbol{x}) + \lambda(g(\boldsymbol{x}) - c)$$

## Step 2: Solve

$$\min_{\boldsymbol{x}} \quad \max_{\lambda \geq 0} \quad \mathcal{L}(\boldsymbol{x}, \lambda)$$

### Find saddle point:

$$\nabla \mathcal{L}(\boldsymbol{x}, \lambda) \quad \text{s.t. } \lambda \geq 0$$

### Equivalent to solving:

$$\nabla f(x) = \lambda \nabla g(x) \quad \text{s.t. } \lambda \geq 0$$

# SVM Primal vs Dual

## Construct Lagrangian

## Primal

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

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

Lagrange Multipliers

$$\min_{\boldsymbol{x}} \quad f(\boldsymbol{x}) \quad \text{s.t.} \quad g(\boldsymbol{x}) \leq c$$

Construct Lagrangian

$$\mathcal{L}(\boldsymbol{x},\lambda) = f(\boldsymbol{x}) + \lambda(g(\boldsymbol{x}) - c)$$

Solve: $\min_{\boldsymbol{x}} \quad \max_{\lambda \geq 0} \quad \mathcal{L}(\boldsymbol{x},\lambda)$

# SVM Dual Optimization

$$\mathcal{L}(w, b, \boldsymbol{\alpha}) = \frac{1}{2}\boldsymbol{w}^T\boldsymbol{w} - \sum_i^N \alpha_i[\, y^{(i)}(\boldsymbol{w}^T\boldsymbol{x}^{(i)} + b) - 1]$$

# SVM Dual Optimization

Dual

$$\max_{\boldsymbol{\alpha}} \quad \sum_i^N \alpha_i - \frac{1}{2}\sum_i^N \sum_j^N \alpha_i \alpha_j y^{(i)} y^{(j)} \boldsymbol{x}^{(i)T} \boldsymbol{x}^{(j)}$$

$$\text{s.t.} \quad \alpha_i \geq 0 \quad \forall i$$

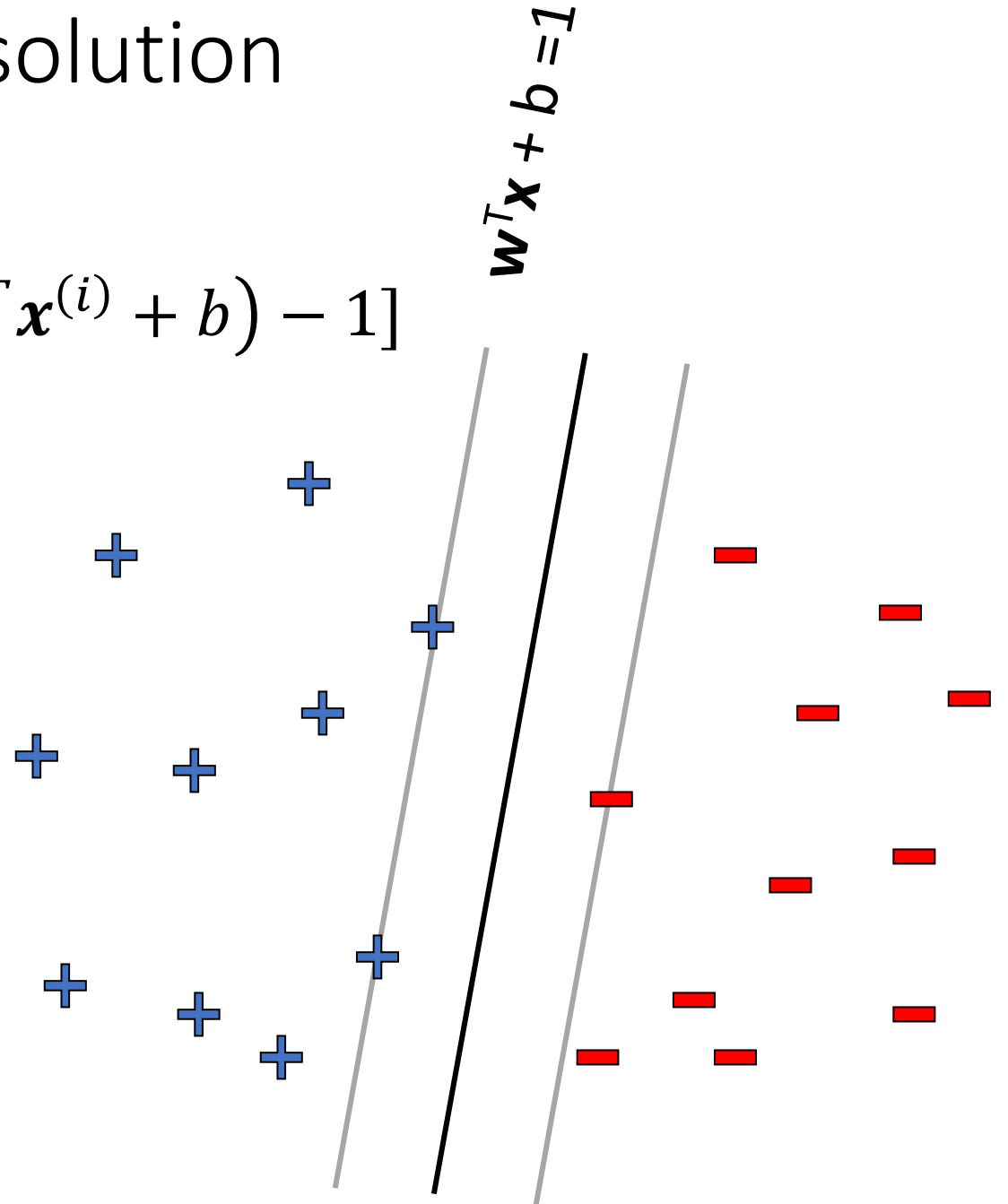$$\boldsymbol{w} = \sum_i^N \alpha_i y^{(i)} \boldsymbol{x}^{(i)}$$

$$b = y^{(k)} - \boldsymbol{w}^T \boldsymbol{x}^{(k)} \text{ for any } k \text{ where } \alpha_k > 0$$

Prediction

# Dual SVM: Sparsity of dual solution

$$\min_{\boldsymbol{w},b} \max_{\boldsymbol{\alpha} \geq 0} \quad \mathcal{L}(\boldsymbol{w}, b, \boldsymbol{\alpha})$$

$$\min_{\boldsymbol{w},b} \max_{\boldsymbol{\alpha} \geq 0} \quad \frac{1}{2}\boldsymbol{w}^T\boldsymbol{w} - \sum_i^N \alpha_i[\ y^{(i)}(\boldsymbol{w}^T\boldsymbol{x}^{(i)} + b) - 1]$$

$\boldsymbol{w}^T\boldsymbol{x} + b = 1$

# Dual SVM: Sparsity of dual solution

$$\min_{\boldsymbol{w},b} \ \max_{\boldsymbol{\alpha} \geq 0} \quad \mathcal{L}(\boldsymbol{w}, b, \boldsymbol{\alpha})$$

$$\min_{\boldsymbol{w},b} \ \max_{\boldsymbol{\alpha} \geq 0} \quad \frac{1}{2}\boldsymbol{w}^T\boldsymbol{w} - \sum_i^N \alpha_i[\ y^{(i)}\big(\boldsymbol{w}^T\boldsymbol{x}^{(i)} + b\big) - 1]$$

$\boldsymbol{w}^T\boldsymbol{x} + b = 1$

# Dual SVM: Sparsity of dual solution



$$\mathbf{w} = \sum_j \alpha_j y_j \mathbf{x}_j$$

Only few $\alpha_j$s can be non-zero : where constraint is active and tight

$$(\mathbf{w}.\mathbf{x}_j + b)y_j = 1$$

**Support vectors** – training points j whose $\alpha_j$s are non-zero

29

# Support Vector Machines

Next steps

- Different optimization formulation
  - Primal → dual
  - "Support vectors"

- Support non-linear classification
  - Feature maps
  - Kernel trick

- Support non-separable data
  - Hard-margin SVM → soft-margin SVM

# Kernels: Motivation

Most real-world problems exhibit data that is not linearly separable.
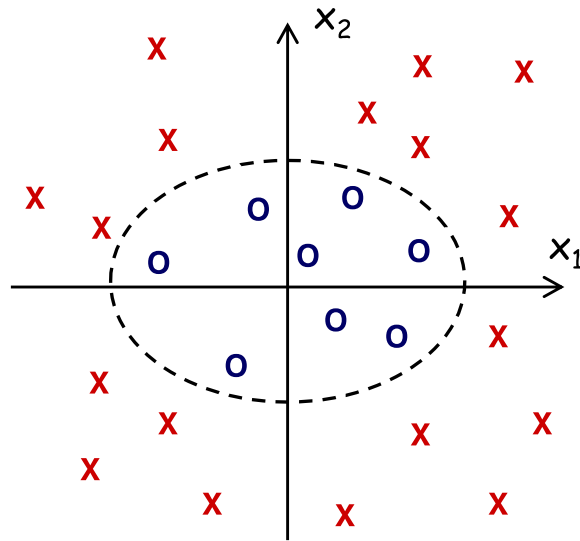
Example: pixel representation for Facial Recognition:



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

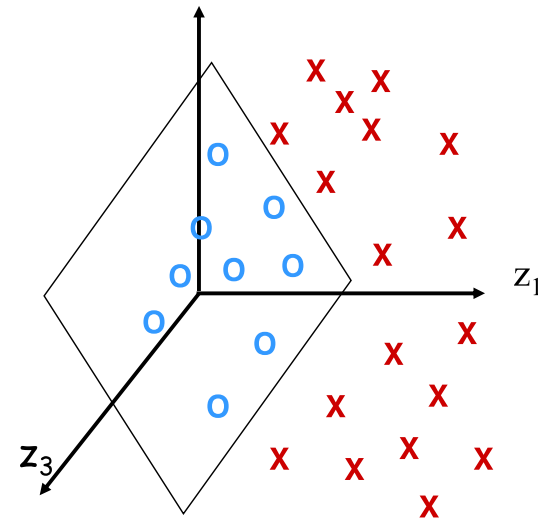**A:** Preprocess the data to produce **nonlinear features**

# Example: Polynomial Kernel

https://www.youtube.com/watch?v=3liCbRZPrZA

# 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

- **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:

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

  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)**

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

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

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

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

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

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

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

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

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

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

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

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

# 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 $k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \phi(\mathbf{x}^{(i)}) \cdot \phi(\mathbf{x}^{(j)})$

# 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:

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

  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 φ 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

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

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

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

# Kernel Examples

**Side Note:** The feature space might not be unique!

**Explicit representation #1:**

$$\phi: R^2 \rightarrow R^3, (x_1, x_2) \rightarrow \Phi(x) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

$$\phi(x) \cdot \phi(z) = \left(x_1^2, x_2^2, \sqrt{2}x_1x_2\right) \cdot (z_1^2, z_2^2, \sqrt{2}z_1z_2)$$

$$= (x_1z_1 + x_2z_2)^2 = (x \cdot z)^2 = K(x, z)$$

**Explicit representation #2:**

$$\phi: R^2 \rightarrow R^4, (x_1, x_2) \rightarrow \Phi(x) = (x_1^2, x_2^2, x_1x_2, x_2x_1)$$

$$\phi(x) \cdot \phi(z) = (x_1^2, x_2^2, x_1x_2, x_2x_1) \cdot (z_1^2, z_2^2, z_1z_2, z_2z_1)$$

$$= (x \cdot z)^2 = K(x, z)$$

**These two different feature representations correspond to the same kernel function!**

# 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 (RBF) | $K(\mathbf{x}, \mathbf{z}) = \exp(-\dfrac{\|\mathbf{x} - \mathbf{z}\|_2^2}{2\sigma^2})$ | 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) |

# Kernels: Mercer's Theorem

What functions are valid kernels that correspond to feature vectors $\varphi(\mathbf{x})$?

Answer: Mercer kernels for $k(\mathbf{x}, \mathbf{z})$ and

matrix K, where $K_{i,j} = k(x^{(i)}, x^{(j)})$

- $k(\mathbf{x}, \mathbf{z})$ is continuous

- K is symmetric

- K is positive semi-definite, i.e. $\mathbf{z}^T K \mathbf{z} \geq 0$ for all $\mathbf{z}$

# SVMs with Kernels

- Choose a set of features and kernel function
- Solve dual problem to obtain support vectors $\alpha_i$
- At classification time, compute:

$$\mathbf{w} \cdot \Phi(\mathbf{x}) = \sum_i \alpha_i y_i K(\mathbf{x}, \mathbf{x}_i)$$

$$b = y_k - \sum_i \alpha_i y_i K(\mathbf{x}_k, \mathbf{x}_i)$$

for any $k$ where $C > \alpha_k > 0$

**Classify as** → $sign\left(\mathbf{w} \cdot \Phi(\mathbf{x}) + b\right)$
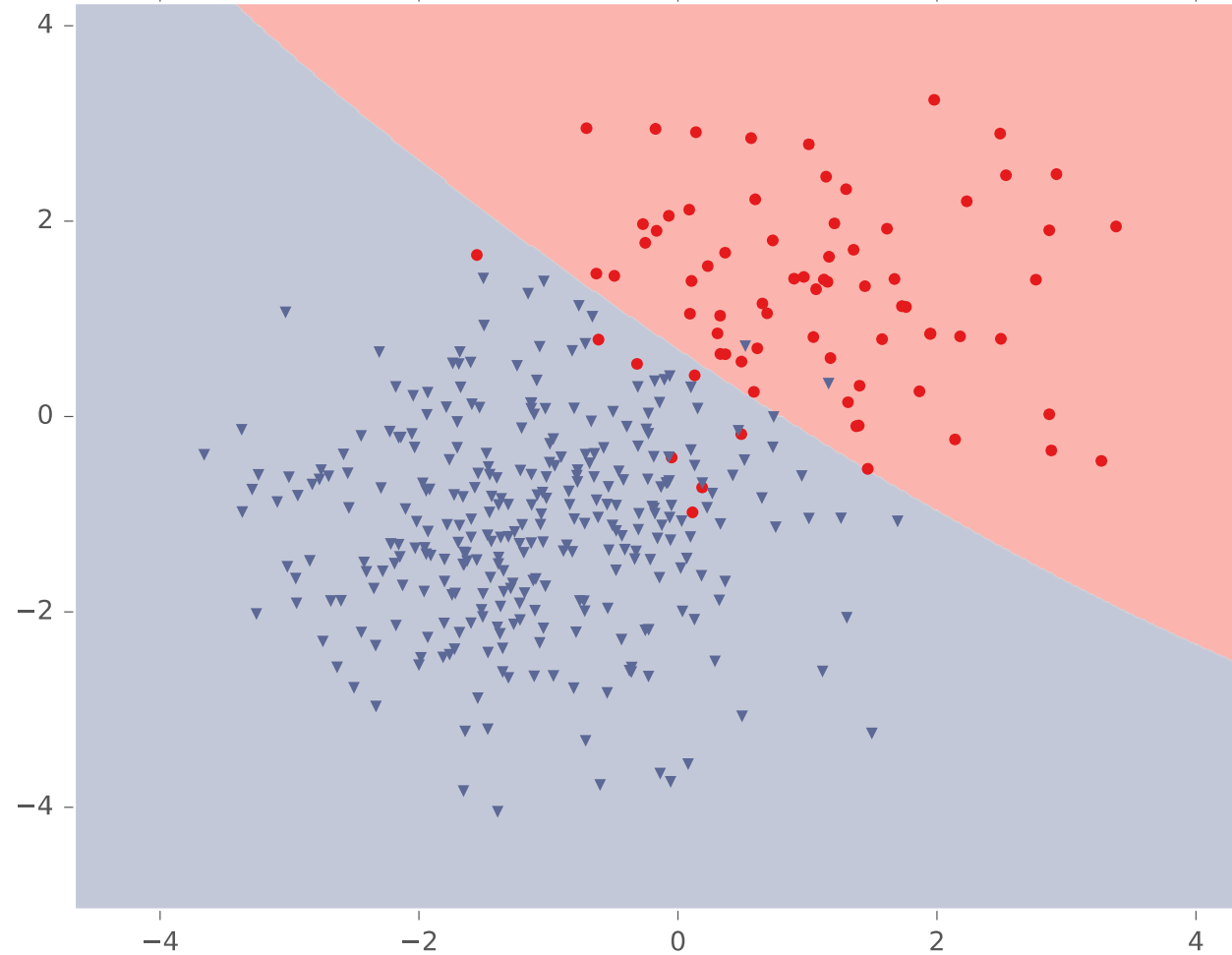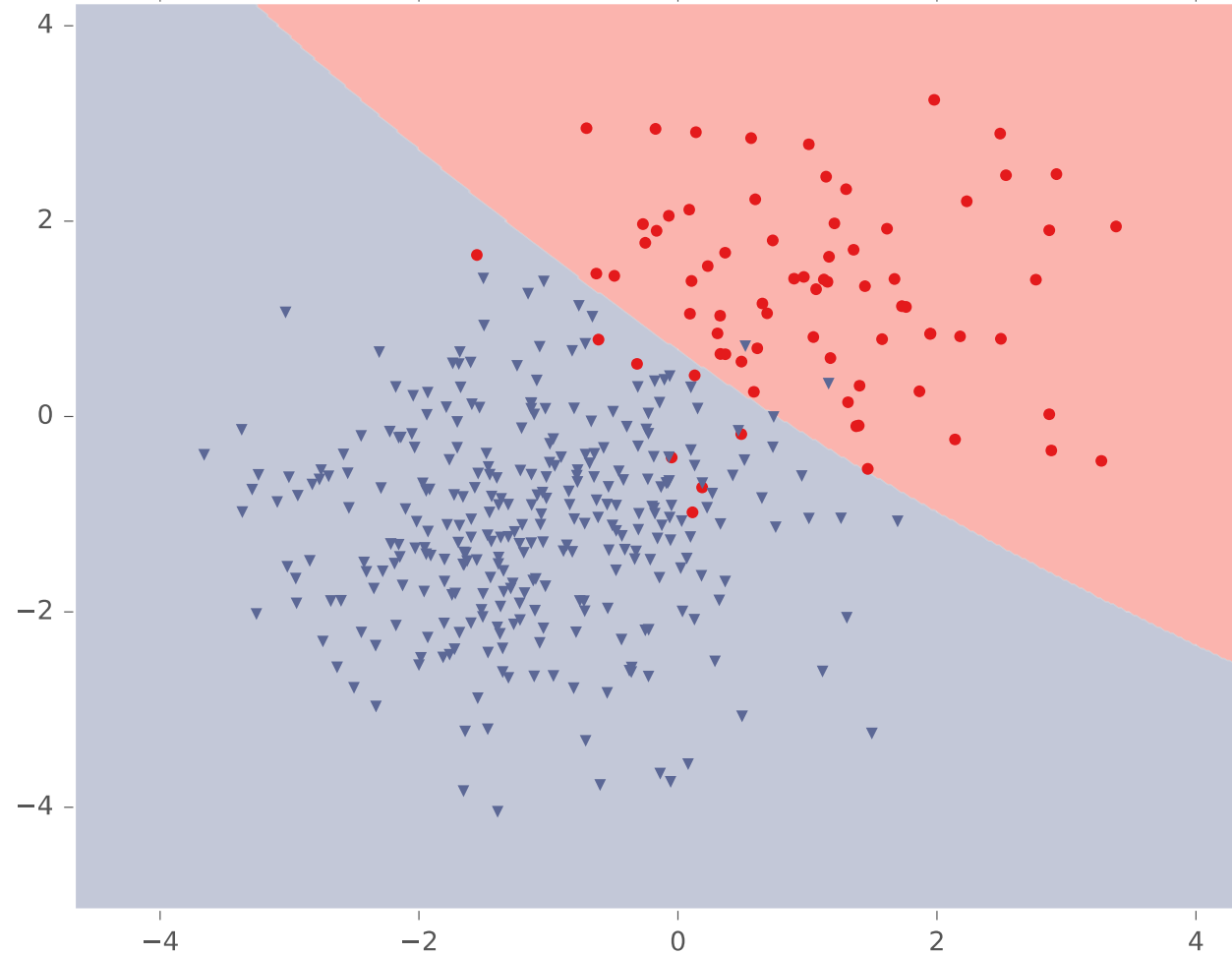
# RBF Kernel Example

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



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

45

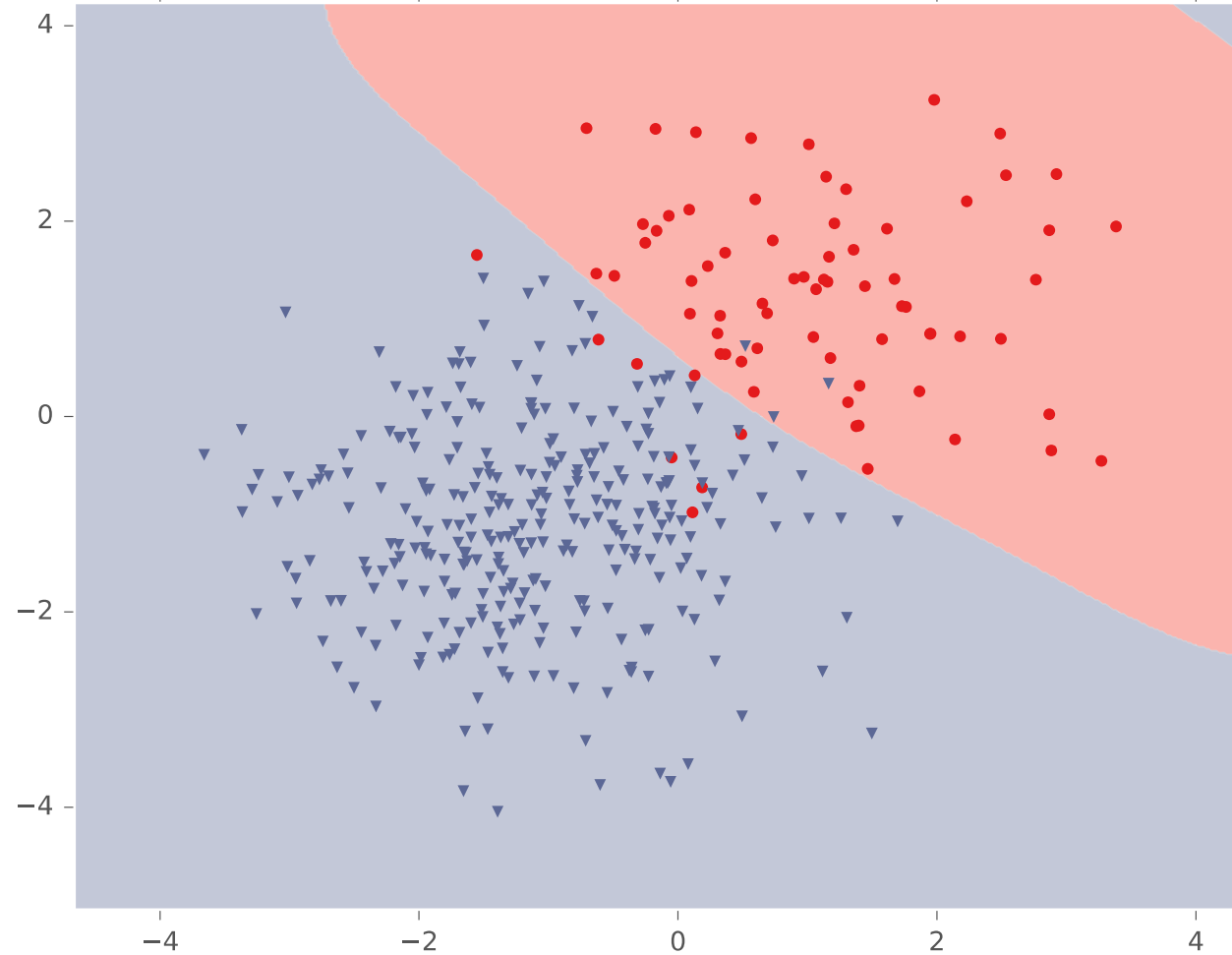# RBF Kernel Example



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

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

# RBF Kernel Example
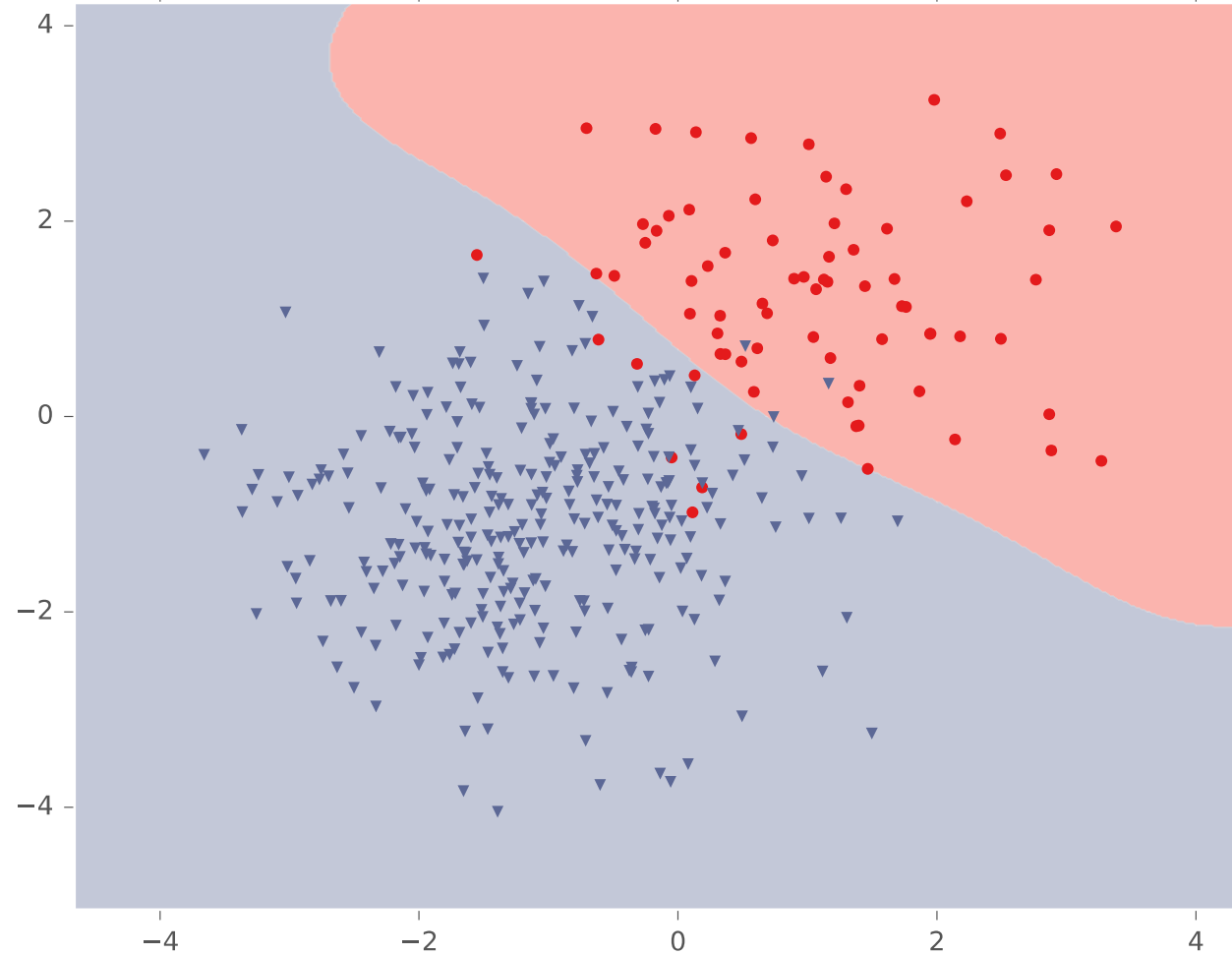


Classification with SVM (kernel=rbf, gamma=0.020000)

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

47

# RBF Kernel Example



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

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

48

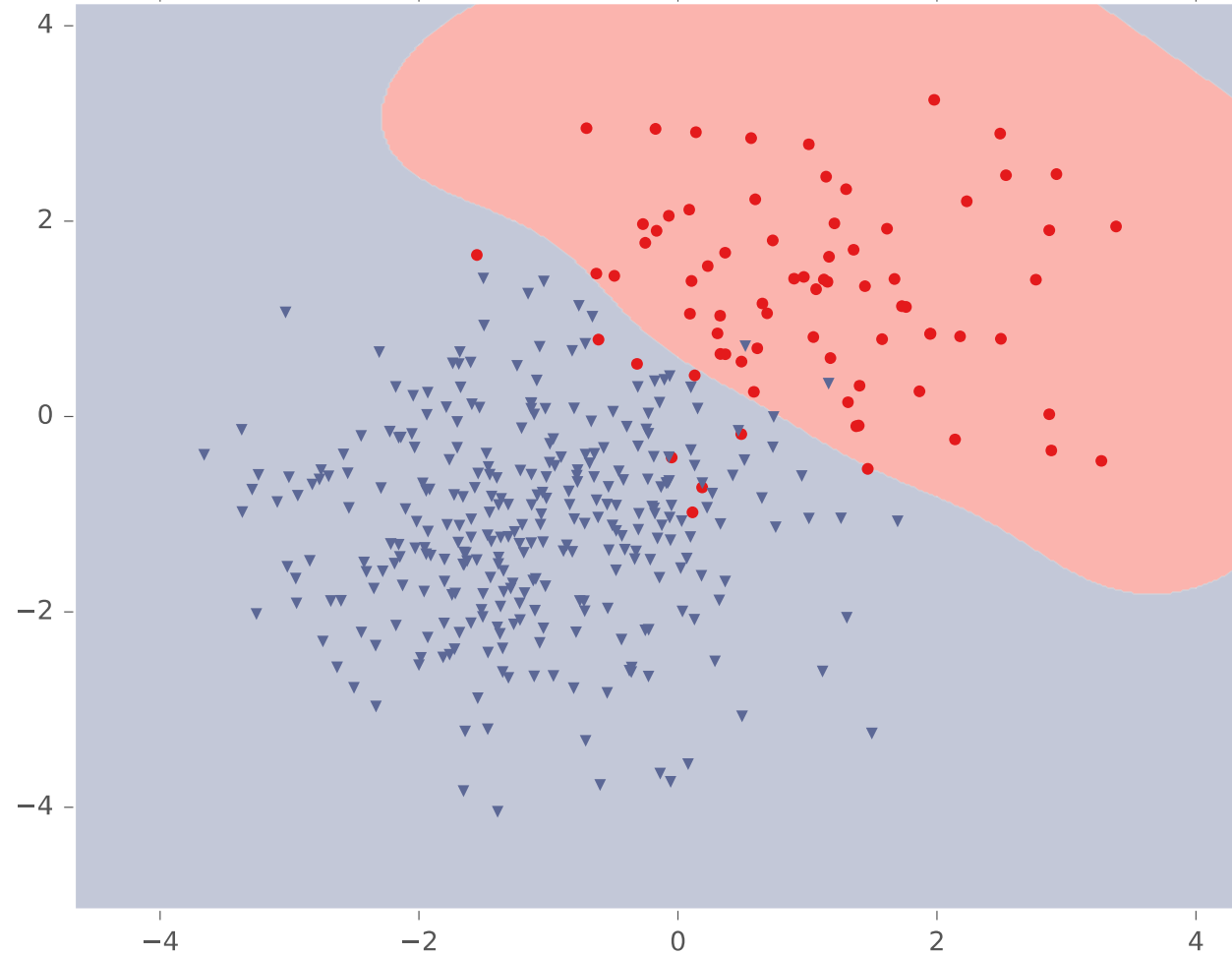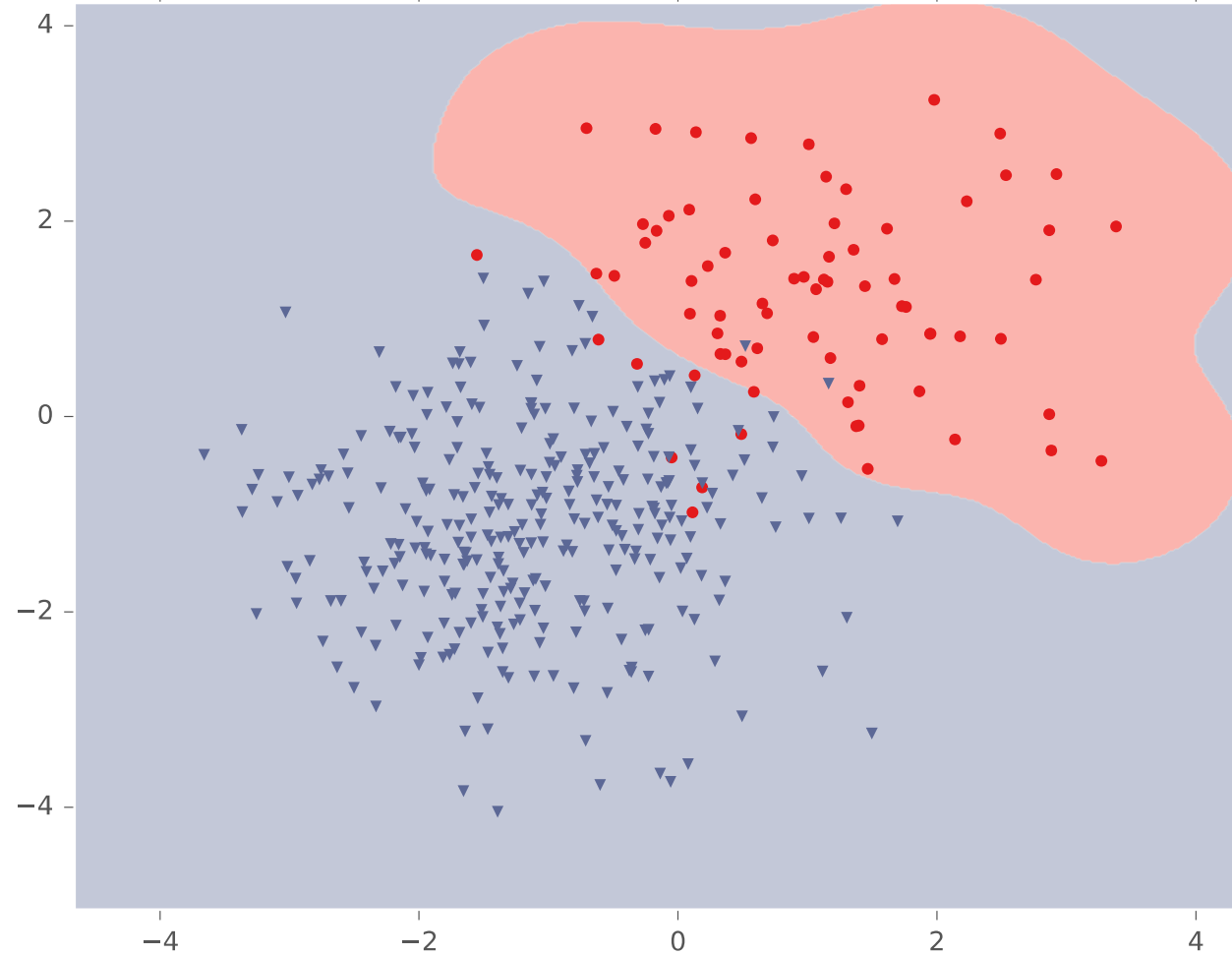# RBF Kernel Example

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



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

49

# RBF Kernel Example

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

Slide credit: CMU MLD Matt Gormley
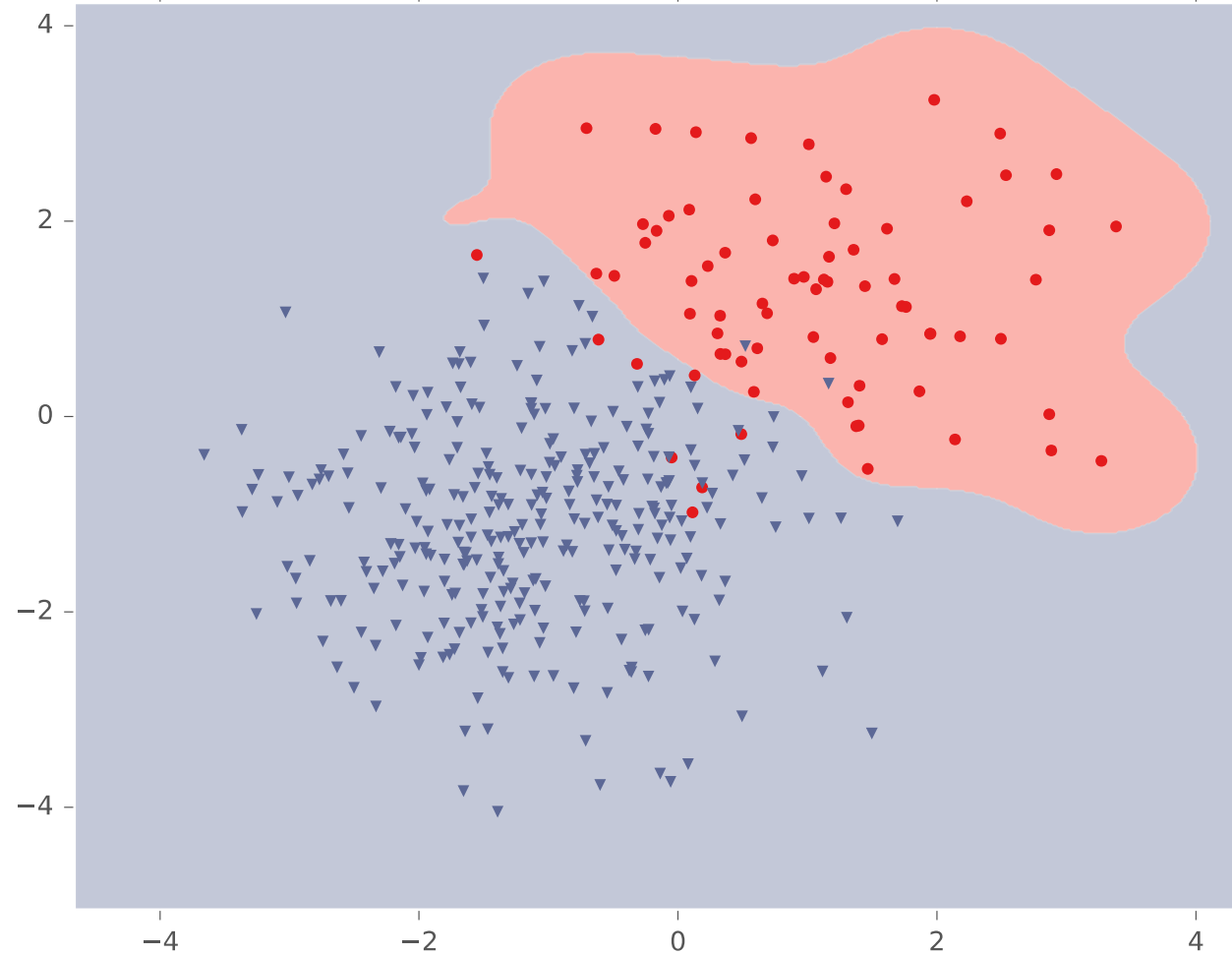
50

# RBF Kernel Example



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

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

51

# RBF Kernel Example

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



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

52

# RBF Kernel Example



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

**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$
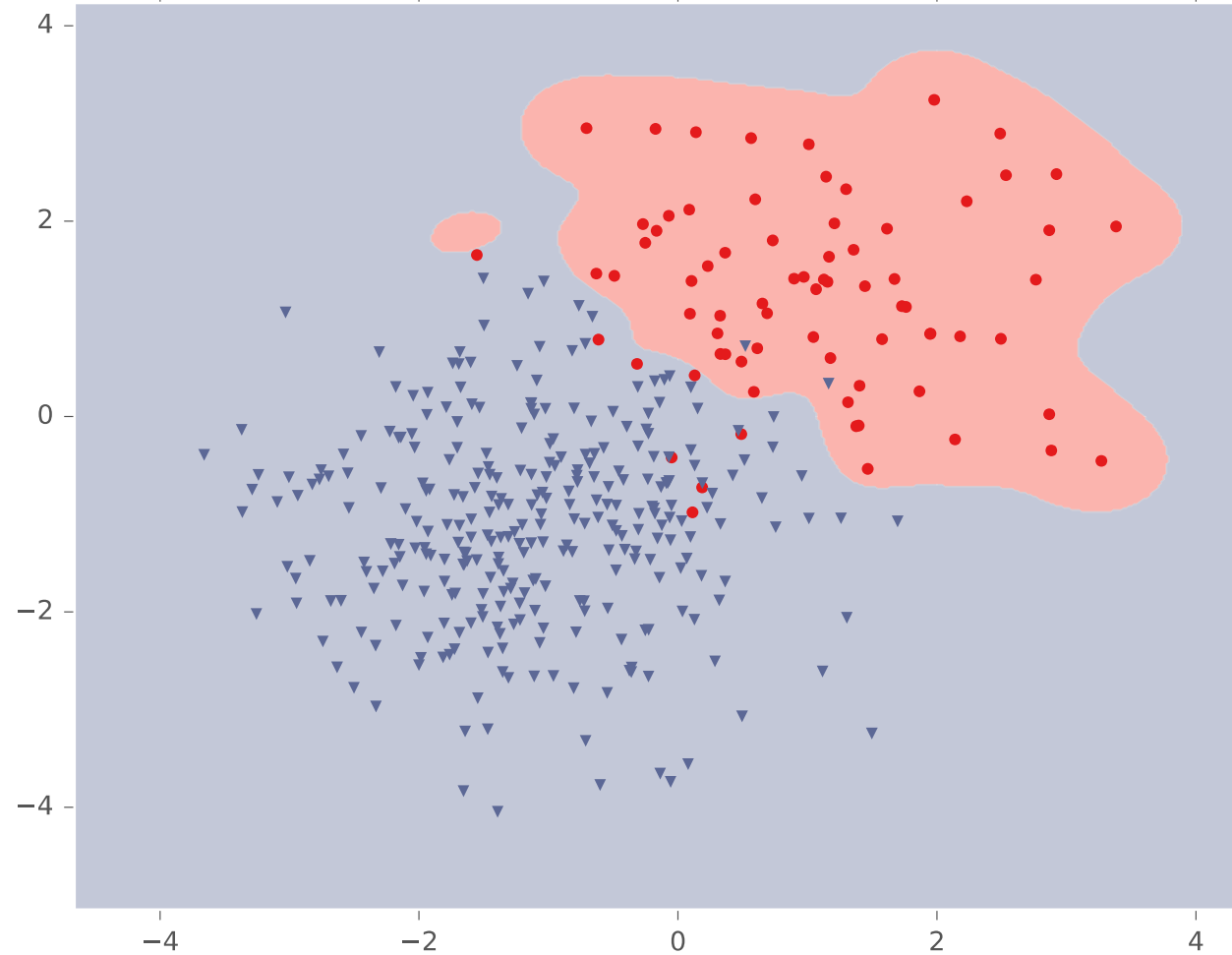
53

# RBF Kernel Example

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

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

54

# RBF Kernel Example

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

**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$
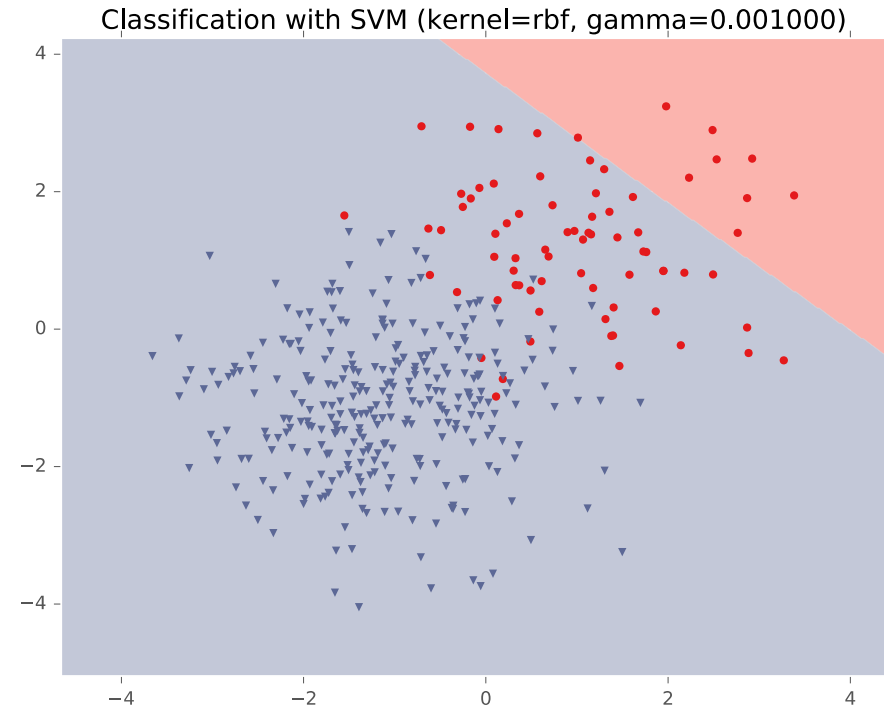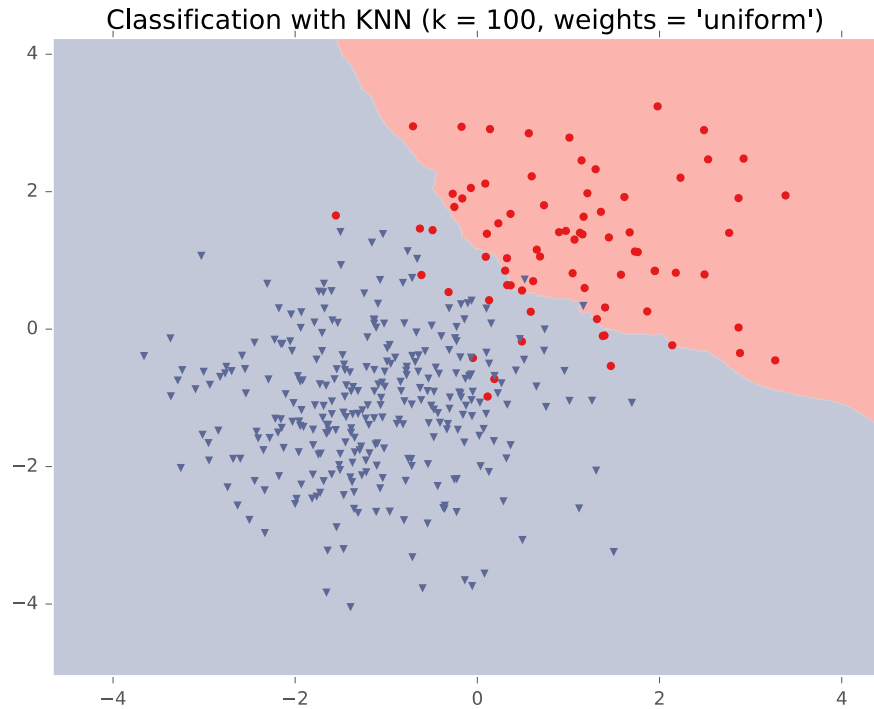
55

# 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:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$
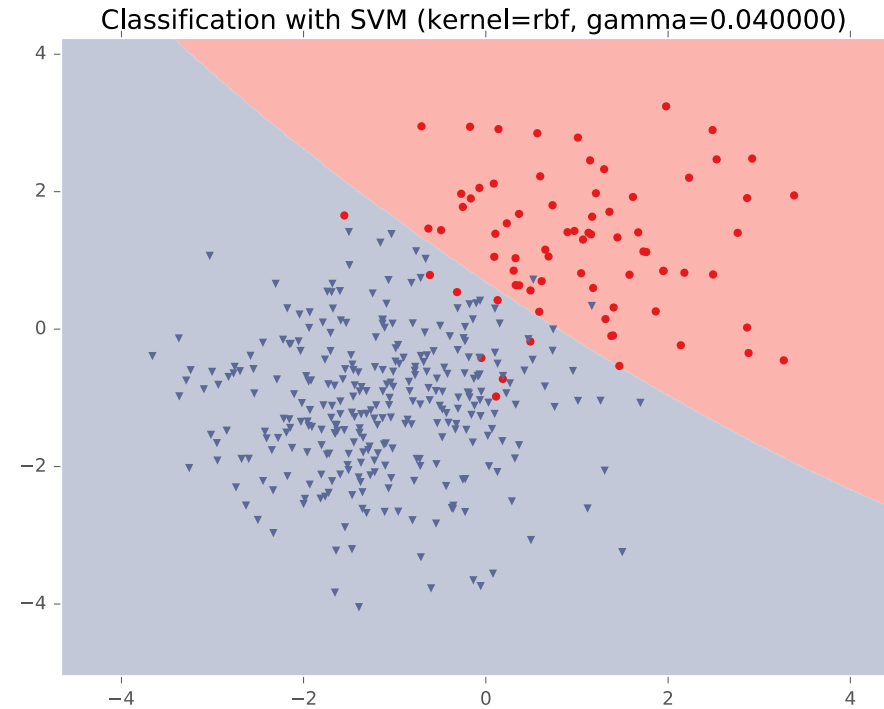
57

# RBF Kernel Example

## KNN vs. SVM



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

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

**RBF Kernel:** $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \exp(-\gamma||\mathbf{x}^{(i)} - \mathbf{x}^{(j)}||_2^2)$
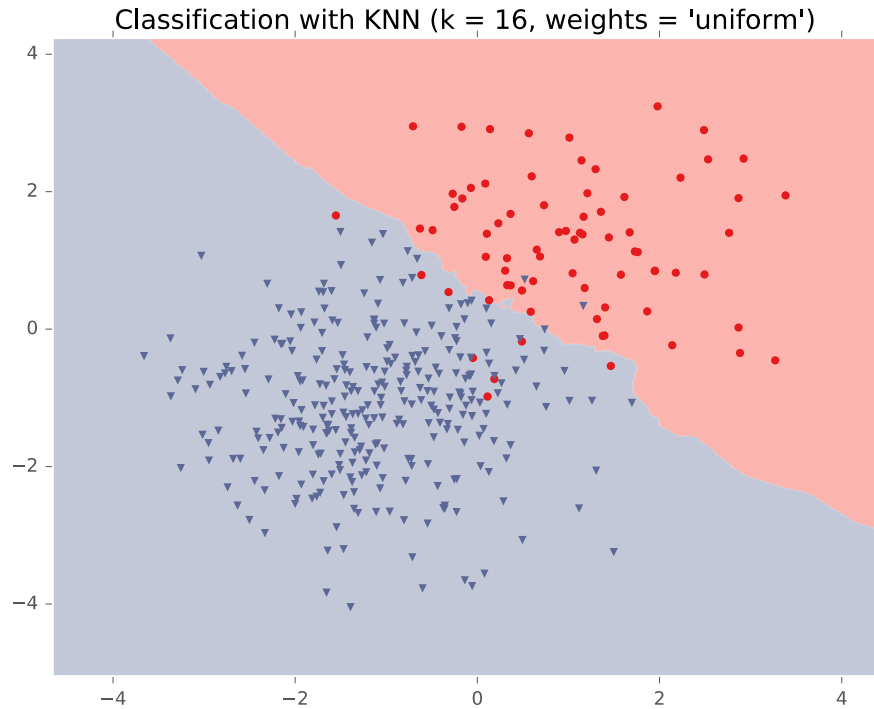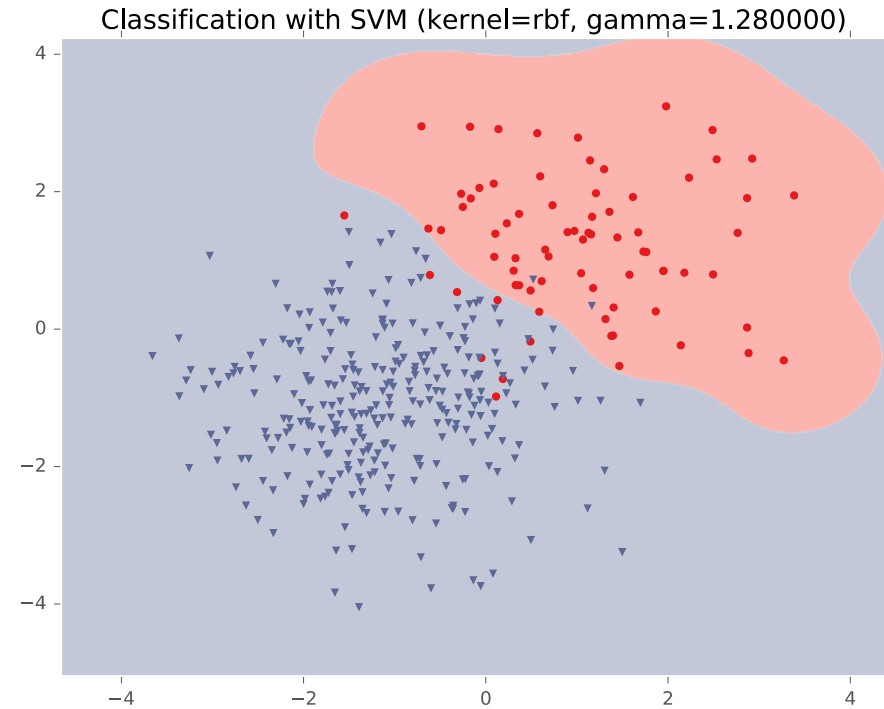
58

# RBF Kernel Example

**KNN vs. SVM**



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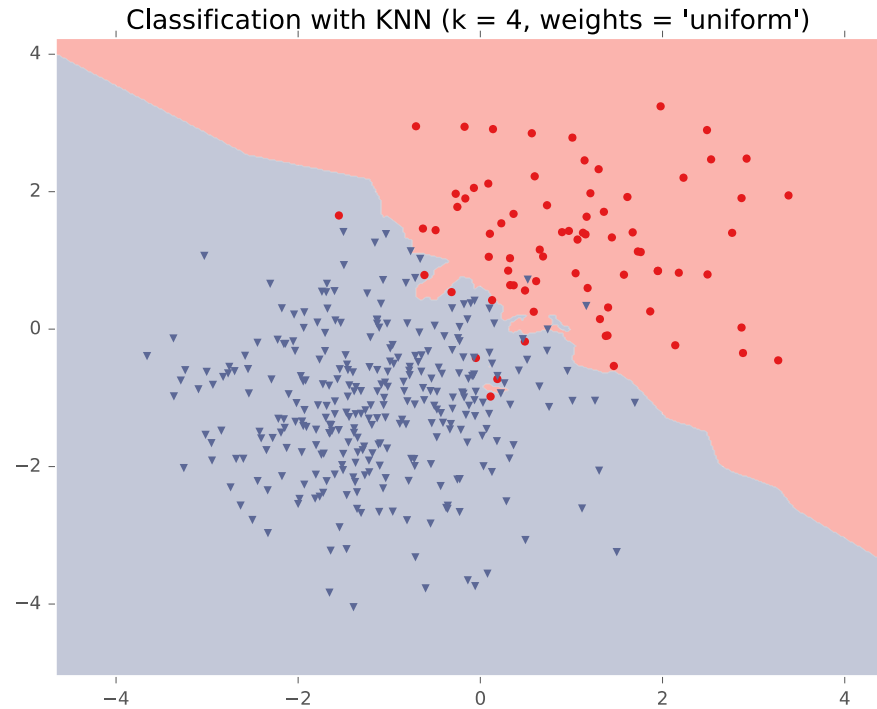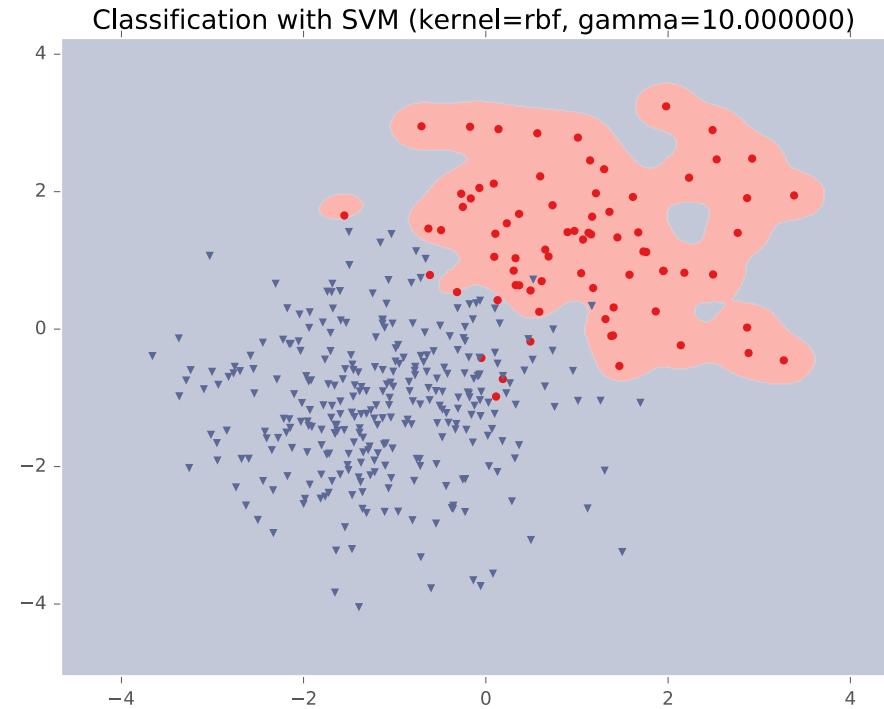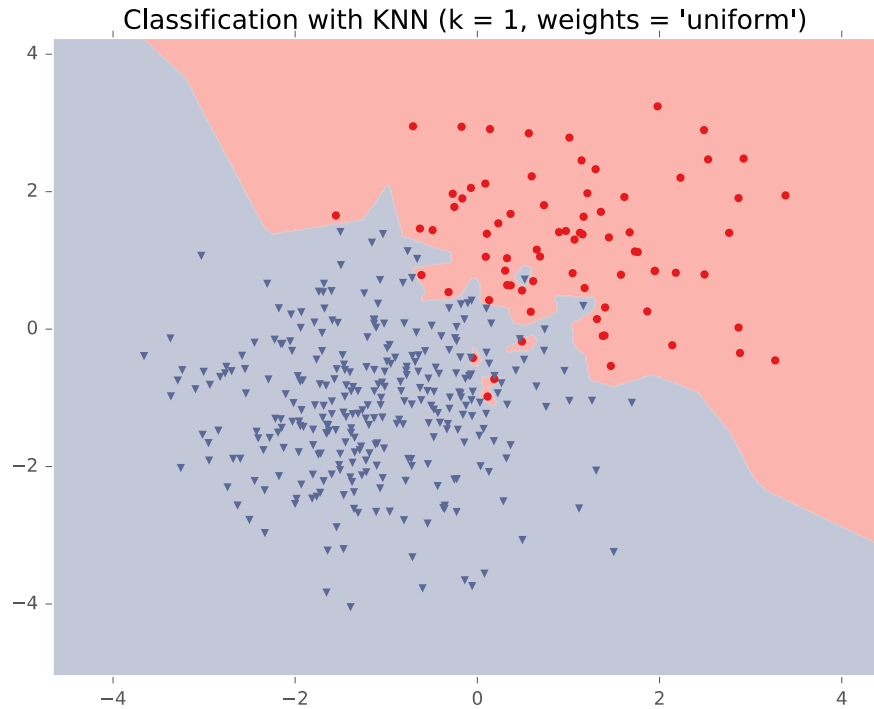
Classification with SVM (kernel=rbf, gamma=1.280000)

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

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
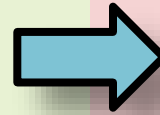
# Support Vector Machines

Next steps

- Different optimization formulation
    - Primal → dual
    - "Support vectors"

- Support non-linear classification
    - Feature maps
    - Kernel trick

- Support non-separable data
    - Hard-margin SVM → soft-margin SVM

# 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)

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

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

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

- 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 $\alpha^{(i)} \neq 0$

# Soft-Margin SVM

**Hard-margin SVM (Primal)**

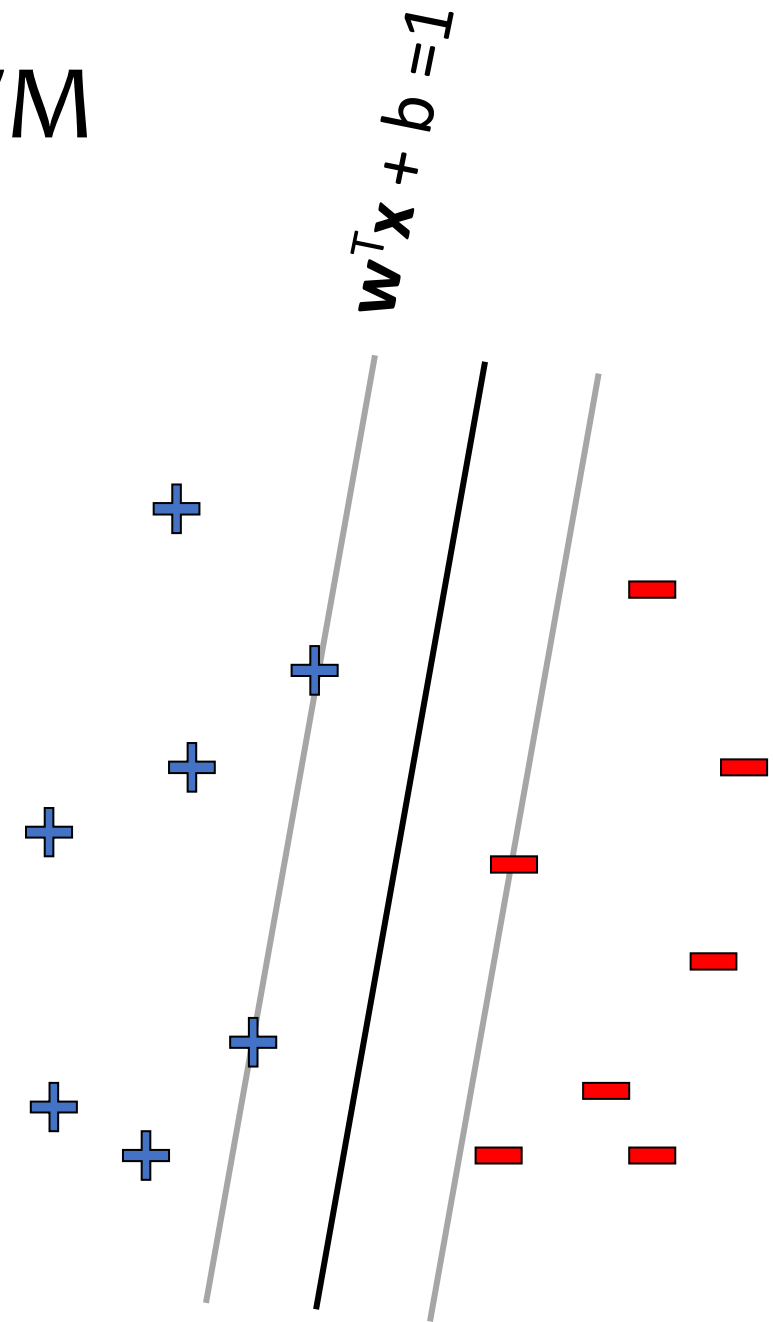$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|_2^2$$

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

**Soft-margin SVM (Primal)**

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

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

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

$\mathbf{w}^T\mathbf{x} + b = 1$

# Soft-Margin SVM

Hard-margin SVM (Primal)

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

$$\textbf{s.t. } y^{(i)}(\mathbf{w}^T\mathbf{x}^{(i)} + b) \geq 1, \quad \forall i = 1,\ldots,N$$
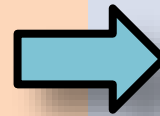
Hard-margin SVM (Lagrangian Dual)

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

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

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

Soft-margin SVM (Primal)

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

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

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

Soft-margin SVM (Lagrangian Dual)

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

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

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

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