

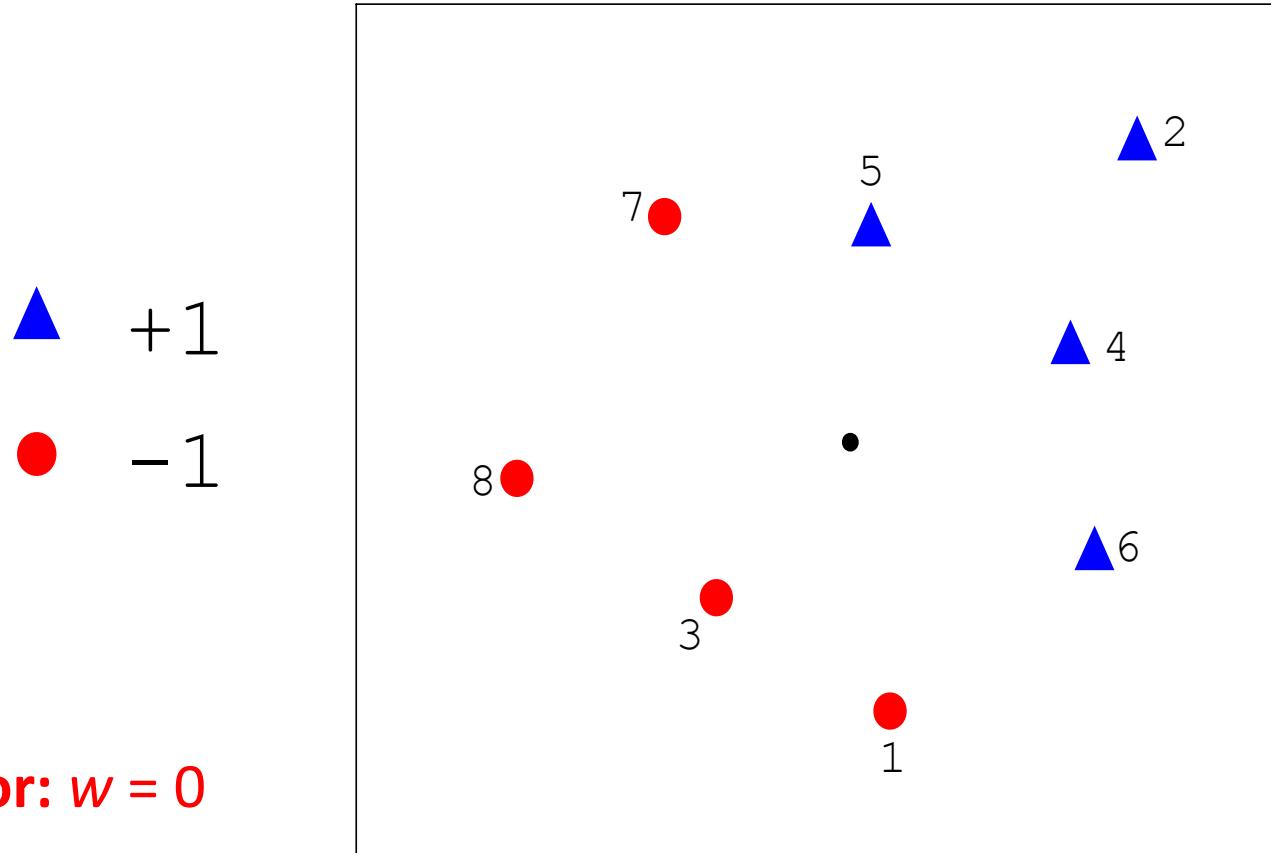
Kernels

DSE 220

Recall: Perceptron

Input space $X = \mathbb{R}^p$, label space $Y = \{-1, 1\}$

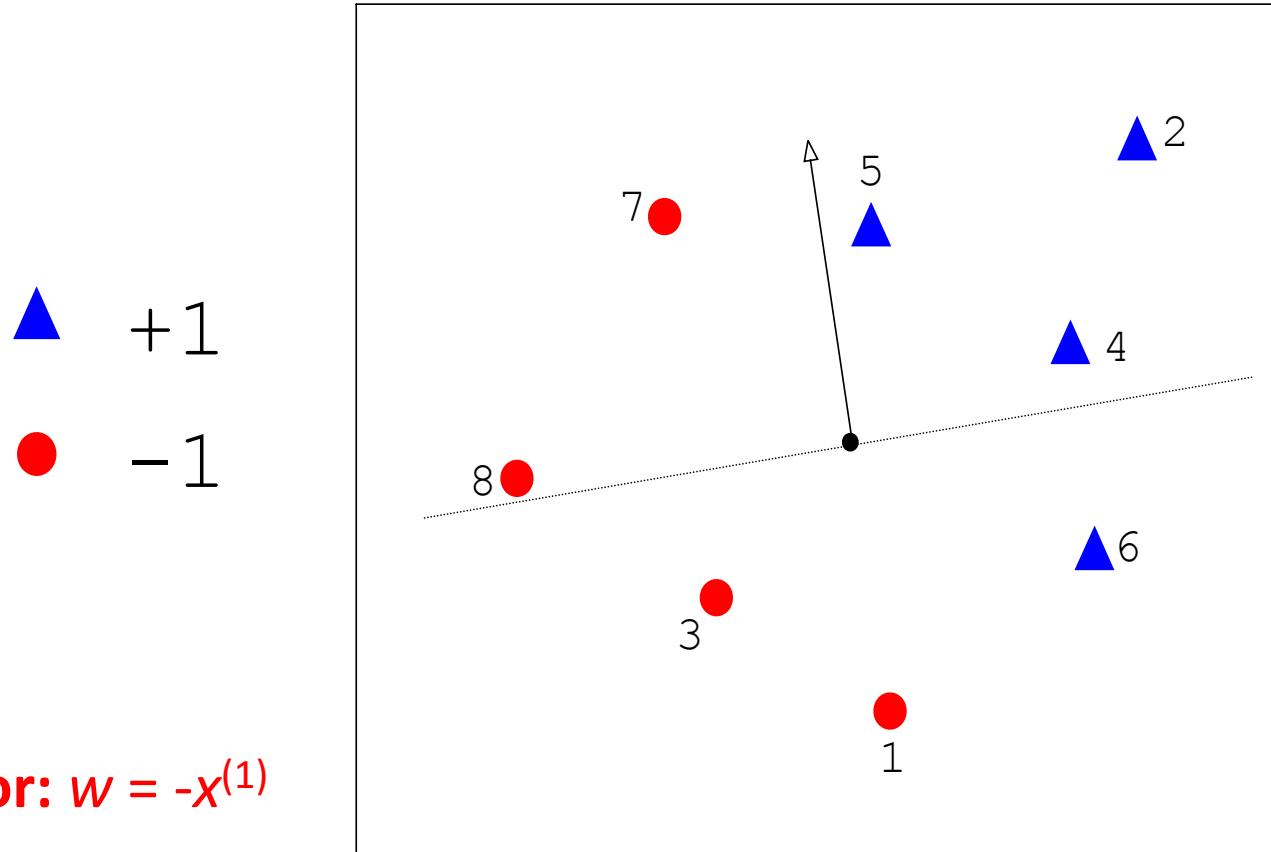
- $w = 0$
- while some (x, y) is misclassified:
 - $w = w + yx$



Recall: Perceptron

Input space $X = \mathbb{R}^p$, label space $Y = \{-1, 1\}$

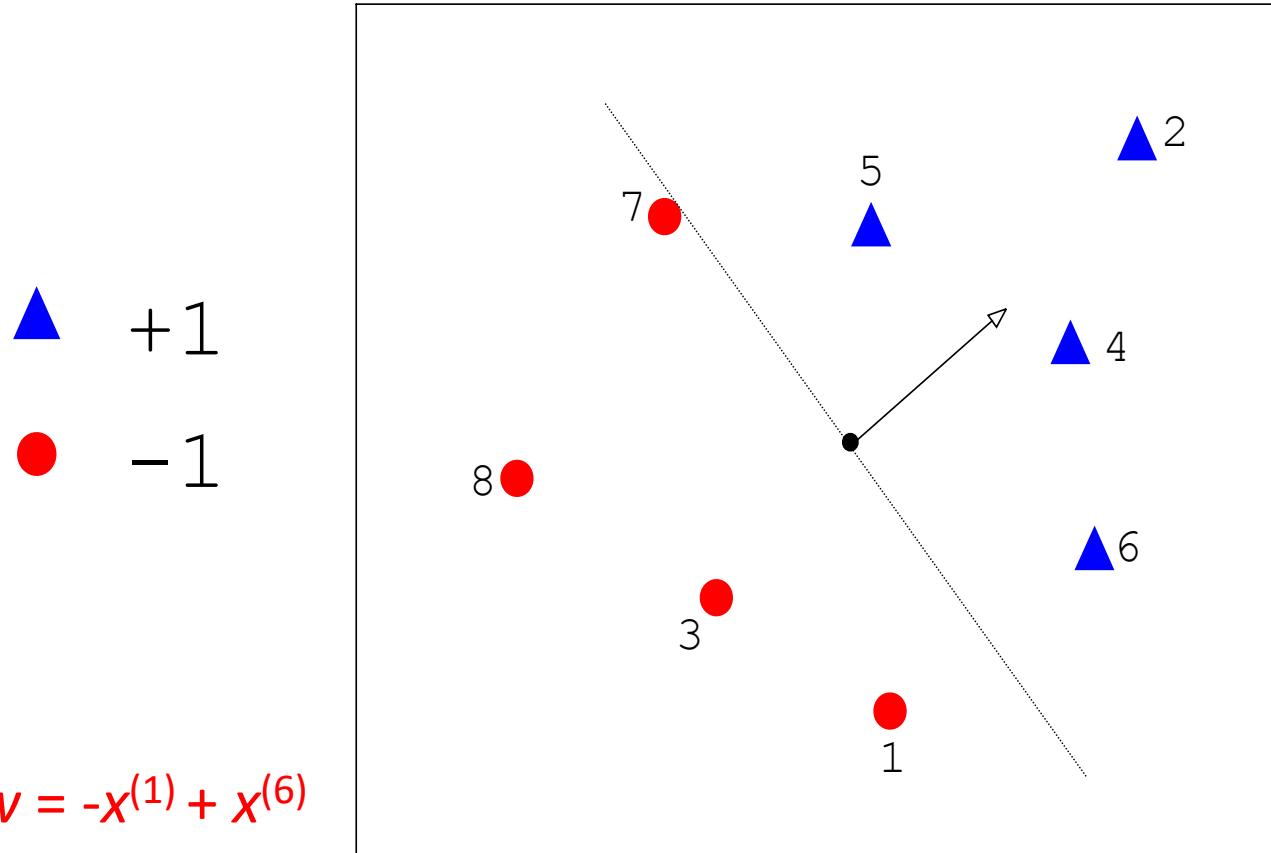
- $w = 0$
- while some (x, y) is misclassified:
 - $w = w + yx$



Recall: Perceptron

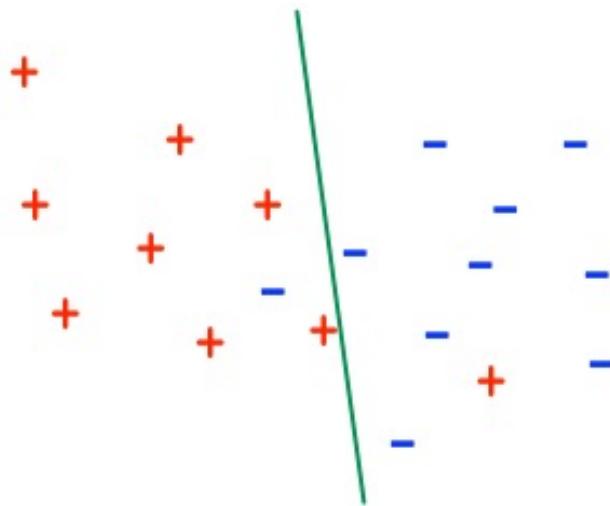
Input space $X = \mathbb{R}^p$, label space $Y = \{-1, 1\}$

- $w = 0$
- while some (x, y) is misclassified:
 - $w = w + yx$



Deviations from linear separability

Noise

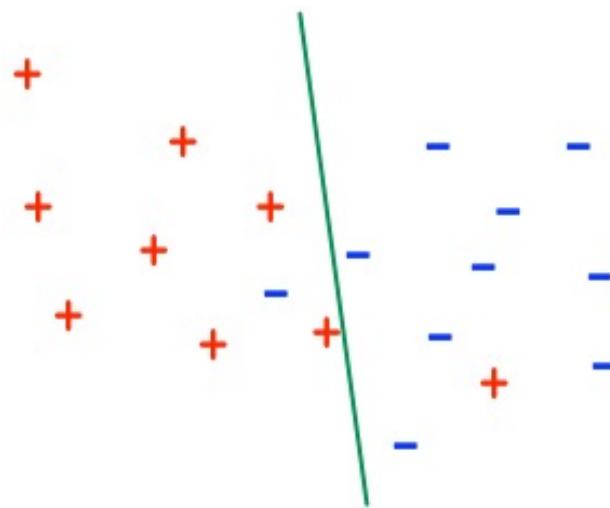


Find a separator that minimizes a convex loss function related to the number of mistakes.

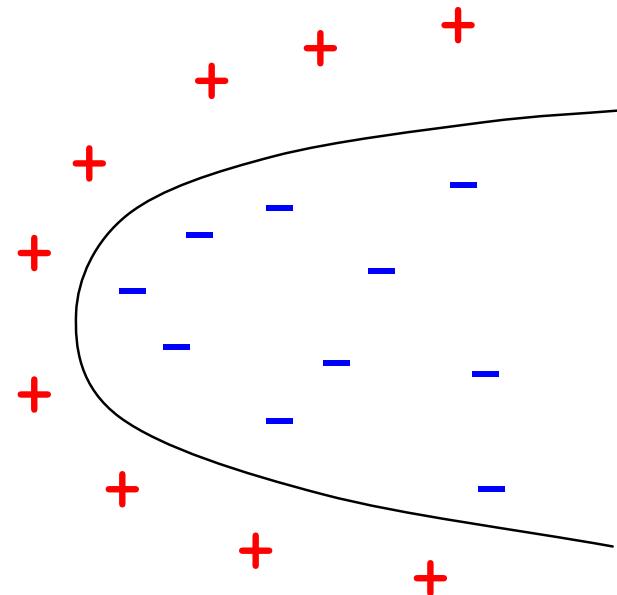
e.g. SVM, logistic regression.

Deviations from linear separability

Noise



Systematic deviation



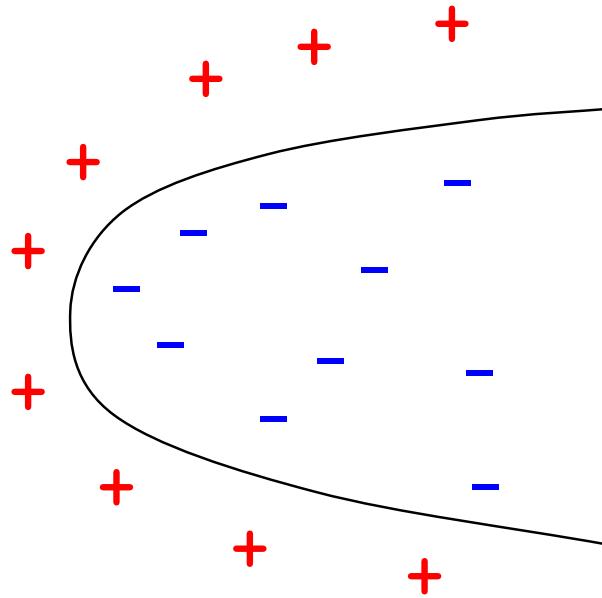
Find a separator that minimizes a convex loss function related to the number of mistakes.

e.g. SVM, logistic regression.

What to do with this?

Systematic inseparability

In this case, the actual boundary looks quadratic.



Quick fix: in addition to the regular features $x = (x_1, x_2, \dots, x_p)$, add in extra features:

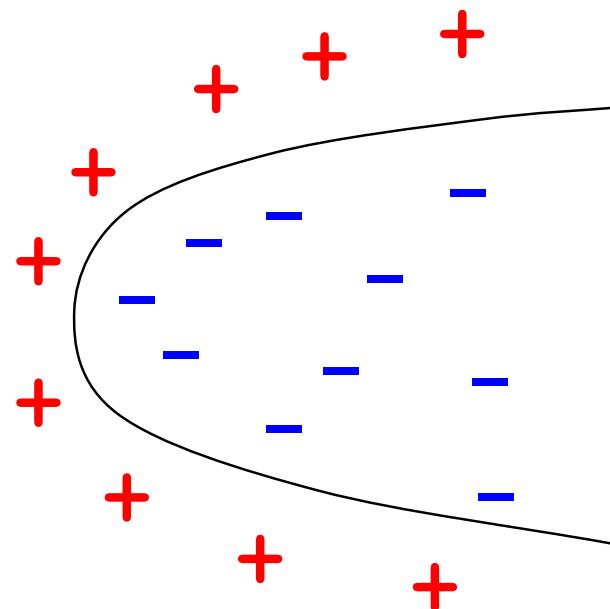
$$x_1^2, x_2^2, \dots, x_p^2$$

$$x_1x_2, x_1x_3, \dots, x_{p-1}x_p$$

The new, enhanced data vectors are of the form:

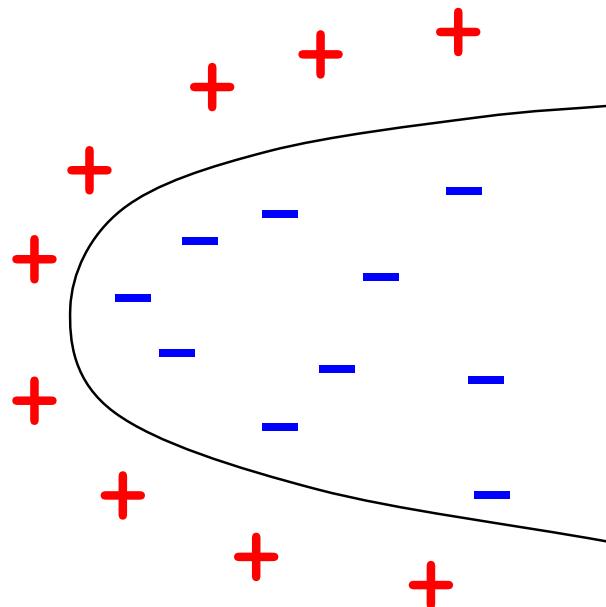
$$\Phi(x) = (1, \sqrt{2}x_1, \dots, \sqrt{2}x_p, x_1^2, \dots, x_p^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_{p-1}x_p).$$

Adding new features



Actual boundary is something like $x_1 = x_2^2 + 5$.

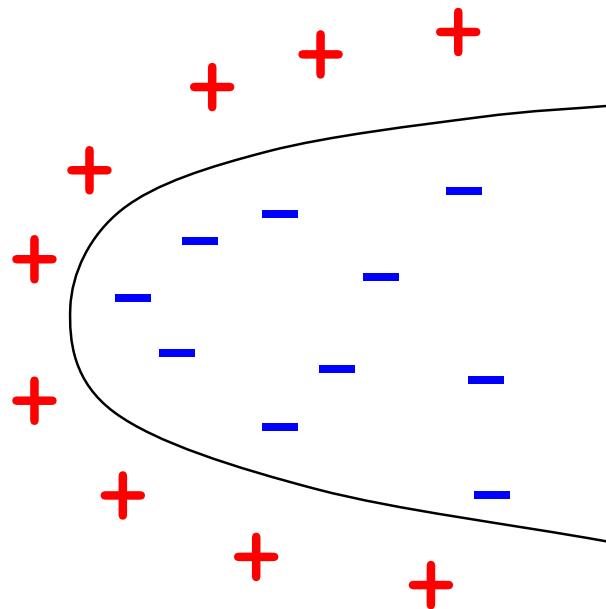
Adding new features



Actual boundary is something like $x_1 = x_2^2 + 5$.

- This is quadratic in $x = (1, x_1, x_2)$
- But it is linear in $\Phi(x) = (1, x_1, x_2, x_1^2, x_2^2, x_1 x_2)$

Adding new features



Actual boundary is something like $x_1 = x_2^2 + 5$.

- This is quadratic in $x = (1, x_1, x_2)$
- But it is linear in $\Phi(x) = (1, x_1, x_2, x_1^2, x_2^2, x_1 x_2)$

By embedding the data in a higher-dimensional feature space, we can keep using a linear classifier!

Quick quiz

- 1 Suppose $x = (1, x_1, x_2, x_3)$. What is the dimension of $\Phi(x)$?

Quick quiz

- 1 Suppose $x = (1, x_1, x_2, x_3)$. What is the dimension of $\Phi(x)$?

Ten

Features: $1, x_1, x_2, x_3, x_1^2, x_2^2, x_3^2, x_1x_2, x_2x_3, x_1x_3$

Quick quiz

- ① Suppose $x = (1, x_1, x_2, x_3)$. What is the dimension of $\Phi(x)$?

Ten

Features: $1, x_1, x_2, x_3, x_1^2, x_2^2, x_3^2, x_1x_2, x_2x_3, x_1x_3$

- ② What if $x = (1, x_1, \dots, x_p)$?

Quick quiz

- 1 Suppose $x = (1, x_1, x_2, x_3)$. What is the dimension of $\Phi(x)$?

Ten

Features: 1, x_1 , x_2 , x_3 , x_1^2 , x_2^2 , x_3^2 , x_1x_2 , x_2x_3 , x_1x_3

- 2 What if $x = (1, x_1, \dots, x_p)$?

Use the formula – $(1 + 2p + (p^2))$

Perceptron revisited

Learning in the higher-dimensional feature space:

- $w = 0$
- while some $y (w \cdot \Phi(x)) < 0$:
 - $w = w + y \Phi(x)$

Final w is a weighted linear sum of various $\Phi(x)$.

Perceptron revisited

Learning in the higher-dimensional feature space:

- $w = 0$
- while some $y (w \cdot \Phi(x)) < 0$:
 - $w = w + y \Phi(x)$

Final w is a weighted linear sum of various $\Phi(x)$.

Problem: number of features has now increased dramatically.

For MNIST, from 784 to 307720!

Perceptron revisited

Learning in the higher-dimensional feature space:

- $w = 0$
- while some $y (w \cdot \Phi(x)) < 0$:
 - $w = w + y \Phi(x)$

Final w is a weighted linear sum of various $\Phi(x)$.

Problem: number of features has now increased dramatically.

For MNIST, from 784 to 307720!

The **kernel trick** (Aizenman, Braverman, Rozonoer (1964)):

- The only time we ever access $\Phi(x)$ is to compute $w \cdot \Phi(x)$. If, say,

$$w = a_1 \Phi(x^{(1)}) + a_2 \Phi(x^{(2)}) + a_3 \Phi(x^{(3)}),$$

then $w \cdot \Phi(x)$ is a weighted sum of dot products $\Phi(x) \cdot \Phi(x^{(i)})$.

Perceptron revisited

Learning in the higher-dimensional feature space:

- $w = 0$
- while some $y (w \cdot \Phi(x)) < 0$:
 - $w = w + y \Phi(x)$

Final w is a weighted linear sum of various $\Phi(x)$.

Problem: number of features has now increased dramatically.

For MNIST, from 784 to 307720!

The **kernel trick** (Aizenman, Braverman, Rozonoer (1964)):

- The only time we ever access $\Phi(x)$ is to compute $w \cdot \Phi(x)$. If, say,

$$w = a_1 \Phi(x^{(1)}) + a_2 \Phi(x^{(2)}) + a_3 \Phi(x^{(3)}),$$

then $w \cdot \Phi(x)$ is a weighted sum of dot products $\Phi(x) \cdot \Phi(x^{(i)})$.

- Can we compute such dot products without writing out the $\Phi(x)$'s?

Computing dot products

In 2-d:

$$\phi(x) \cdot \phi(z)$$

$$= (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2) \cdot (1, \sqrt{2}z_1, \sqrt{2}z_2, z_1^2, z_2^2, \sqrt{2}z_1z_2)$$

$$= 1 + 2x_1z_1 + 2x_2z_2 + x_1^2z_1^2 + x_2^2z_2^2 + 2x_1x_2z_1z_2$$

$$= (1 + x_1z_1 + x_2z_2)^2 = (1 + x \cdot z)^2$$

Computing dot products

In 2-d:

$$\begin{aligned}\Phi(x) \cdot \Phi(z) &= (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2) \cdot (1, \sqrt{2}z_1, \sqrt{2}z_2, z_1^2, z_2^2, \sqrt{2}z_1z_2) \\ &= 1 + 2x_1z_1 + 2x_2z_2 + x_1^2z_1^2 + x_2^2z_2^2 + 2x_1x_2z_1z_2 \\ &= (1 + x_1z_1 + x_2z_2)^2 = (1 + x \cdot z)^2\end{aligned}$$

In p dimensions:

$$\begin{aligned}\Phi(x) \cdot \Phi(z) &= (1, \sqrt{2}x_1, \dots, \sqrt{2}x_p, x_1^2, \dots, x_p^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_{p-1}x_p) \cdot \\ &\quad (1, \sqrt{2}z_1, \dots, \sqrt{2}z_p, z_1^2, \dots, z_p^2, \sqrt{2}z_1z_2, \dots, \sqrt{2}z_{p-1}z_p) \\ &= 1 + 2 \sum_i x_i z_i + \sum_i x_i^2 z_i^2 + 2 \sum_{i \neq j} x_i x_j z_i z_j \\ &= (1 + x_1 z_1 + \dots + x_p z_p)^2 = (1 + x \cdot z)^2\end{aligned}$$

Computing dot products

In 2-d:

$$\begin{aligned}\Phi(x) \cdot \Phi(z) &= (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2) \cdot (1, \sqrt{2}z_1, \sqrt{2}z_2, z_1^2, z_2^2, \sqrt{2}z_1z_2) \\ &= 1 + 2x_1z_1 + 2x_2z_2 + x_1^2z_1^2 + x_2^2z_2^2 + 2x_1x_2z_1z_2 \\ &= (1 + x_1z_1 + x_2z_2)^2 = (1 + x \cdot z)^2\end{aligned}$$

In p dimensions:

$$\begin{aligned}\Phi(x) \cdot \Phi(z) &= (1, \sqrt{2}x_1, \dots, \sqrt{2}x_p, x_1^2, \dots, x_p^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_{p-1}x_p) \cdot \\ &\quad (1, \sqrt{2}z_1, \dots, \sqrt{2}z_p, z_1^2, \dots, z_p^2, \sqrt{2}z_1z_2, \dots, \sqrt{2}z_{p-1}z_p) \\ &= 1 + 2 \sum_i x_i z_i + \sum_i x_i^2 z_i^2 + 2 \sum_{i \neq j} x_i x_j z_i z_j \\ &= (1 + x_1 z_1 + \dots + x_p z_p)^2 = (1 + x \cdot z)^2\end{aligned}$$

For MNIST: computing dot products in the 307720-dimensional quadratic feature space takes time proportional to just 784, the original dimension!

The kernel trick

Why does it work?

The kernel trick

Why does it work?

- ➊ The only time we ever look at data – during training or subsequent classification – is to compute dot products $w \cdot \Phi(x)$.

The kernel trick

Why does it work?

- ➊ The only time we ever look at data – during training or subsequent classification – is to compute dot products $w \cdot \Phi(x)$.
- ➋ And w itself is a linear combination of $\Phi(x)$'s. If, say,

$$w = \alpha_1 \Phi(x^{(1)}) + \alpha_{22} \Phi(x^{(22)}) + \alpha_{37} \Phi(x^{(37)}),$$

we can store w as $[(1, \alpha_1), (22, \alpha_{22}), (37, \alpha_{37})]$.

The kernel trick

Why does it work?

- ➊ The only time we ever look at data – during training or subsequent classification – is to compute dot products $w \cdot \Phi(x)$.
- ➋ And w itself is a linear combination of $\Phi(x)$'s. If, say,

$$w = \alpha_1 \Phi(x^{(1)}) + \alpha_{22} \Phi(x^{(22)}) + \alpha_{37} \Phi(x^{(37)}),$$

we can store w as $[(1, \alpha_1), (22, \alpha_{22}), (37, \alpha_{37})]$.

- ➌ Dot products $\Phi(x) \cdot \Phi(x')$ can be computed efficiently, without ever writing out the high-dimensional embedding $\Phi(\cdot)$.

Kernel Perceptron

Learning from data $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}) \in \mathcal{X} \times \{-1, 1\}$

Primal form:

- $w = 0$
- while there is some i with $y^{(i)}(w \cdot \Phi(x^{(i)})) < 0$:
 - $w = w + y^{(i)} \Phi(x^{(i)})$

Kernel Perceptron

Learning from data $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}) \in \mathcal{X} \times \{-1, 1\}$

Primal form:

- $w = 0$
- while there is some i with $y^{(i)}(w \cdot \Phi(x^{(i)})) < 0$:
 - $w = w + y^{(i)} \Phi(x^{(i)})$

Dual form: $w = \sum_j \alpha_j y^{(j)} \Phi(x^{(j)})$, where $\alpha \in \mathbb{R}^n$.

- $\alpha = 0$
- while there is some i with $y^{(i)} \left(\underbrace{\sum_j \alpha_j y^{(j)} \Phi(x^{(j)}) \cdot \Phi(x^{(i)})}_{\text{efficient}} \right) < 0$:
 - $\alpha_i = \alpha_i + 1$

Kernel Perceptron

Learning from data $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}) \in \mathcal{X} \times \{-1, 1\}$

Primal form:

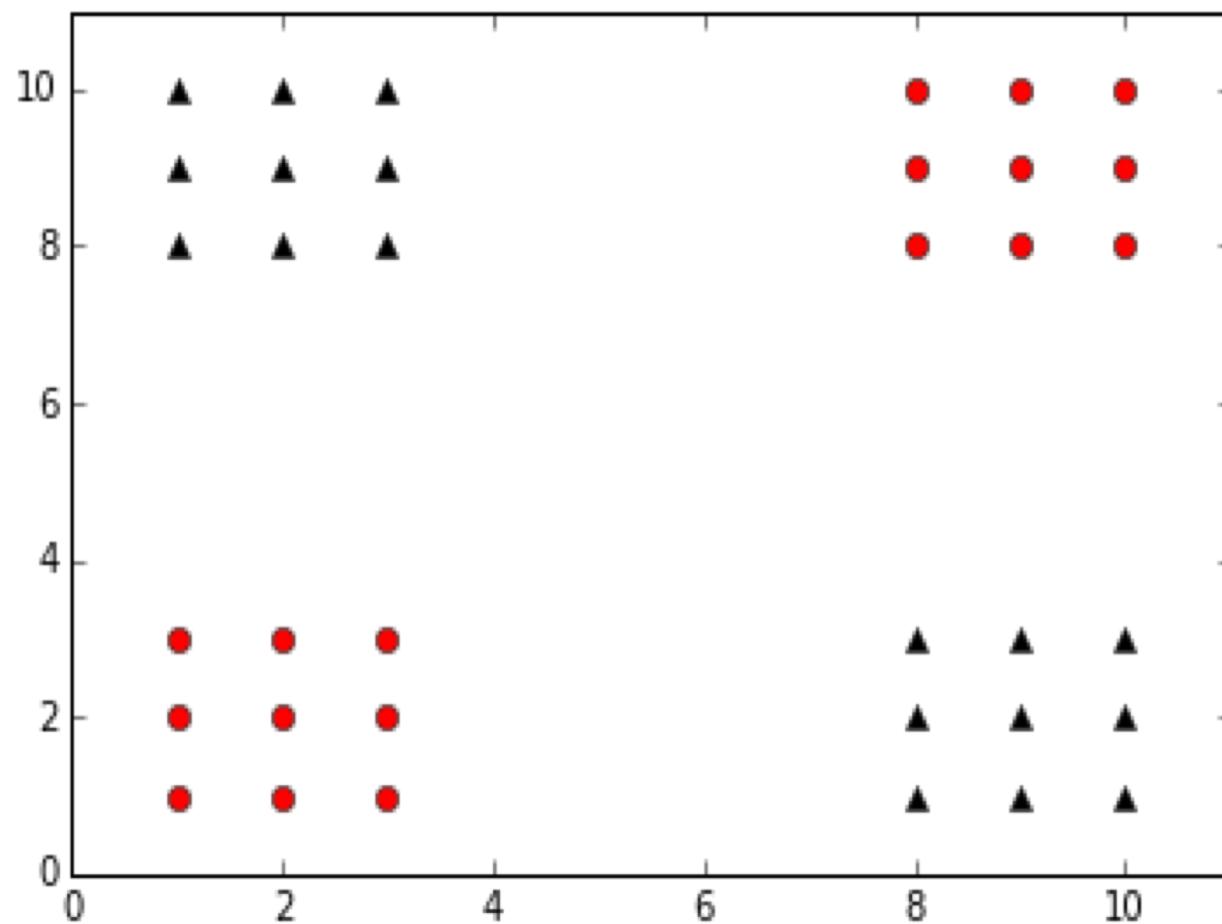
- $w = 0$
- while there is some i with $y^{(i)}(w \cdot \Phi(x^{(i)})) < 0$:
 - $w = w + y^{(i)} \Phi(x^{(i)})$

Dual form: $w = \sum_j \alpha_j y^{(j)} \Phi(x^{(j)})$, where $\alpha \in \mathbb{R}^n$.

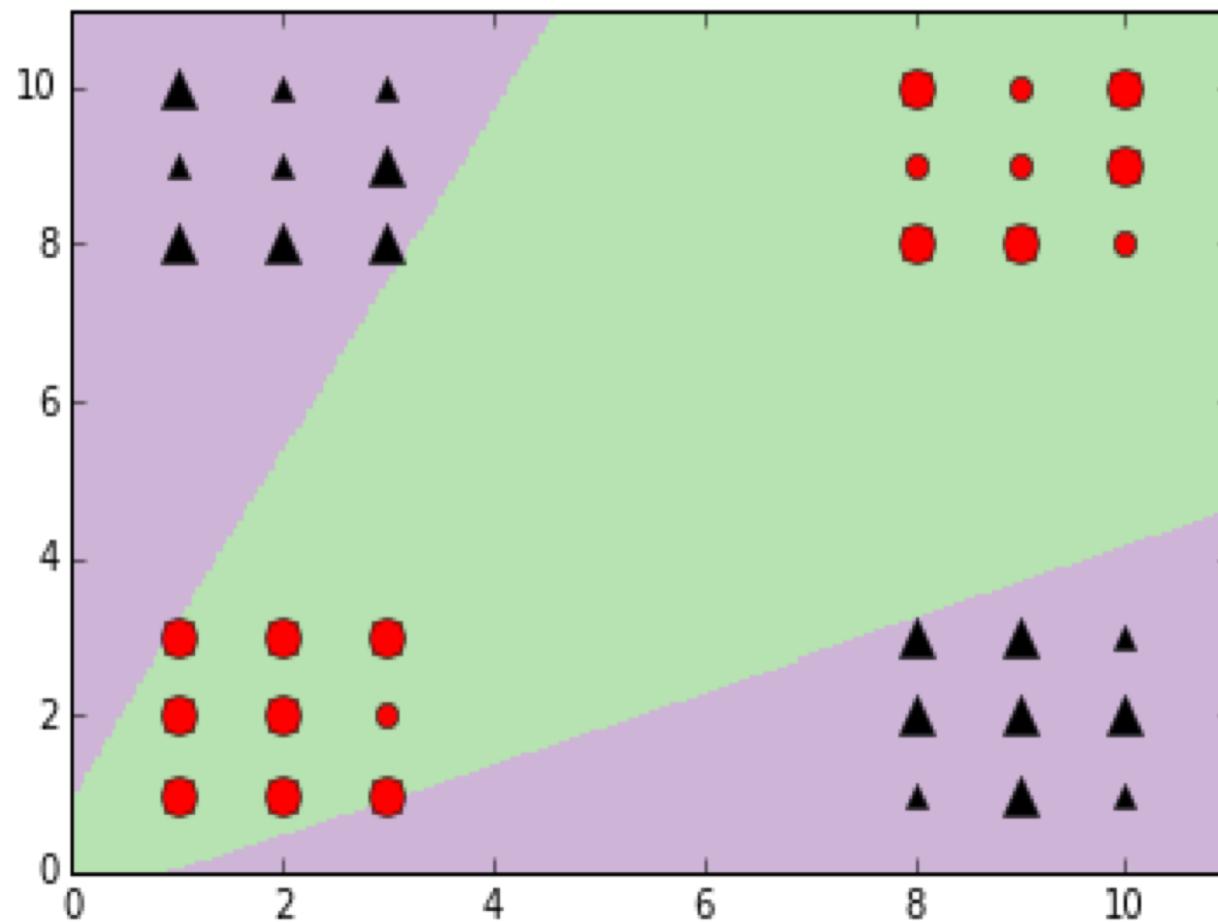
- $\alpha = 0$
- while there is some i with $y^{(i)} \left(\underbrace{\sum_j \alpha_j y^{(j)} \Phi(x^{(j)}) \cdot \Phi(x^{(i)})}_{\text{efficient}} \right) < 0$:
 - $\alpha_i = \alpha_i + 1$

To classify a new point x : Return $\text{sign} \left(\sum_j \alpha_j y^{(j)} \Phi(x^{(j)}) \cdot \Phi(x) \right)$.

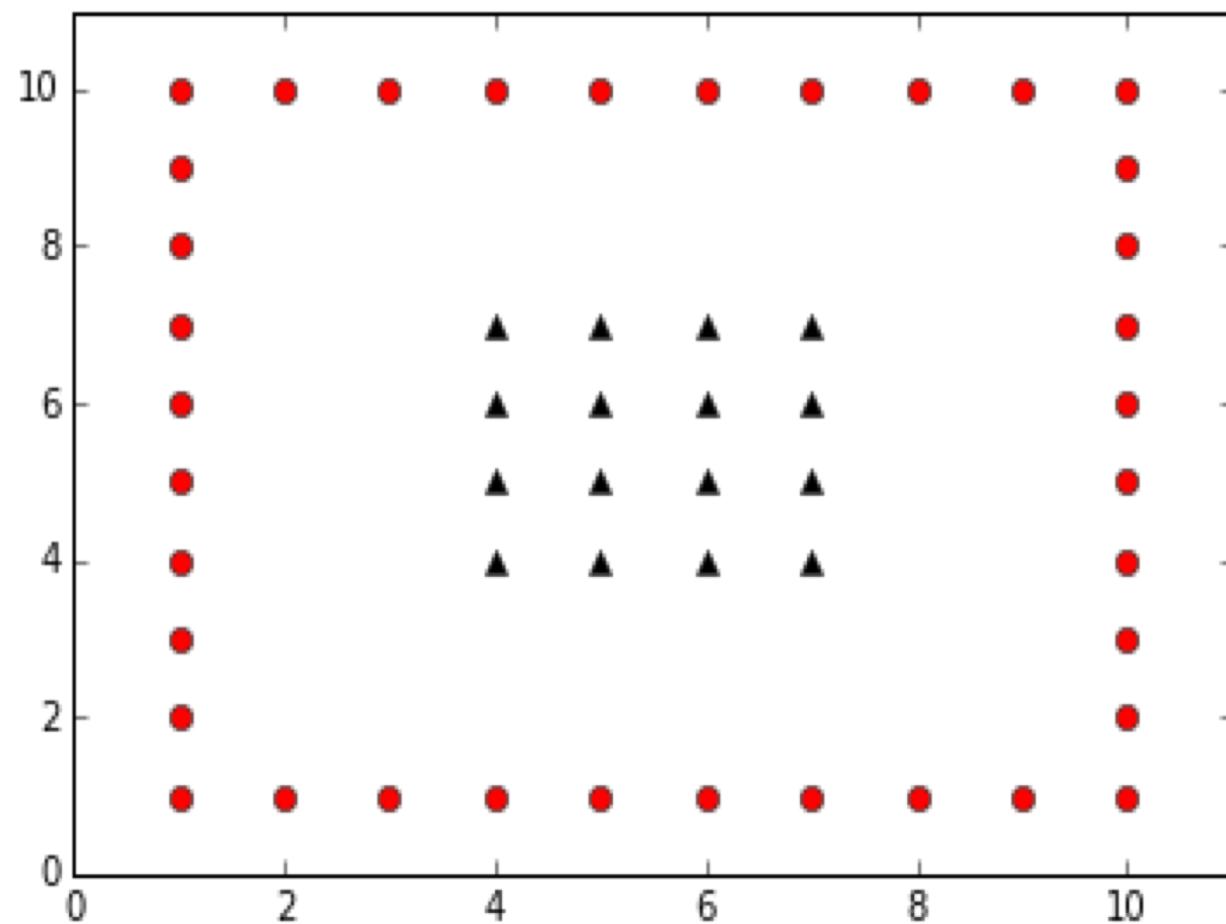
Kernel Perceptron: Examples



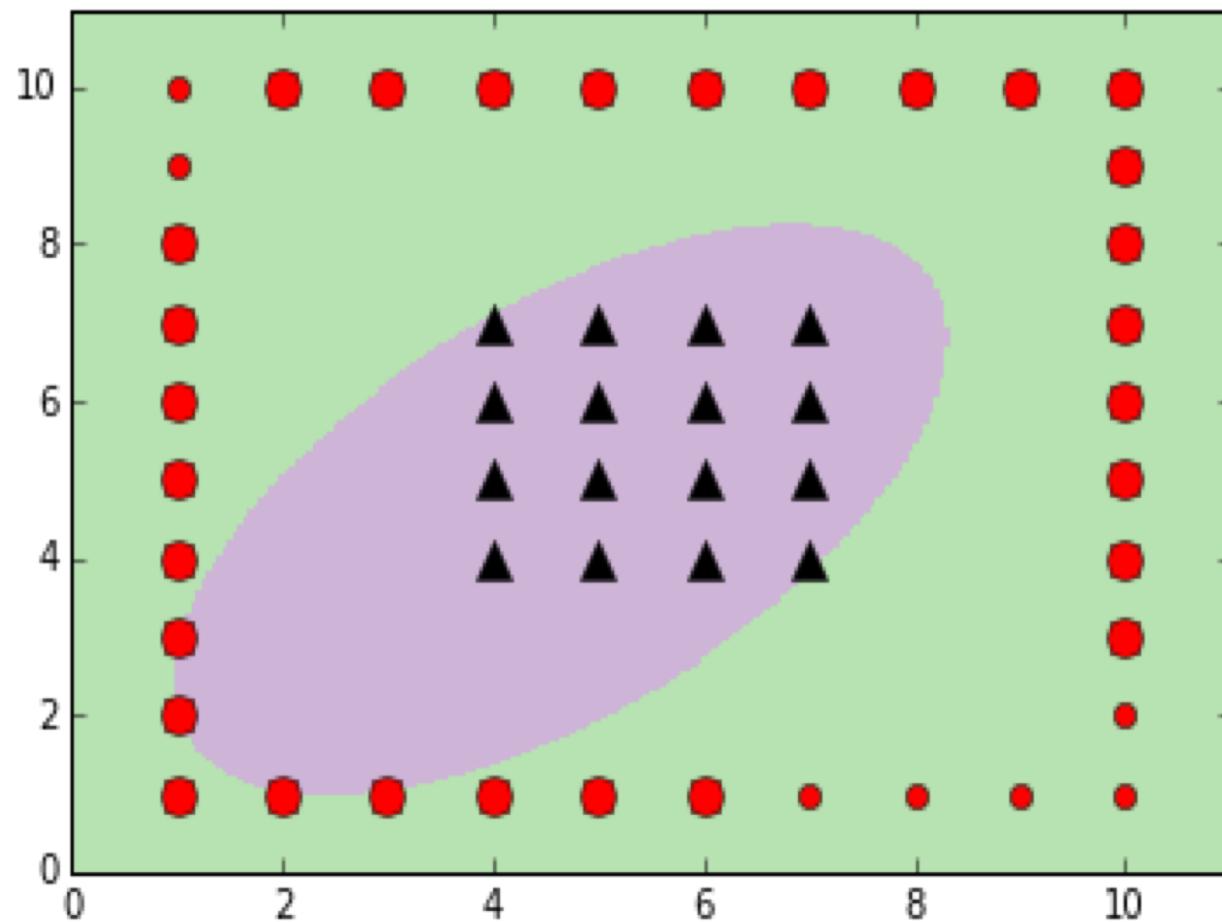
Kernel Perceptron: Examples



Kernel Perceptron: Examples



Kernel Perceptron: Examples



Quick quiz

Recall the kernel perceptron algorithm:

- $\alpha = 0$
- while there is some i with $y^{(i)}(\sum_j \alpha_j y^{(j)} \Phi(x^{(j)}) \cdot \Phi(x^{(i)})) < 0$:
 - $\alpha_i = \alpha_i + 1$

Suppose we run it on a data set and find that it converges after k updates.

- ➊ How many dot product computations are needed to classify a new point?

Quick quiz

Recall the kernel perceptron algorithm:

- $\alpha = 0$
- while there is some i with $y^{(i)}(\sum_j \alpha_j y^{(j)} \Phi(x^{(j)}) \cdot \Phi(x^{(i)})) < 0$:
 - $\alpha_i = \alpha_i + 1$

Suppose we run it on a data set and find that it converges after k updates.

- ① How many dot product computations are needed to classify a new point?

Solution: At most k

Quick quiz

Recall the kernel perceptron algorithm:

- $\alpha = 0$
- while there is some i with $y^{(i)}(\sum_j \alpha_j y^{(j)} \Phi(x^{(j)}) \cdot \Phi(x^{(i)})) < 0$:
 - $\alpha_i = \alpha_i + 1$

Suppose we run it on a data set and find that it converges after k updates.

- ➊ How many dot product computations are needed to classify a new point?

Solution: At most k

- ➋ What is the total number of dot product computations during learning?

Quick quiz

Recall the kernel perceptron algorithm:

- $\alpha = 0$
- while there is some i with $y^{(i)}(\sum_j \alpha_j y^{(j)} \Phi(x^{(j)}) \cdot \Phi(x^{(i)})) < 0$:
 - $\alpha_i = \alpha_i + 1$

Suppose we run it on a data set and find that it converges after k updates.

- ① How many dot product computations are needed to classify a new point?

Solution: At most k

- ② What is the total number of dot product computations during learning?

Solution: 1 term after first update, 2 terms after second update and so on till k updates. Therefore, in order of $(1 + 2 + \dots + k) = \text{at least } \binom{k}{2}$ updates.

Does this work with SVMs?

$$\begin{aligned} \text{(PRIMAL)} \quad & \min_{w \in \mathbb{R}^p, b \in \mathbb{R}, \xi \in \mathbb{R}^n} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.: } & y^{(i)}(w \cdot x^{(i)} + b) \geq 1 - \xi_i \quad \text{for all } i = 1, 2, \dots, n \\ & \xi \geq 0 \end{aligned}$$

$$\begin{aligned} \text{(DUAL)} \quad & \max_{\alpha \in \mathbb{R}^n} \sum_{i=1}^n \alpha_i - \sum_{i,j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} (x^{(i)} \cdot x^{(j)}) \\ \text{s.t.: } & \sum_{i=1}^n \alpha_i y^{(i)} = 0 \\ & 0 \leq \alpha_i \leq C \end{aligned}$$

At optimality, $w = \sum_i \alpha_i y^{(i)} x^{(i)}$, with

$$0 < \alpha_i < C \Rightarrow y^{(i)}(w \cdot x^{(i)} + b) = 1$$

$$\alpha_i = C \Rightarrow y^{(i)}(w \cdot x^{(i)} + b) = 1 - \xi_i$$

Kernel SVM

- ① **Embedding.** Pick a mapping $x \mapsto \Phi(x)$.

Kernel SVM

① **Embedding.** Pick a mapping $x \mapsto \Phi(x)$.

② **Learning.** Solve the dual problem:

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^n} \quad & \sum_{i=1}^n \alpha_i - \sum_{i,j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} (\underbrace{\Phi(x^{(i)}) \cdot \Phi(x^{(j)})}_{\text{efficient}}) \\ \text{s.t.:} \quad & \sum_{i=1}^n \alpha_i y^{(i)} = 0 \\ & 0 \leq \alpha_i \leq C \end{aligned}$$

This yields $w = \sum_i \alpha_i y^{(i)} \Phi(x^{(i)})$.

Offset b is obtained from the complementary slackness conditions.

Kernel SVM

① **Embedding.** Pick a mapping $x \mapsto \Phi(x)$.

② **Learning.** Solve the dual problem:

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^n} \quad & \sum_{i=1}^n \alpha_i - \sum_{i,j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} (\underbrace{\Phi(x^{(i)}) \cdot \Phi(x^{(j)})}_{\text{efficient}}) \\ \text{s.t.:} \quad & \sum_{i=1}^n \alpha_i y^{(i)} = 0 \\ & 0 \leq \alpha_i \leq C \end{aligned}$$

This yields $w = \sum_i \alpha_i y^{(i)} \Phi(x^{(i)})$.

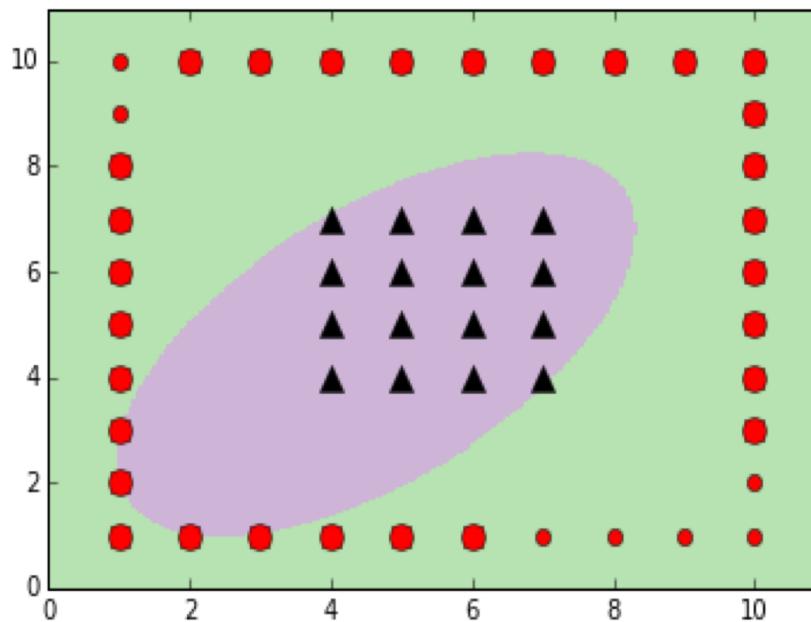
Offset b is obtained from the complementary slackness conditions.

③ **Classification.** Given a new point x , classify as

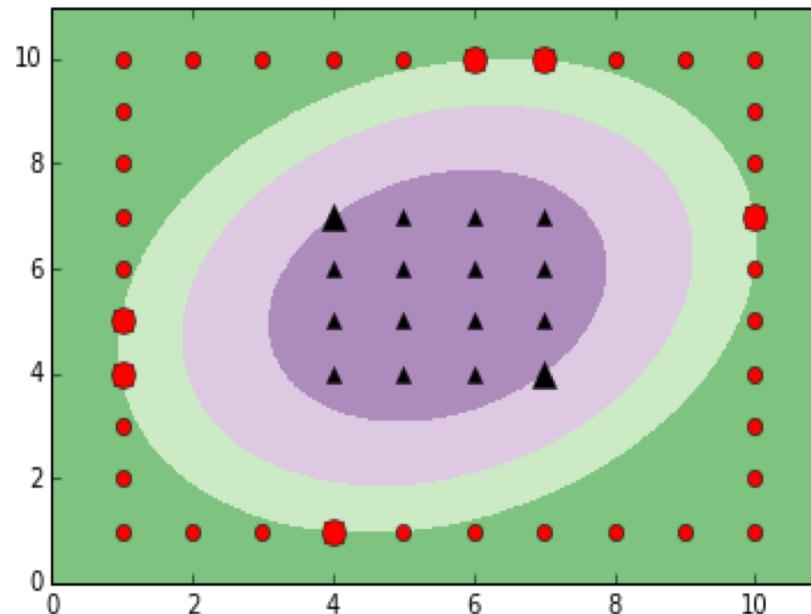
$$\text{sign} \left(\sum_i \alpha_i y^{(i)} (\underbrace{\Phi(x^{(i)}) \cdot \Phi(x)}_{\text{again}}) + b \right).$$

Kernel Perceptron vs. Kernel SVM: examples

Perceptron:

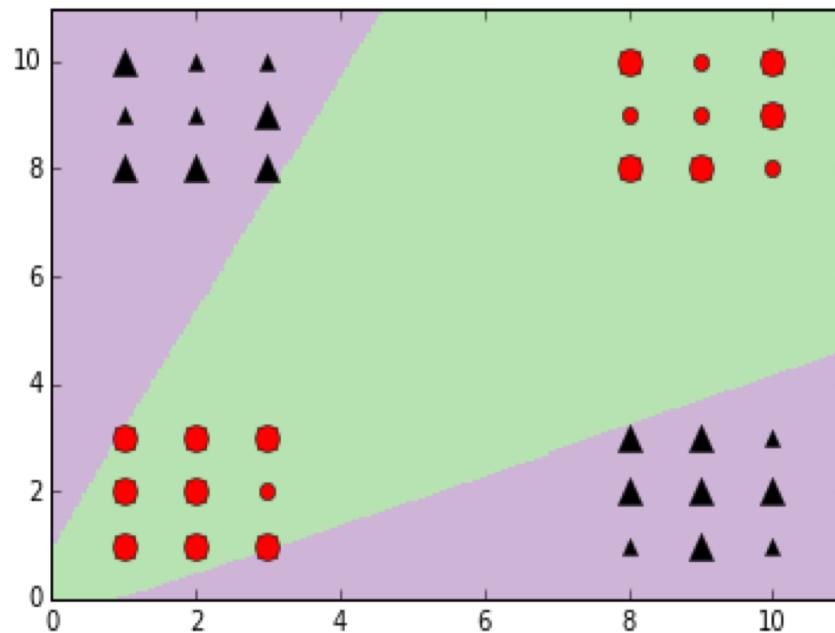


SVM:

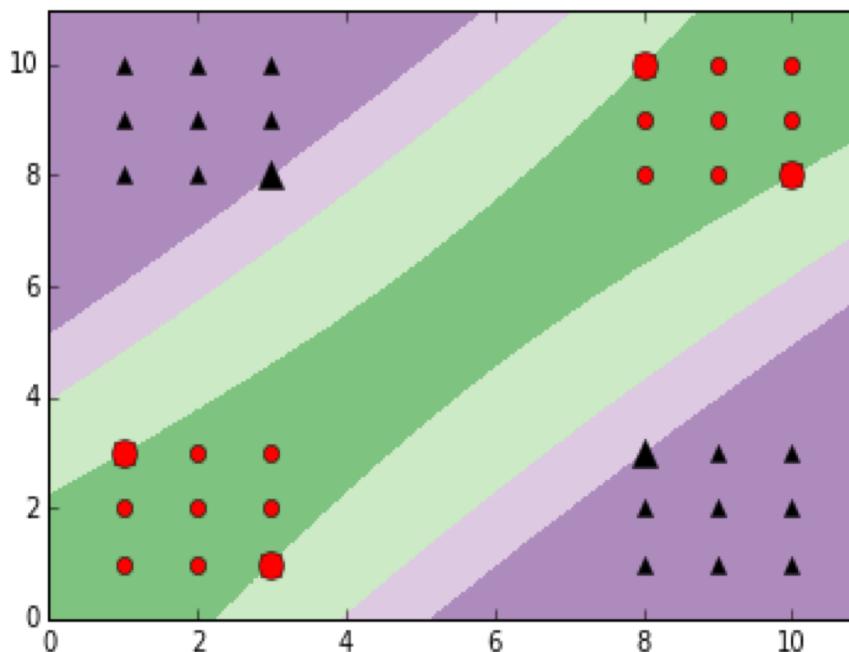


Kernel Perceptron vs. Kernel SVM: examples

Perceptron:

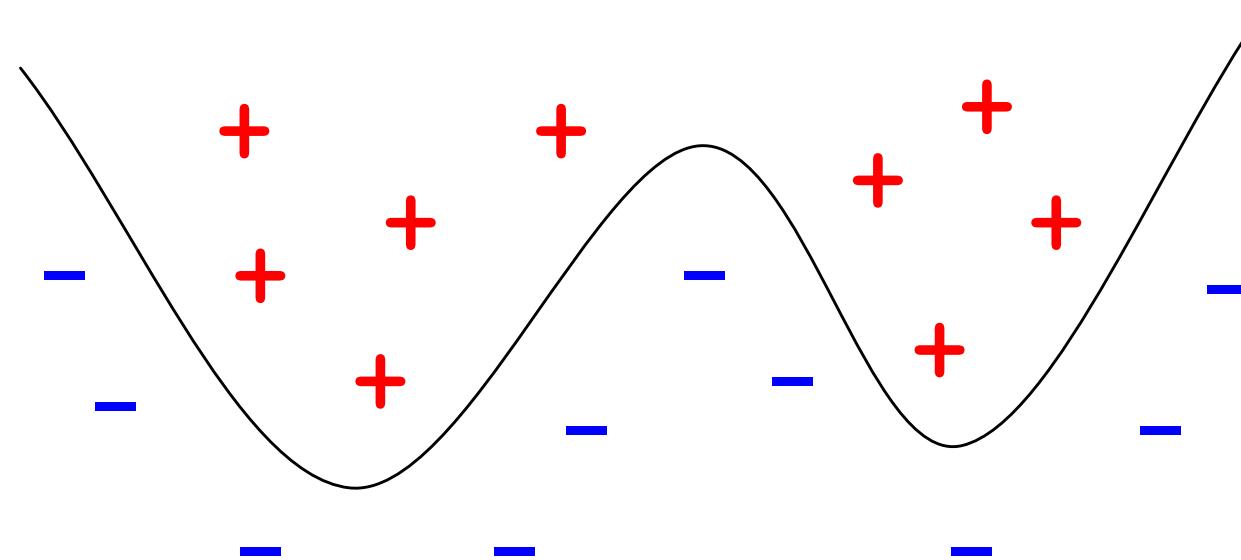


SVM:



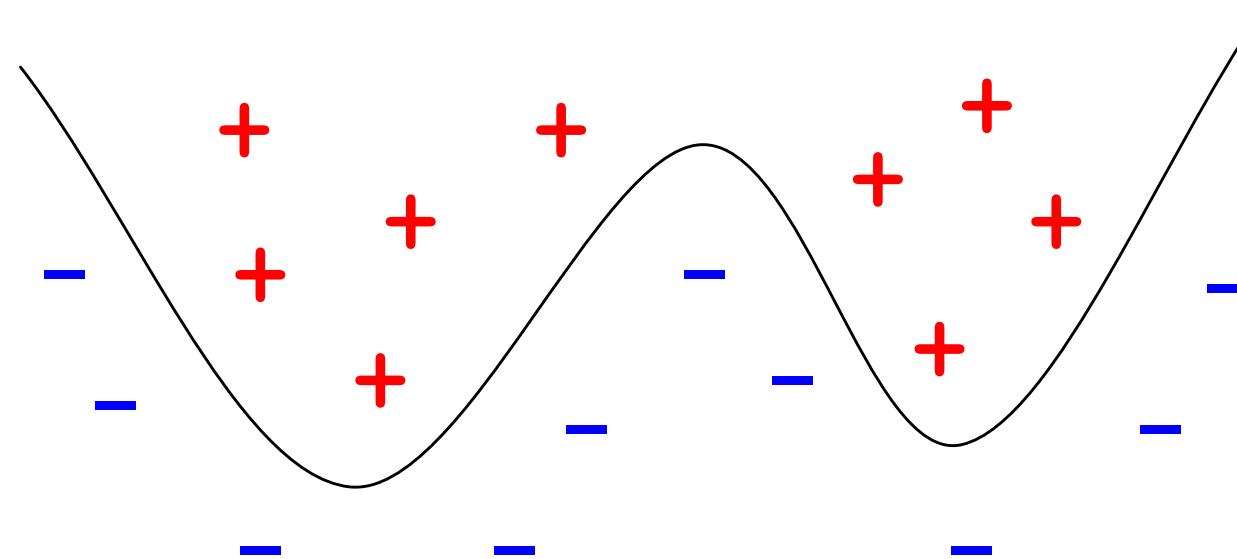
Polynomial decision boundaries

To get a decision surface which is an arbitrary polynomial of order d :



Polynomial decision boundaries

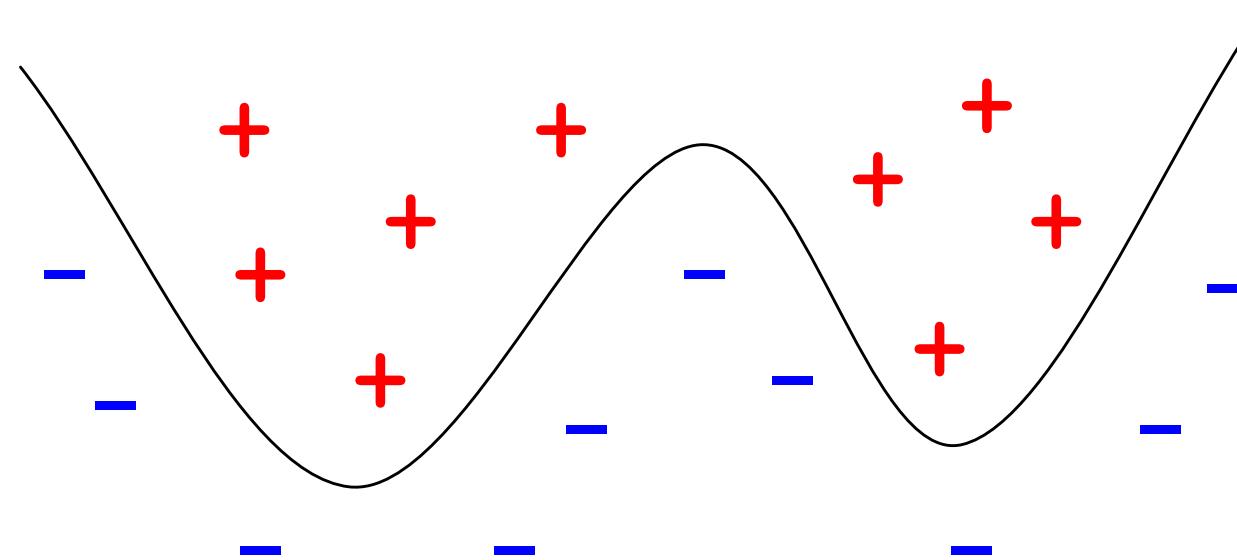
To get a decision surface which is an arbitrary polynomial of order d :



- Let $\Phi(x)$ consist of all terms of order $\leq d$, such as $x_1x_2^2x_3^{d-3}$.
(How many such terms are there, roughly?)

Polynomial decision boundaries

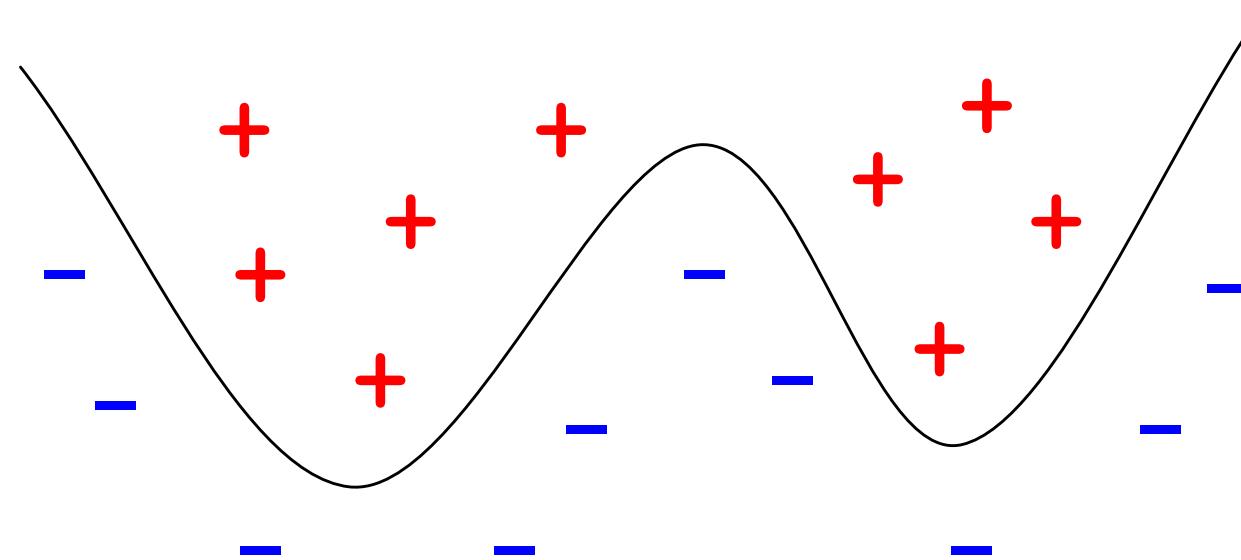
To get a decision surface which is an arbitrary polynomial of order d :



- Let $\Phi(x)$ consist of all terms of order $\leq d$, such as $x_1x_2^2x_3^{d-3}$.
(How many such terms are there, roughly?)
- Same trick works: $\Phi(x) \cdot \Phi(z) = (1 + x \cdot z)^d$.

Polynomial decision boundaries

To get a decision surface which is an arbitrary polynomial of order d :



- Let $\Phi(x)$ consist of all terms of order $\leq d$, such as $x_1x_2^2x_3^{d-3}$.
(How many such terms are there, roughly?)
- Same trick works: $\Phi(x) \cdot \Phi(z) = (1 + x \cdot z)^d$.

The **kernel function**, a measure of similarity:

$$k(x, z) = \Phi(x) \cdot \Phi(z).$$

The kernel function

Now shift attention:

- away from the embedding $\Phi(x)$, which we never explicitly construct,
- towards the thing we actually use: the **similarity measure** $k(x, z)$

Rewrite learning algorithm and final classifier in terms of k .

The kernel function

Now shift attention:

- away from the embedding $\Phi(x)$, which we never explicitly construct,
- towards the thing we actually use: the **similarity measure** $k(x, z)$

Rewrite learning algorithm and final classifier in terms of k .

The classifier $w \cdot \Phi(x)$, for

$$w = \alpha_1 y^{(1)} \Phi(x^{(1)}) + \cdots + \alpha_n y^{(n)} \Phi(x^{(n)}),$$

becomes a **similarity-weighted vote**,

$$F(x) = \alpha_1 y^{(1)} k(x^{(1)}, x) + \cdots + \alpha_n y^{(n)} k(x^{(n)}, x).$$

The kernel function

Now shift attention:

- away from the embedding $\Phi(x)$, which we never explicitly construct,
- towards the thing we actually use: the **similarity measure** $k(x, z)$

Rewrite learning algorithm and final classifier in terms of k .

The classifier $w \cdot \Phi(x)$, for

$$w = \alpha_1 y^{(1)} \Phi(x^{(1)}) + \cdots + \alpha_n y^{(n)} \Phi(x^{(n)}),$$

becomes a **similarity-weighted vote**,

$$F(x) = \alpha_1 y^{(1)} k(x^{(1)}, x) + \cdots + \alpha_n y^{(n)} k(x^{(n)}, x).$$

Can we choose k to be any similarity function suitable for the application domain?

Quick quiz

For each of the following kernel functions, find a corresponding embedding $x \rightarrow \Phi(x)$.

- ① $k(x, z) = \cos(\text{angle between } x \text{ and } z)$.

Quick quiz

For each of the following kernel functions, find a corresponding embedding $x \rightarrow \Phi(x)$.

① $k(x, z) = \cos(\text{angle between } x \text{ and } z)$.

1. We know, $x.z = |x||z|\cos\theta \Rightarrow \cos\theta = \frac{x.z}{|x||z|} \Rightarrow \phi(x) = \frac{x}{|x|}$

Quick quiz

For each of the following kernel functions, find a corresponding embedding $x \rightarrow \Phi(x)$.

① $k(x, z) = \cos(\text{angle between } x \text{ and } z)$.

1. We know, $x.z = |x||z|\cos\theta \Rightarrow \cos\theta = \frac{x.z}{|x||z|} \Rightarrow \phi(x) = \frac{x}{|x|}$

② $k(x, z) = x^T A z$, where A is a symmetric positive semidefinite matrix.

Quick quiz

For each of the following kernel functions, find a corresponding embedding $x \rightarrow \Phi(x)$.

① $k(x, z) = \cos(\text{angle between } x \text{ and } z)$.

1. We know, $x.z = |x||z|\cos\theta \Rightarrow \cos\theta = \frac{x.z}{|x||z|} \Rightarrow \phi(x) = \frac{x}{|x|}$

② $k(x, z) = x^T A z$, where A is a symmetric positive semidefinite matrix.

2. Since A is PSD, we can re-write it as UU^T

$$\begin{aligned} x^T A z &= x^T U U^T z = (U^T x)^T (z U)^T \\ &\Rightarrow \phi(x) = U^T x \end{aligned}$$

Mercer's condition

A function $k : X \times X \rightarrow \mathbb{R}$ is a valid kernel function if it corresponds to some embedding: that is, if there exists Φ defined on X such that

$$k(x, z) = \Phi(x) \cdot \Phi(z).$$

Mercer's condition

A function $k : X \times X \rightarrow \mathbb{R}$ is a valid kernel function if it corresponds to some embedding: that is, if there exists Φ defined on X such that

$$k(x, z) = \Phi(x) \cdot \Phi(z).$$

Mercer (1909): This is equivalent to requiring that for any finite subset $\{x^{(1)}, \dots, x^{(m)}\} \subset \mathcal{X}$, the $m \times m$ similarity matrix

$$K_{ij} = k(x^{(i)}, x^{(j)})$$

is positive semidefinite.

Mercer's condition

A function $k : X \times X \rightarrow \mathbb{R}$ is a valid kernel function if it corresponds to some embedding: that is, if there exists Φ defined on X such that

$$k(x, z) = \Phi(x) \cdot \Phi(z).$$

Mercer (1909): This is equivalent to requiring that for any finite subset $\{x^{(1)}, \dots, x^{(m)}\} \subset \mathcal{X}$, the $m \times m$ similarity matrix

$$K_{ij} = k(x^{(i)}, x^{(j)})$$

is positive semidefinite.

Do you see why this matrix is p.s.d. if $k(x, z) = \Phi(x) \cdot \Phi(z)$?

Mercer's condition

A function $k : X \times X \rightarrow \mathbb{R}$ is a valid kernel function if it corresponds to some embedding: that is, if there exists Φ defined on X such that

$$k(x, z) = \Phi(x) \cdot \Phi(z).$$

Mercer (1909): This is equivalent to requiring that for any finite subset $\{x^{(1)}, \dots, x^{(m)}\} \subset \mathcal{X}$, the $m \times m$ similarity matrix

$$K_{ij} = k(x^{(i)}, x^{(j)})$$

is positive semidefinite.

Do you see why this matrix is p.s.d. if $k(x, z) = \Phi(x) \cdot \Phi(z)$? UU^T form matrix

Mercer's condition

A function $k : X \times X \rightarrow \mathbb{R}$ is a valid kernel function if it corresponds to some embedding: that is, if there exists Φ defined on X such that

$$k(x, z) = \Phi(x) \cdot \Phi(z).$$

Mercer (1909): This is equivalent to requiring that for any finite subset $\{x^{(1)}, \dots, x^{(m)}\} \subset \mathcal{X}$, the $m \times m$ similarity matrix

$$K_{ij} = k(x^{(i)}, x^{(j)})$$

is positive semidefinite.

Do you see why this matrix is p.s.d. if $k(x, z) = \Phi(x) \cdot \Phi(z)$? UU^T form matrix

A popular similarity function: the **Gaussian kernel** or **RBF kernel**

$$k(x, z) = e^{-\|x-z\|^2/2\sigma^2},$$

where σ is an adjustable scale parameter.

RBF Kernel Intuition

The feature space of Radial Basis Function (RBF) kernel or the Gaussian kernel has an **infinite** number of dimensions.

$$k(x, z) = e^{-\|x-z\|^2/2\sigma^2}$$

The numerator may be recognized as the squared Euclidean distance between two feature vectors, and σ is a free parameter.

RBF Kernel Intuition

The feature space of Radial Basis Function (RBF) kernel or the Gaussian kernel has an **infinite** number of dimensions.

$$k(x, z) = e^{-\|x-z\|^2/2\sigma^2}$$

The numerator may be recognized as the squared Euclidean distance between two feature vectors, and σ is a free parameter.

Since the value of the RBF kernel decreases as distance increases and ranges between zero and one (when $x = z$), it has a ready interpretation as a similarity measure – as it has a scale of [0,1].

RBF Kernel Intuition

The feature space of Radial Basis Function (RBF) kernel or the Gaussian kernel has an **infinite** number of dimensions.

$$k(x, z) = e^{-\|x-z\|^2/2\sigma^2}$$

The numerator may be recognized as the squared Euclidean distance between two feature vectors, and σ is a free parameter.

Since the value of the RBF kernel decreases as distance increases and ranges between zero and one (when $x = z$), it has a ready interpretation as a similarity measure – as it has a scale of [0,1].

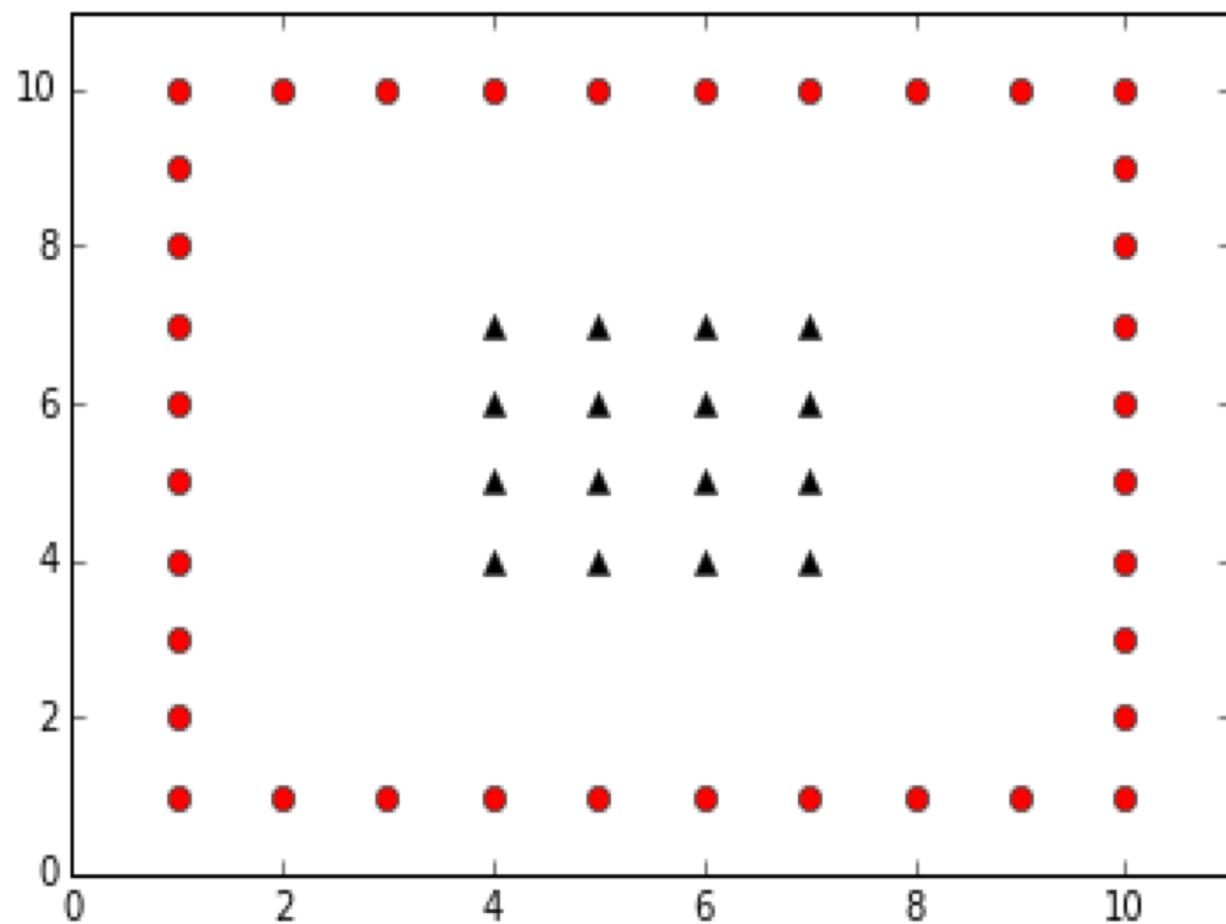
Idea behind infinite dimensions – Taylor series expansion of e^x is given by:

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots$$

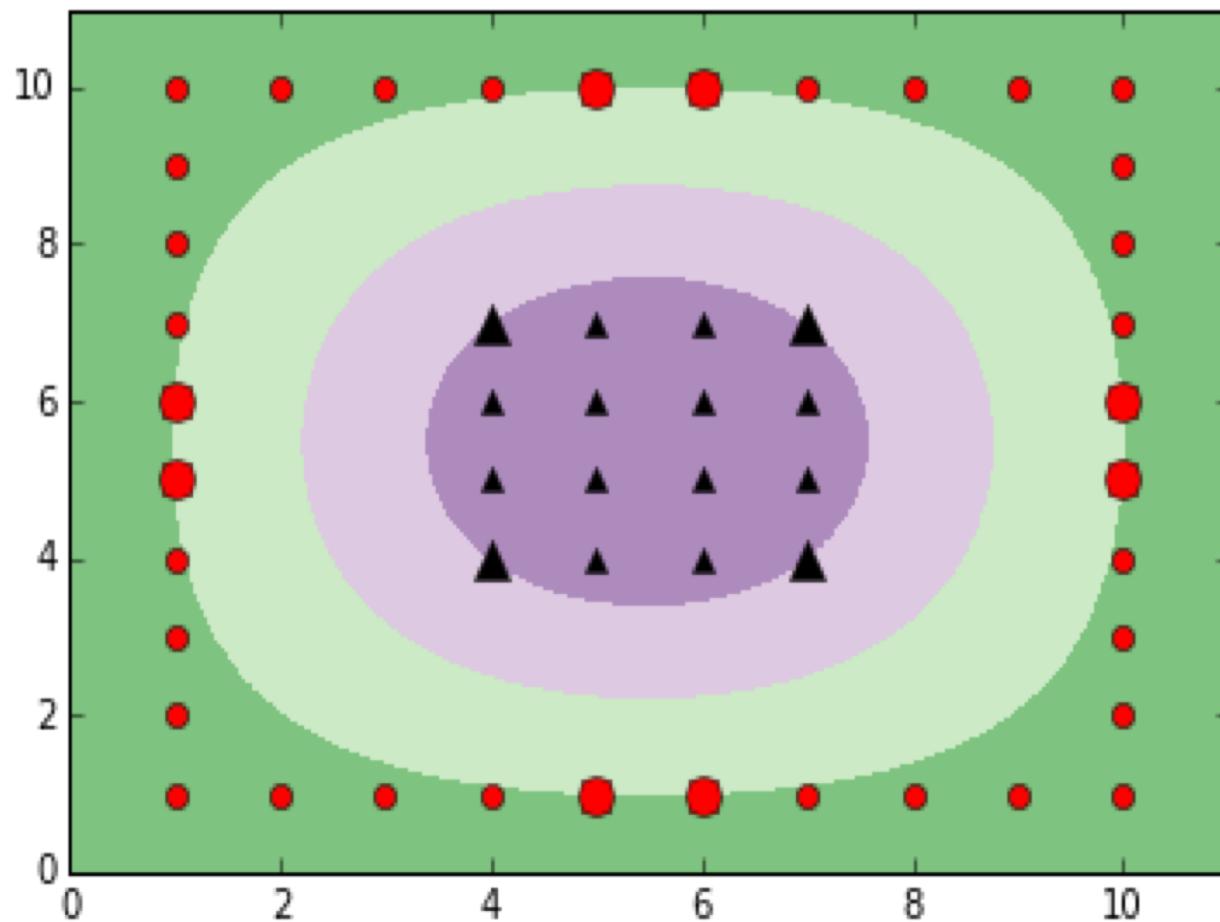
Infinite terms!

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

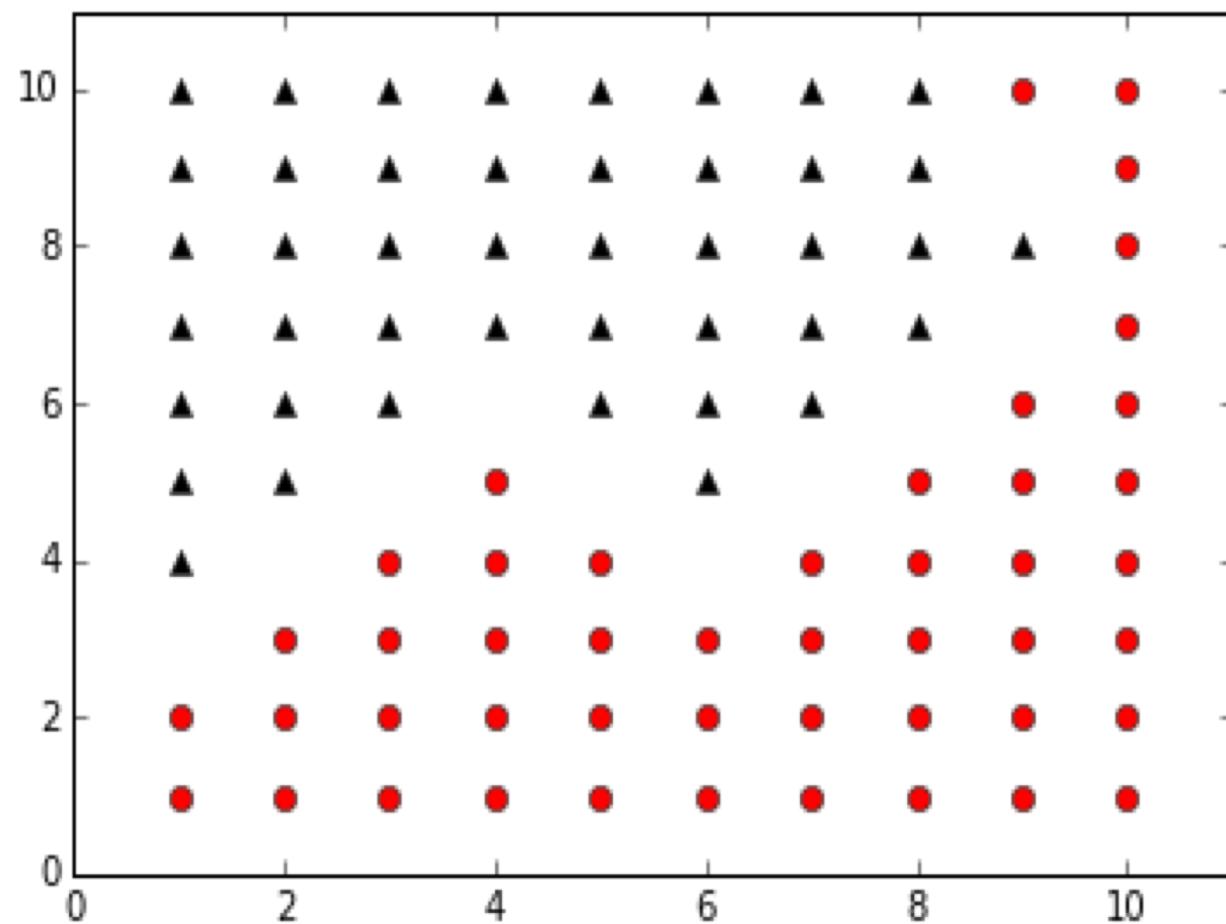
RBF kernel: examples



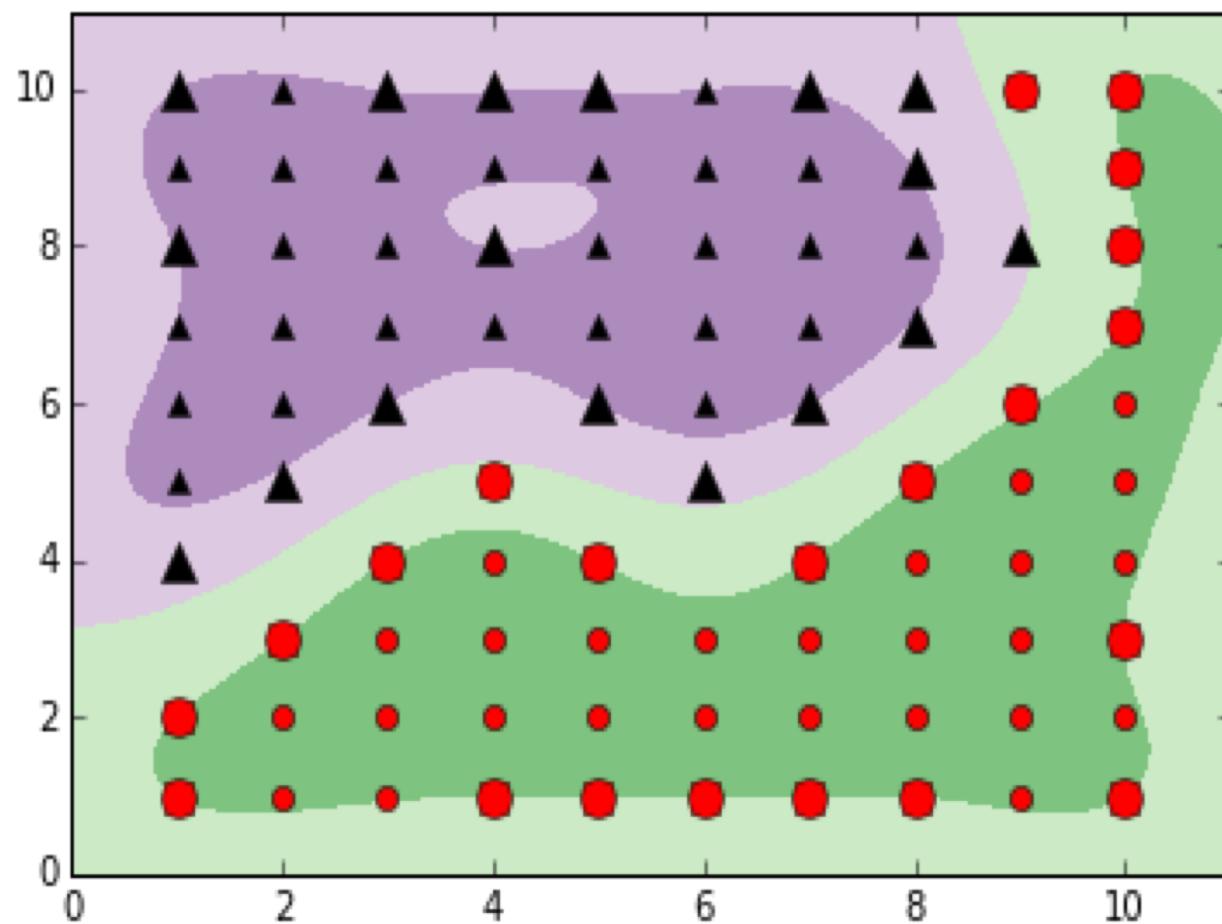
RBF kernel: examples



RBF kernel: examples



RBF kernel: examples



Quick quiz

Recall that the RBF kernel is given by $k(x, z) = e^{-\|x-z\|^2/2\sigma^2}$.

- ① How does this function behave as $\sigma \uparrow \infty$?

Quick quiz

Recall that the RBF kernel is given by $k(x, z) = e^{-\|x-z\|^2/2\sigma^2}$.

- ➊ How does this function behave as $\sigma \uparrow \infty$?
 1. As σ reaches infinity, $k(x, z)$ reaches one. This means everything is similar (called global/alike).

Quick quiz

Recall that the RBF kernel is given by $k(x, z) = e^{-\|x-z\|^2/2\sigma^2}$.

- ➊ How does this function behave as $\sigma \uparrow \infty$?
 1. As σ reaches infinity, $k(x, z)$ reaches one. This means everything is similar (called global/alike).
- ➋ How does this function behave as $\sigma \downarrow 0$?

Quick quiz

Recall that the RBF kernel is given by $k(x, z) = e^{-\|x-z\|^2/2\sigma^2}$.

① How does this function behave as $\sigma \uparrow \infty$?

1. As σ reaches infinity, $k(x, z)$ reaches one. This means everything is similar (called global/alike).

② How does this function behave as $\sigma \downarrow 0$?

2. As σ reaches zero, $k(x, z)$ reaches zero. This means everything is dissimilar and we have a local notion of similarity (as similarity measure of x, z is zero).

Quick quiz

Recall that the RBF kernel is given by $k(x, z) = e^{-\|x-z\|^2/2\sigma^2}$.

① How does this function behave as $\sigma \uparrow \infty$?

1. As σ reaches infinity, $k(x, z)$ reaches one. This means everything is similar (called global/alike).

② How does this function behave as $\sigma \downarrow 0$?

2. As σ reaches zero, $k(x, z)$ reaches zero. This means everything is dissimilar and we have a local notion of similarity (as similarity measure of x, z is zero).

③ As the amount of data increases, would it make sense to increase σ or to decrease it?

Quick quiz

Recall that the RBF kernel is given by $k(x, z) = e^{-\|x-z\|^2/2\sigma^2}$.

① How does this function behave as $\sigma \uparrow \infty$?

- As σ reaches infinity, $k(x, z)$ reaches one. This means everything is similar (called global/alike).

② How does this function behave as $\sigma \downarrow 0$?

- As σ reaches zero, $k(x, z)$ reaches zero. This means everything is dissimilar and we have a local notion of similarity (as similarity measure of x, z is zero).

③ As the amount of data increases, would it make sense to increase σ or to decrease it?

- Decrease σ as you get more data to ensure that you have less smoothing and everything doesn't seem the same to your model (rely on local values).

Kernels: postscript

① Customized kernels

- For many different domains (NLP, biology, speech, ...)
- Over many different structures (sequences, sets, trees, graphs, ...)

Kernels: postscript

① Customized kernels

- For many different domains (NLP, biology, speech, ...)
- Over many different structures (sequences, sets, trees, graphs, ...)

② Learning the kernel function

Given a set of plausible kernels, find a linear combination of them that works well.

Kernels: postscript

① Customized kernels

- For many different domains (NLP, biology, speech, ...)
- Over many different structures (sequences, sets, trees, graphs, ...)

② Learning the kernel function

Given a set of plausible kernels, find a linear combination of them that works well.

③ Speeding up learning and prediction

The $n \times n$ kernel matrix $k(x^{(i)}, x^{(j)})$ is a bottleneck for large n .

One idea:

- Go back to the primal space!
- Replace the embedding Φ by a low-dimensional mapping $\tilde{\Phi}$ such that

$$\tilde{\Phi}(x) \cdot \tilde{\Phi}(z) \approx \Phi(x) \cdot \Phi(z).$$