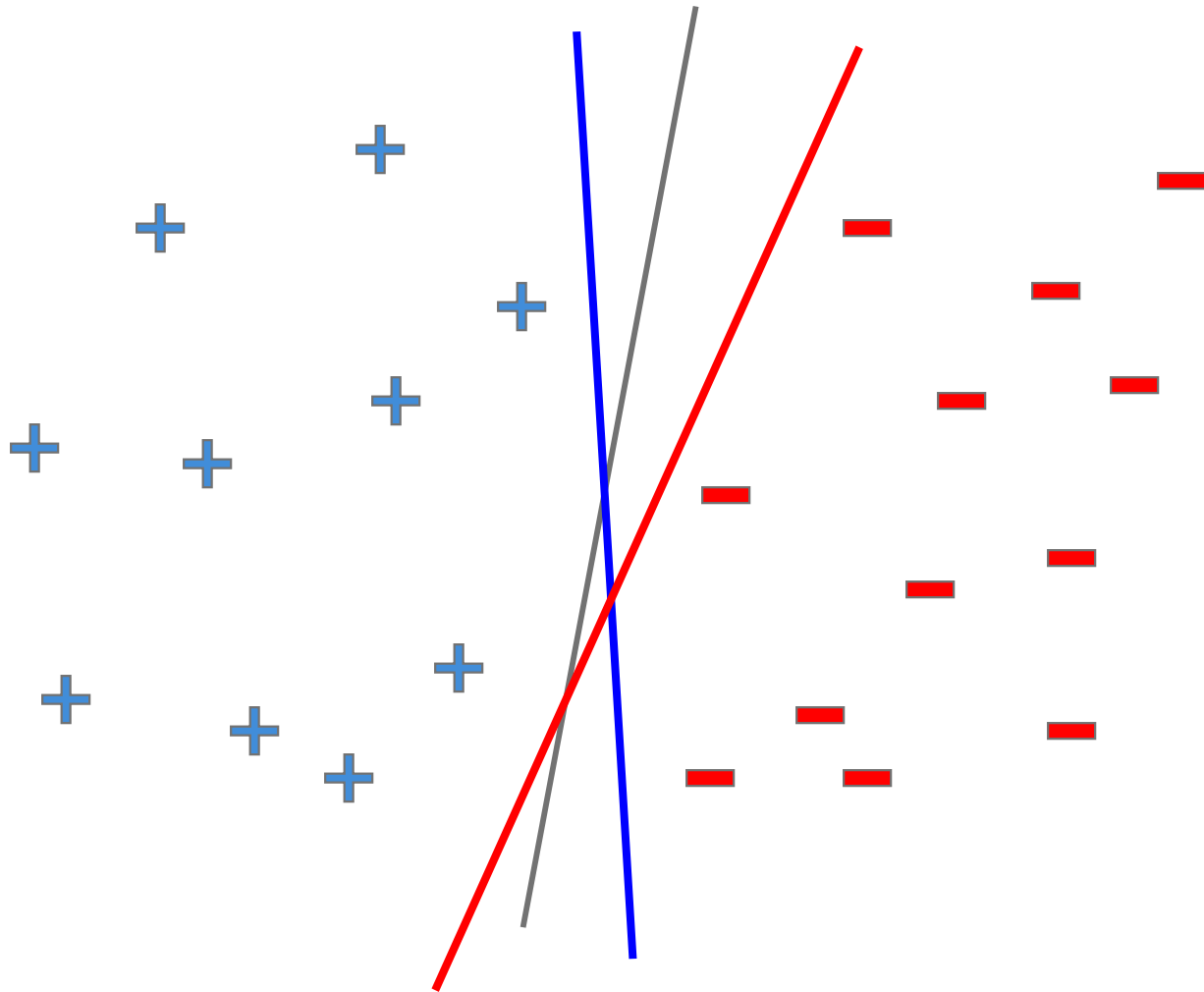


Advanced Introduction to Machine Learning CMU-10715

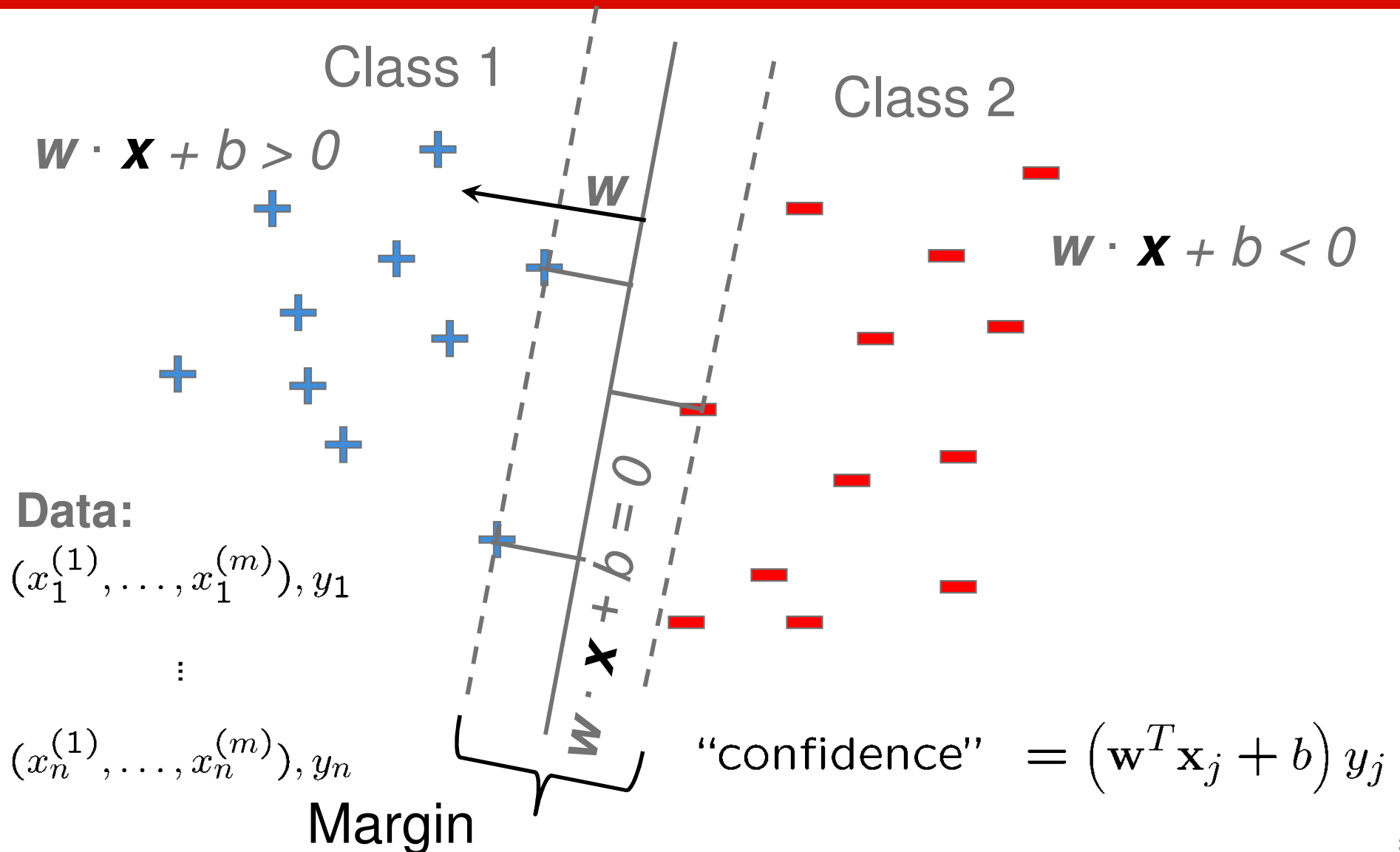
Support Vector Machines

Barnabás Póczos, 2014 Fall

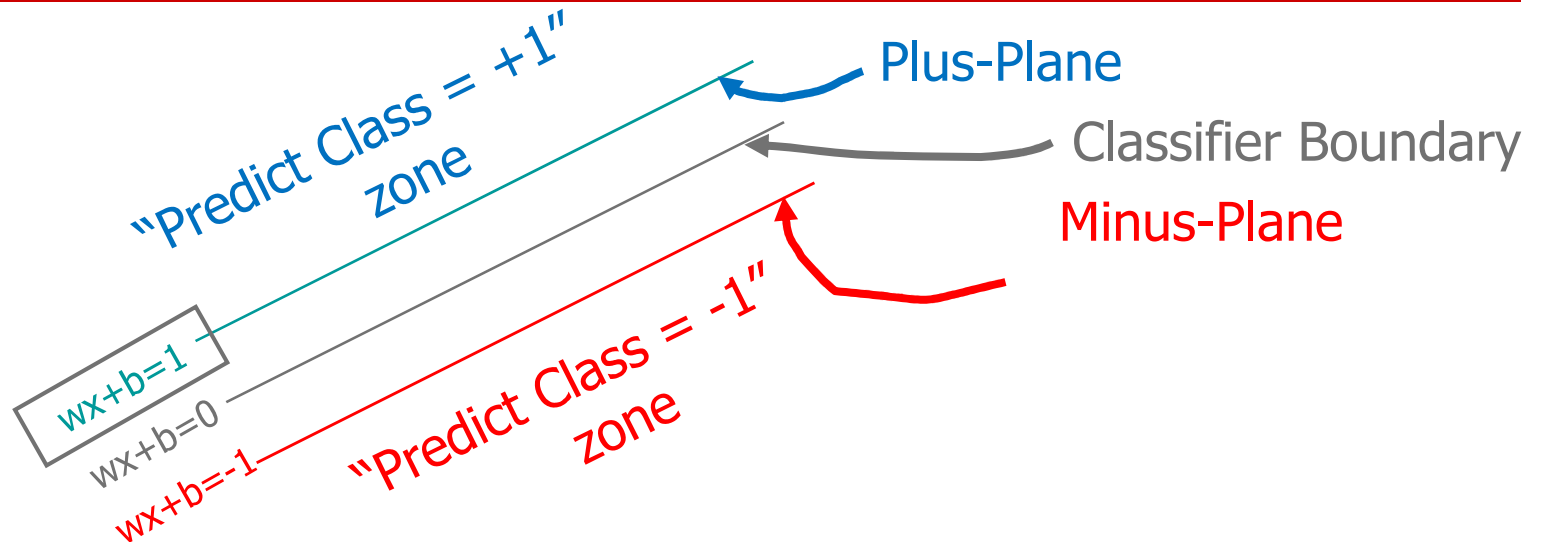
Linear classifiers which line is better?



Pick the one with the largest margin!



Scaling



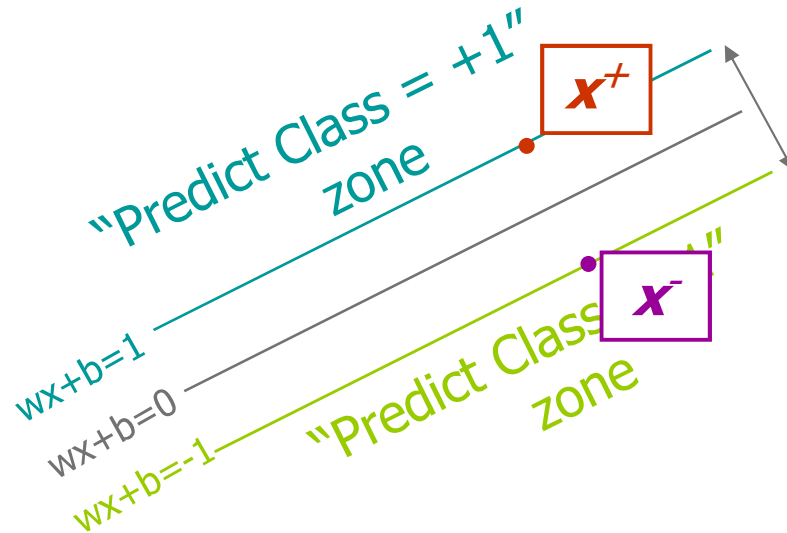
Classification rule:

Classify as..	+1	if	$w \cdot x + b \geq 1$
	-1	if	$w \cdot x + b \leq -1$
	Universe explodes	if	$-1 < w \cdot x + b < 1$

How large is the margin of this classifier?

Goal: Find the maximum margin classifier

Computing the margin width



$$M = \text{Margin Width} = \frac{2}{\sqrt{\mathbf{w} \cdot \mathbf{w}}}$$

$$\omega^T (\underbrace{x^+ - x^-}_{\lambda \omega}) = 2$$

$$\lambda = \frac{2}{\omega^T \omega}$$

Let \mathbf{x}^+ and \mathbf{x}^- be such that

- $\mathbf{w} \cdot \mathbf{x}^+ + b = +1$
- $\mathbf{w} \cdot \mathbf{x}^- + b = -1$
- $\mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$
- $|\mathbf{x}^+ - \mathbf{x}^-| = M = ?$ (Margin)

$$M = |\mathbf{x}^+ - \mathbf{x}^-| = |\lambda \omega| = \frac{|\lambda \omega|}{\frac{\omega^T \omega}{\omega^T \omega}}$$

$$= \frac{2}{\sqrt{\omega^T \omega}}$$

Maximize $M \equiv$ minimize $\mathbf{w} \cdot \mathbf{w}$!

The Primal Hard SVM

- Given $D = \{(\mathbf{x}_i, y_i), i = 1, \dots, n\}$ training data set.
- Assume that D is **linearly separable**.

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^m} \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{subject to } y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1, \forall i = 1, \dots, n$$

Prediction: $f_{\hat{\mathbf{w}}}(\mathbf{x}) = \text{sign}(\langle \hat{\mathbf{w}}, \mathbf{x} \rangle)$

**This is a QP problem (m-dimensional)
(Quadratic cost function, linear constraints)**

Quadratic Programming

Find

$$\text{ARG MIN}_{\omega \in \mathbb{R}^m} \omega^T H \omega + \omega^T q + e$$



Subject to

$$A \omega \leq b$$

$$A \in \mathbb{R}^{n \times m} \quad \omega \in \mathbb{R}^m \quad b \in \mathbb{R}^n$$

and to

$$C \omega = d$$

$$C \in \mathbb{R}^{s \times m} \quad d \in \mathbb{R}^s$$

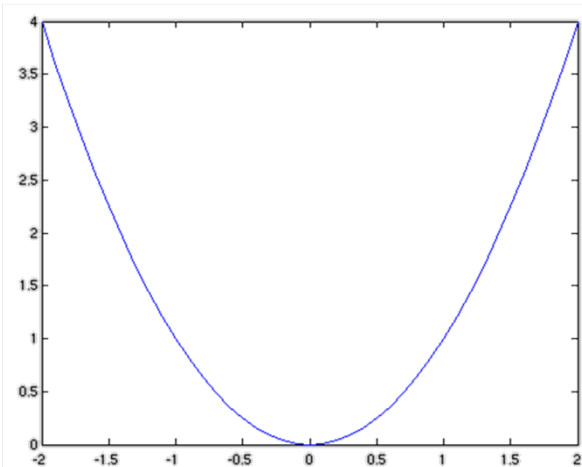
Efficient Algorithms exist for QP.

They often solve the dual problem instead of the primal.

Constrained Optimization

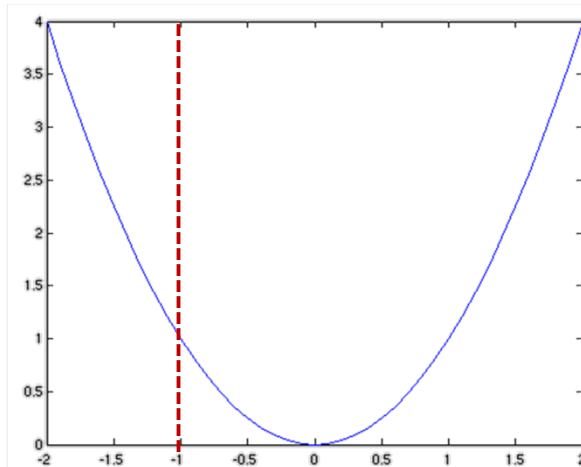
$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq b \end{aligned}$$

$$\min_x x^2$$



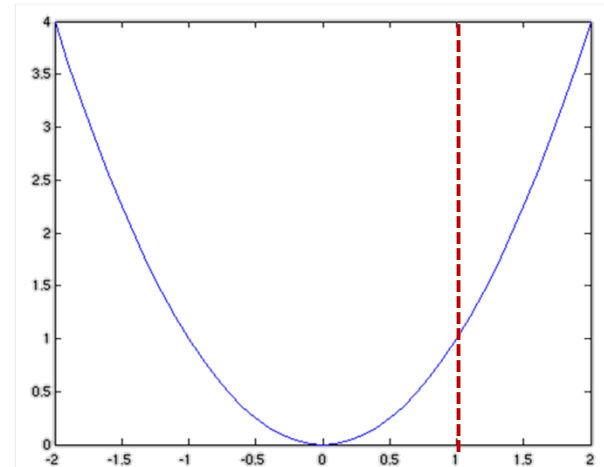
$$x^* = 0$$

$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq -1 \end{aligned}$$



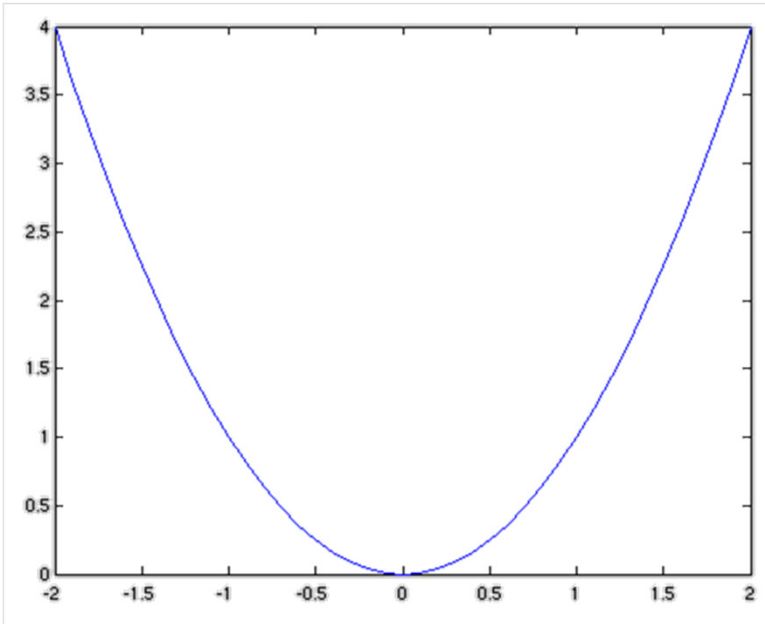
$$x^* = 0$$

$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq 1 \end{aligned}$$



$$x^* = 1$$

Lagrange Multiplier



$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq b \end{aligned}$$

Moving the constraint to objective function
Lagrangian:

$$\begin{aligned} L(x, \alpha) &= x^2 - \alpha(x - b) \\ \text{s.t.} \quad & \alpha \geq 0 \end{aligned}$$

Solve:

$$\begin{aligned} \min_x \max_{\alpha} \quad & L(x, \alpha) \\ \text{s.t.} \quad & \alpha \geq 0 \end{aligned}$$

Constraint is active when $\alpha > 0$

Lagrange Multiplier – Dual Variables

Solving:

$$\min_x \max_{\alpha} \overbrace{x^2 - \alpha(x - b)}^{L(x, \alpha)}$$

s.t. $\alpha \geq 0$

$$\frac{\partial L}{\partial x} = 0 \Rightarrow x^* = \frac{\alpha}{2}$$

$$\frac{\partial L}{\partial \alpha} = 0 \Rightarrow \alpha^* = \max(2b, 0)$$

$$\begin{aligned} 2x - \alpha = 0 &\Rightarrow x^* = \frac{\alpha}{2} \\ L(x^*, \alpha) &= \frac{\alpha^2}{4} - \alpha\left(\frac{\alpha}{2} - b\right) \\ &= -\frac{\alpha^2}{4} + b\alpha \\ \frac{\partial L(x^*, \alpha)}{\partial \alpha} &= -\frac{\alpha}{2} + b \\ \alpha^* &= 2b \end{aligned}$$

When $\alpha > 0$, constraint is tight


From Primal to Dual

Primal problem:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^m} \frac{1}{2} \|\mathbf{w}\|^2$$

subject to $y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1, \forall i = 1, \dots, n$

Lagrange function:

$\alpha = (\alpha_1, \dots, \alpha_n)^T \geq 0$ Lagrange multipliers 

$$L(\mathbf{w}, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i \langle \mathbf{x}_i, \mathbf{w} \rangle - 1)$$

The Lagrange Problem

$$L(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i \langle \mathbf{x}_i, \mathbf{w} \rangle - 1)$$

The Lagrange problem:

$$(\hat{\mathbf{w}}, \hat{\boldsymbol{\alpha}}) = \arg \min_{\mathbf{w} \in \mathbb{R}^m} \max_{0 \leq \boldsymbol{\alpha} \in \mathbb{R}^n} L(\mathbf{w}, \boldsymbol{\alpha})$$

$$0 = \left. \frac{\partial L(\mathbf{w}, \boldsymbol{\alpha})}{\partial \mathbf{w}} \right|_{\mathbf{w}=\hat{\mathbf{w}}} = \hat{\mathbf{w}} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$$

$$\Rightarrow \hat{\mathbf{w}} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$$

The Dual Problem

$$L(\mathbf{w}, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i \langle \mathbf{x}_i, \mathbf{w} \rangle - 1)$$

$$\Rightarrow \hat{\mathbf{w}} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$$

$$\Rightarrow L(\hat{\mathbf{w}}, \boldsymbol{\alpha}) = \frac{1}{2} \|\hat{\mathbf{w}}\|^2 - \sum_{i=1}^n \alpha_i (y_i \langle \mathbf{x}_i, \hat{\mathbf{w}} \rangle - 1)$$

$$= \frac{1}{2} \underbrace{\left\| \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \right\|^2}_{\boldsymbol{\alpha}^T \mathbf{Y} \mathbf{G} \mathbf{Y} \boldsymbol{\alpha}} + \boldsymbol{\alpha}^T \mathbf{1}_n - \underbrace{\sum_{i=1}^n \alpha_i y_i \langle \mathbf{x}_i, \sum_{j=1}^n \alpha_j y_j \mathbf{x}_j \rangle}_{\boldsymbol{\alpha}^T \mathbf{Y} \mathbf{G} \mathbf{Y} \boldsymbol{\alpha}}$$

$$= \boldsymbol{\alpha}^T \mathbf{1}_n - \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{Y} \mathbf{G} \mathbf{Y} \boldsymbol{\alpha}$$

$$\mathbf{Y} \doteq \text{diag}(y_1, \dots, y_n), \quad y_i \in \{-1, 1\}^n$$

$$\mathbf{G} \in \mathbb{R}^{n \times n} \doteq \{G_{ij}\}_{i,j}^{n,n}, \quad \text{where } G_{ij} \doteq \langle \mathbf{x}_i, \mathbf{x}_j \rangle \text{ Gram matrix.}$$

The Dual Hard SVM

$$Y \doteq \text{diag}(y_1, \dots, y_n), \quad y_i \in \{-1, 1\}^n$$

$$G \in \mathbb{R}^{n \times n} \doteq \{G_{ij}\}_{i,j}^{n,n}, \quad \text{where } G_{ij} \doteq \langle \mathbf{x}_i, \mathbf{x}_j \rangle \text{ Gram matrix.}$$

$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \alpha^T \mathbf{1}_n - \frac{1}{2} \alpha^T Y G Y \alpha$$

subject to $\alpha_i \geq 0, \forall i = 1, \dots, n$

Quadratic Programming (n-dimensional)

Lemma $\hat{\mathbf{w}} = \sum_{i=1}^n \hat{\alpha}_i y_i \mathbf{x}_i$

Prediction: $f_{\hat{\mathbf{w}}}(x) = \text{sign}(\langle \hat{\mathbf{w}}, \mathbf{x} \rangle) = \text{sign}\left(\sum_{i=1}^n \hat{\alpha}_i y_i \underbrace{\langle \mathbf{x}_i, \mathbf{x} \rangle}_{k(\mathbf{x}_i, \mathbf{x})}\right)$

The Problem with Hard SVM

It assumes samples are linearly separable...

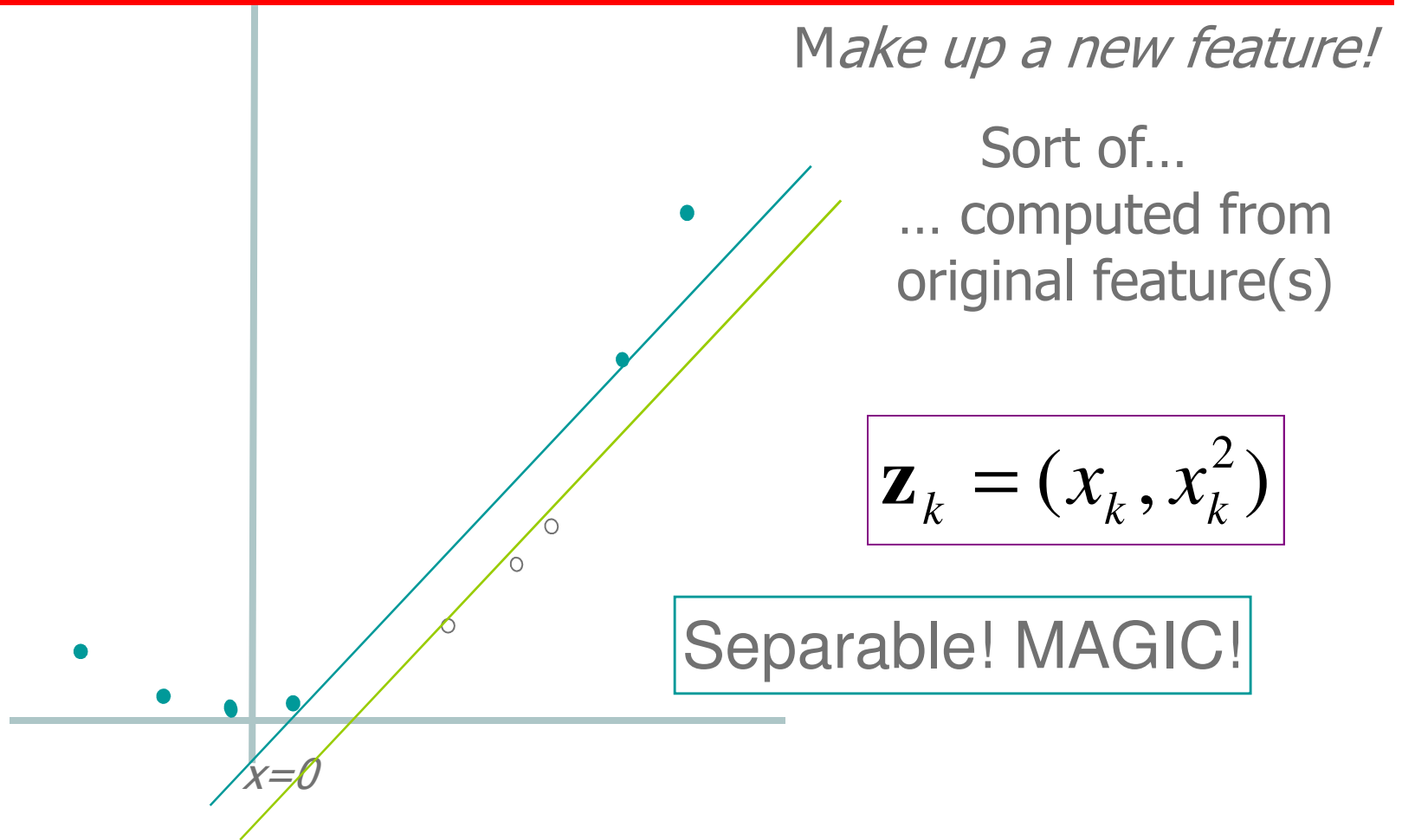
**What can we do if data is
not linearly separable???**

Hard 1-dimensional Dataset

If the data set is **not** linearly separable, then adding new features (mapping the data to a larger feature space) the data might become linearly separable



Hard 1-dimensional Dataset



Now drop this “augmented” data into our linear SVM.

Feature mapping

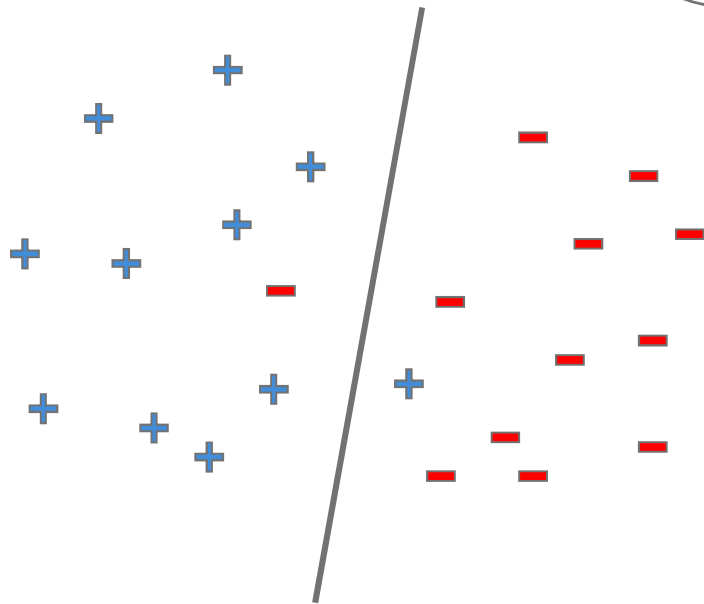
- n general! points in an $n-1$ dimensional space is always linearly separable by a hyperspace!
⇒ it is good to map the data to high dimensional spaces
- Having n training data, is it always enough to map the data into a feature space with dimension $n-1$?
 - *Nope... We have to think about the test data as well!
Even if we don't know how many test data we have and what they are...*
 - *We might want to map our data to a huge (∞) dimensional feature space*
 - *Overfitting? Generalization error?...*
We don't care now...

How to do feature mapping?

Let us have n training objects: $\vec{x}_i = [\vec{x}_{i,1}, \vec{x}_{i,2}] \in \mathbb{R}^2$, $i = 1, \dots, n$

The possible test objects are denoted by $\vec{x} = [\vec{x}_1, \vec{x}_2] \in \mathbb{R}^2$

Let $\phi(\vec{x}) \doteq [\sin(\vec{x}_2), \exp(\vec{x}_2 + \vec{x}_1), \vec{x}_1, \vec{x}_2^{\tan(\vec{x}_1)}, \dots]$



**Use features of features
of features of features....**

The Problem with Hard SVM

It assumes samples are linearly separable...

Solutions:

1. Use feature transformation to a larger space
 - ⇒ each training samples are linearly separable in the feature space
 - ⇒ Hard SVM can be applied 😊
 - ⇒ overfitting... 😞
2. **Soft margin SVM** instead of Hard SVM
 - Slack variables... We will discuss them now

Hard SVM

The Hard SVM problem can be rewritten:

$$\hat{\mathbf{w}}_{hard} = \arg \min_{\mathbf{w} \in \mathbb{R}^m} \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{subject to } y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1, \forall i = 1, \dots, n$$



$$\hat{\mathbf{w}}_{hard} = \arg \min_{\mathbf{w} \in \mathbb{R}^m} \sum_{i=1}^n l_{0-\infty}(\langle \mathbf{x}_i, \mathbf{w} \rangle, y_i) + \frac{1}{2} \|\mathbf{w}\|^2$$

where

$$l_{0-\infty}(a, b) \doteq \begin{cases} \infty : ab < 1 & \text{Misclassification, or inside the margin} \\ 0 : ab \geq 1 & \text{Correct classification and outside of the margin} \end{cases}$$

From Hard to Soft constraints

Instead of using hard constraints (points are linearly separable)

$$\hat{\mathbf{w}}_{hard} = \arg \min_{\mathbf{w} \in \mathbb{R}^m} \sum_{i=1}^n l_{0-\infty}(\langle \mathbf{x}_i, \mathbf{w} \rangle, y_i) + \frac{1}{2} \|\mathbf{w}\|^2$$

We can try solve the soft version of it: . Introduce a λ parameter!

(Your loss is only 1 instead of ∞ if you misclassify an instance)

$$\hat{\mathbf{w}}_{soft} = \arg \min_{\mathbf{w} \in \mathbb{R}^m} \sum_{i=1}^n l_{0-1}(\langle \mathbf{x}_i, \mathbf{w} \rangle, y_i) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

where

$$l_{0-1}(y, f(\mathbf{x})) = \begin{cases} 1 & : yf(\mathbf{x}) < 0 & \text{Misclassification} \\ 0 & : yf(\mathbf{x}) > 0 & \text{Correct classification} \end{cases}$$

Problems with l_{0-1} loss

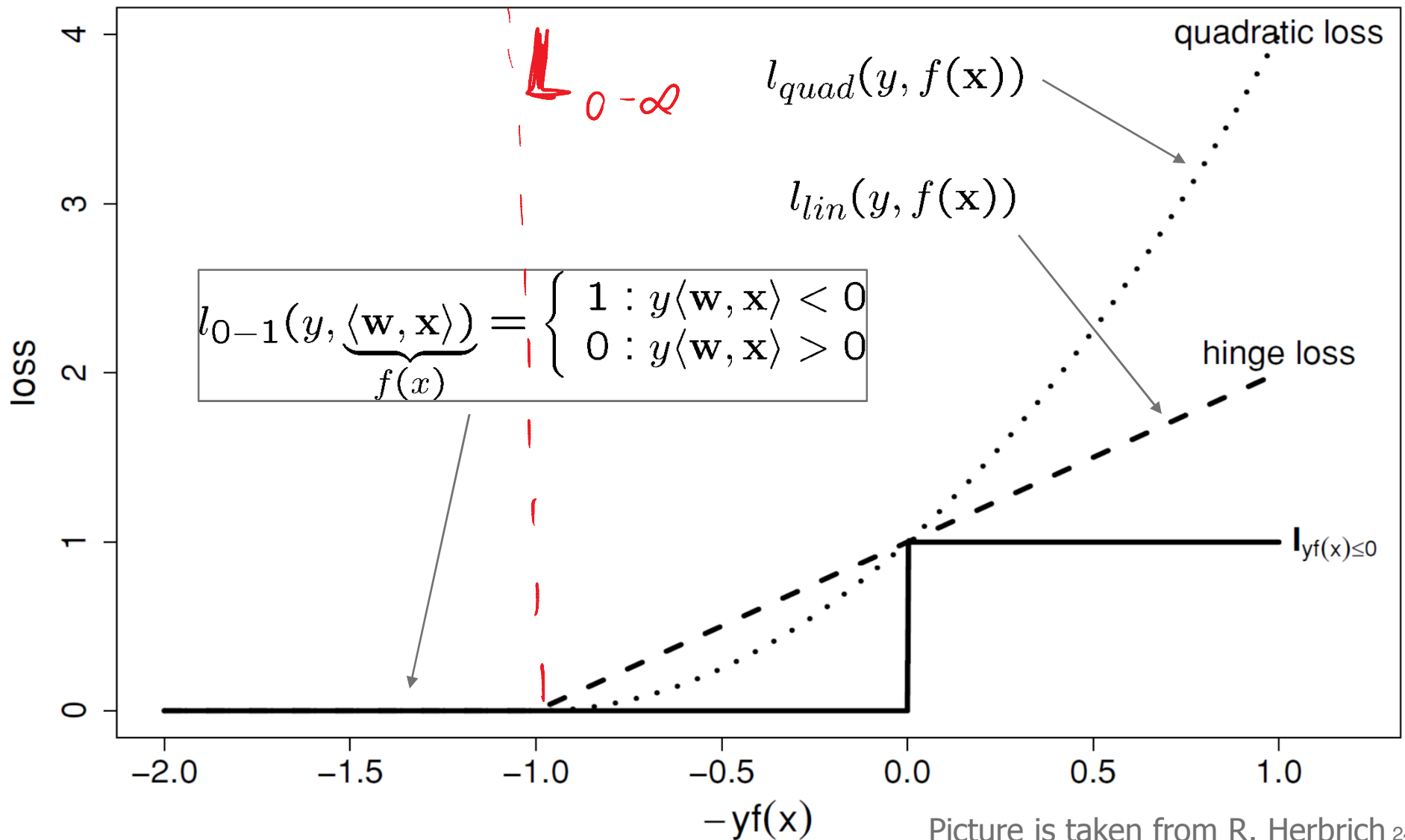
$$\hat{\mathbf{w}}_{soft} = \arg \min_{\mathbf{w} \in \mathbb{R}^m} \sum_{i=1}^n l_{0-1}(\langle \mathbf{x}_i, \mathbf{w} \rangle, y_i) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

$$l_{0-1}(y, f(\mathbf{x})) = \begin{cases} 1 & : yf(\mathbf{x}) < 0 \\ 0 & : yf(\mathbf{x}) > 0 \end{cases}$$

It is not convex in $yf(\mathbf{x}) \Rightarrow$ It is not convex in \mathbf{w} , either...
... and we only like convex functions...

Let us approximate it with convex functions!

Approximation of the Heaviside step function



Approximations of l_{0-1} loss

- Piecewise linear approximations (hinge loss, l_{lin})

$$l_{lin}(f(\mathbf{x}), y) = \max\{1 - yf(\mathbf{x}), 0\}$$

[We want $yf(\mathbf{x}) > 1$]

- Quadratic approximation (l_{quad})

$$l_{quad}(f(\mathbf{x}), y) = \max\{1 - yf(\mathbf{x}), 0\}^2$$

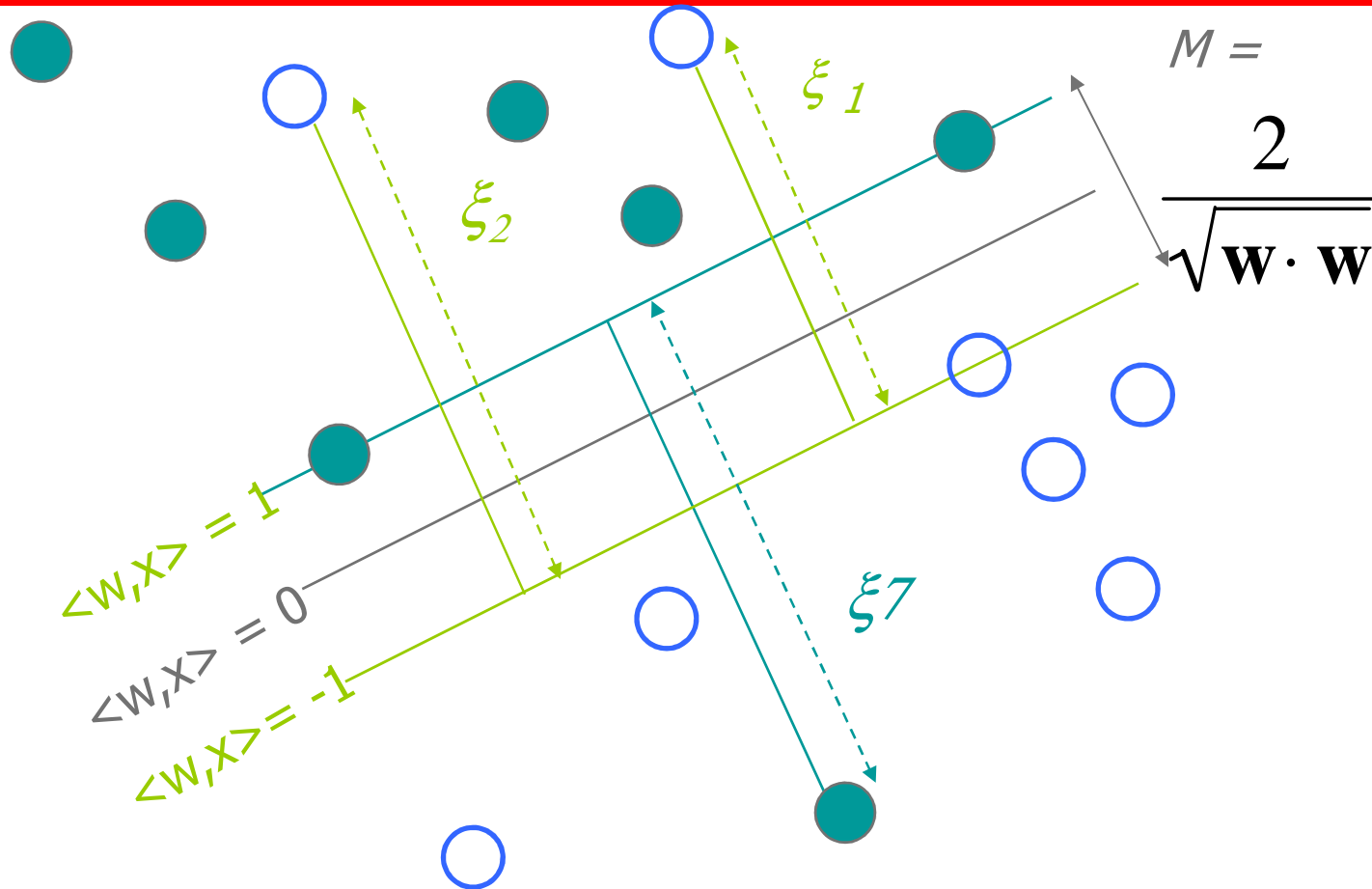
The hinge loss approximation of l_{0-1}

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^m} \sum_{i=1}^n \underbrace{l_{lin}(\langle \mathbf{x}_i, \mathbf{w} \rangle, y_i)}_{\xi_i \geq 0} + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

Where,

$$\begin{aligned} \xi_i &\doteq l_{lin}(f(\mathbf{x}_i), y_i) = \max\{1 - y_i f(\mathbf{x}_i), 0\} \\ &\geq 1 - y_i \underbrace{\langle \mathbf{w}, \mathbf{x}_i \rangle}_{f(\mathbf{x}_i)} \geq l_{0-1}(y_i, f(\mathbf{x}_i)) \end{aligned}$$

The Slack Variables



$$\xi_i \doteq l_{lin}(f(\mathbf{x}_i), y_i) = \max\{1 - y_i f(\mathbf{x}_i), 0\}$$

$$= \max\{1 - y_i (\mathbf{w}^T \mathbf{x}_i), 0\}$$

The Primal Soft SVM problem

$$\hat{\mathbf{w}}_{soft} = \arg \min_{\mathbf{w} \in \mathbb{R}^m} \sum_{i=1}^n \underbrace{l_{lin}(\langle \mathbf{x}_i, \mathbf{w} \rangle, y_i)}_{\xi_i \geq 0} + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

where

$$\xi_i \doteq l_{lin}(f(\mathbf{x}_i), y_i) = \max\{1 - y_i(\mathbf{w}^T \mathbf{x}_i), 0\}$$

Equivalently,

$$\hat{\mathbf{w}}_{soft} = \arg \min_{\mathbf{w} \in \mathbb{R}^m, \boldsymbol{\xi} \in \mathbb{R}^n} \sum_{i=1}^n \xi_i + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

subject to $y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1 - \xi_i, \forall i = 1, \dots, n$

$$\xi_i \geq 0, \forall i = 1, \dots, n$$

ξ_i : Slack variables

The Primal Soft SVM problem

$$\hat{\mathbf{w}}_{soft} = \arg \min_{\mathbf{w} \in \mathbb{R}^m, \boldsymbol{\xi} \in \mathbb{R}^n} \sum_{i=1}^n \xi_i + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

$$\text{subject to } y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1 - \xi_i, \quad \forall i = 1, \dots, n$$

Equivalently,

$$\xi_i \geq 0, \quad \forall i = 1, \dots, n$$

We can use this form, too... where $C = \frac{1}{\lambda}$

$$\hat{\mathbf{w}}_{soft} = \arg \min_{\mathbf{w} \in \mathbb{R}^m, \boldsymbol{\xi} \in \mathbb{R}^n} C \underbrace{\sum_{i=1}^n \xi_i}_{\boldsymbol{\xi}^T \mathbf{1}_n} + \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{subject to } y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1 - \xi_i, \quad \forall i = 1, \dots, n$$

$$\xi_i \geq 0, \quad \forall i = 1, \dots, n$$

What is the dual form of primal soft SVM?

The Dual Soft SVM (using hinge loss)

$$\hat{\mathbf{w}}_{soft} = \arg \min_{\mathbf{w} \in \mathcal{R}^m, \boldsymbol{\xi} \in \mathbb{R}^n} C \sum_{i=1}^n \xi_i + \frac{1}{2} \|\mathbf{w}\|^2$$

subject to $y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1 - \xi_i, \forall i = 1, \dots, n$

$$\xi_i \geq 0, \forall i = 1, \dots, n$$

$\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n)^T \geq 0$ Lagrange multipliers

$\boldsymbol{\beta} = (\beta_1, \dots, \beta_n)^T \geq 0$ Lagrange multipliers

$$(\hat{\mathbf{w}}, \hat{\boldsymbol{\xi}}, \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}) = \arg \min_{\substack{\mathbf{w} \in \mathbb{R}^m \\ \mathbf{0} \leq \boldsymbol{\xi} \in \mathbb{R}^n}} \max_{\substack{0 \leq \boldsymbol{\alpha} \\ 0 \leq \boldsymbol{\beta}}} L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta})$$

where

$$L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i (y_i \langle \mathbf{x}_i, \mathbf{w} \rangle - 1 + \xi_i) - \sum_{i=1}^n \beta_i \xi_i$$

The Dual Soft SVM (using hinge loss)

$$L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \boldsymbol{\xi}^T \mathbf{1}_n - \sum_{i=1}^n \alpha_i y_i \langle \mathbf{x}_i, \mathbf{w} \rangle + \boldsymbol{\alpha}^T \mathbf{1}_n - \boldsymbol{\xi}^T (\boldsymbol{\alpha} + \boldsymbol{\beta})$$

$$0 = \left. \frac{\partial L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta})}{\partial \mathbf{w}} \right|_{\mathbf{w}=\hat{\mathbf{w}}} = \hat{\mathbf{w}} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \Rightarrow \boxed{\hat{\mathbf{w}} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i}$$

$$0 = \left. \frac{\partial L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta})}{\partial \boldsymbol{\xi}} \right|_{\boldsymbol{\xi}=\hat{\boldsymbol{\xi}}} = C \mathbf{1}_n - \boldsymbol{\alpha} - \boldsymbol{\beta} \Rightarrow \boldsymbol{\beta} = C \mathbf{1}_n - \boldsymbol{\alpha} \geq 0$$

$$\Rightarrow 0 \leq \boldsymbol{\alpha} \leq C$$

$$\Rightarrow (\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}) = \arg \max_{\substack{0 \leq \boldsymbol{\alpha} \leq C \\ 0 \leq \boldsymbol{\beta}}} L(\hat{\mathbf{w}}, \hat{\boldsymbol{\xi}}, \boldsymbol{\alpha}, \boldsymbol{\beta})$$

$$\Rightarrow \hat{\boldsymbol{\alpha}} = \arg \max_{0 \leq \boldsymbol{\alpha} \leq C} \boldsymbol{\alpha}^T \mathbf{1}_n - \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{Y} \mathbf{G} \mathbf{Y} \boldsymbol{\alpha}$$

The Dual Soft SVM (using hinge loss)

$$Y \doteq \text{diag}(y_1, \dots, y_n) \in \{-1, 1\}^n$$

$$G \in \mathbb{R}^{n \times n} \doteq \{G_{ij}\}_{i,j}^{n,n}, \text{ where } G_{ij} \doteq \overbrace{\langle \mathbf{x}_i, \mathbf{x}_j \rangle}^{k(\mathbf{x}_i, \mathbf{x}_j)}, \text{ Gram matrix.}$$

$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \alpha^T \mathbf{1}_n - \frac{1}{2} \alpha^T Y G Y \alpha$$

subject to $0 \leq \alpha_i \leq C$

where $C = \frac{1}{\lambda}$

If $\lambda \rightarrow 0 \Rightarrow \text{soft-SVM} \rightarrow \text{hard-SVM}$

This is the same as the dual hard-SVM problem, but now we have the additional $0 \leq \alpha_i \leq C$ constraints.

SVM classification in the dual space

Solve the dual problem

$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \alpha^T \mathbf{1}_n - \frac{1}{2} \alpha^T Y G Y \alpha$$

subject to $0 \leq \alpha_i \leq C$

where $C = \frac{1}{\lambda}$. Let $\hat{\mathbf{w}} = \sum_{i=1}^n \hat{\alpha}_i y_i \mathbf{x}_i$.

On test data \mathbf{x} : $f_{\hat{\mathbf{w}}}(\mathbf{x}) = \langle \hat{\mathbf{w}}, \mathbf{x} \rangle = \sum_{i=1}^n \hat{\alpha}_i y_i \underbrace{\langle \mathbf{x}_i, \mathbf{x} \rangle}_{k(\mathbf{x}_i, \mathbf{x})}$

Why is it called Support Vector Machine?

$\alpha = (\alpha_1, \dots, \alpha_n)^T \geq 0$ Lagrange multipliers

$$L(\mathbf{w}, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i \underbrace{(y_i \langle \mathbf{x}_i, \mathbf{w} \rangle - 1)}_0$$

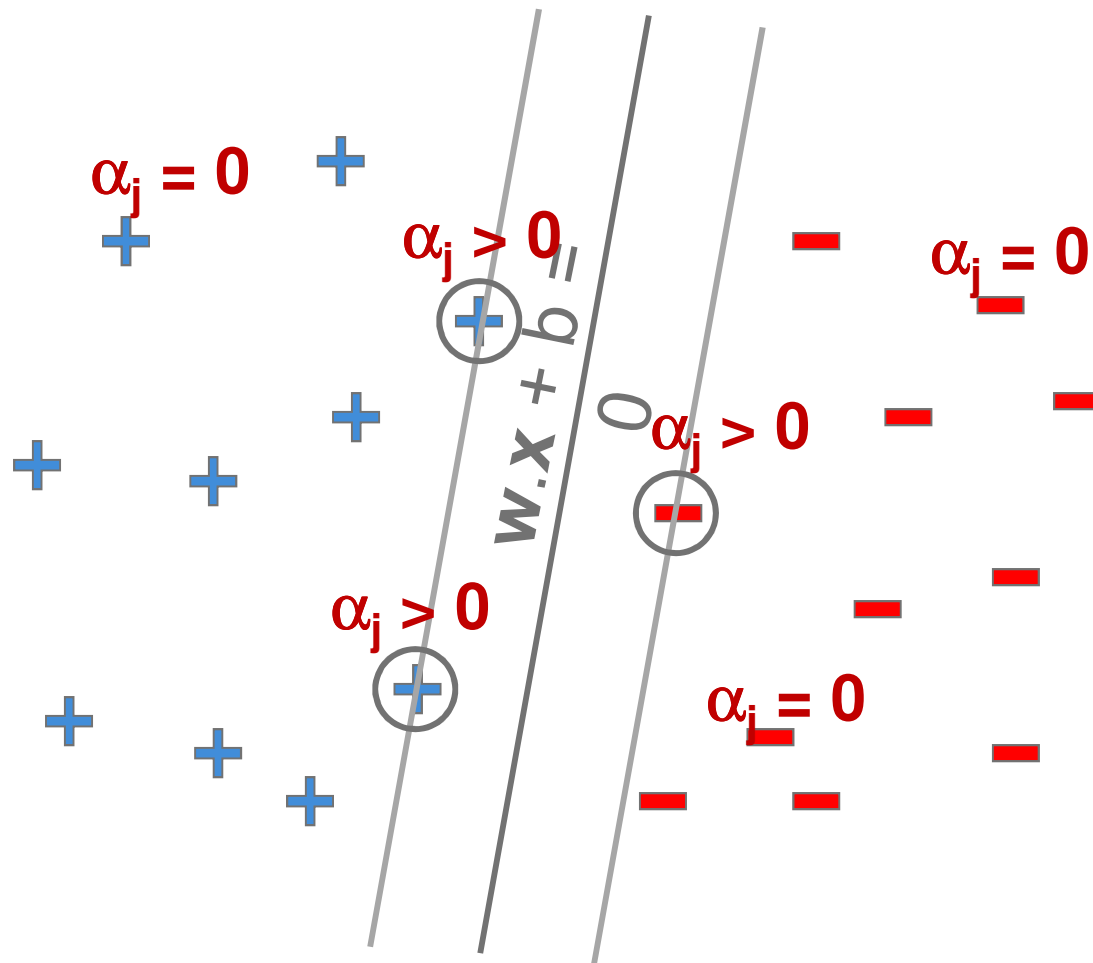
KKT conditions

COMPLEMENTARY SLACKNESS CONDITION

$$d_i > 0 \Rightarrow y_i \langle \mathbf{x}_i, \mathbf{w} \rangle = 1$$
$$\langle \mathbf{x}_i, \mathbf{w} \rangle = \begin{matrix} +1 \\ -1 \end{matrix}$$

EITHER $d_i = 0$ OR $(y_i \langle \mathbf{x}_i, \mathbf{w} \rangle - 1) = 0$ AND $d_i > 0$
 $(y_i \langle \mathbf{x}_i, \mathbf{w} \rangle - 1) \neq 0$
 \Downarrow
 \mathbf{x}_i IS ON THE MARGIN LINES
SUPPORT VECTORS

Dual SVM Interpretation: Sparsity



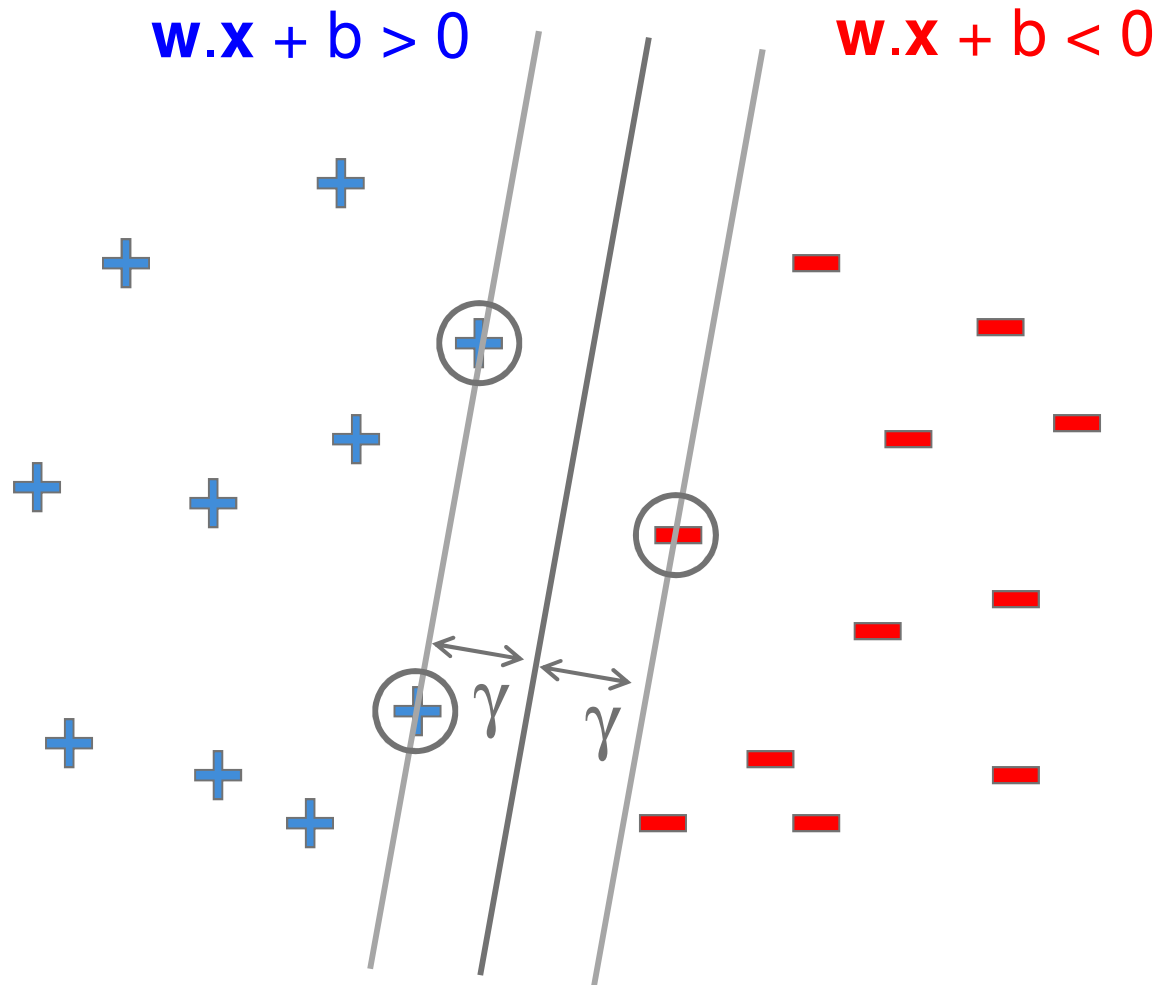
$$\mathbf{w} = \sum_j \alpha_j y_j \mathbf{x}_j$$

Only few α_j s can be non-zero : where constraint is tight

$$(\langle \mathbf{w}, \mathbf{x}_j \rangle + b) y_j = 1$$

Support vectors – training points j whose α_j s are non-zero

Support Vectors



Linear hyperplane defined by “support vectors”

Moving other points a little doesn't effect the decision boundary

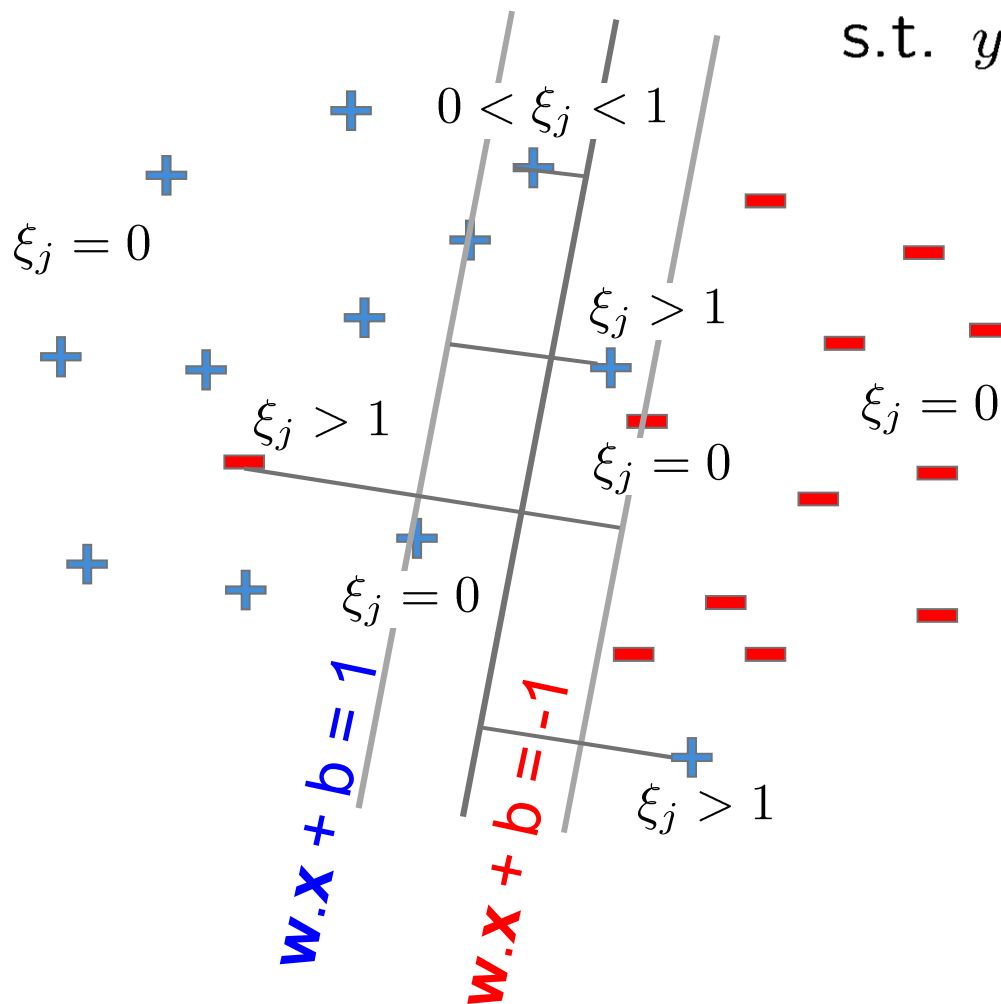
only need to store the support vectors to predict labels of new points

Support vectors in Soft SVM

$$\hat{\mathbf{w}}_{soft} = \arg \min_{\mathbf{w} \in \mathbb{R}^m, \boldsymbol{\xi} \in \mathbb{R}^n} C \sum_{i=1}^n \xi_i + \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{s.t. } y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1 - \xi_i, \quad \forall i = 1, \dots, n$$

$$\xi_i \geq 0, \quad \forall i = 1, \dots, n$$

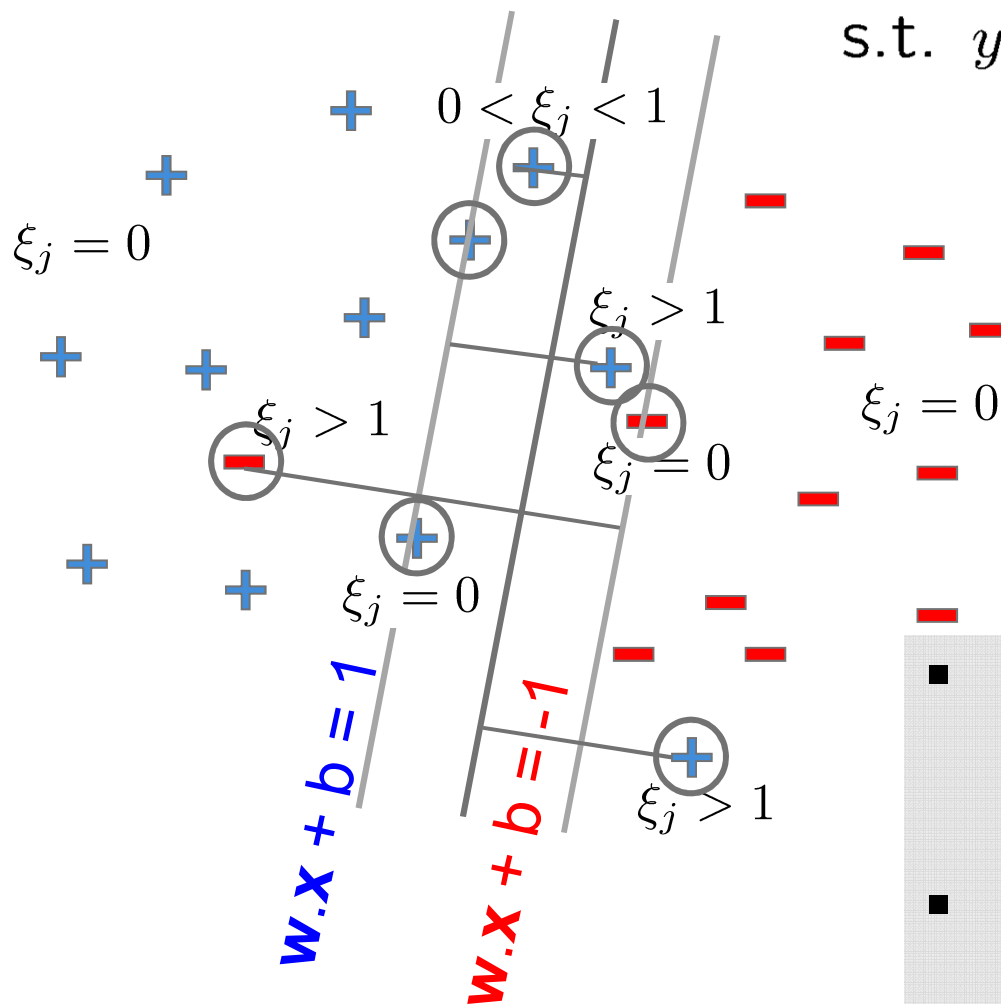


Support vectors in Soft SVM

$$\hat{\mathbf{w}}_{soft} = \arg \min_{\mathbf{w} \in \mathbb{R}^m, \boldsymbol{\xi} \in \mathbb{R}^n} C \sum_{i=1}^n \xi_i + \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{s.t. } y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1 - \xi_i, \quad \forall i = 1, \dots, n$$

$$\xi_i \geq 0, \quad \forall i = 1, \dots, n$$



- **Margin support vectors**
 $y_i \langle \mathbf{x}_i, \mathbf{w} \rangle = 1$
- **Nonmargin support vectors**
 $\xi_i > 0$

SVM classification in the dual space

“Without b”

$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^m} \alpha^T \mathbf{1}_m - \frac{1}{2} \alpha^T YGY \alpha$$

subject to $0 \leq \alpha_i \leq C$

“With b”

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i (\mathbf{x}_i \cdot \mathbf{w} + b) - 1)$$

$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \alpha^T \mathbf{1}_n - \frac{1}{2} \alpha^T YGY \alpha$$

subject to $0 \leq \alpha_i \leq C$

$$\sum_i \alpha_i y_i = 0$$

SVM with Linear Programs

QP:

$$\min_{\mathbf{w} \in \mathbb{R}^m, \boldsymbol{\xi} \in \mathbb{R}^n} C \sum_{i=1}^n \xi_i + \frac{1}{2} \|\mathbf{w}\|^2$$

Max margin

subject to $y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1 - \xi_i, \forall i = 1, \dots, n$

$$\xi_i \geq 0, \forall i = 1, \dots, n$$

LP:

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^n, \boldsymbol{\xi} \in \mathbb{R}^n} C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i$$

Min support vectors

subject to $y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1 - \xi_i, \forall i = 1, \dots, n$

$$\xi_i \geq 0, \forall i = 1, \dots, n$$

$$\alpha_i \geq 0, \forall i = 1, \dots, n$$

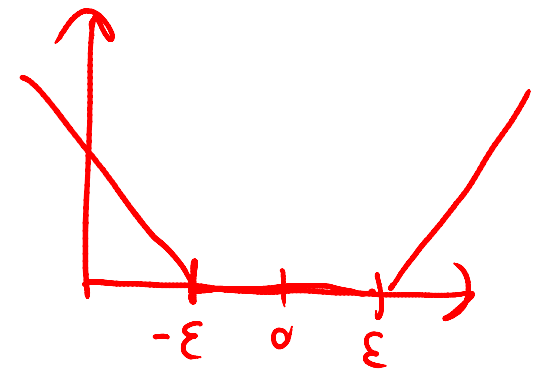
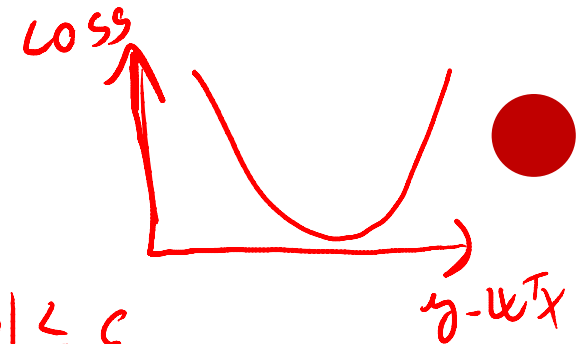
$$\mathbf{w} = \sum_{j=1}^n \alpha_j y_j \mathbf{x}_j$$

SVM for Regression

$$\text{Loss}(\gamma, \omega^T x) = (\gamma - \omega^T x)^2$$

$$0 \text{ if } |\gamma - \omega^T x| \leq \epsilon$$

$$|\gamma - \omega^T x| \text{ OTHERWISE}$$



Ridge Regression

Linear regression: $f(x) = \langle \mathbf{w}, \phi(x) \rangle$

Primal:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathcal{K}} \sum_{i=1}^n \xi_i^2$$

subject to $y_i - \underbrace{\langle \phi(x_i), \mathbf{w} \rangle}_{x_i} = \xi_i, \forall i = 1, \dots, n$

and $\|\mathbf{w}\| \leq B$

$$L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}, \beta) = \sum_{i=1}^n \xi_i^2 + \sum_{i=1}^n \alpha_i (y_i - \langle \phi(x_i), \mathbf{w} \rangle - \xi_i) + \lambda (\|\mathbf{w}\|^2 - B^2)$$

$\lambda > 0$

Dual for a given λ : ...after some calculations...

$$\hat{\boldsymbol{\alpha}} = \arg \max_{\boldsymbol{\alpha} \in \mathbb{R}^n} \lambda \sum_{i=1}^n \alpha_i^2 - 2 \sum_{i=1}^n \alpha_i y_i + \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(x_i, x_j)$$

This can be solved in closed form:

Kernel Ridge Regression Algorithm

Given $D = \{(x_i, y_i), i = 1, \dots, n\}$ training data set.

$k(\cdot, \cdot)$ kernel, $\lambda > 0$ parameter. $\mathbf{y} \doteq (y_1, \dots, y_n)^T \in \{-1, 1\}^n$

- $\mathbf{G} \in \mathbb{R}^{n \times n} \doteq \{G_{ij}\}_{i,j}^{n,n}$,

where $G_{ij} \doteq \overbrace{\langle \underbrace{\mathbf{x}_i}_{\phi(x_i)}, \underbrace{\mathbf{x}_j}_{\phi(x_j)} \rangle_{\mathcal{K}}}_{k(x_i, x_j)}$, Gram matrix.

- $\hat{\boldsymbol{\alpha}} = (\mathbf{G} + \lambda \mathbf{I}_n)^{-1} \mathbf{y}$

- $\hat{\mathbf{w}} = \sum_{i=1}^n \hat{\alpha}_i \phi(x_i).$

- $f(x) = \langle \hat{\mathbf{w}}, \phi(x) \rangle = \sum_{i=1}^n \hat{\alpha}_i k(x_i, x)$

SVM vs. Logistic Regression

SVM : **Hinge loss**

$$\text{loss}(f(x_j), y_j) = (1 - (\mathbf{w} \cdot x_j + b)y_j)_+$$

Logistic Regression : **Log loss** (log conditional likelihood)

$$\text{loss}(f(x_j), y_j) = -\log P(y_j | x_j, \mathbf{w}, b) = \log(1 + e^{-(\mathbf{w} \cdot x_j + b)y_j})$$



Difference between SVMs and Logistic Regression

	SVMs	Logistic Regression
Loss function	Hinge loss	Log-loss
High dimensional features with kernels	Yes!	No (but there is kernel logistic regression too)
Solution sparse	Often yes!	Almost always no!
Semantics of output	“Margin”	“Real probabilities”

Constructing Kernels

Common Kernels

- Polynomials of degree d

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v})^d$$

- Polynomials of degree up to d

$$K(\mathbf{u}, \mathbf{v}) = (\mathbf{u} \cdot \mathbf{v} + 1)^d$$

- Gaussian/Radial kernels

$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma^2}\right)$$

- Sigmoid

$$K(\mathbf{u}, \mathbf{v}) = \tanh(\eta \mathbf{u} \cdot \mathbf{v} + \nu)$$

Designing new kernels from kernels

$k_1 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, $k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ are kernels \Rightarrow

1. $k(x, \tilde{x}) = k_1(x, \tilde{x}) + k_2(x, \tilde{x})$,
2. $k(x, \tilde{x}) = c \cdot k_1(x, \tilde{x})$, for all $c \in \mathbb{R}^+$,
3. $k(x, \tilde{x}) = k_1(x, \tilde{x}) + c$, for all $c \in \mathbb{R}^+$,
4. $k(x, \tilde{x}) = k_1(x, \tilde{x}) \cdot k_2(x, \tilde{x})$,
5. $k(x, \tilde{x}) = f(x) \cdot f(\tilde{x})$, for any function $f : \mathcal{X} \rightarrow \mathbb{R}$

are also kernels.

Designing new kernels from kernels

1. $k(x, \tilde{x}) = (k_1(x, \tilde{x}) + \theta_1)^{\theta_2}$, for all $\theta_1 \in \mathbb{R}^+$ and $\theta_2 \in \mathbb{N}$.
2. $k(x, \tilde{x}) = \exp\left(\frac{k_1(x, \tilde{x})}{\sigma^2}\right)$, for all $\sigma \in \mathbb{R}^+$,
3. $k(x, \tilde{x}) = \exp\left(-\frac{k_1(x, x) - 2k_1(x, \tilde{x}) + k_1(\tilde{x}, \tilde{x})}{2\sigma^2}\right)$, for all $\sigma \in \mathbb{R}^+$
4. $k(x, \tilde{x}) = \frac{k_1(x, \tilde{x})}{\sqrt{k_1(x, x) \cdot k_1(\tilde{x}, \tilde{x})}}$

Designing new kernels from kernels

The meaning of

$$k(x, \tilde{x}) = \frac{k_1(x, \tilde{x})}{\sqrt{k_1(x, x)k_1(\tilde{x}, \tilde{x})}}$$

is that we can normalize the data in the feature space without performing the explicit mapping.

Use the normalized kernel k_{norm} :

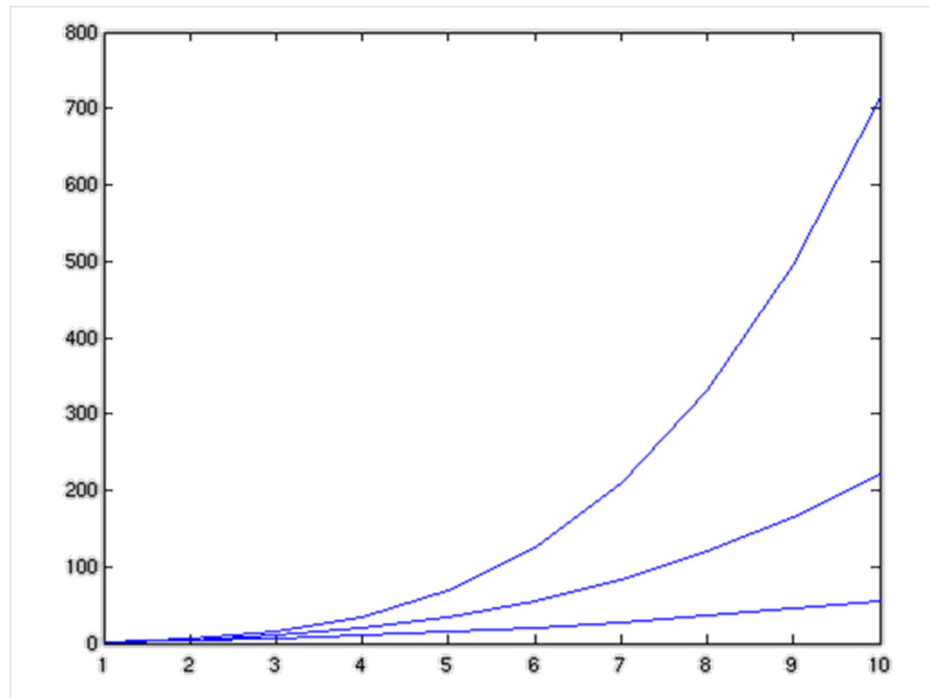
$$k_{norm}(x, \tilde{x}) \doteq \frac{k(x, \tilde{x})}{\sqrt{k(x, x)k(\tilde{x}, \tilde{x})}} = \frac{\langle x, \tilde{x} \rangle}{\sqrt{\|x\|^2\|\tilde{x}\|^2}} = \left\langle \frac{x}{\|x\|}, \frac{\tilde{x}}{\|\tilde{x}\|} \right\rangle$$

Higher Order Polynomials

m – input features

d – degree of polynomial

$$\text{num. terms} = \binom{d + m - 1}{d} = \frac{(d + m - 1)!}{d!(m - 1)!} \sim m^d$$



grows fast!
d = 6, m = 100
about 1.6 billion terms

Dot Product of Polynomials

$\Phi(\mathbf{x}) =$ polynomials of degree exactly d

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

$$d=1 \quad \Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \cdot \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = x_1 z_1 + x_2 z_2 = \mathbf{x} \cdot \mathbf{z}$$

$$\begin{aligned} d=2 \quad \Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) &= \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{bmatrix} \cdot \begin{bmatrix} z_1^2 \\ z_1 z_2 \\ z_2^2 \end{bmatrix} = x_1^2 z_1^2 + x_2^2 z_2^2 + x_1 x_2 z_1 z_2 \\ &= (x_1 z_1 + x_2 z_2)^2 \\ &= (\mathbf{x} \cdot \mathbf{z})^2 \end{aligned}$$

$$d \quad \Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) = K(\mathbf{x}, \mathbf{z}) = (\mathbf{x} \cdot \mathbf{z})^d$$

$$\dim(\mathcal{X}) = N$$

Name	Kernel function	$\dim(\mathcal{K})$
p th degree polynomial	$k(\vec{u}, \vec{v}) = (\langle \vec{u}, \vec{v} \rangle_{\mathcal{X}})^p$ $p \in \mathbb{N}^+$	$\binom{N+p-1}{p}$
complete polynomial	$k(\vec{u}, \vec{v}) = (\langle \vec{u}, \vec{v} \rangle_{\mathcal{X}} + c)^p$ $c \in \mathbb{R}^+, p \in \mathbb{N}^+$	$\binom{N+p}{p}$
RBF kernel	$k(\vec{u}, \vec{v}) = \exp\left(-\frac{\ \vec{u} - \vec{v}\ _{\mathcal{X}}^2}{2\sigma^2}\right)$ $\sigma \in \mathbb{R}^+$	∞
Mahalanobis kernel	$k(\vec{u}, \vec{v}) = \exp\left(-(\vec{u} - \vec{v})' \mathbf{\Sigma} (\vec{u} - \vec{v})\right)$ $\mathbf{\Sigma} = \text{diag}\left(\sigma_1^{-2}, \dots, \sigma_N^{-2}\right),$ $\sigma_1, \dots, \sigma_N \in \mathbb{R}^+$	∞

Picture is taken from R. Herbrich

The RBF kernel

Note:

The RBF kernel maps the input space \mathcal{X} onto the surface of an infinite dimensional hypersphere.

Proof:

$$\|\phi(x)\| = \sqrt{k(x, x)} = \sqrt{\exp(0)} = 1$$

Note:

The RBF kernel is shift invariant:

$$k(u + a, v + a) = k(u, v), \quad \forall a$$

Overfitting

- Huge feature space with kernels, what about overfitting???
- Maximizing margin leads to sparse set of support vectors
- Some interesting theory says that SVMs search for simple hypothesis with large margin
- Often robust to overfitting

String kernels

P-spectrum kernel:

$P=3$: s ="statistics" t ="computation"

They contain the following substrings of length 3

"sta", "tat", "ati", "tis", "ist", "sti", "tic", "ics"

"com", "omp", "mpu", "put", "uta", "tat", "ati", "tio", "ion"

Common substrings: "tat", "ati"

$k(s,t)=2$

Distribution kernels

Euclidean:

$$K(p, q) = \int p(x) q(x) dx$$

Bhattacharyya's affinity:

$$K(p, q) = \int \sqrt{p(x)} \sqrt{q(x)} dx$$

Mean map:

$$\Phi(p) = \mathbb{E}_{x \sim p} K(\cdot, x)$$

$$K(p, q) = \mathbb{E}_{x \sim p} \mathbb{E}_{y \sim q} K(x, y)$$

Set kernels

Mean map:

$$\frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m \ell_2(x_i, y_j) \stackrel{\text{def}}{=} \langle \{x_1, \dots, x_n\}, \{y_1, \dots, y_m\} \rangle$$

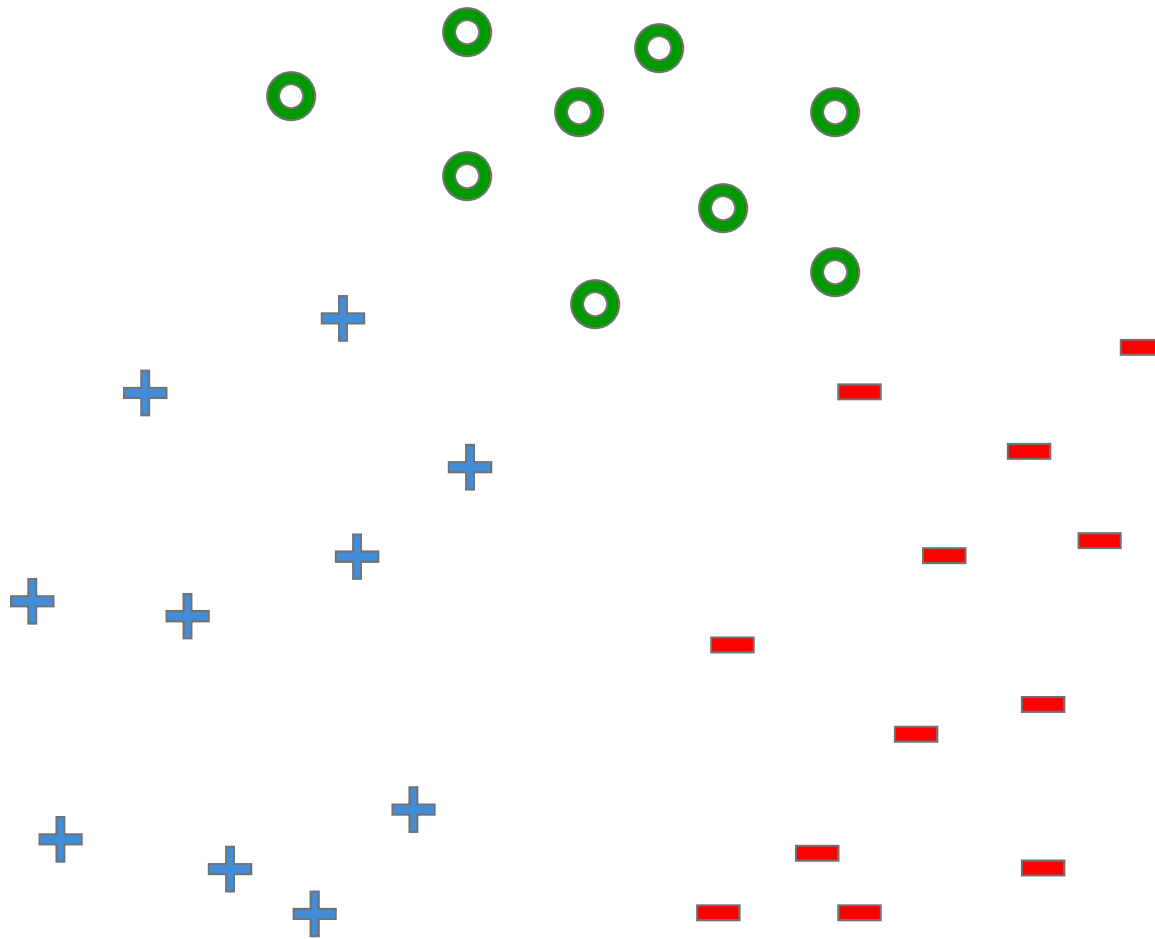
Intersection kernel:

$$\begin{aligned} \ell_2(A_1, A_2) &= \int_{A_1 \cap A_2} I(x) d\mu(x) \\ &= \mu(A_1 \cap A_2) \end{aligned}$$

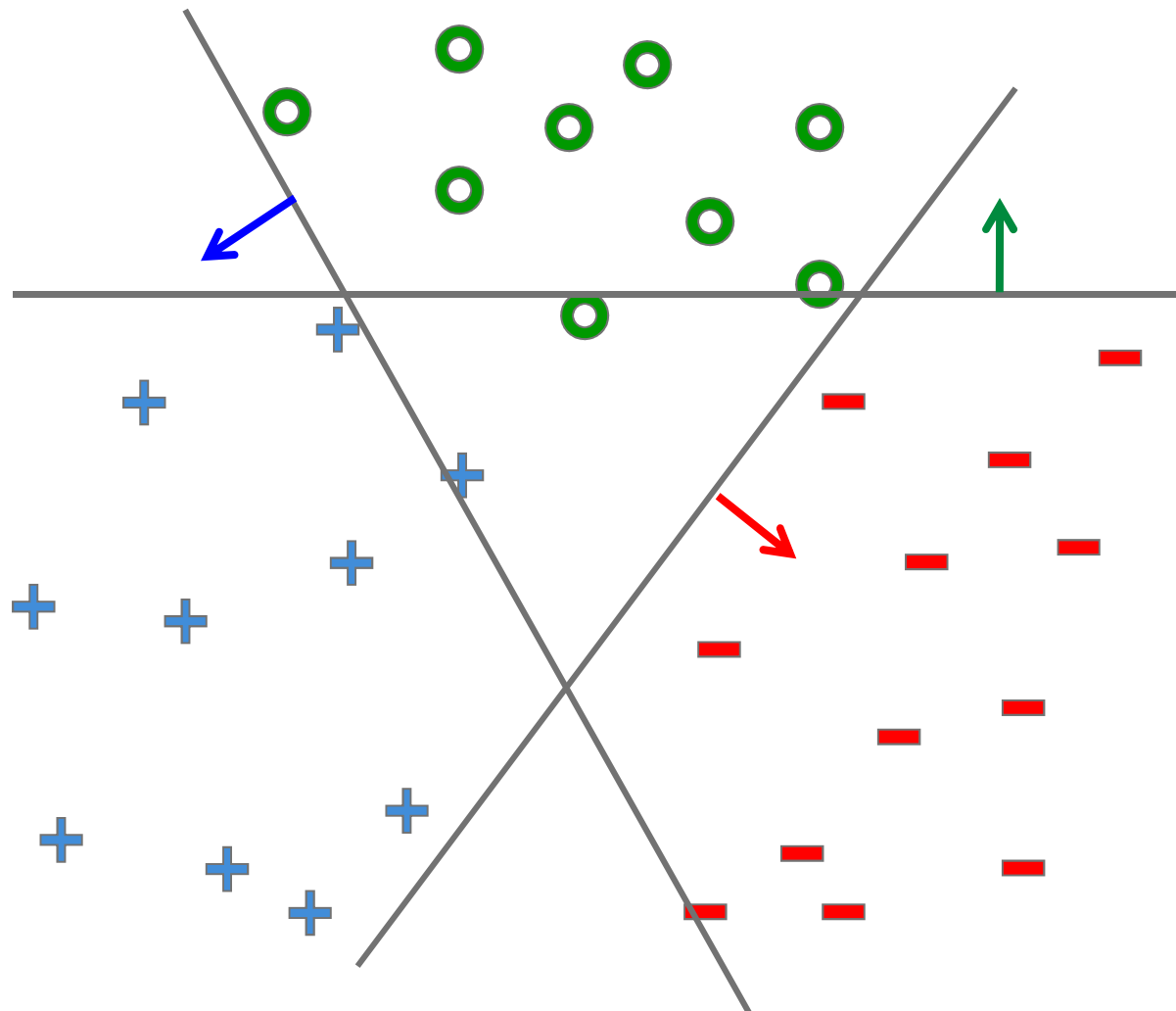
Union complement kernel:

$$1 - \mu(A_1 \cup A_2) \quad \mu(\mathcal{Q}) = 1$$

What about multiple classes?



One against all



Learn 3 classifiers
separately:

Class k vs. rest

$$(\mathbf{w}_k, b_k)_{k=1,2,3}$$

$$y = \arg \max_k \mathbf{w}_k \cdot \mathbf{x} + b_k$$

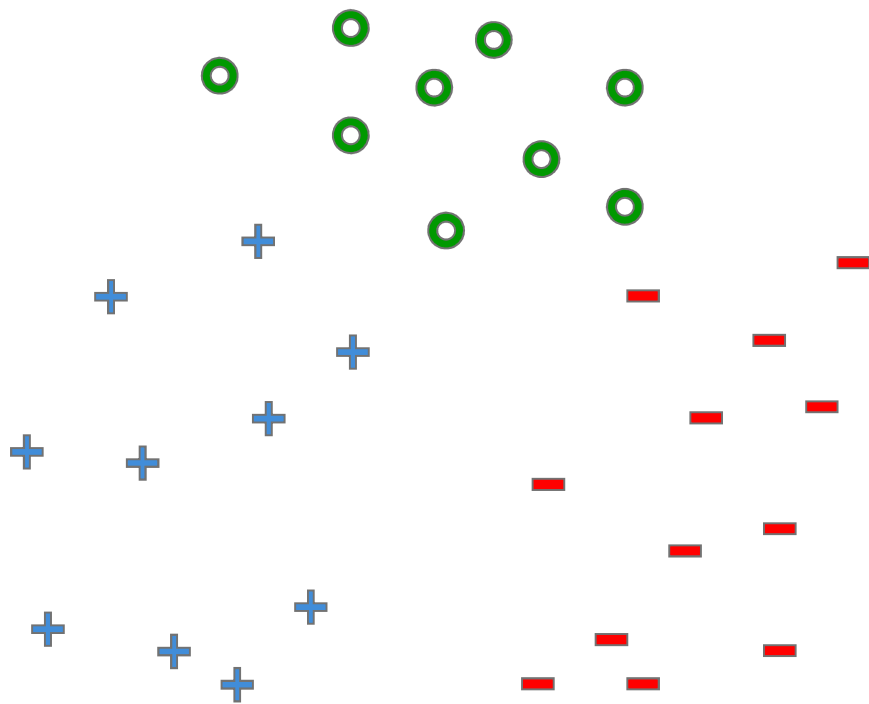
But \mathbf{w}_k s may not be
based on the same
scale.

Note: $(a\mathbf{w}) \cdot \mathbf{x} + (ab)$ is
also a solution

Learn 1 classifier: Multi-class SVM

Simultaneously learn 3 sets of weights

$$\mathbf{w}^{(y_j)} \cdot \mathbf{x}_j + b^{(y_j)} \geq \mathbf{w}^{(y')} \cdot \mathbf{x}_j + b^{(y')} + 1, \quad \forall y' \neq y_j, \quad \forall j$$

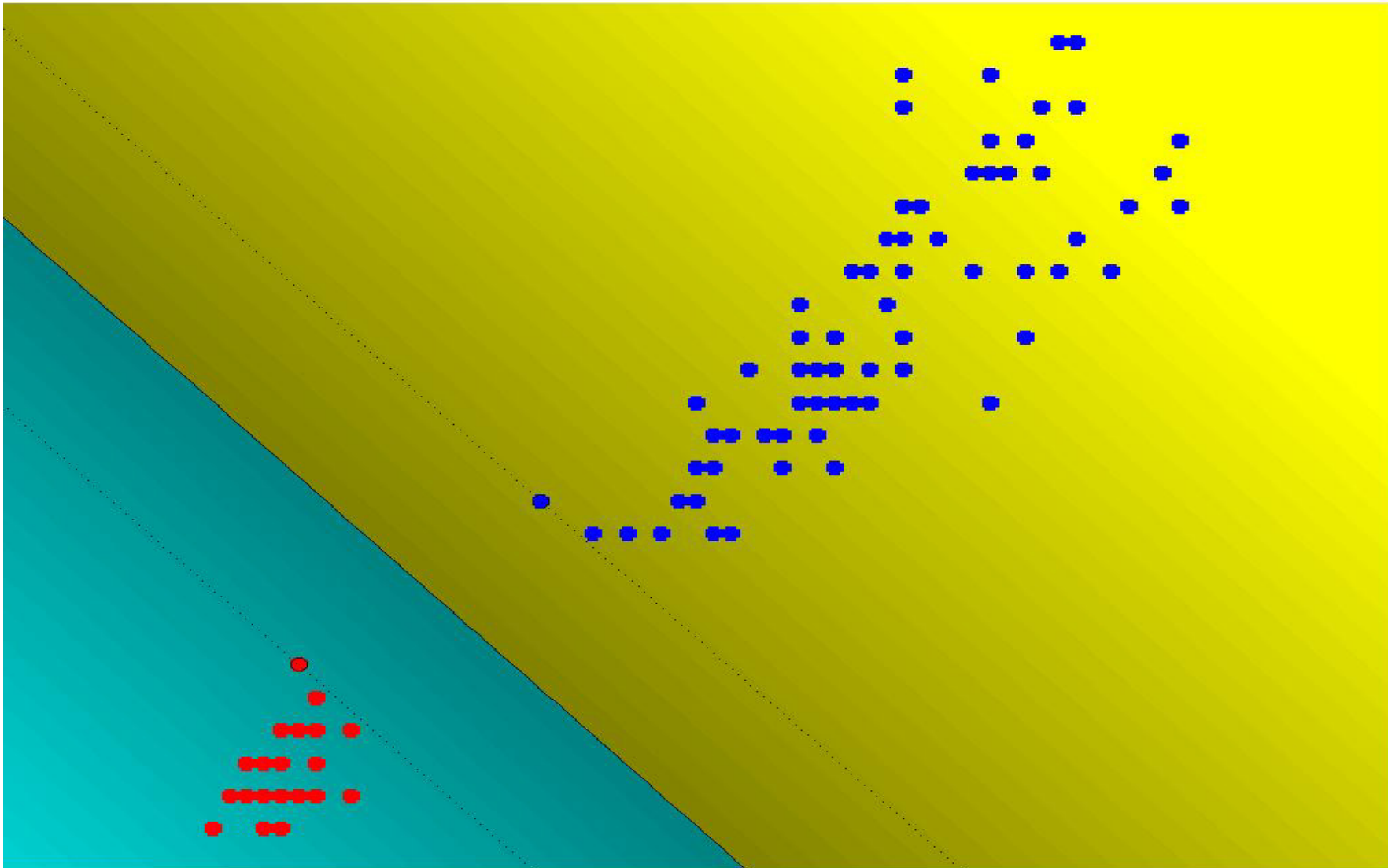


Margin - gap between correct class and nearest other class

$$y = \arg \max \mathbf{w}^{(k)} \cdot \mathbf{x} + b^{(k)}$$

Steve Gunn's svm toolbox

Results, Iris 2vs13, Linear kernel



No. of Support Vectors: 2 (1.7%)

Results, Iris 1vs23, 2nd order kernel

Polynomial



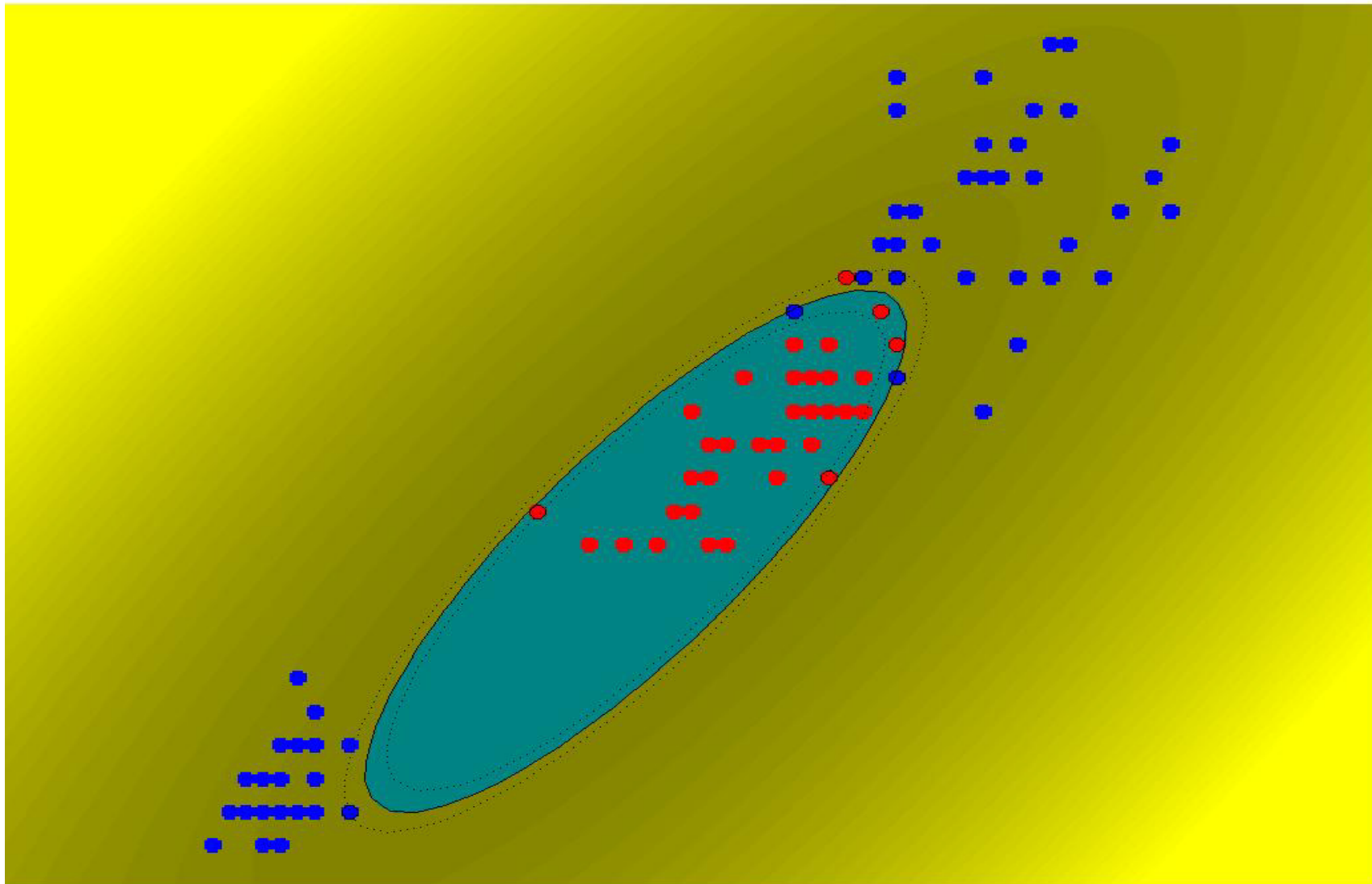
Degree

2

Separable

Bound

Inf



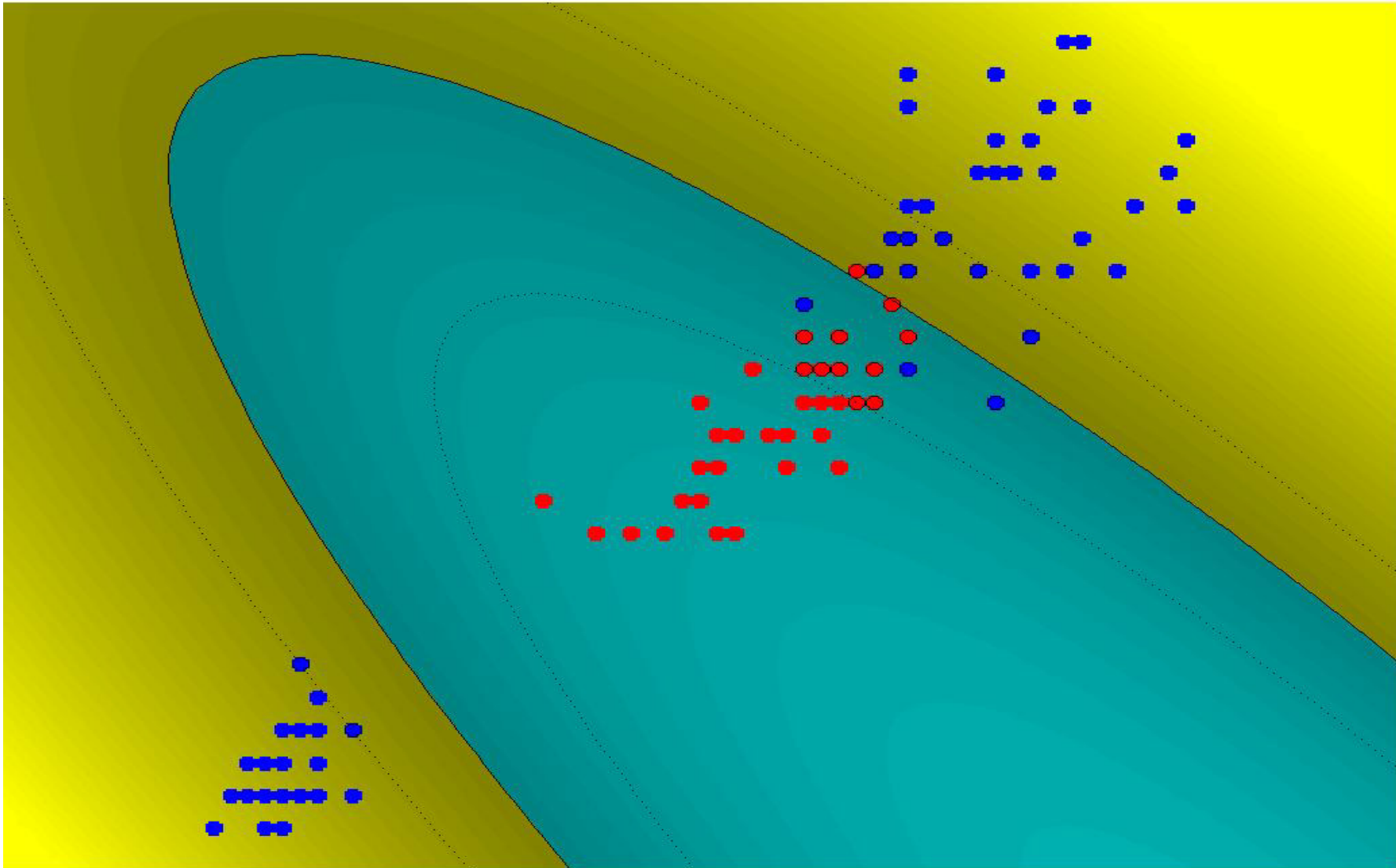
No. of Support Vectors: 12 (10.0%)

Polynomial

Degree 2

Separable

Bound 1



No. of Support Vectors: 30 (25.0%)

$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma^2}\right)$$

$\sigma \rightarrow 0 \Rightarrow$ MORE SUPPORT VECTORS

Gaussian RBF

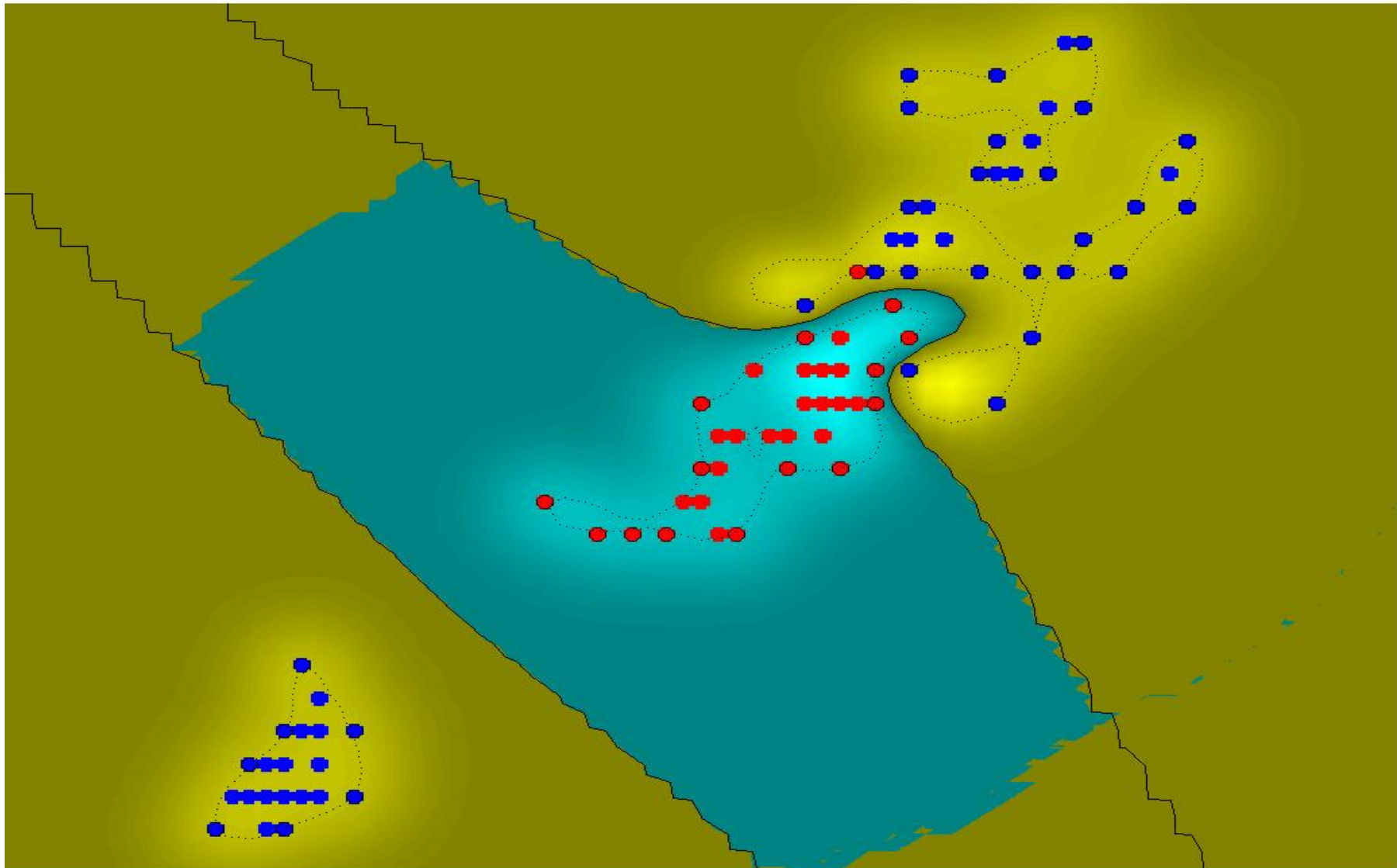
Sigma

.1

Separable

Bound

10



No. of Support Vectors: 55 (45.8%)

Gaussian RBF



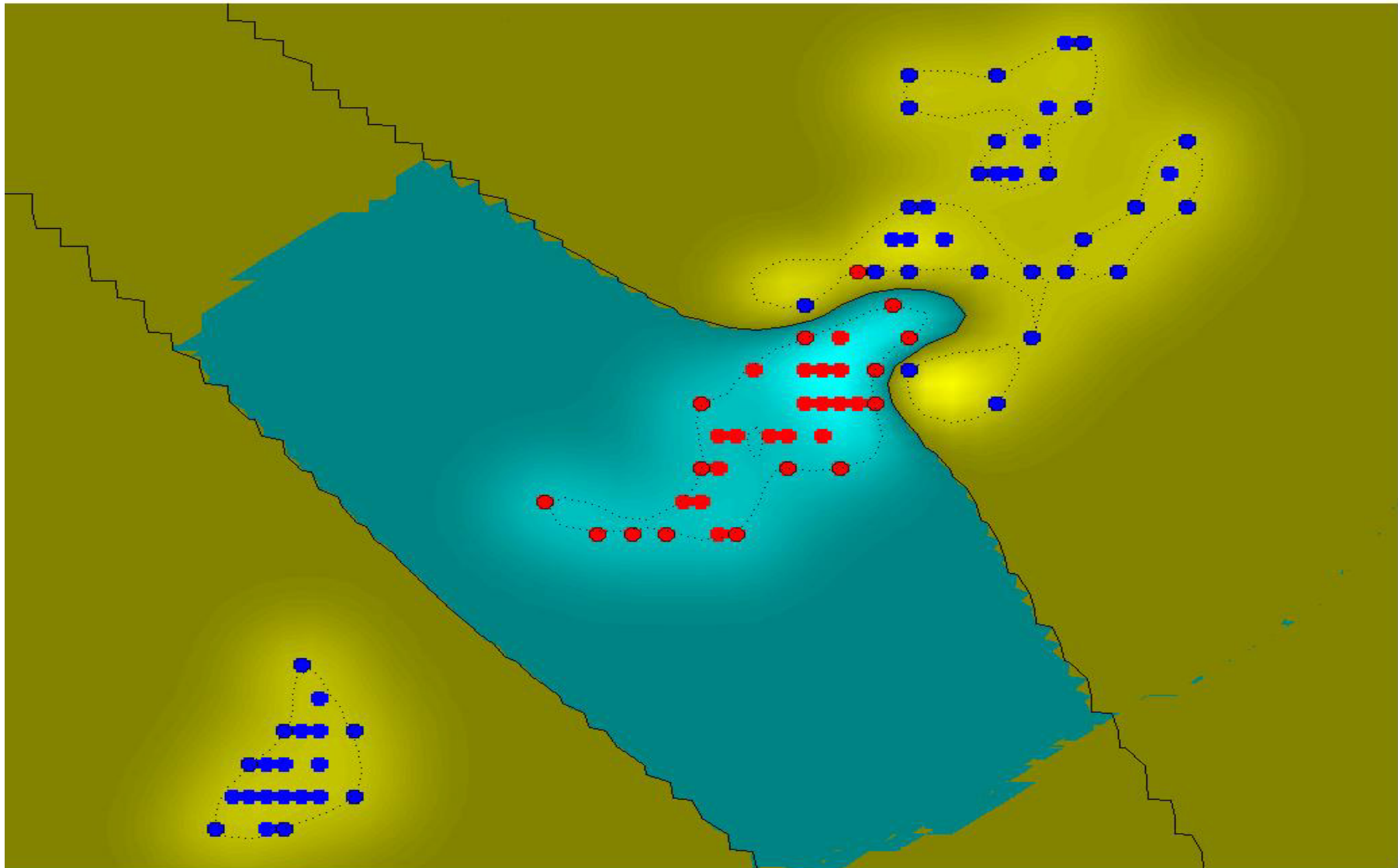
Sigma

.1

Separable

Bound

10



No. of Support Vectors: 55 (45.8%)

Gaussian RBF

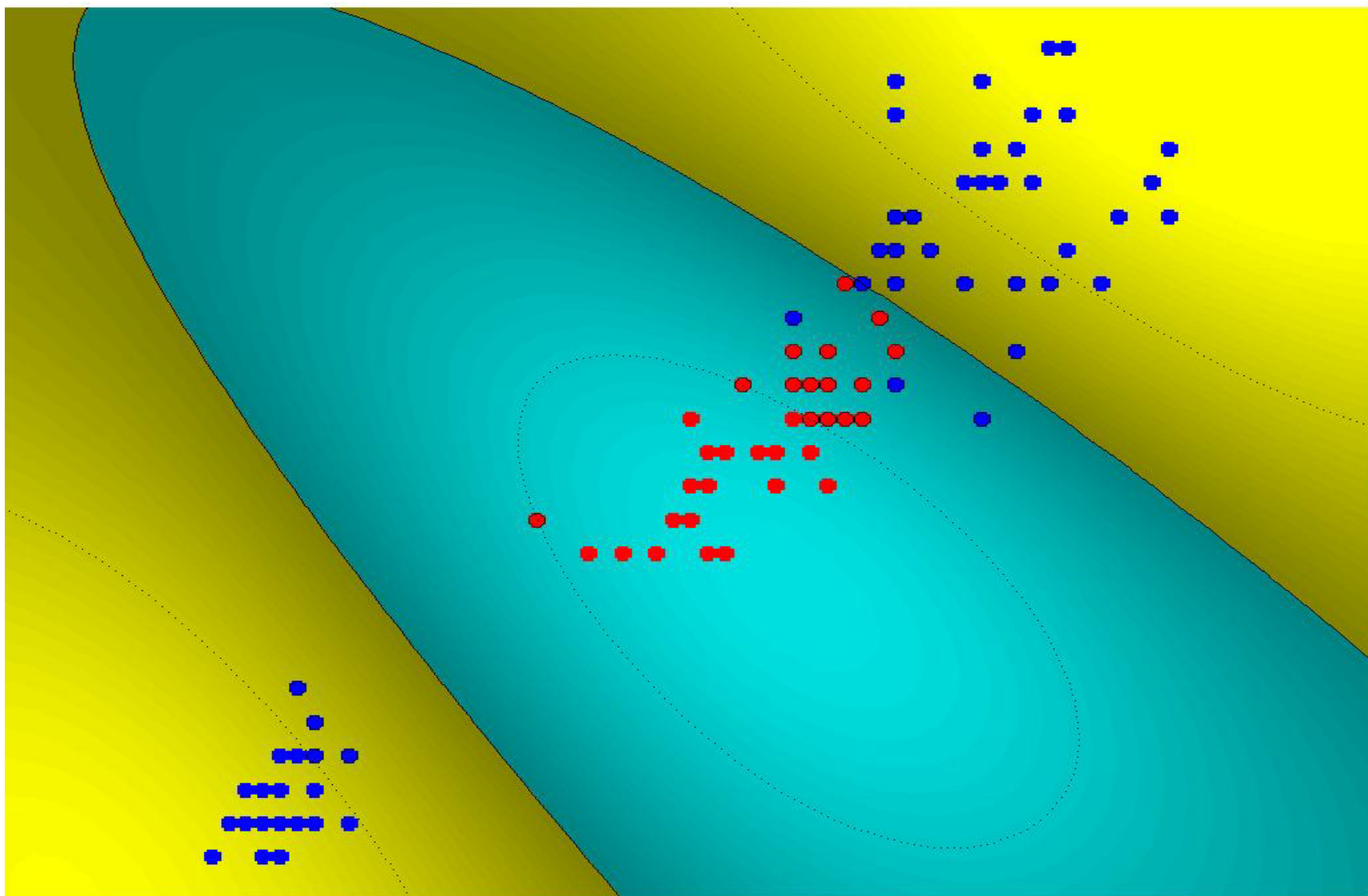
Sigma

1

Separable

Bound

1



No. of Support Vectors: 41 (34.2%)

Polynomial



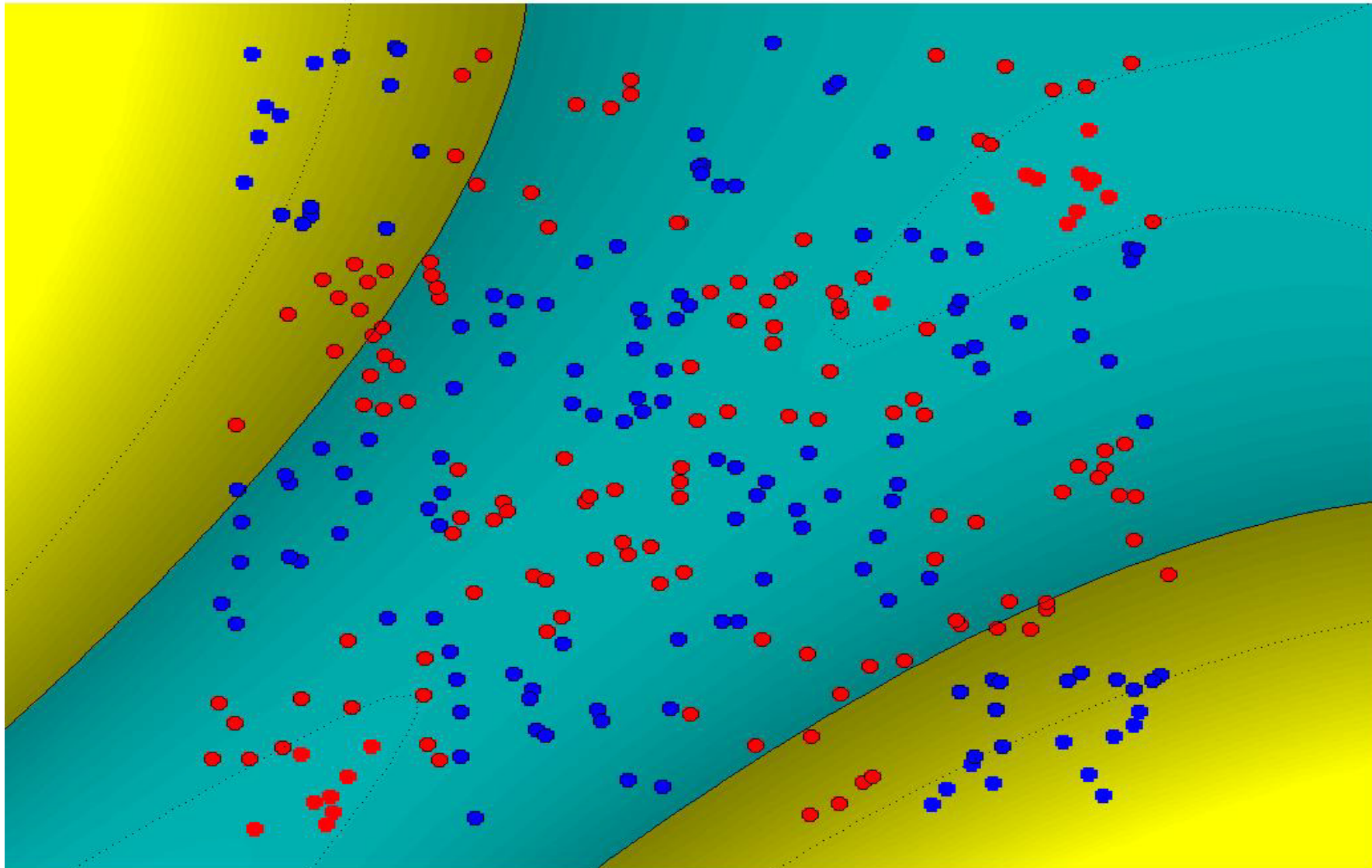
Degree

3

Separable

Bound

1



No. of Support Vectors: 263 (87.7%)

Polynomial



Degree

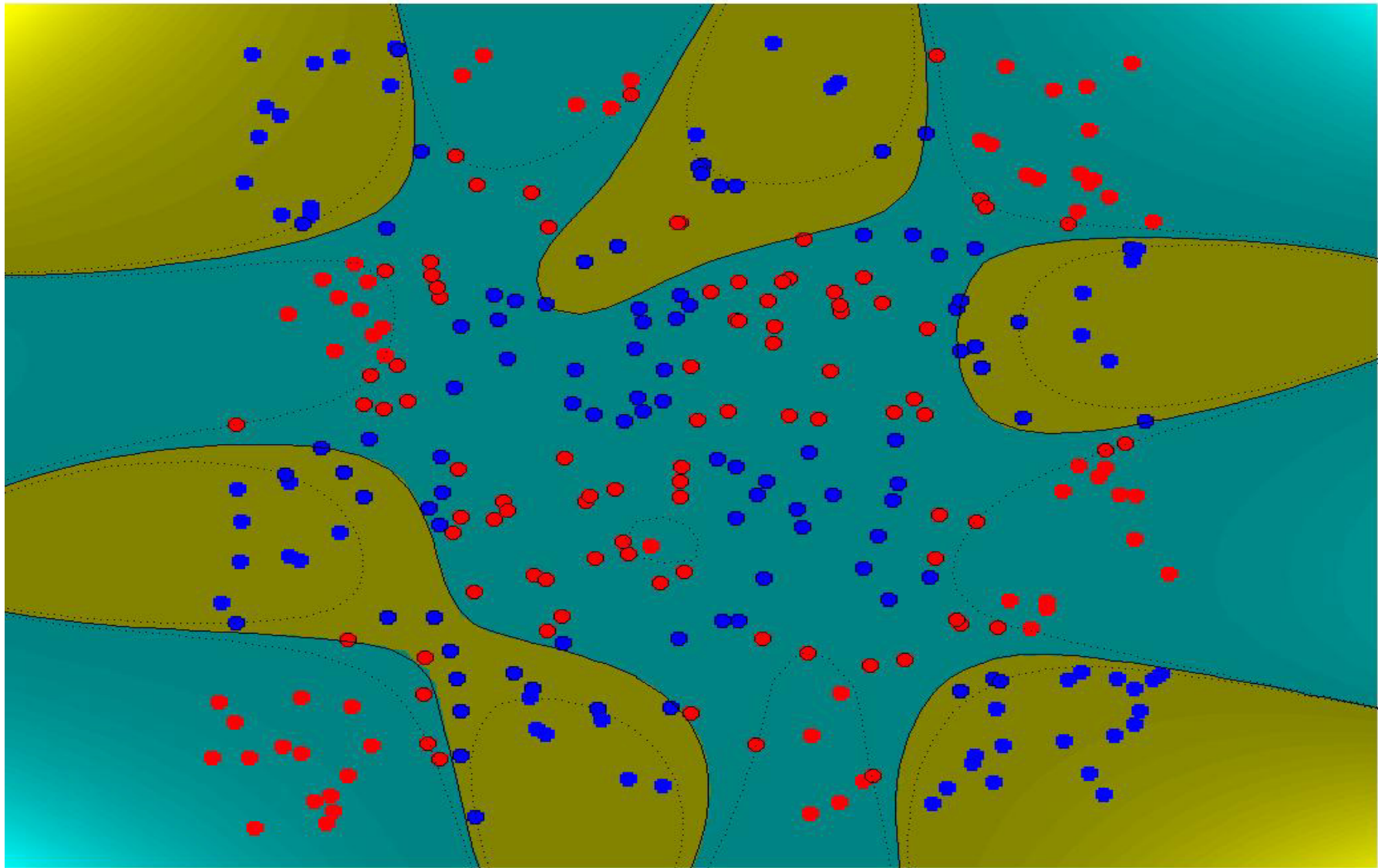
8



Separable

Bound

1



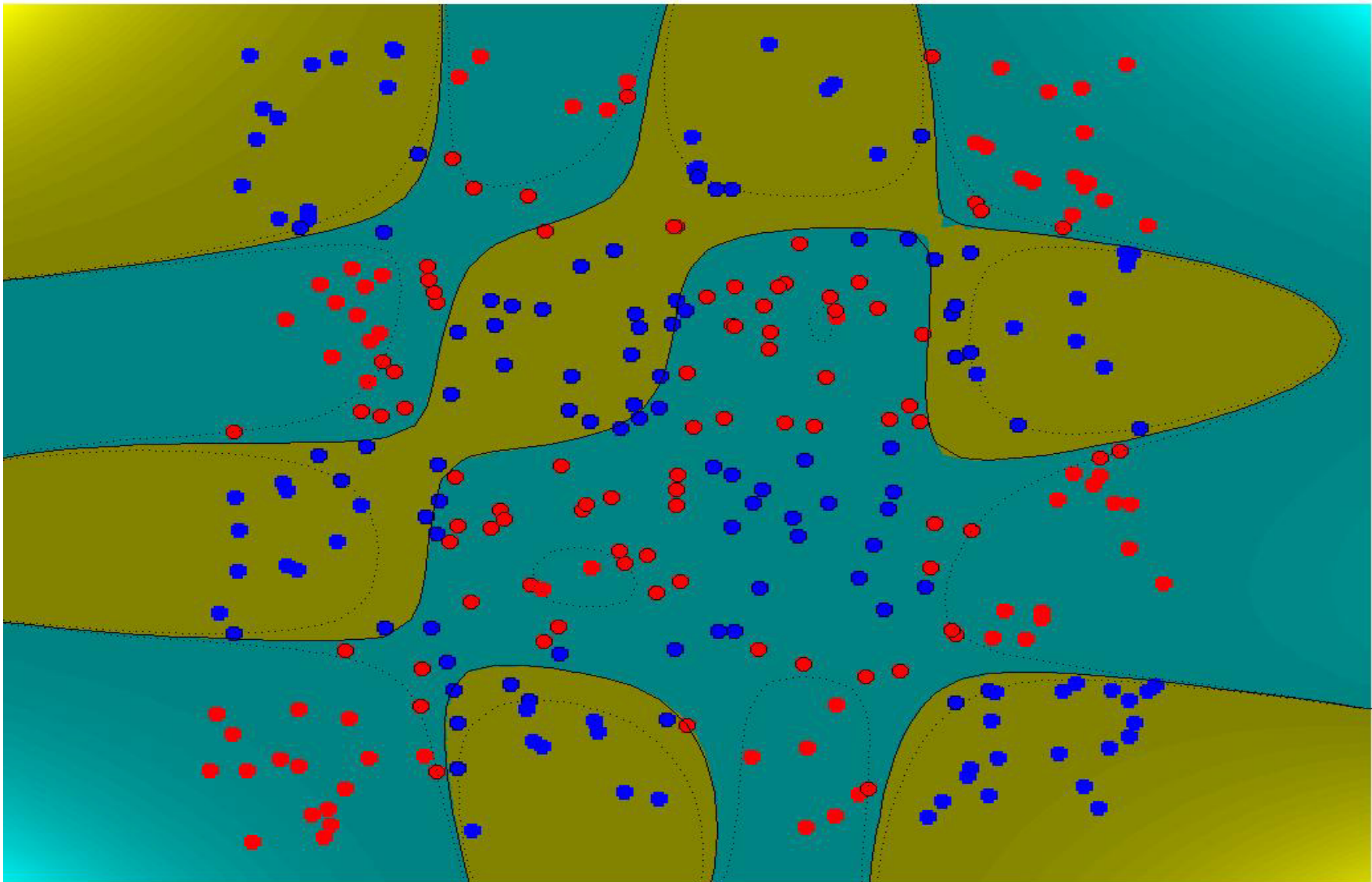
No. of Support Vectors: 183 (61.0%)

Polynomial

Degree 9

Separable

Bound 1



No. of Support Vectors: 164 (54.7%)

Polynomial

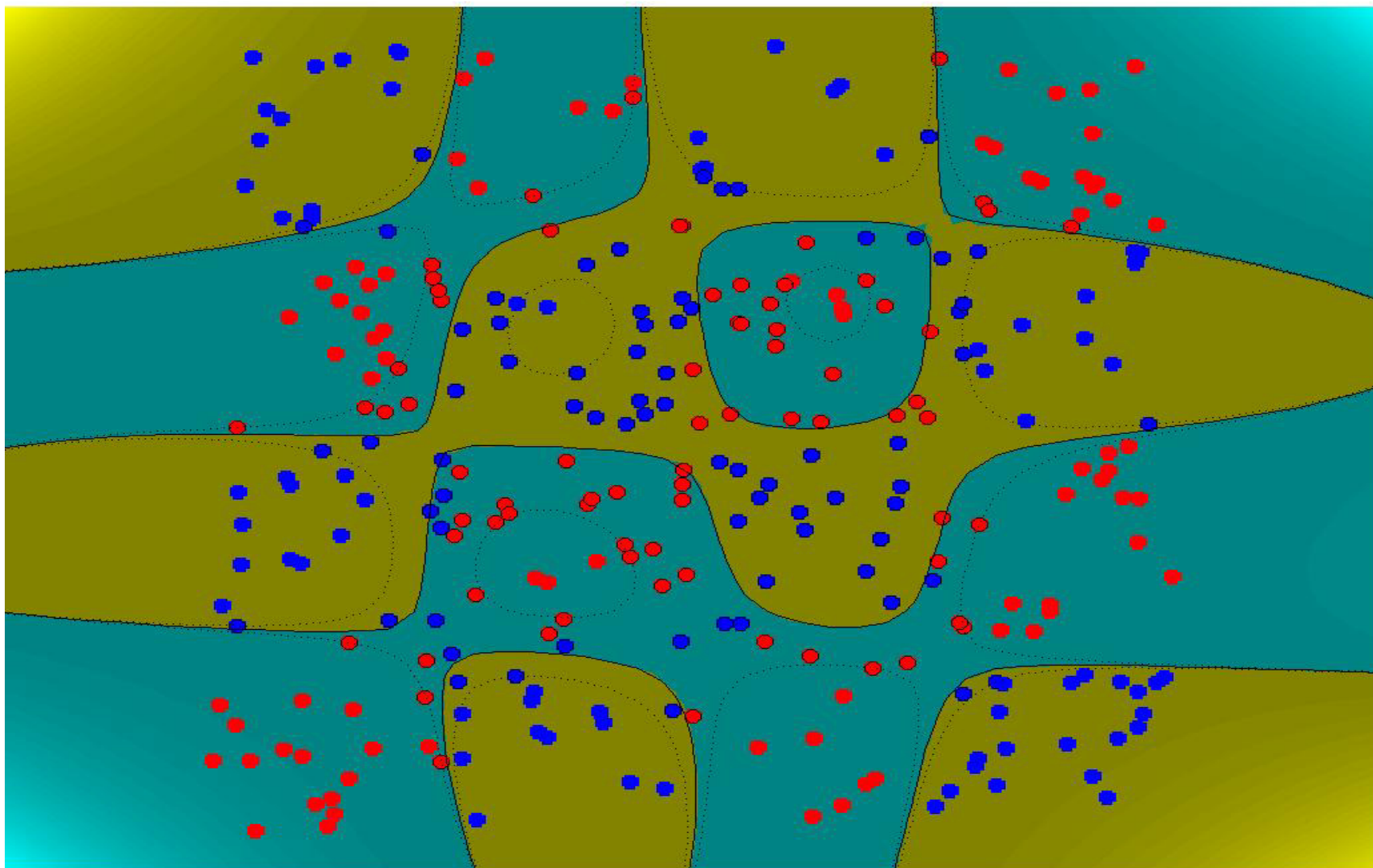
Degree

10

Separable

Bound

1



No. of Support Vectors: 147 (49.0%)

Polynomial



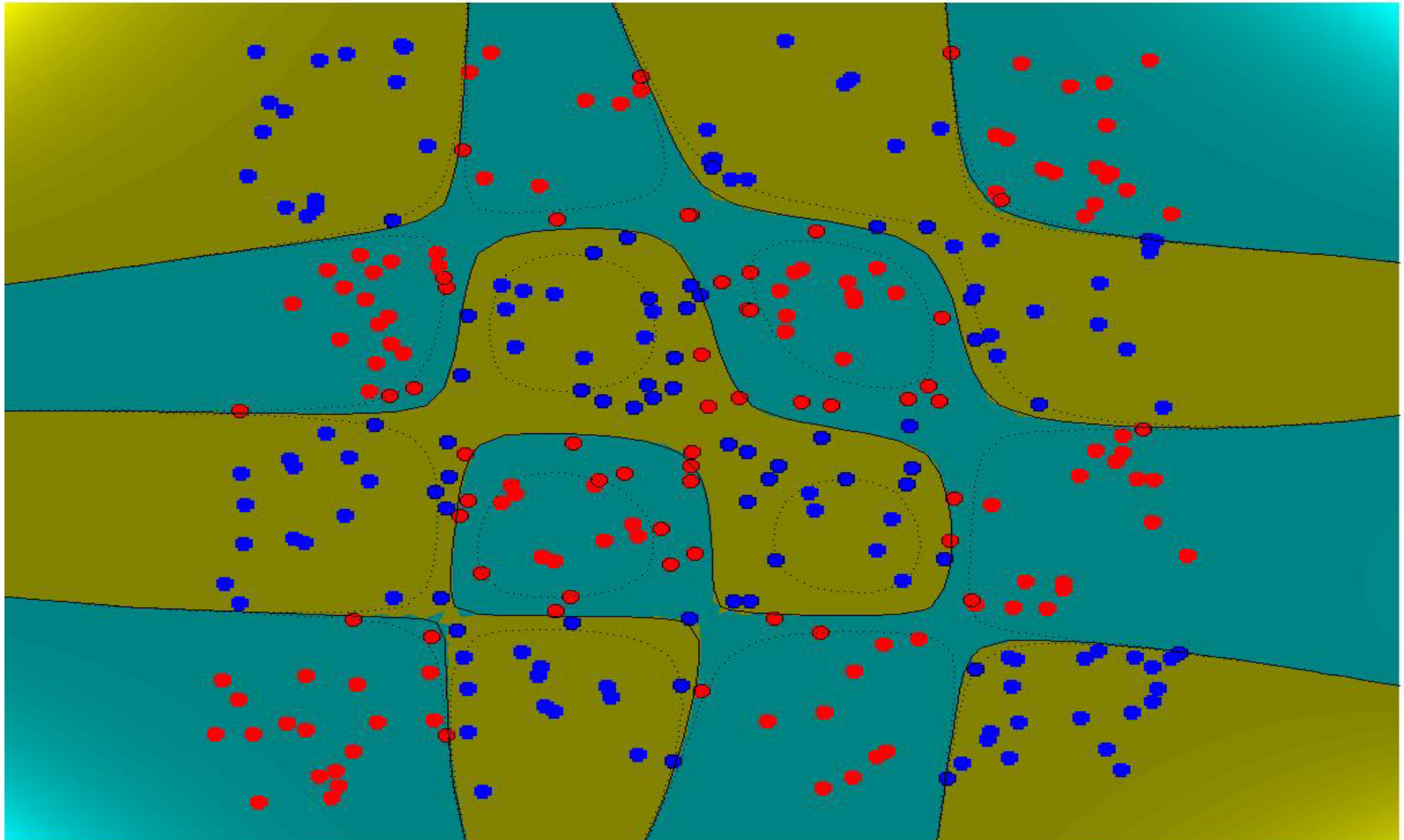
Degree

13

Separable

Bound

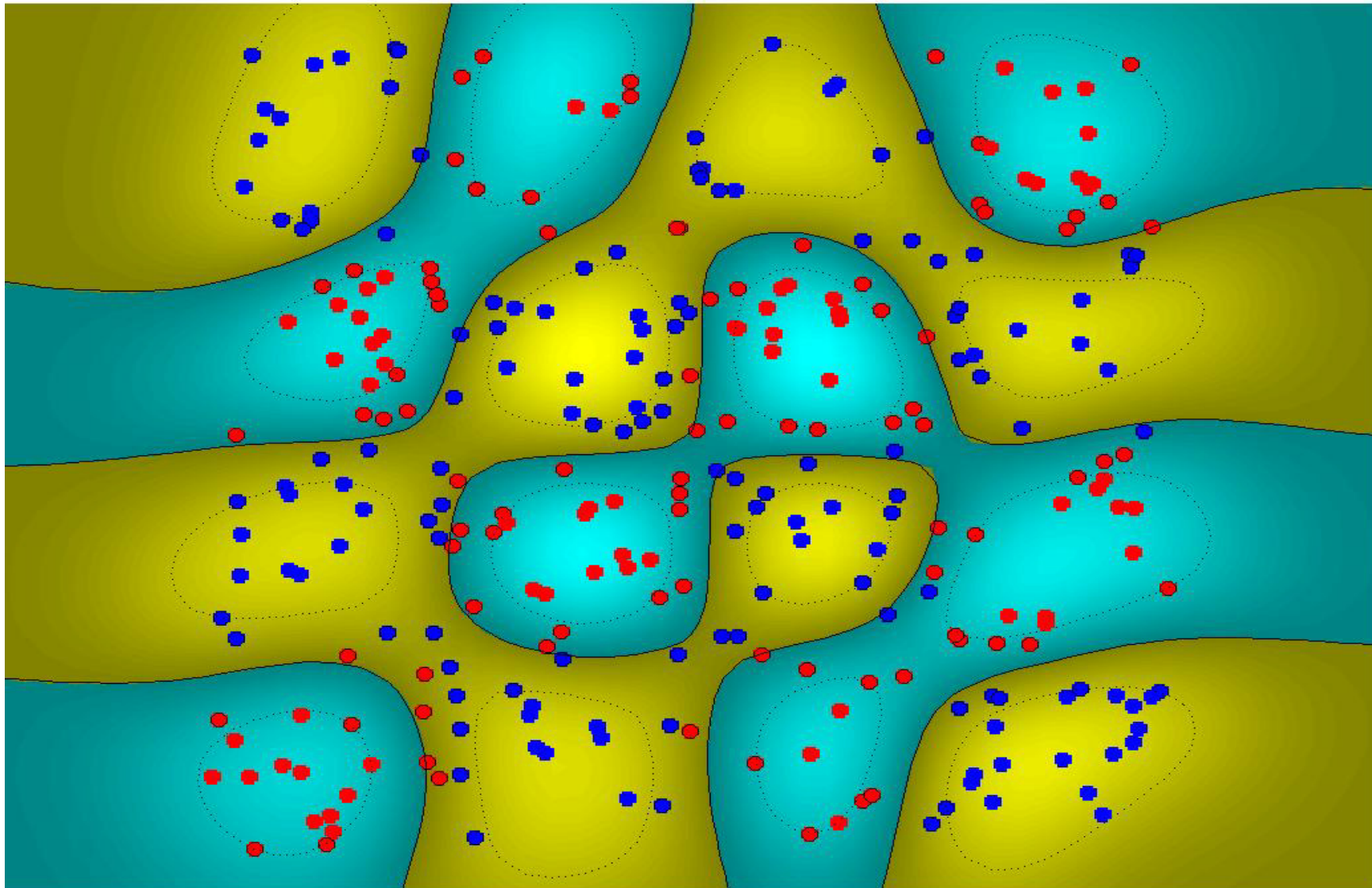
1



No. of Support Vectors: 102 (34.0%)

Results, Chessboard, RBF kernel

Gaussian RBF Sigma .2 Separable Bound 1



No. of Support Vectors: 174 (58.0%)

Sinc= $\sin(\pi x) / (\pi x)$, RBF kernel

Gaussian RBF

Sigma

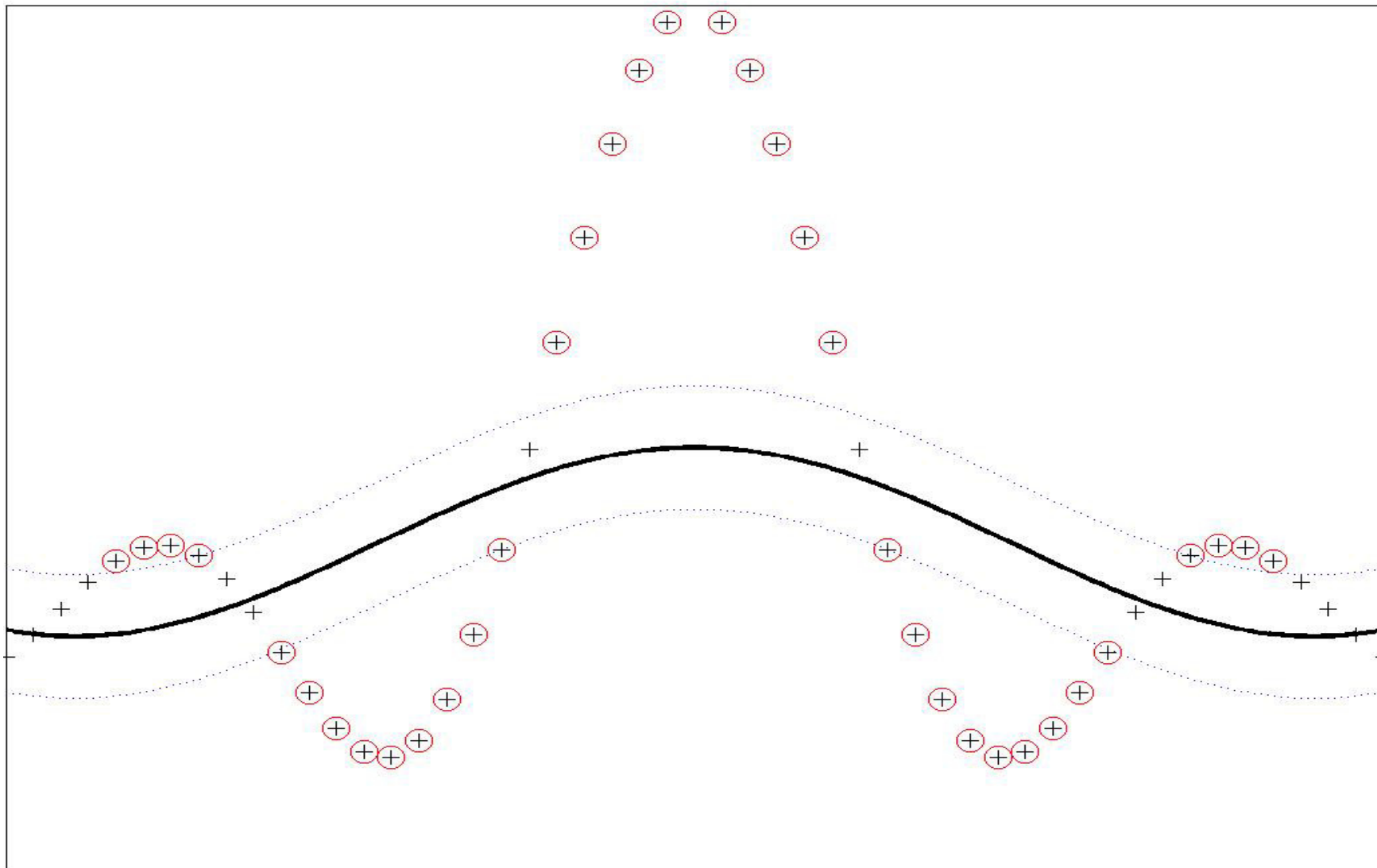
1

Bound

10

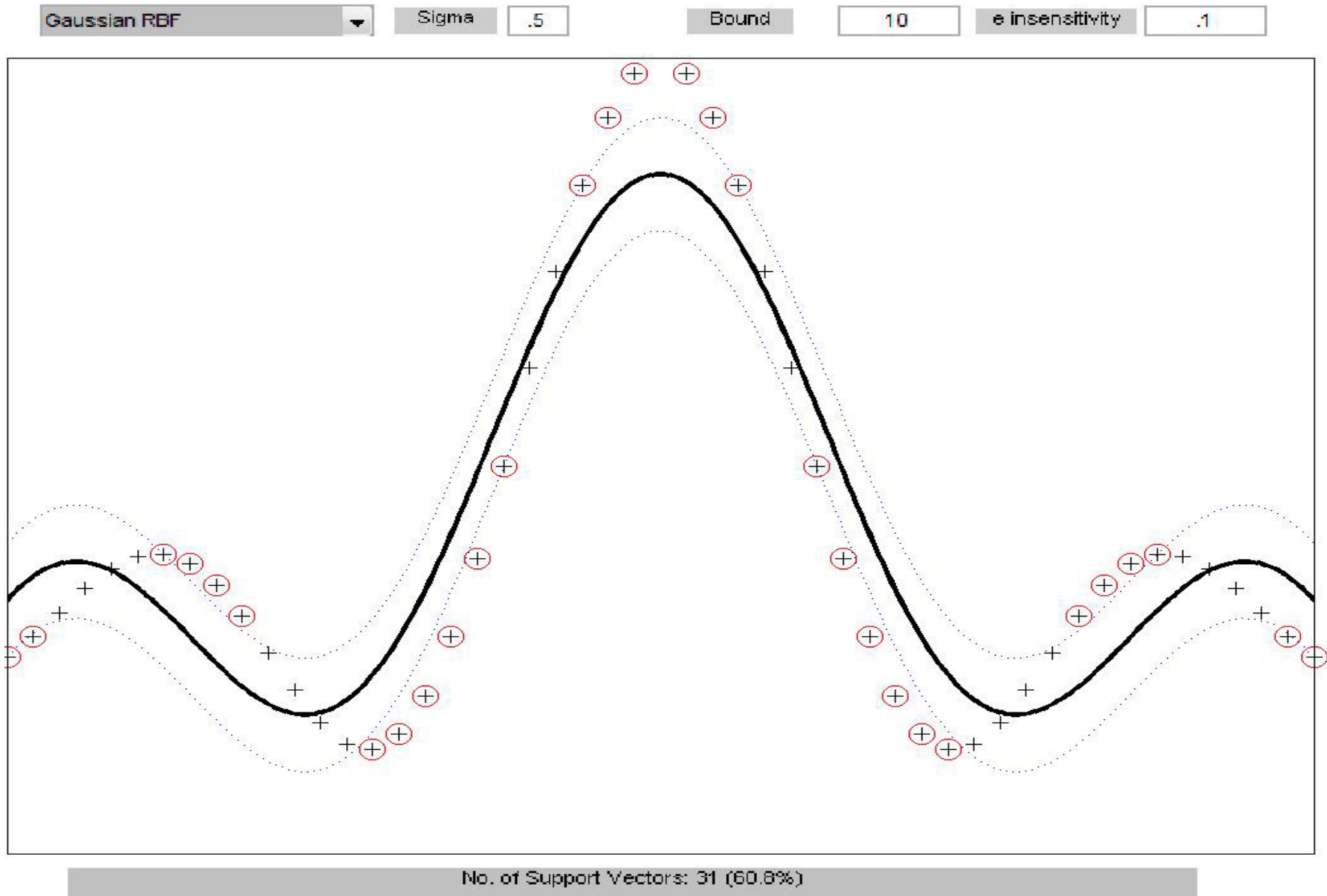
e insensitivity

.1



No. of Support Vectors: 37 (72.5%)

Sinc= $\sin(\pi x) / (\pi x)$, RBF kernel



Sinc= $\sin(\pi x) / (\pi x)$, RBF kernel

Gaussian RBF



Sigma

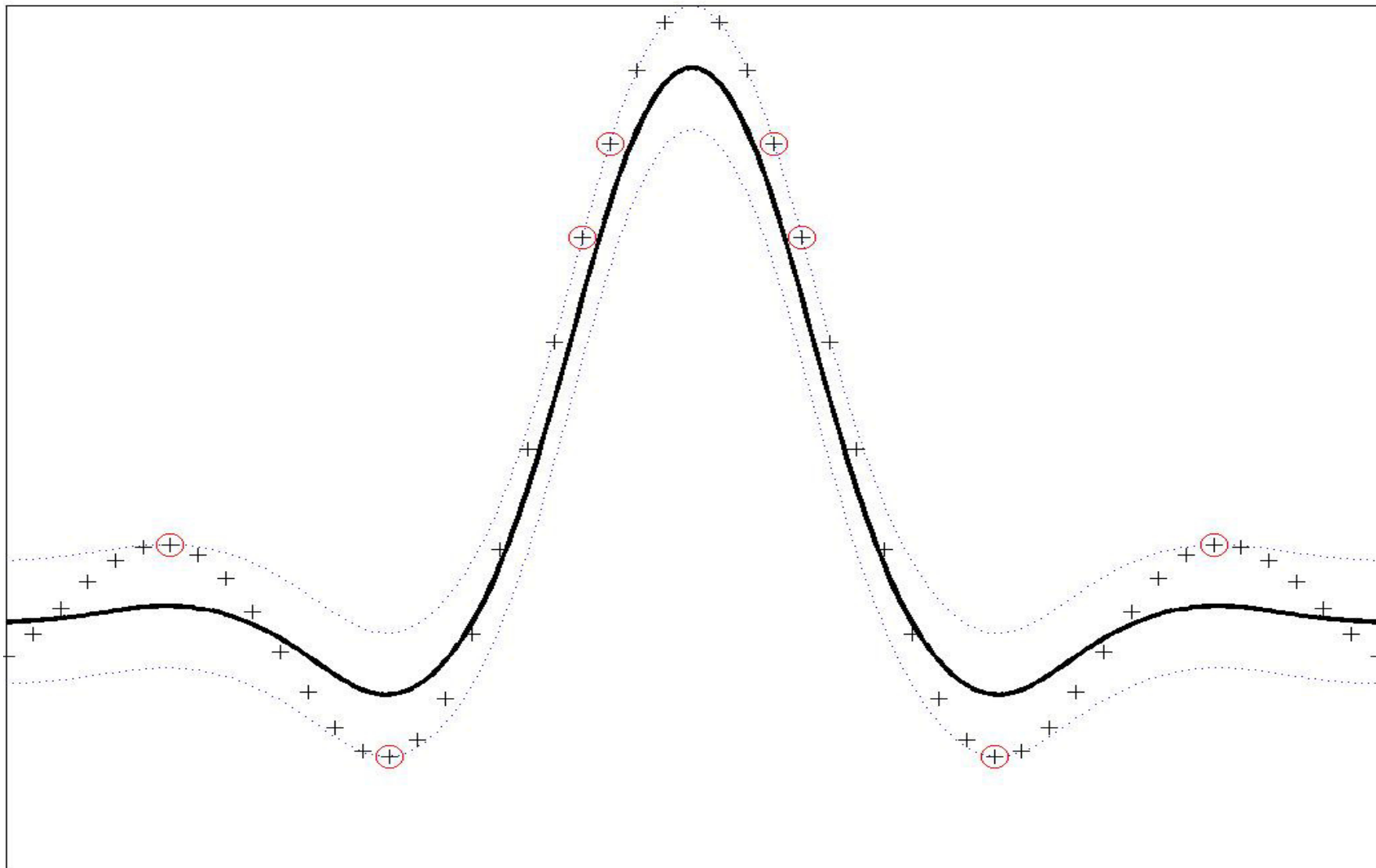
.1

Bound

10

e insensitivity

.1



No. of Support Vectors: 9 (17.6%)

Sinc= $\sin(\pi x) / (\pi x)$, RBF kernel

Gaussian RBF

Sigma

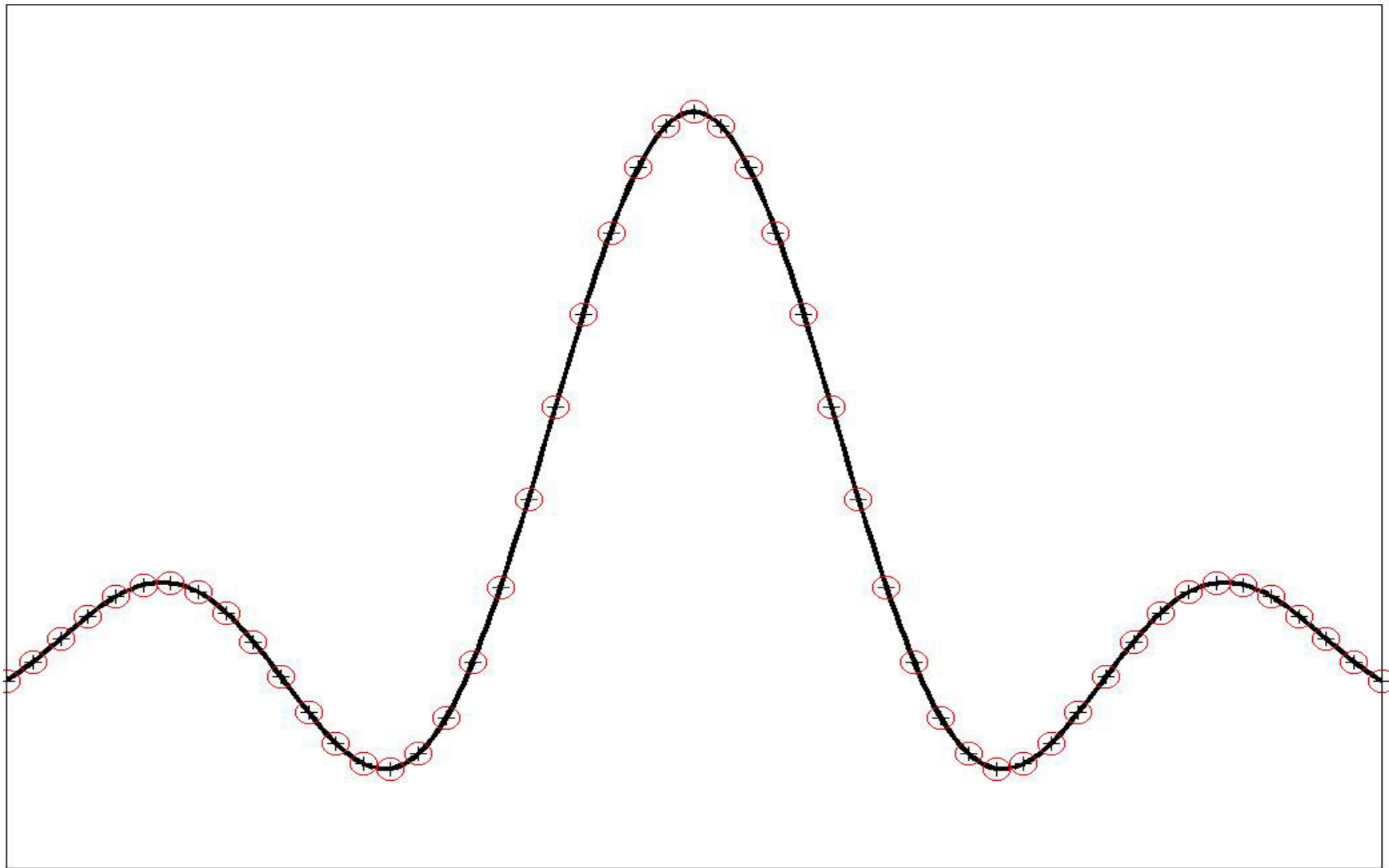
0.1

Bound

Inf

ϵ insensitivity

0



No. of Support Vectors: 51 (100.0%)

Sinc= $\sin(\pi x) / (\pi x)$, RBF kernel

Gaussian RBF

Sigma

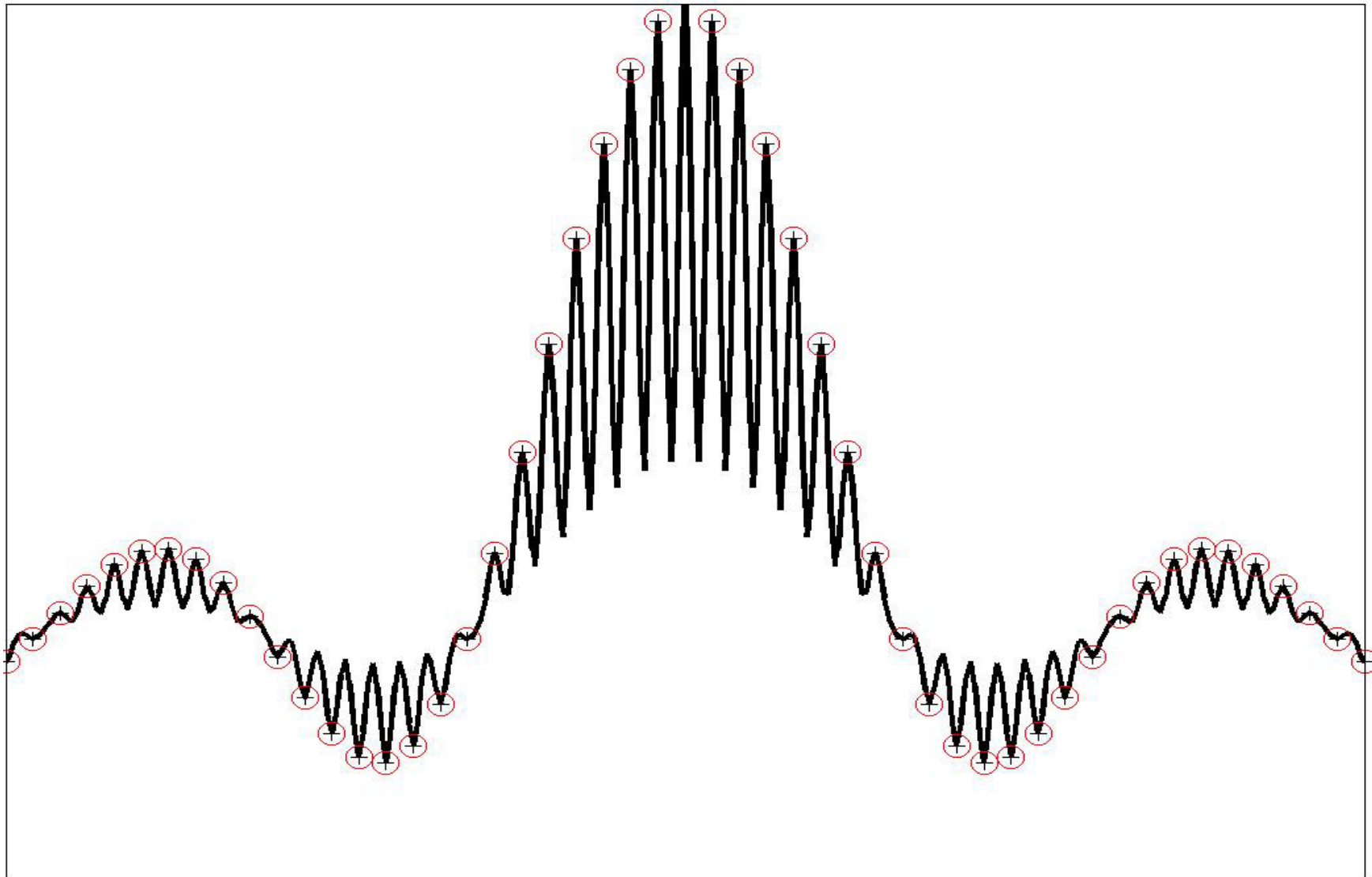
0.01

Bound

Inf

e insensitivity

0



No. of Support Vectors: 51 (100.0%)

Sinc= $\sin(\pi x) / (\pi x)$, RBF kernel

Gaussian RBF

Sigma

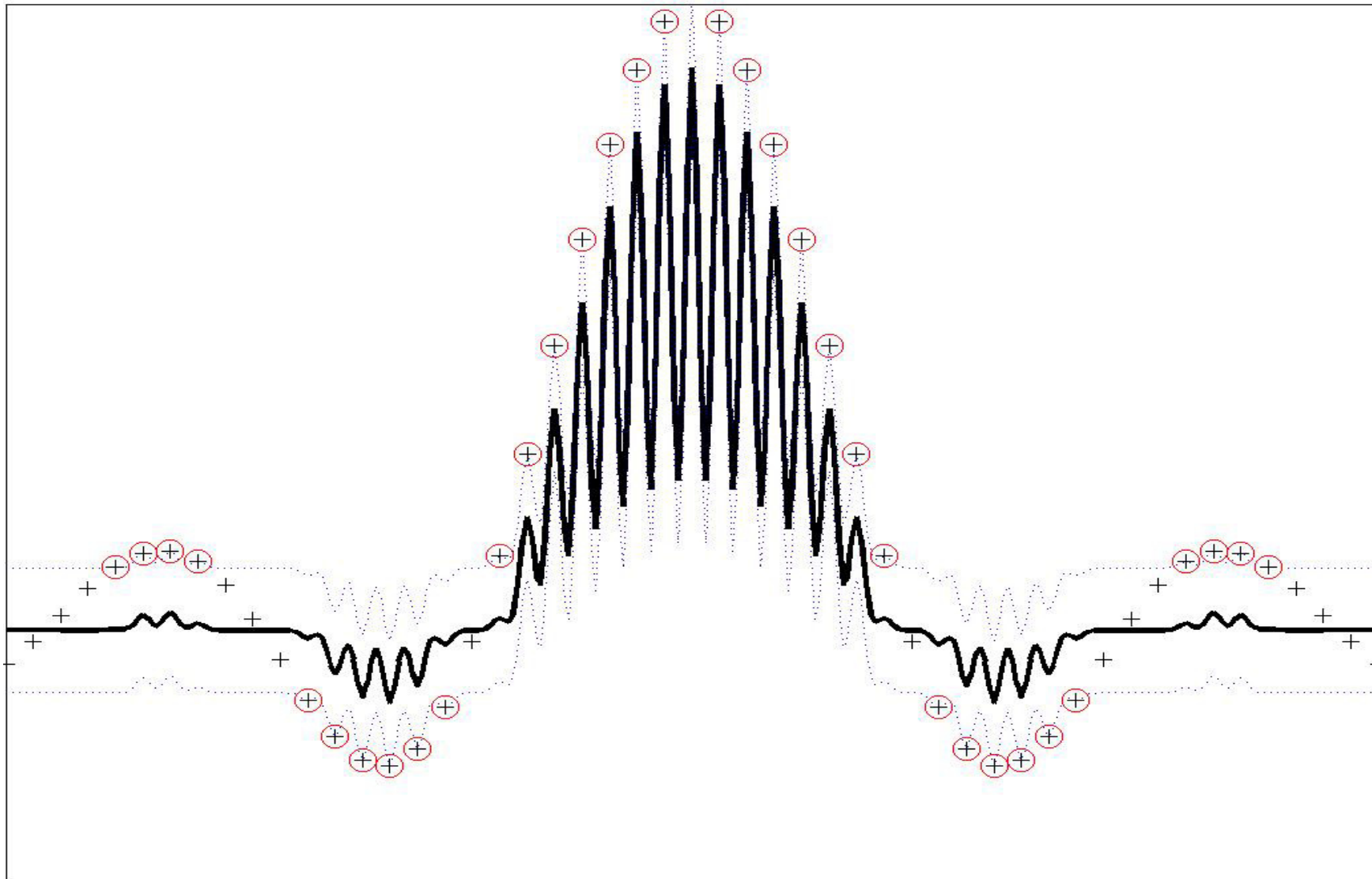
.01

Bound

10

e insensitivity

.1



No. of Support Vectors: 35 (68.6%)

Sinc= $\sin(\pi x) / (\pi x)$, poly kernel

Polynomial

Degree

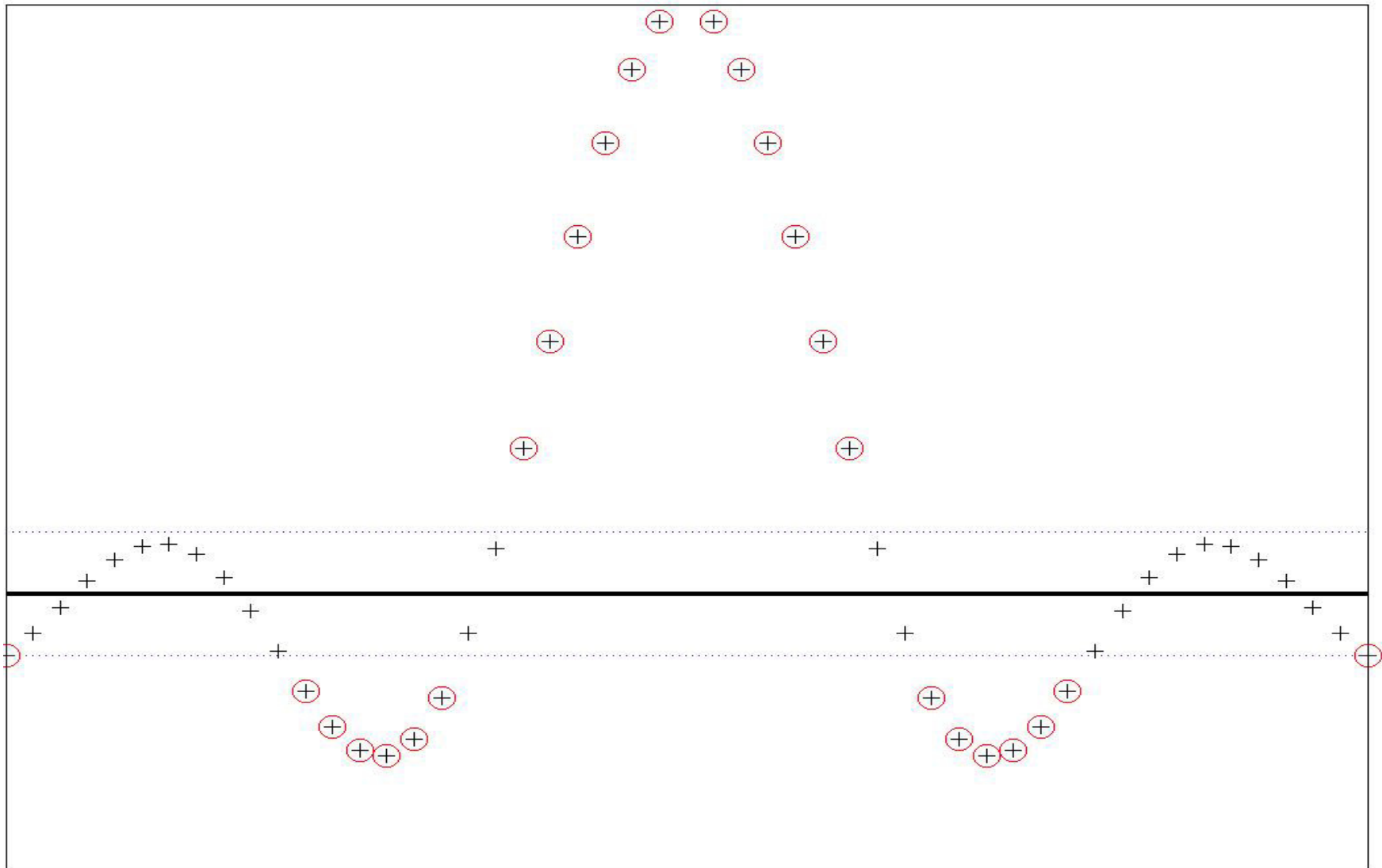
1

Bound

10

ϵ insensitivity

.1



No. of Support Vectors: 27 (52.9%)

Sinc= $\sin(\pi x) / (\pi x)$, poly kernel

Polynomial



Degree

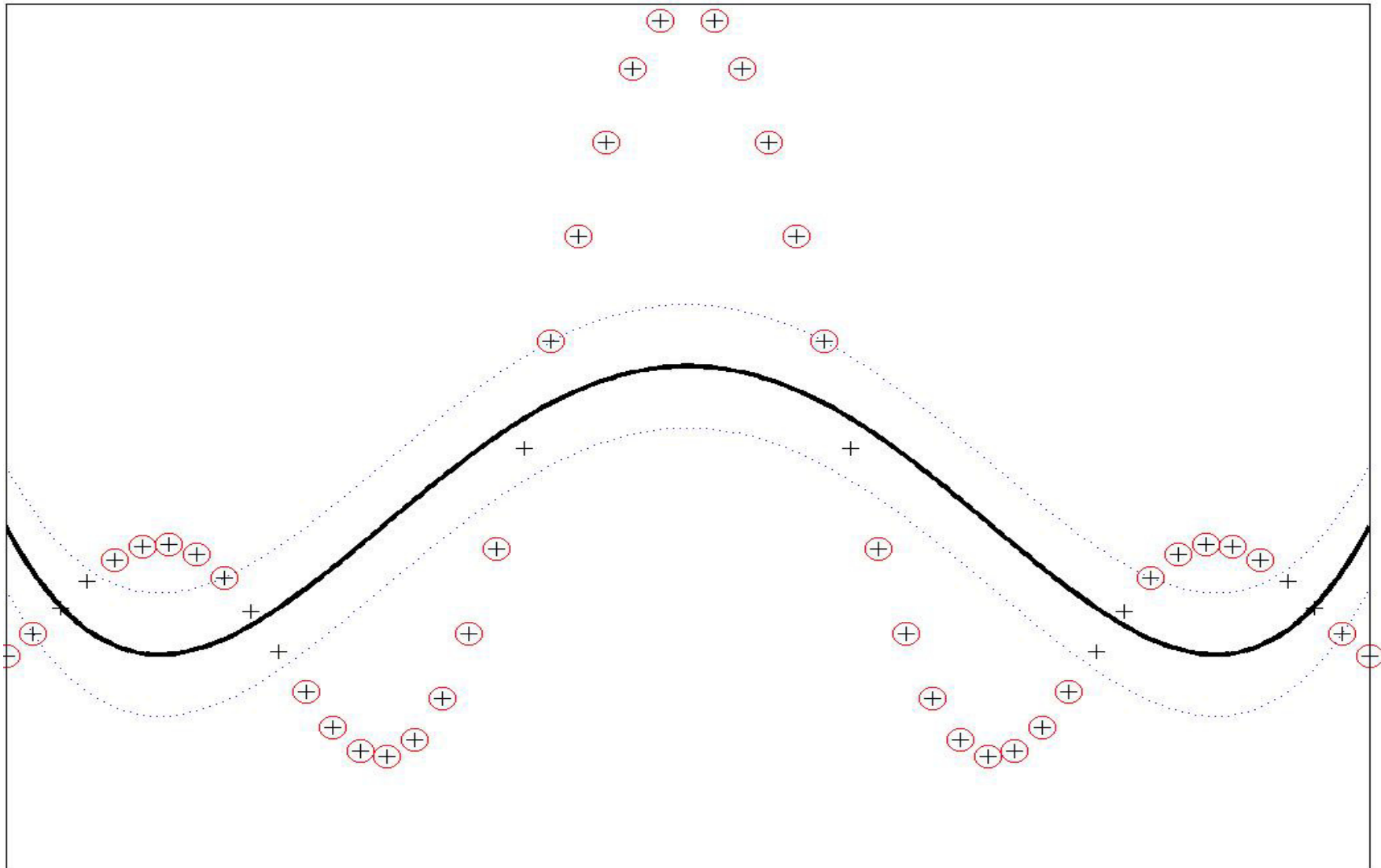
4

Bound

10

e insensitivity

.1



No. of Support Vectors: 41 (80.4%)

Sinc= $\sin(\pi x) / (\pi x)$, poly kernel

Polynomial

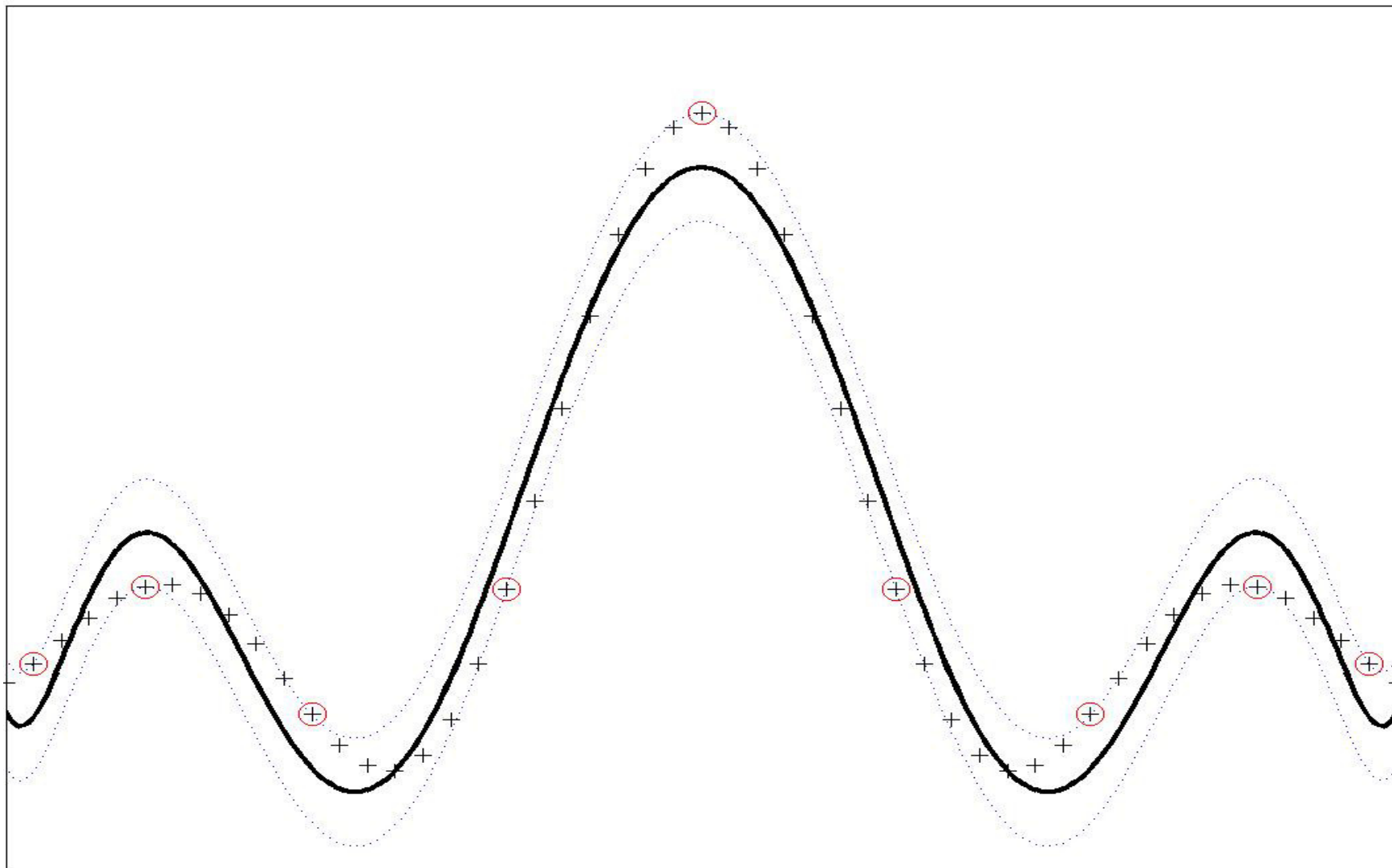
Degree 13

Bound

10

ϵ insensitivity

.1



No. of Support Vectors: 9 (17.6%)

Sinc= $\sin(\pi x) / (\pi x)$, poly kernel

Polynomial



Degree

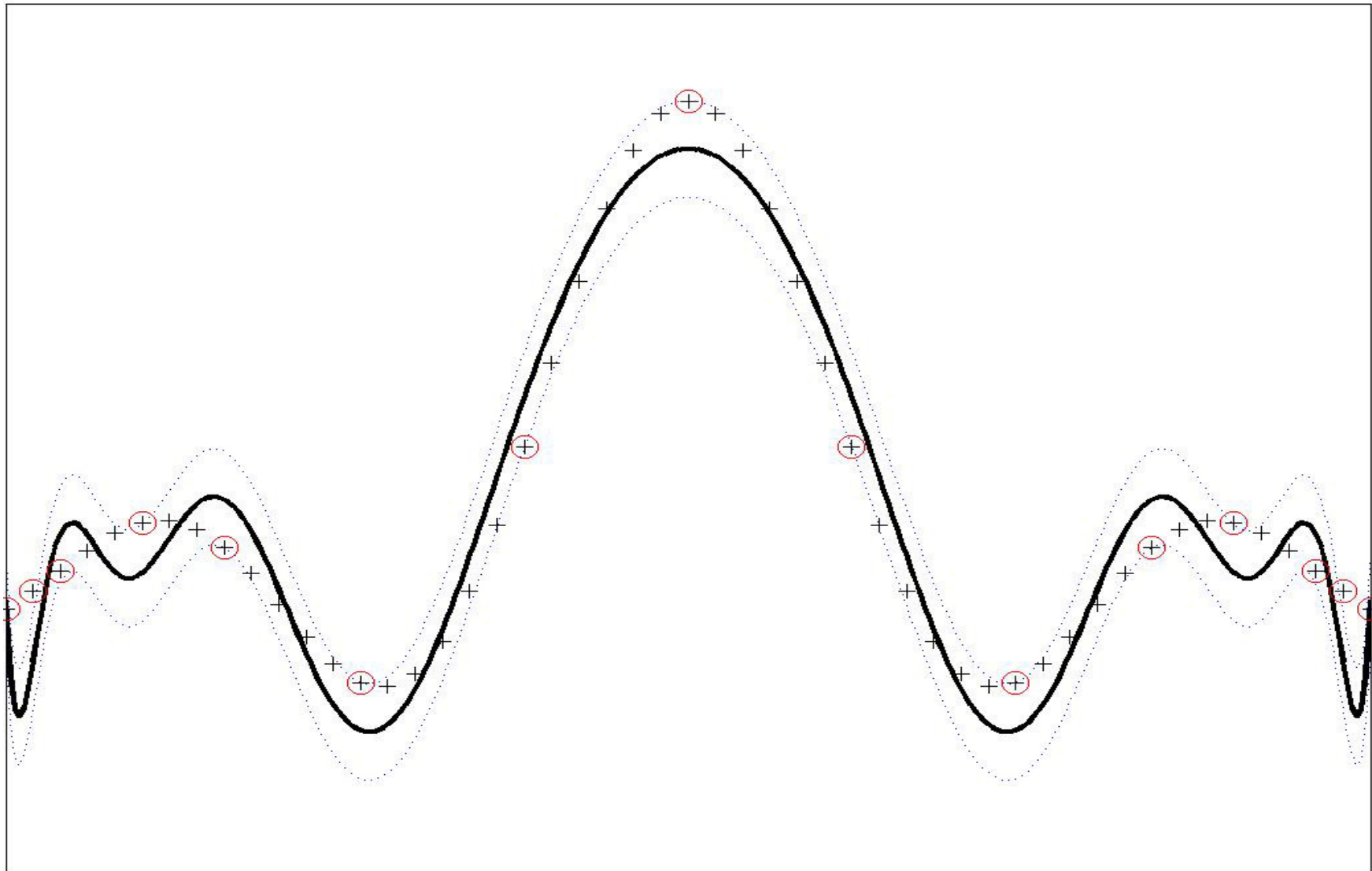
33

Bound

10

ϵ insensitivity

.1



No. of Support Vectors: 15 (29.4%)

Sinc= $\sin(\pi x) / (\pi x)$, poly kernel

Polynomial



Degree

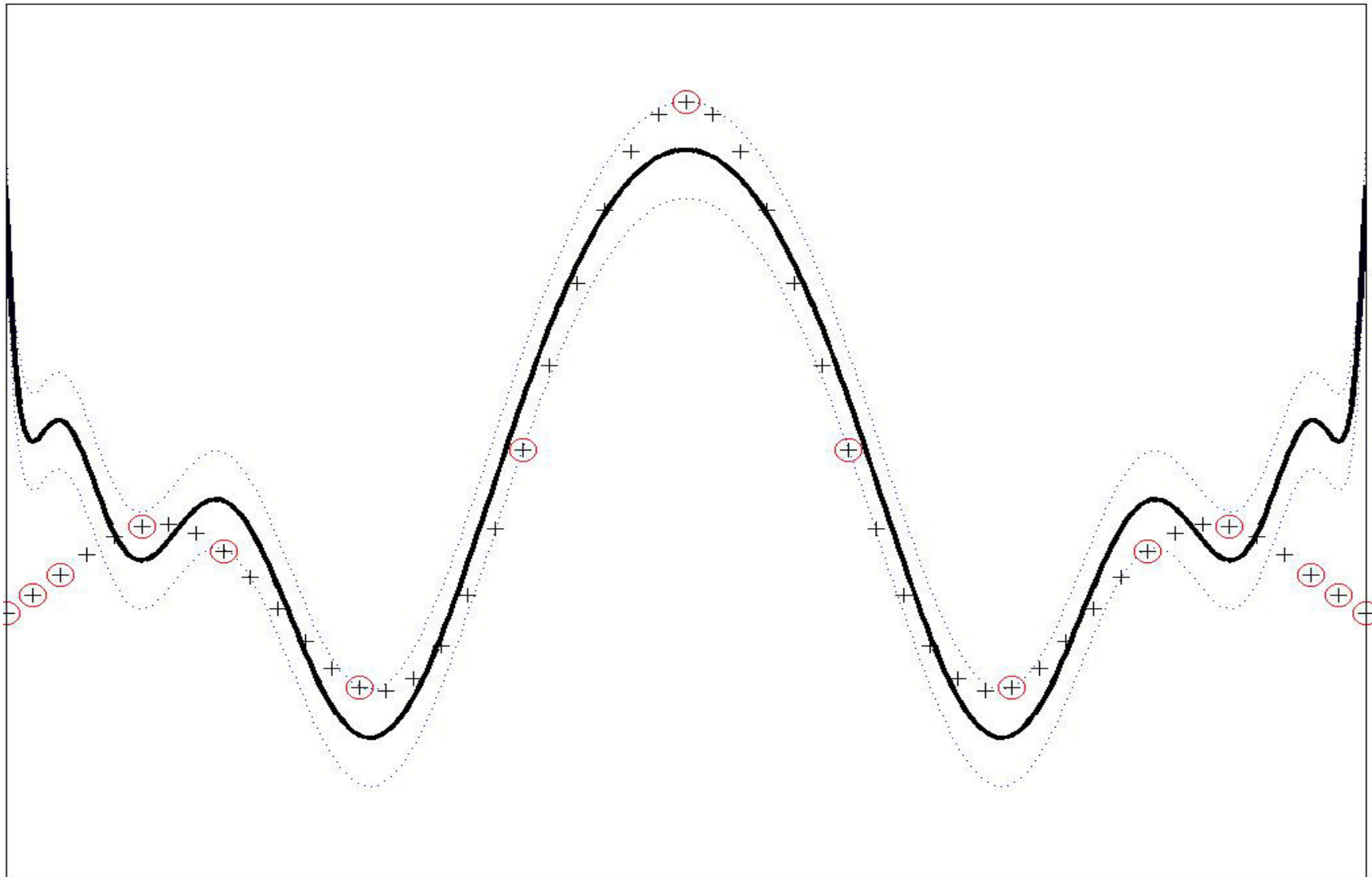
34

Bound

10

e insensitivity

.1



No. of Support Vectors: 15 (29.4%)

Sinc= $\sin(\pi x) / (\pi x)$, poly kernel

Polynomial



Degree

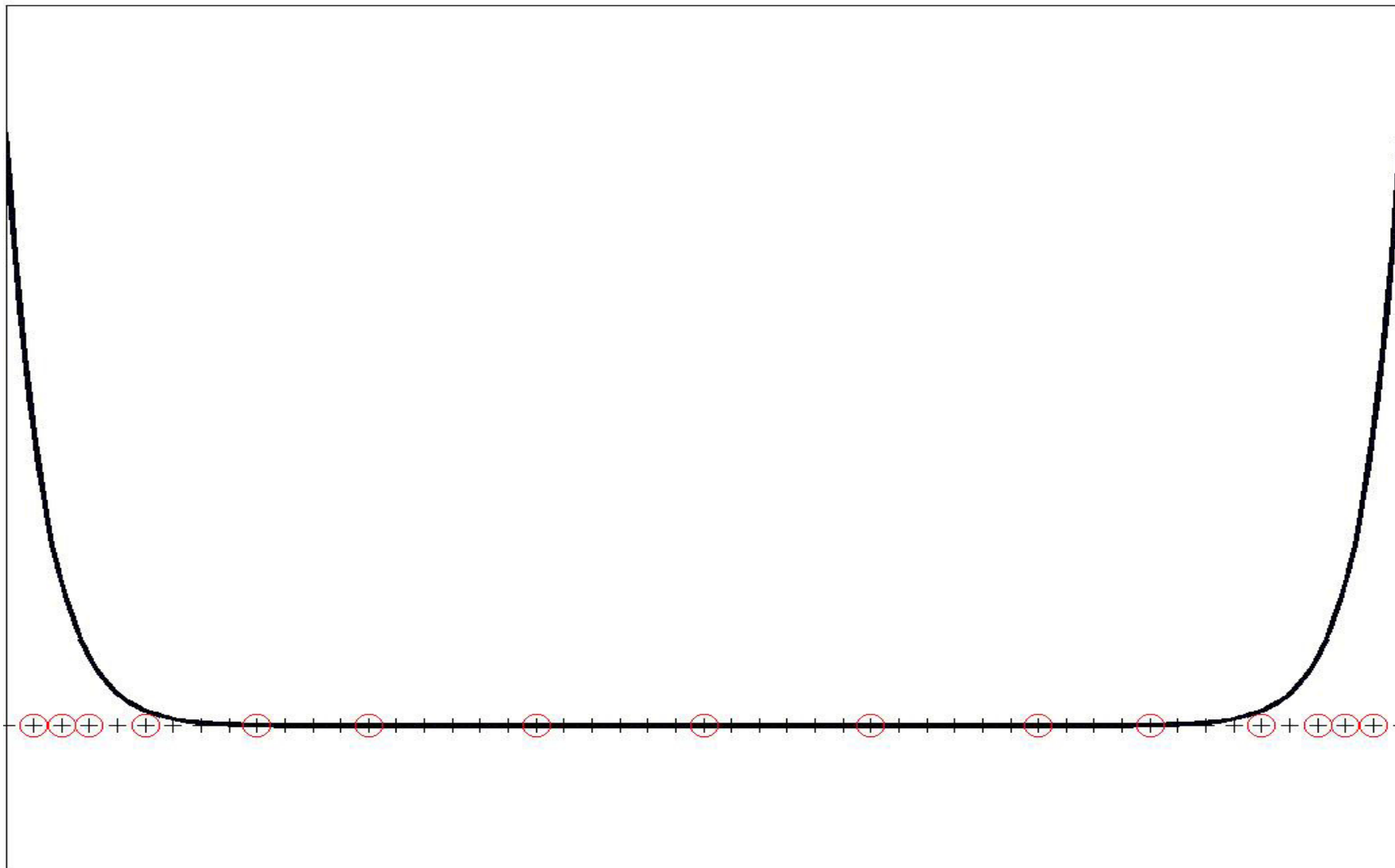
35

Bound

10

ϵ insensitivity

.1



No. of Support Vectors: 17 (33.3%)

What you need to know...

- Dual SVM formulation
 - How it's derived
- Common kernels
- Differences between SVMs and logistic regression

Thanks for your attention 😊