

Introduction to Machine Learning

CMU-10701

Support Vector Machines

Barnabás Póczos & Aarti Singh

2014 Spring

<http://barnabas-cmu-10701.appspot.com/>

Introduction to Machine Learning (10-701), Spring, 2014

Hi there!

**Let me send your Invitation Code to your Andrew email address.
Then please Register with this code.**

Invitation code :	<input type="text"/>
Nick name :	<input type="text"/>
Andrew id :	<input type="text"/>
First name :	<input type="text"/>
Last name :	<input type="text"/>
Email :	<input type="text"/>

Introduction to Machine Learning (10-701), Spring, 2014

Welcome test1! Login successful!

Nick name: test1

Andrew Id: bapoczos_v3

First name: test1

Last name: test1

Email: peergrading.test1@gmail.com

Edit

Log out

Your role:

Author of HW1

Author of HW2

Author of HW3

Author of HW4

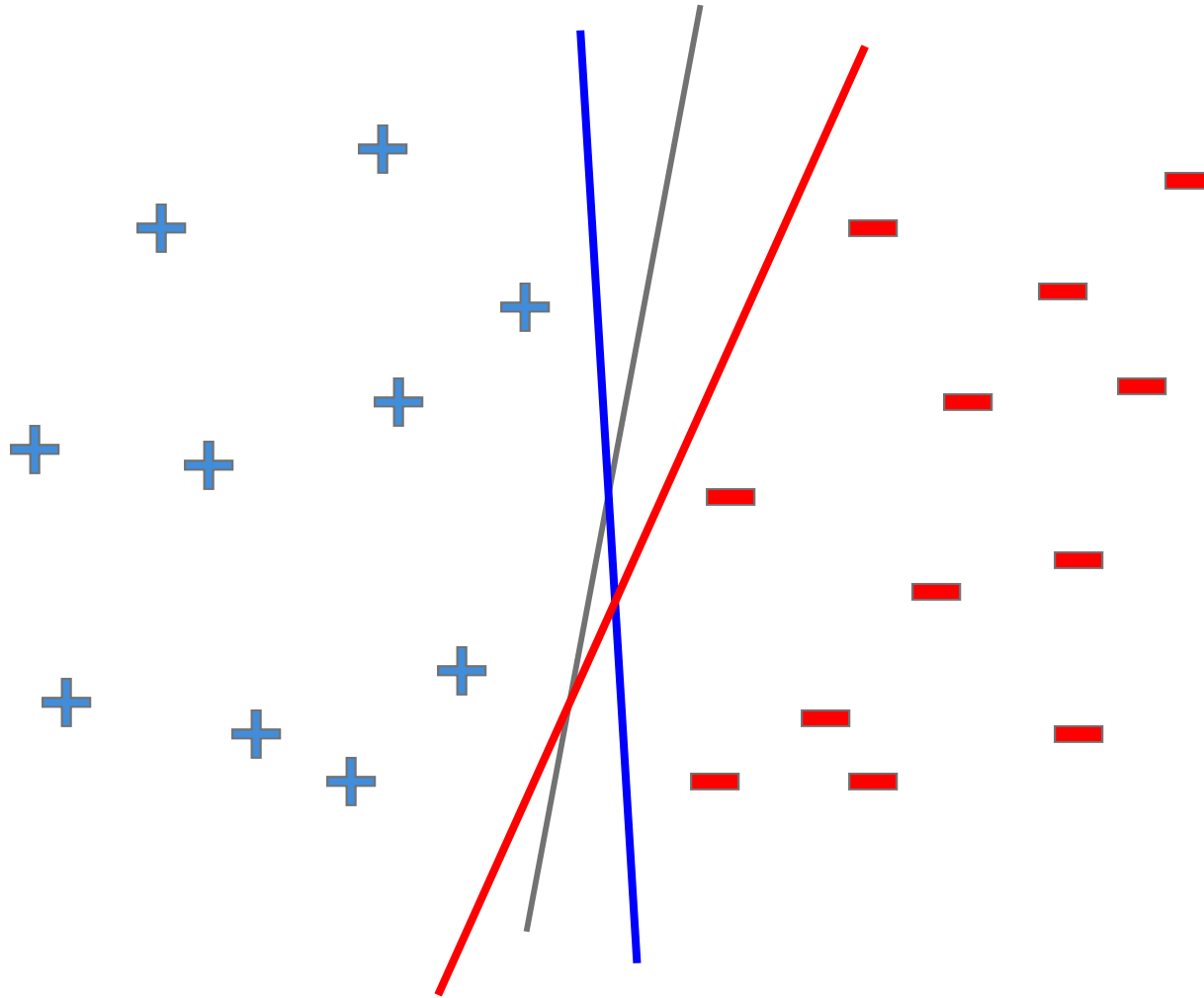
Reviewer of HW1

Reviewer of HW2

Reviewer of HW3

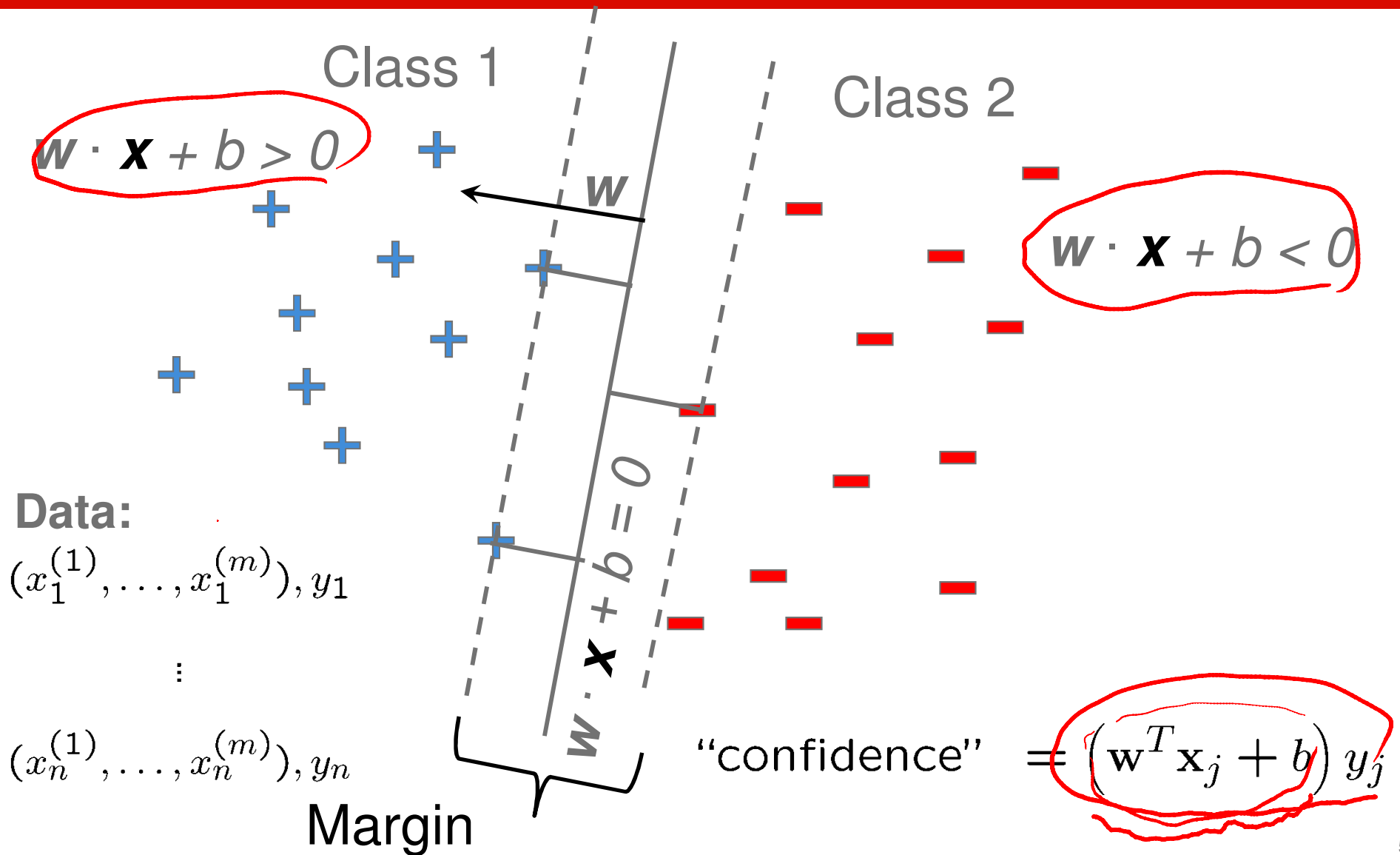
Reviewer of HW4

Linear classifiers which line is better?

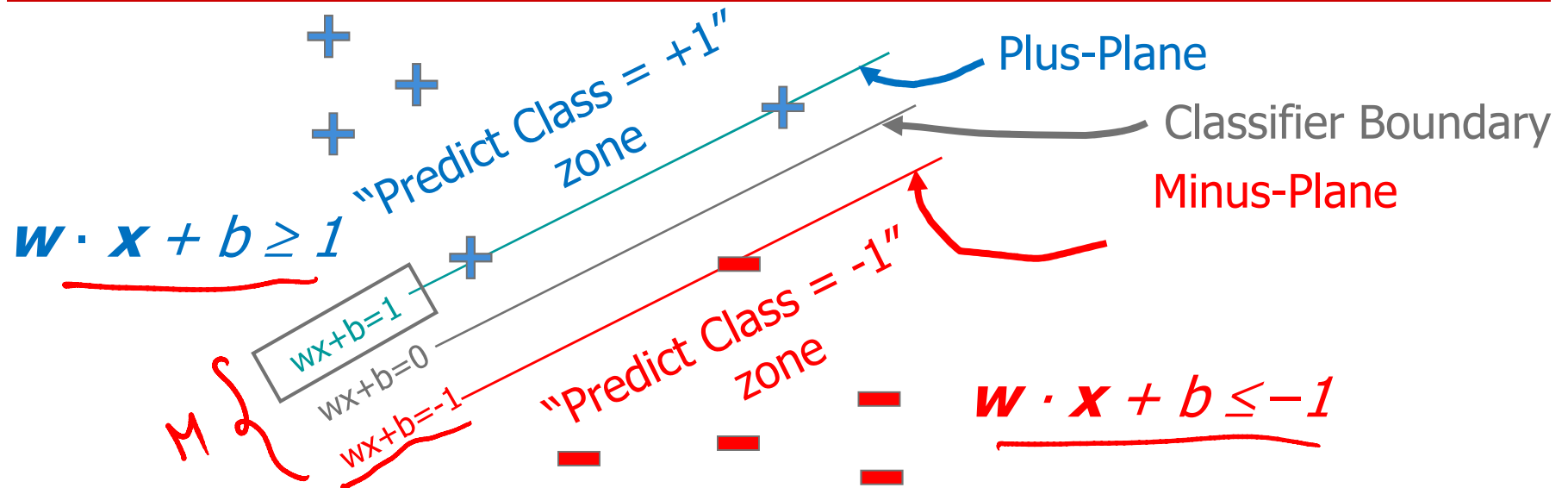


Which decision boundary is better?

Pick the one with the largest margin!



Scaling



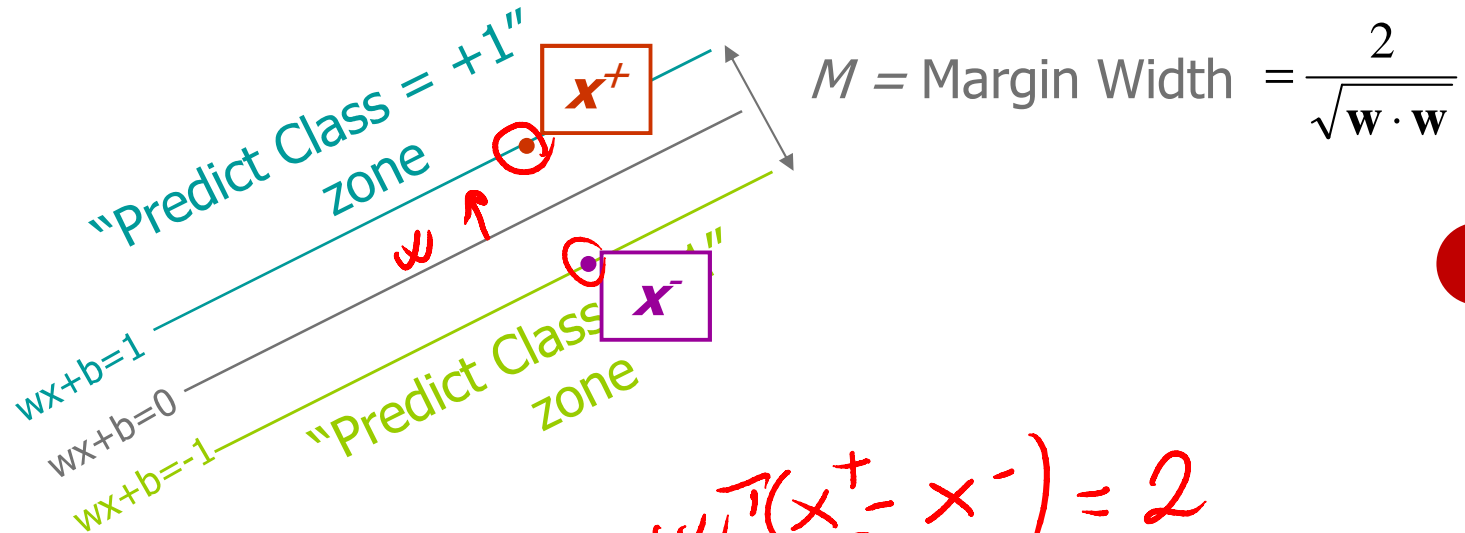
Classification rule:

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



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

- $w \cdot x^+ + b = +1$
- $w \cdot x^- + b = -1$
- $x^+ = x^- + \lambda w$
- $\underline{|x^+ - x^-|} = \underline{M}$

$$w^T(x^+ - x^-) = 2$$

$$\lambda w^T w = 2$$

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

$$M = |x^+ - x^-| = |\lambda w| = \left| \frac{2w}{w^T w} \right| = \frac{2}{\sqrt{w^T w}}$$

Maximize $M \equiv$ minimize $w \cdot w$!

Observations

We can assume $b=0$

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

This is the same as $y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1, \forall i = 1, \dots, n$

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

$$C \omega = d$$

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

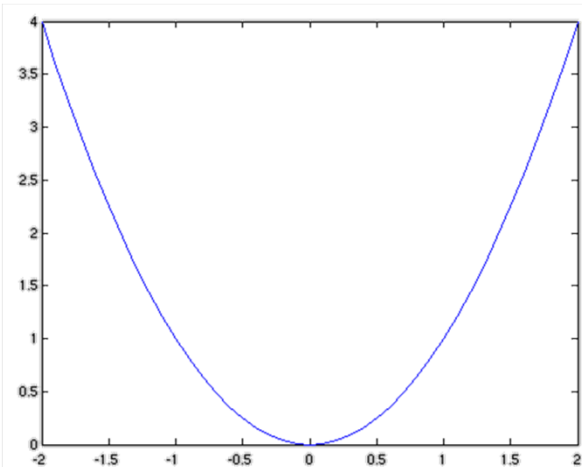
Efficient Algorithms exist for QP.

They often solve the dual problem instead of the primal.

Constrained Optimization

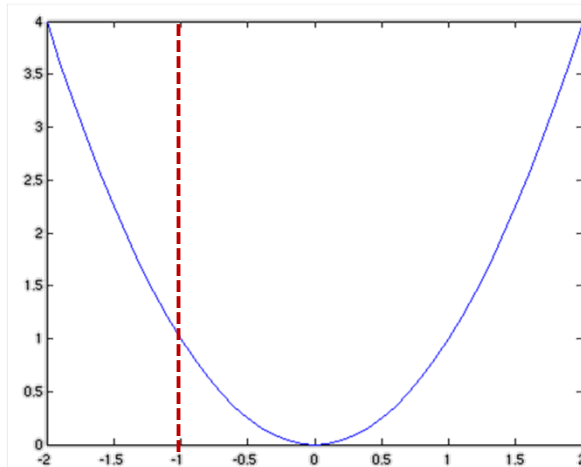
$$\begin{aligned} \min_x \quad & x^2 \\ \text{s.t.} \quad & x \geq b \end{aligned}$$

$$\min_x x^2$$



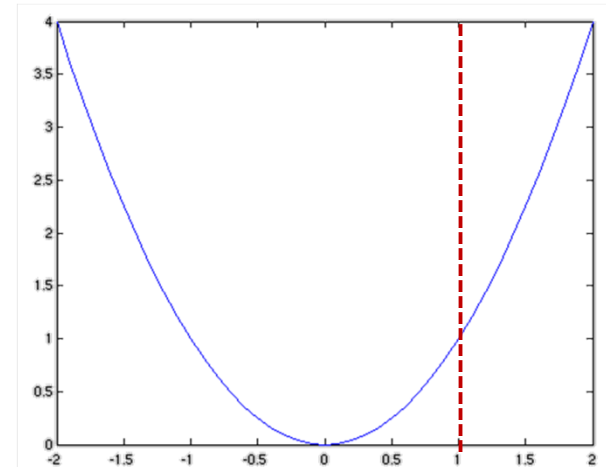
$$x^* = 0$$

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



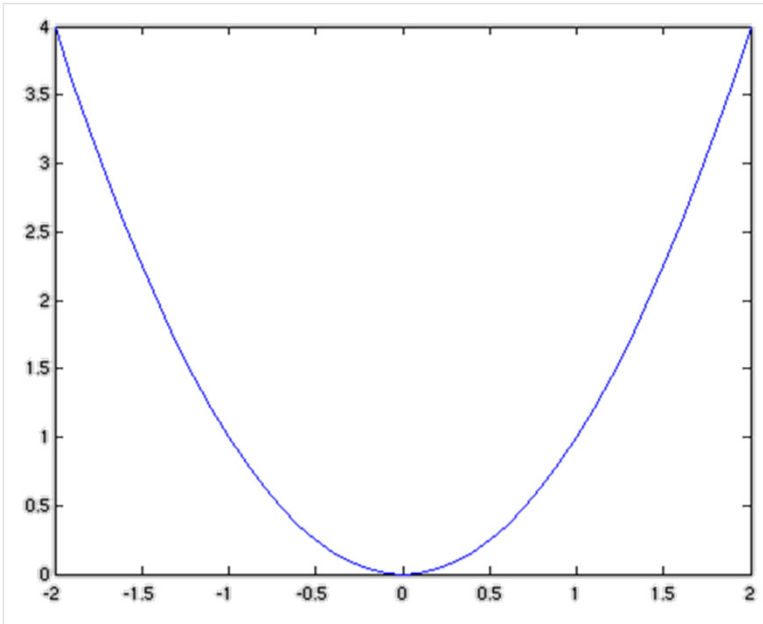
$$x^* = 0$$

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



$$x^* = 1$$

Lagrange Multiplier



$$\min_x x^2$$
$$\text{s.t. } \underline{x \geq b} \quad x - b > 0$$

Moving the constraint to objective function

Lagrangian:

$$L(x, \alpha) = \underline{x^2} - \underline{\alpha}(\underline{x - b})$$
$$\text{s.t. } \underline{\alpha \geq 0}$$

Solve:

$$\left[\begin{array}{l} \min_x \max_{\alpha} L(x, \alpha) \\ \text{s.t. } \alpha \geq 0 \end{array} \right]$$

x^* , α^*

Constraint is active when $\alpha > 0$

Lagrange Multiplier – Dual Variables

Solving:

$$\begin{array}{l} \min_x \max_{\alpha} \underbrace{x^2 - \alpha(x - b)}_{L(x, \alpha)} \\ \text{s.t. } \alpha \geq 0 \end{array}$$

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

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

$$\frac{\partial L(x, \alpha)}{\partial x} = 2x - \alpha = 0$$

$$x^* = \frac{\alpha}{2}$$

$$L(x^*, \alpha) = \frac{\alpha^2}{4} - \alpha \left(\frac{\alpha}{2} - b \right)$$

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

$$x^* = \max(b, 0) \quad \alpha^* = \max(2b, 0)$$

When $\alpha > 0$, constraint is tight

From Primal to Dual

Primal problem:

$$\left[\begin{array}{l} \hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^m} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{subject to } y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1, \forall i = 1, \dots, n \end{array} \right.$$

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

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

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

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

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

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

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

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

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...

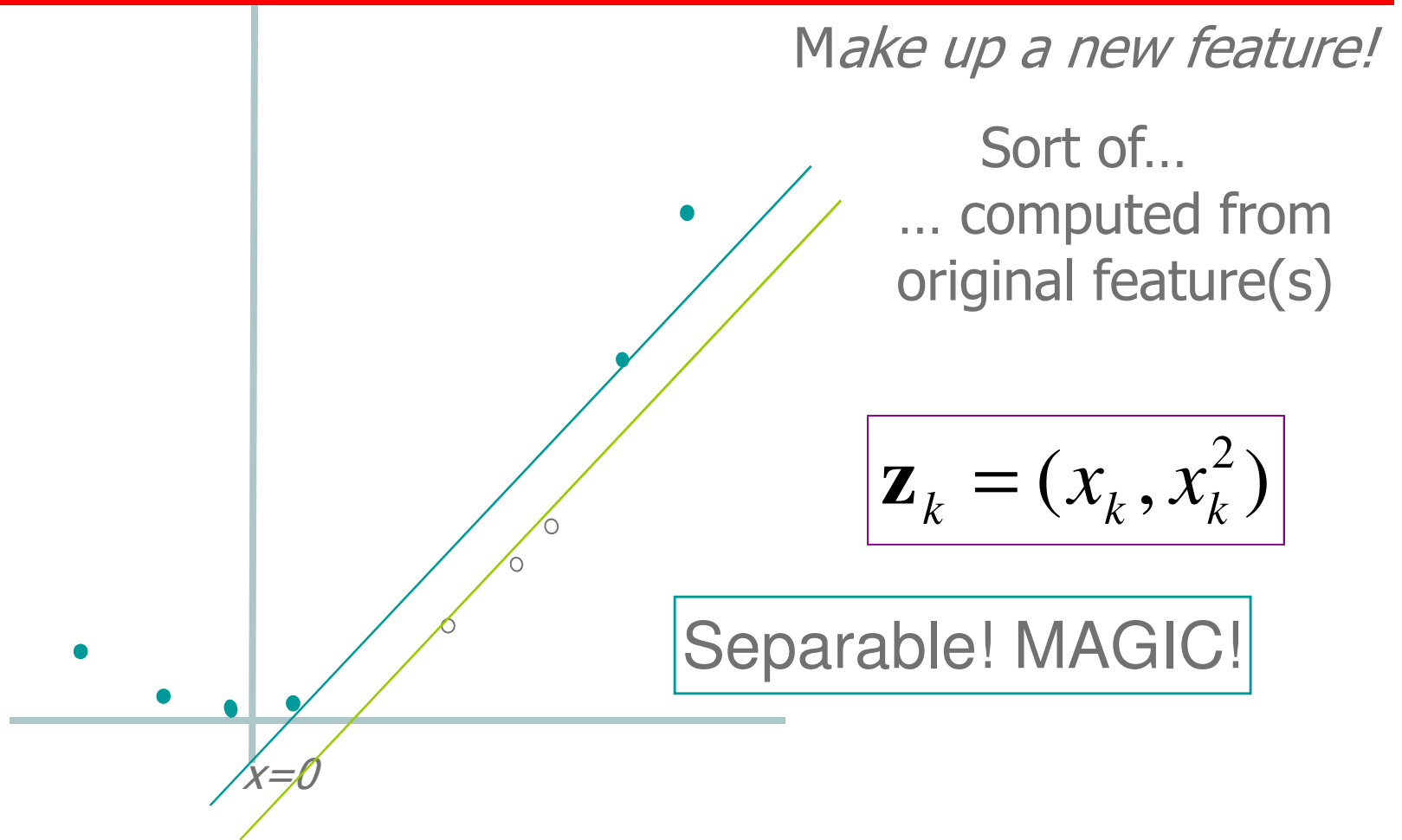
**What can we do if data is
not linearly separable???**

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 good 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...*

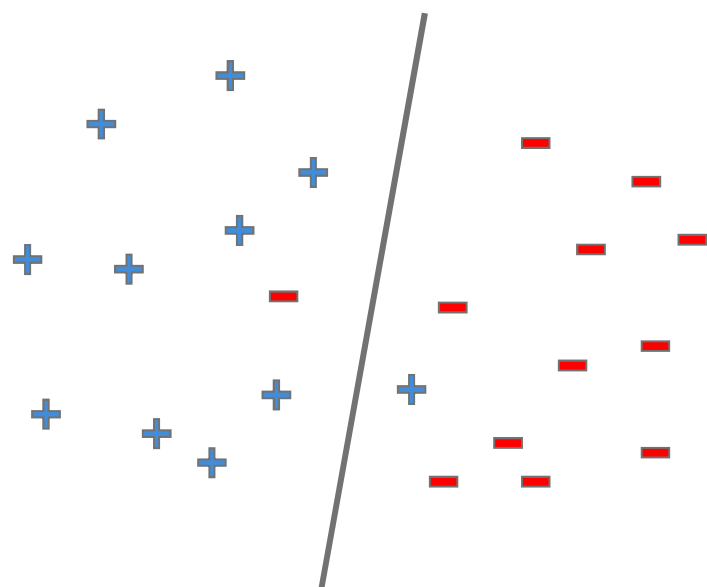
How to do feature mapping?

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 $\vec{x} = [\vec{x}_1, \vec{x}_2] \in \mathbb{R}^2$

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

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



∞

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
 - We will discuss this now

Hard SVM

The Hard SVM problem can be rewritten:

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

subject to $y_i \langle \mathbf{x}_i, \mathbf{w} \rangle > 0, \forall i = 1, \dots, n$



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

where

$$l_{0-\infty}(a, b) \doteq \begin{cases} \infty : ab < 0 & \text{Misclassification} \\ 0 : ab > 0 & \text{Correct classification} \end{cases}$$

From Hard to Soft constraints

Instead of using hard constraints (points are linearly separable)

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

We can try to solve the soft version of it:

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

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

(Handwritten red annotations: $x^T w$ above the first case, and a bracket under the first case)

Problems with l_{0-1} loss

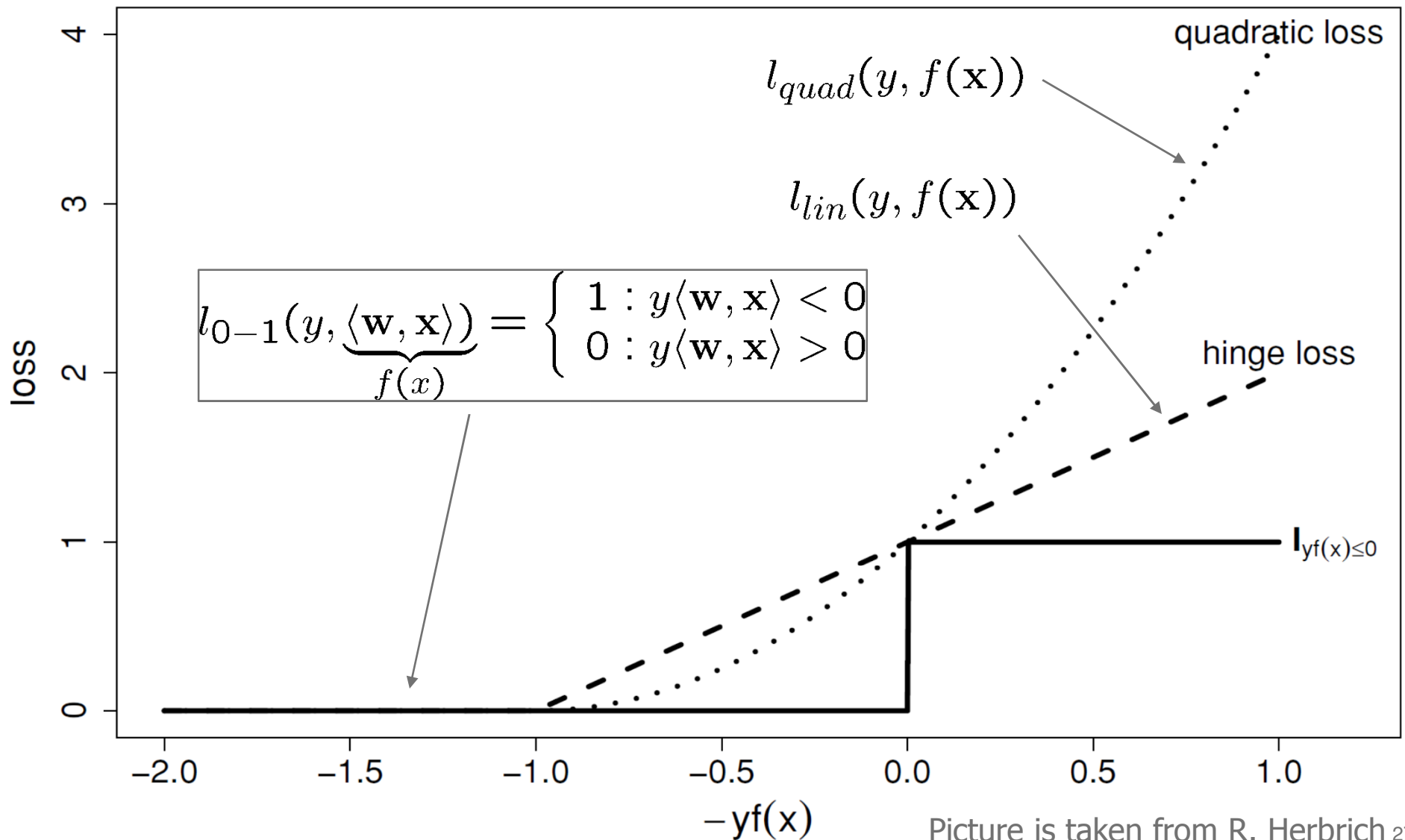
$$\hat{\mathbf{w}}_{soft} = \arg \min_{\mathbf{w} \in \mathbb{R}^m} \sum_{i=1}^n l_{0-1}(\langle \mathbf{x}_i, \mathbf{w} \rangle, y_i) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

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

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

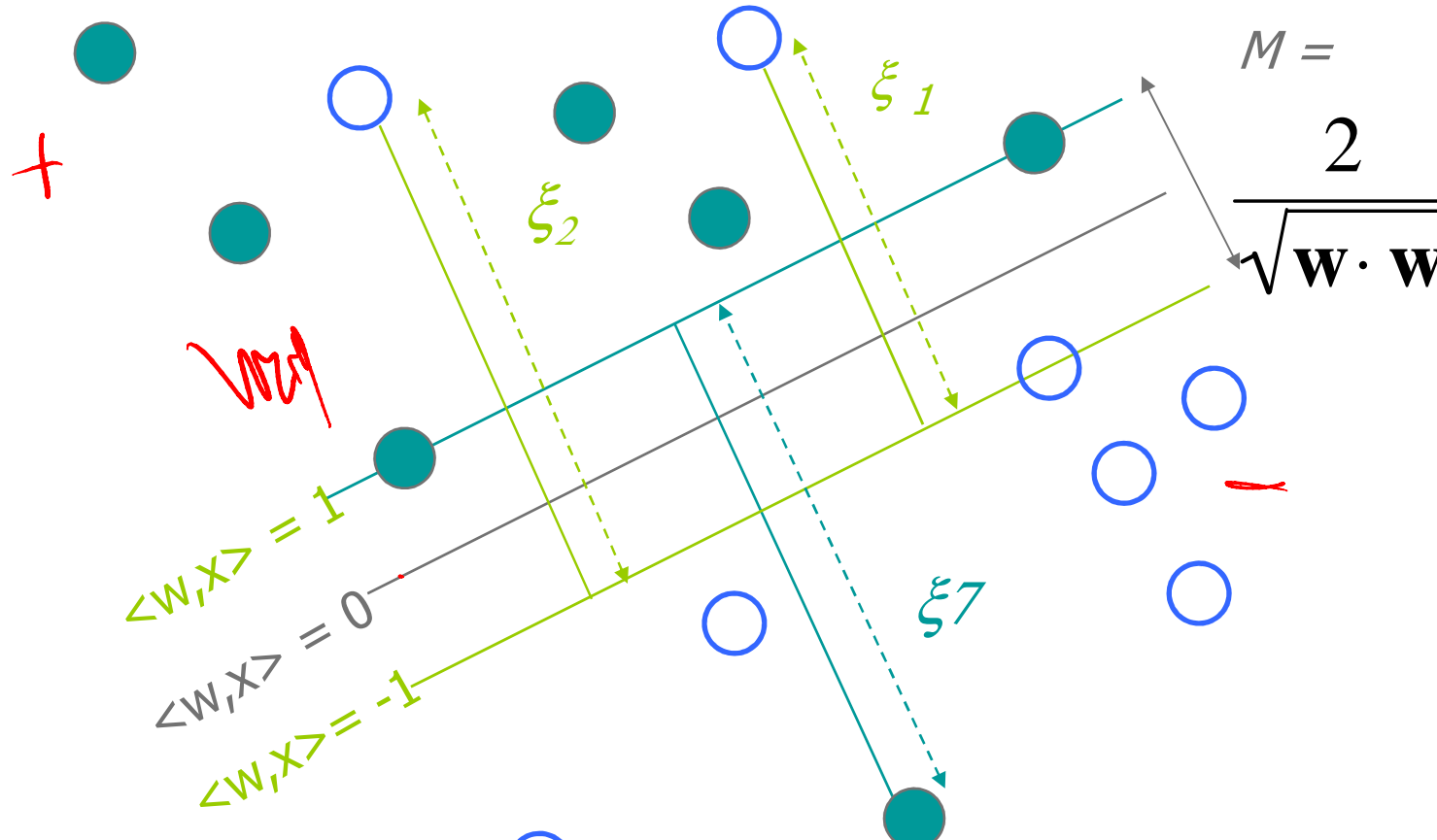
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^m} \sum_{i=1}^n \underbrace{l_{lin}(\langle \mathbf{x}_i, \mathbf{w} \rangle, y_i)}_{\xi_i \geq 0} + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

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 hinge loss upper bounds the 0-1 loss

Geometric interpretation: 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

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

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

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

Equivalently,

We can use this form, too.:

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

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

What is the dual form of primal soft SVM?

The Dual Soft SVM (using hinge loss)

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

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

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

α_i
 β_i

- $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n)^T \geq 0$ Lagrange multipliers
- $\boldsymbol{\beta} = (\beta_1, \dots, \beta_n)^T \geq 0$ Lagrange multipliers

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

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)

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

MIN MAX
 $\mathbf{w}, \boldsymbol{\xi}$ $\boldsymbol{\alpha}, \boldsymbol{\beta}$

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

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

$$\Rightarrow 0 \leq \alpha \leq C$$

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

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

Handwritten:
 $\mathbf{Y} = \begin{pmatrix} y_1 & \dots & 0 \\ \vdots & & \vdots \\ \alpha & \dots & y_n \end{pmatrix}$
 $\mathbf{G}_{ij} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$

The Dual Soft SVM (using hinge loss)

$$Y \doteq \text{diag}(y_1, \dots, y_n) \in \{-1, 1\}^n \quad x_i \mapsto \phi(x_i) \quad \langle \alpha, \epsilon \rangle = \alpha^T \epsilon$$

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

$$k_{ij} = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{\sigma^2}\right)$$

$$\langle \phi(x_i), \phi(x_j) \rangle = k(x_i, x_j)$$

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

$$\text{subject to } 0 \leq \alpha_i \leq C$$

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

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 $0 \leq \alpha_i \leq C \quad \forall i$

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

$\phi(\mathbf{x}_i, \mathbf{x})$
 $\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle$

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

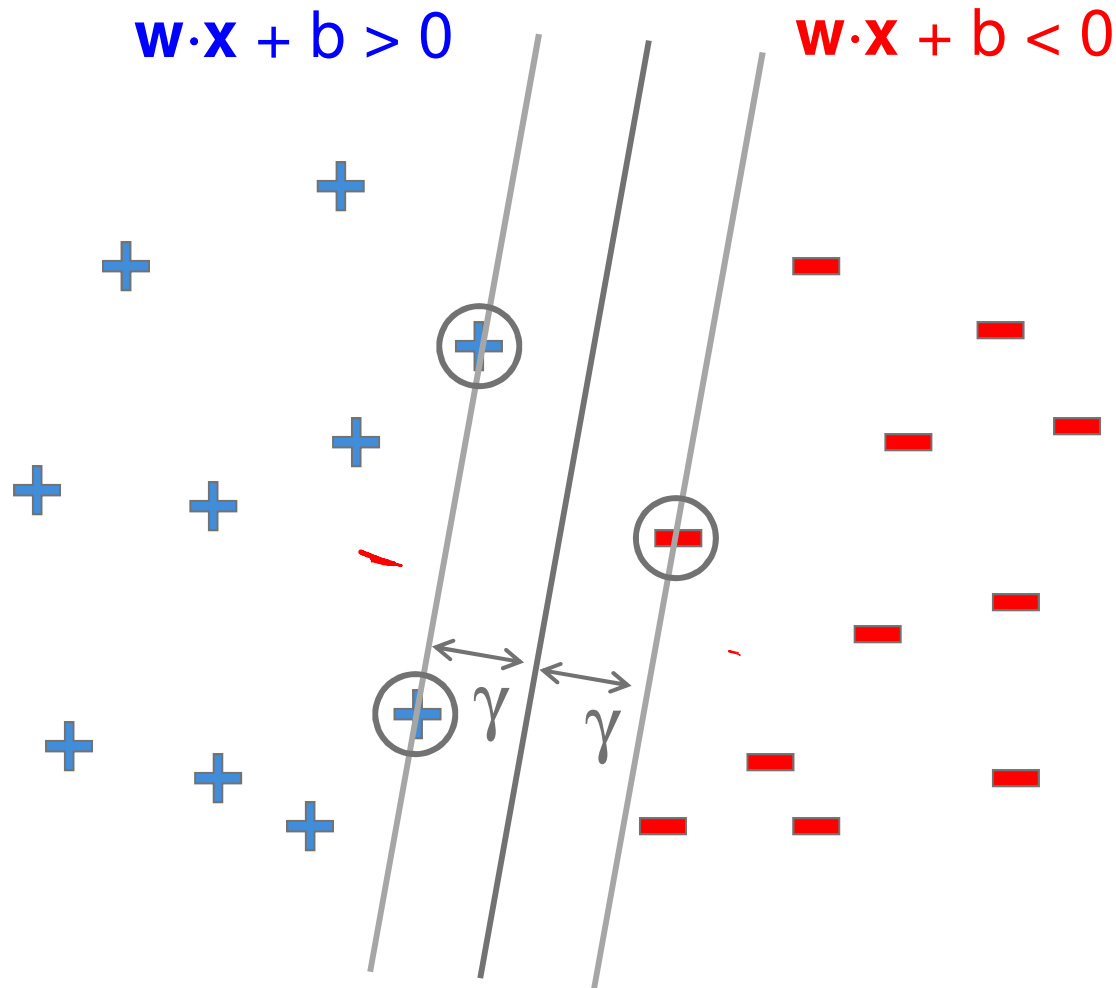
$$d_i > 0 \Rightarrow y_i \langle \mathbf{x}_i, \mathbf{w} \rangle = 1$$
$$\langle \mathbf{x}_i, \mathbf{w} \rangle = \begin{matrix} +1 \\ -1 \end{matrix}$$

EITHER $d_i = 0$ OR $(y_i \langle \mathbf{x}_i, \mathbf{w} \rangle - 1) \neq 0$

$(y_i \langle \mathbf{x}_i, \mathbf{w} \rangle - 1) = 0$ AND $d_i > 0$

\Downarrow
 \mathbf{x}_i IS ON THE MARGIN LINES
SUPPORT VECTORS

Why is it called Support Vector Machine?



Hard SVM:

Linear hyperplane defined by “support vectors”

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

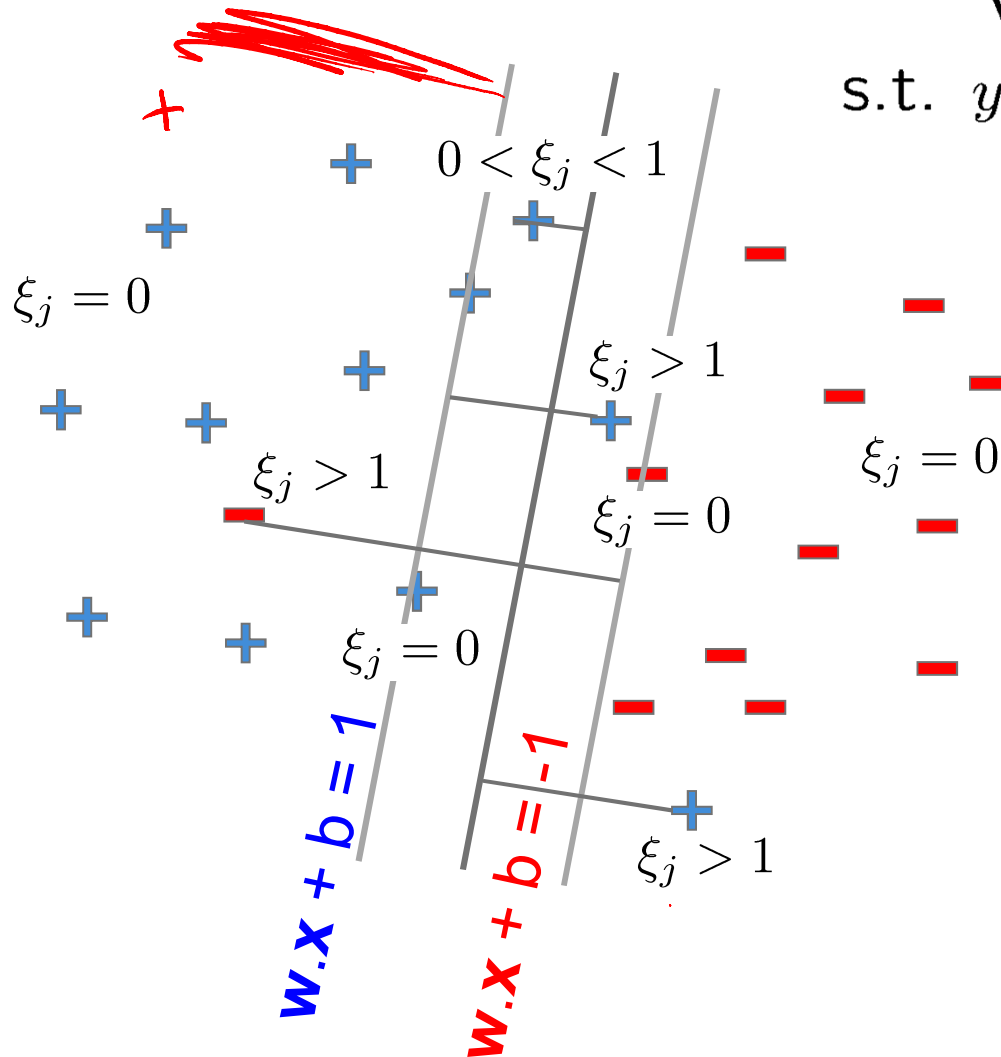
only need to store the support vectors to predict labels of new points

Support vectors in Soft SVM

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

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

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

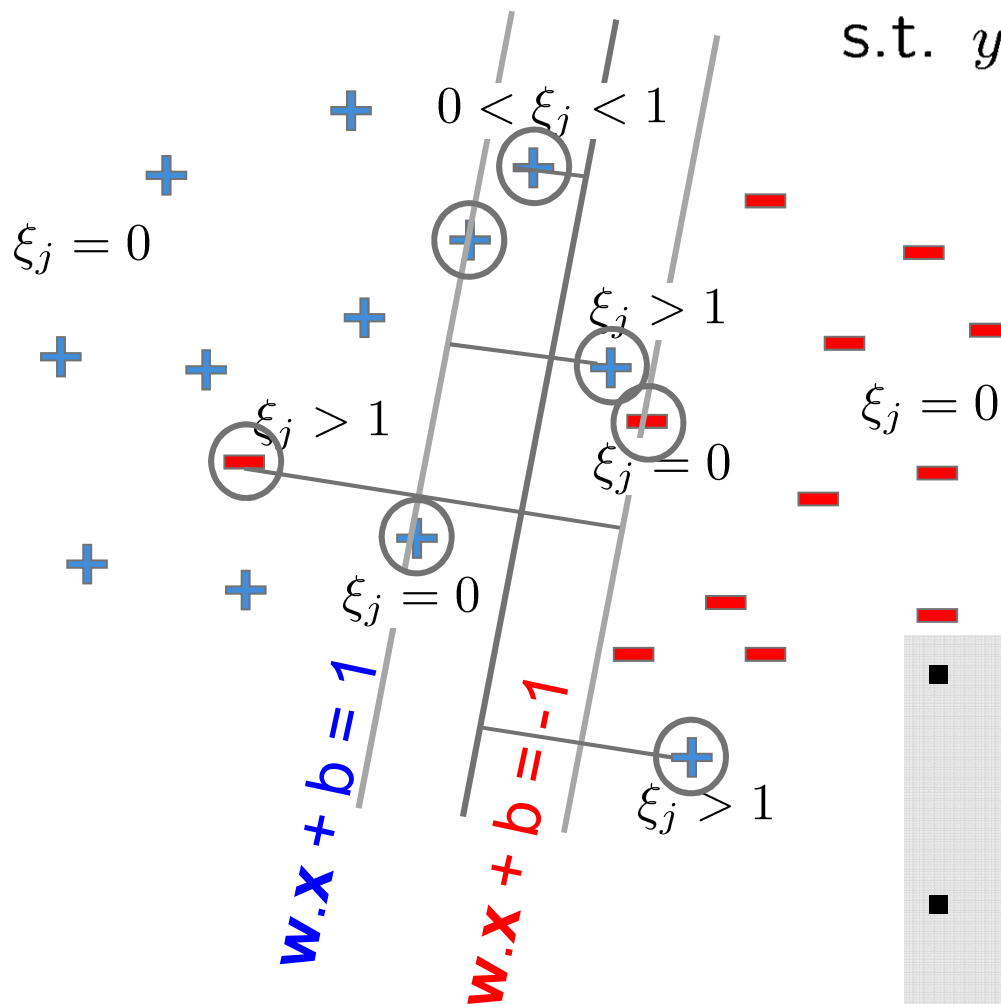


Support vectors in Soft SVM

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

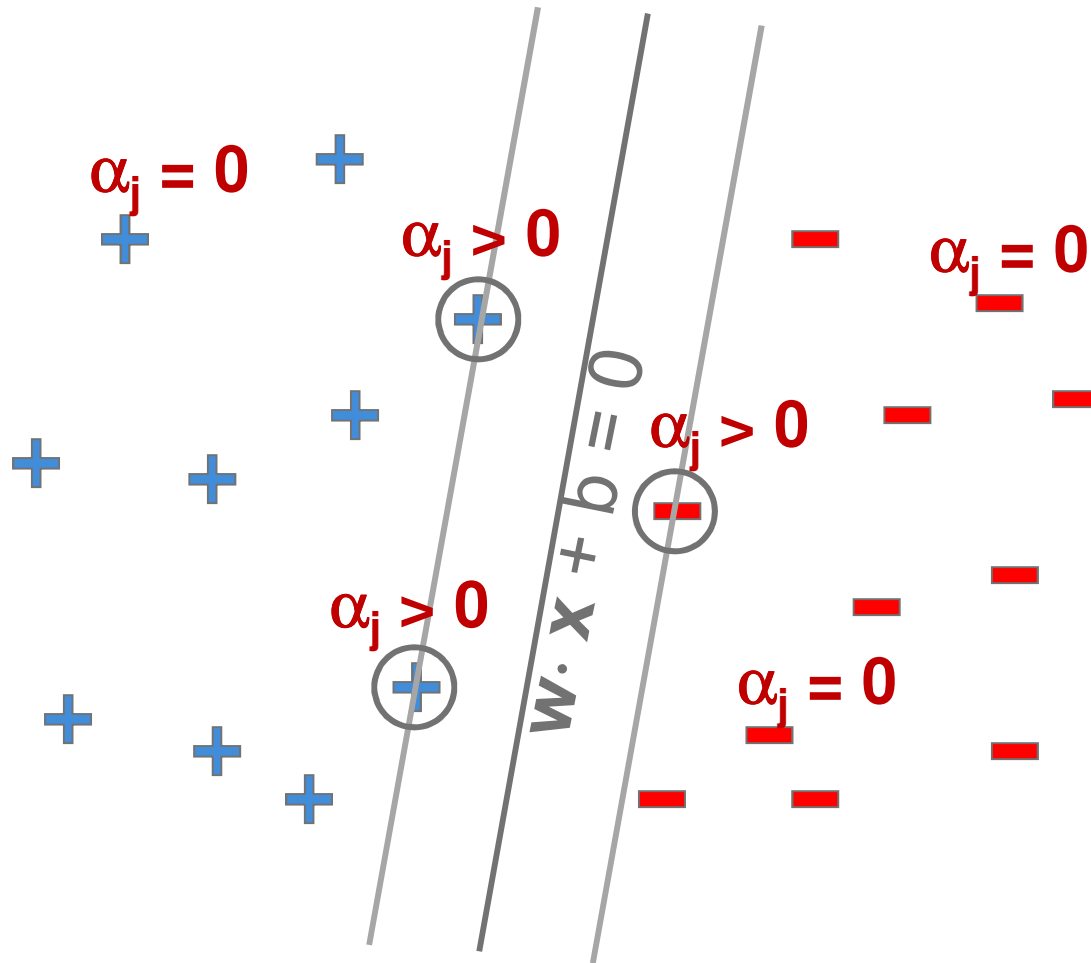
$$\text{s.t. } y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1 - \xi_i, \quad \forall i = 1, \dots, n$$

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



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

Dual SVM Interpretation: Sparsity

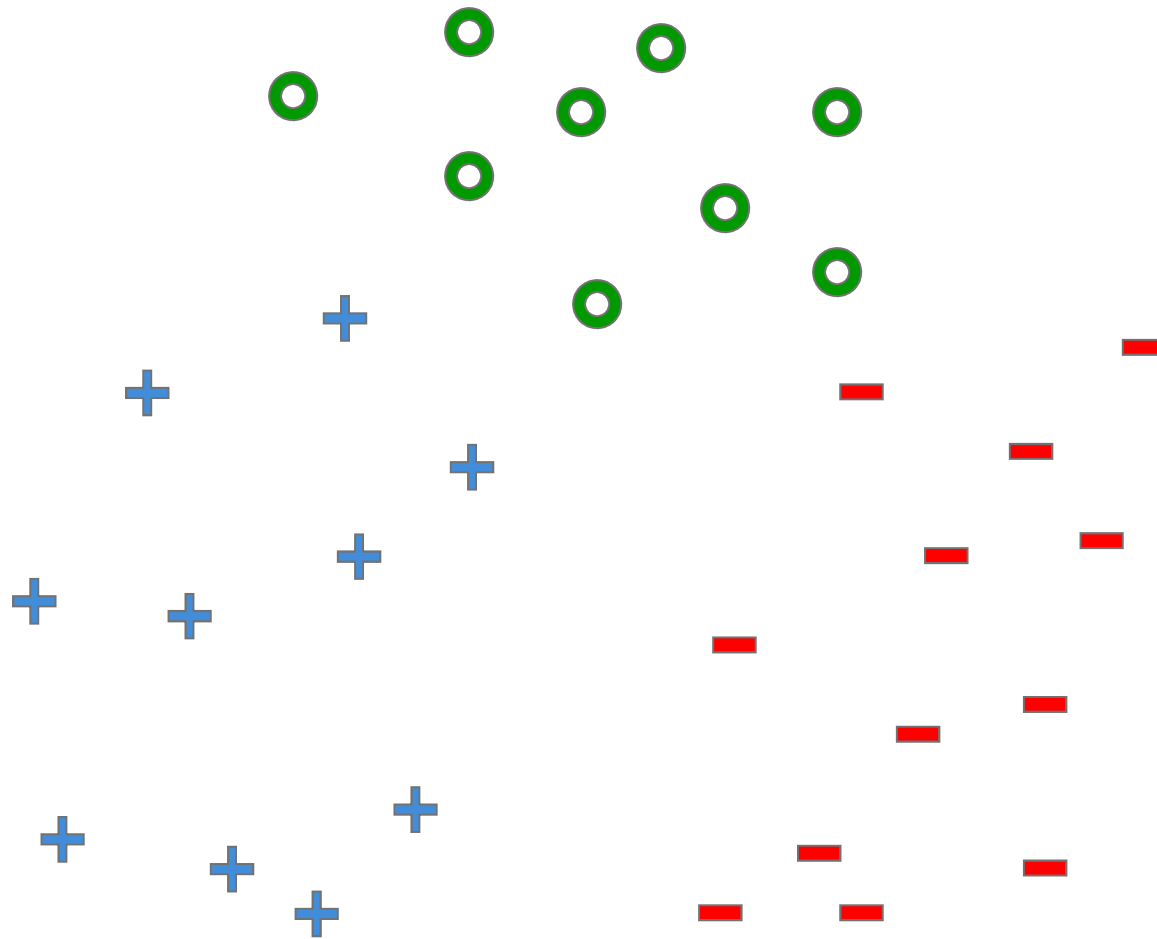


$$\mathbf{w} = \sum_j \alpha_j y_j \mathbf{x}_j$$

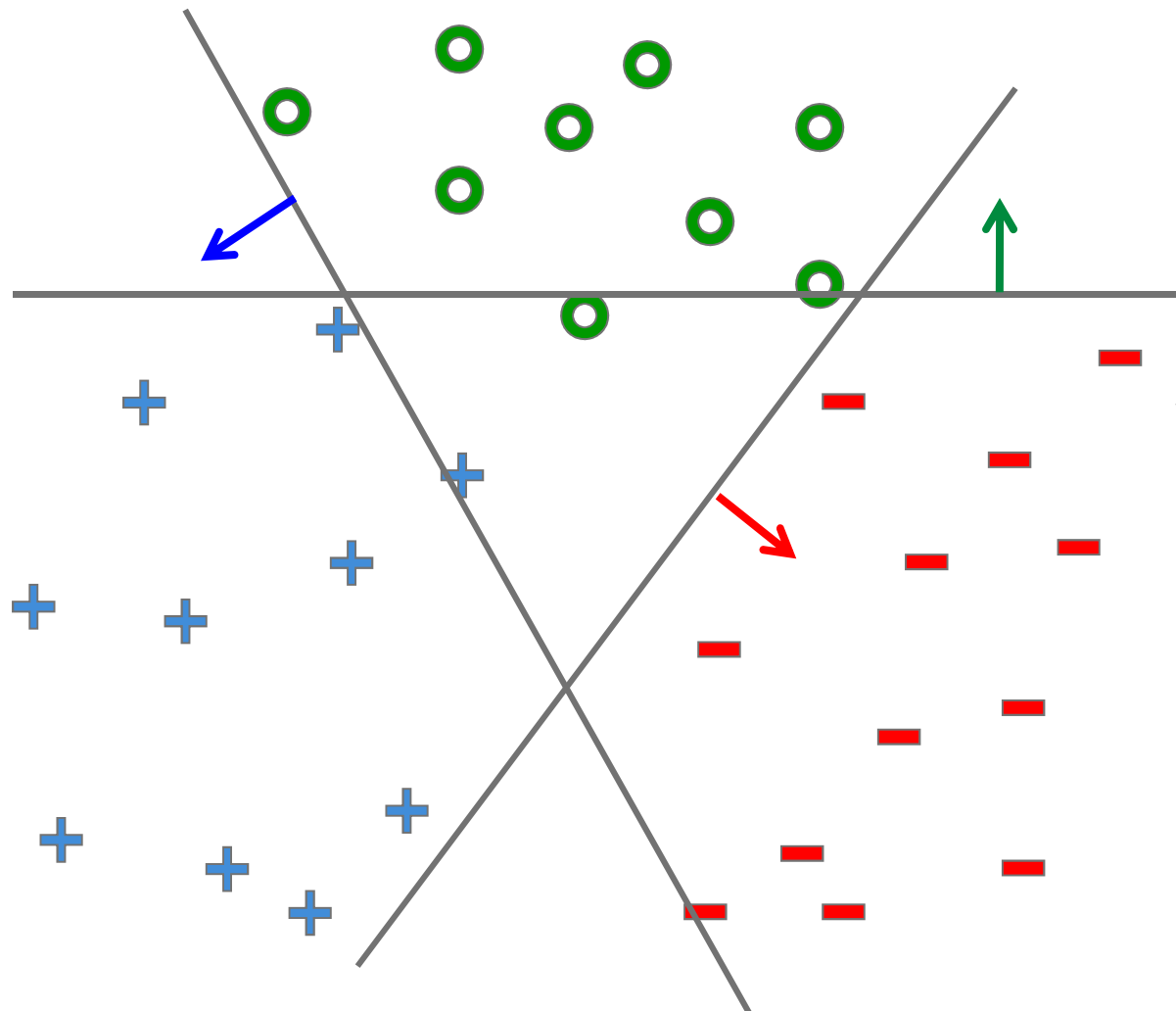
Only few α_j s can be non-zero : where constraint is tight

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

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

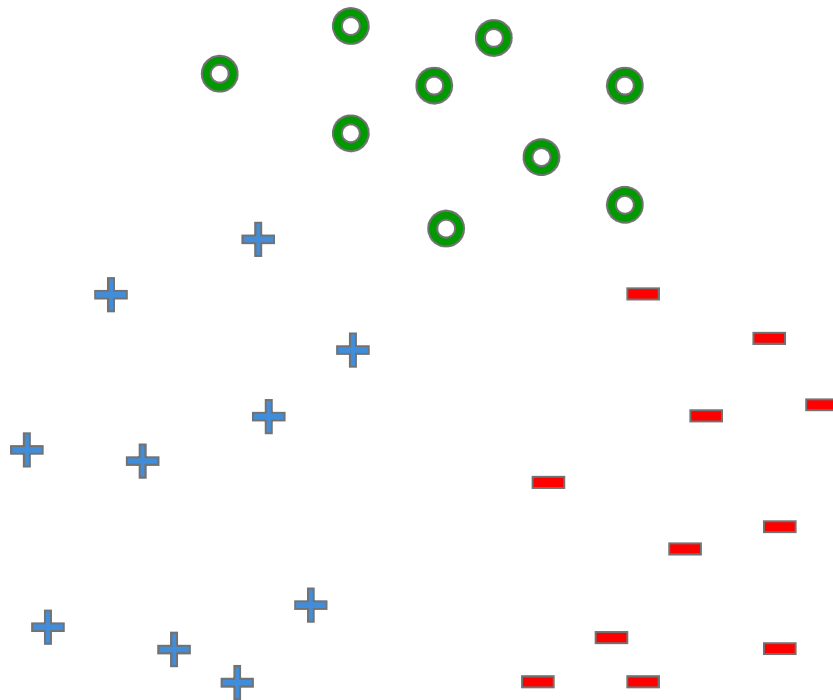
Simultaneously learn 3 sets of weights.

Constraints:

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

Margin - gap between
correct class and nearest
other class

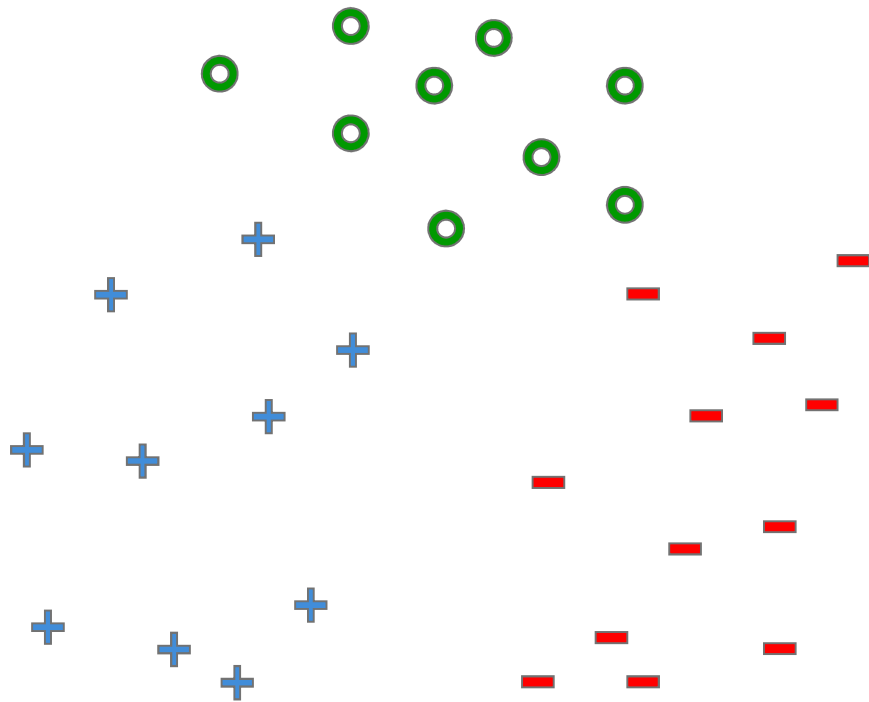
$$y = \arg \max_k \mathbf{w}^{(k)} \cdot \mathbf{x} + b^{(k)}$$



Learn 1 classifier: Multi-class SVM

Simultaneously learn 3 sets of weights

$$\begin{aligned} \text{minimize}_{\mathbf{w}, b} \quad & \sum_y \mathbf{w}^{(y)} \cdot \mathbf{w}^{(y)} + C \sum_j \sum_{y \neq y_j} \xi_j^{(y)} \\ \mathbf{w}^{(y_j)} \cdot \mathbf{x}_j + b^{(y_j)} \geq & \mathbf{w}^{(y)} \cdot \mathbf{x}_j + b^{(y)} + 1 - \xi_j^{(y)}, \quad \forall y \neq y_j, \quad \forall j \\ \xi_j^{(y)} \geq 0 \quad & , \quad \forall y \neq y_j, \quad \forall j \end{aligned}$$



$y = \arg \max_k \mathbf{w}^{(k)} \cdot \mathbf{x} + b^{(k)}$
Joint optimization: \mathbf{w}_k s
have the same scale.

What you need to know

- Maximizing margin
- Derivation of SVM formulation
- Slack variables and hinge loss
- Relationship between SVMs and logistic regression
 - 0/1 loss
 - Hinge loss
 - Log loss
- Tackling multiple class
 - One against All
 - Multiclass SVMs

SVM vs. Logistic Regression

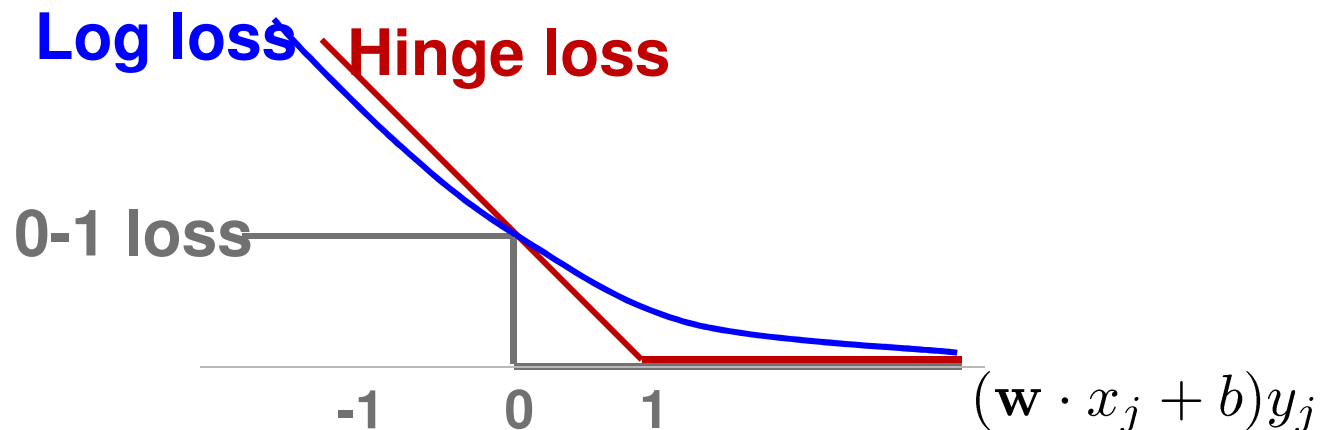
SVM : **Hinge loss**: $\text{loss}(f(\mathbf{x}_j), y_j) = (1 - (\mathbf{w} \cdot \mathbf{x}_j + b)y_j)_+$

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

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

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

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^m} \sum_{i=1}^n \text{loss}(\mathbf{x}_i \cdot \mathbf{w} + b, y_i)$$

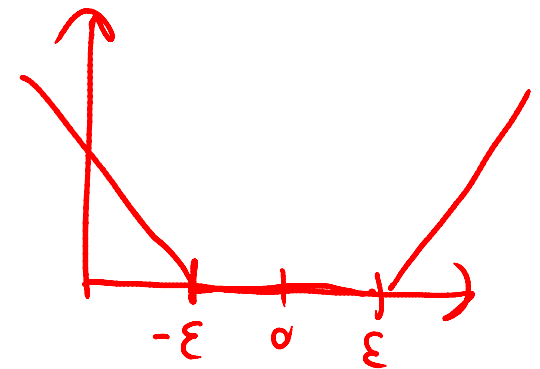
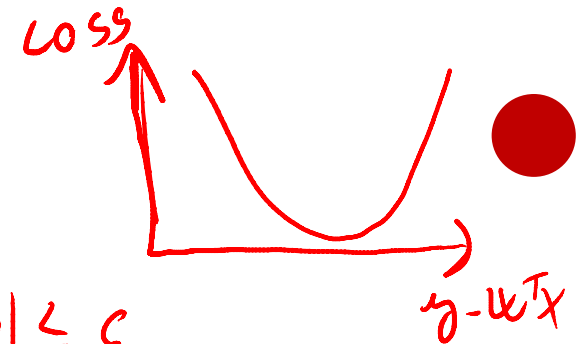


SVM for Regression

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

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

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



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”

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

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

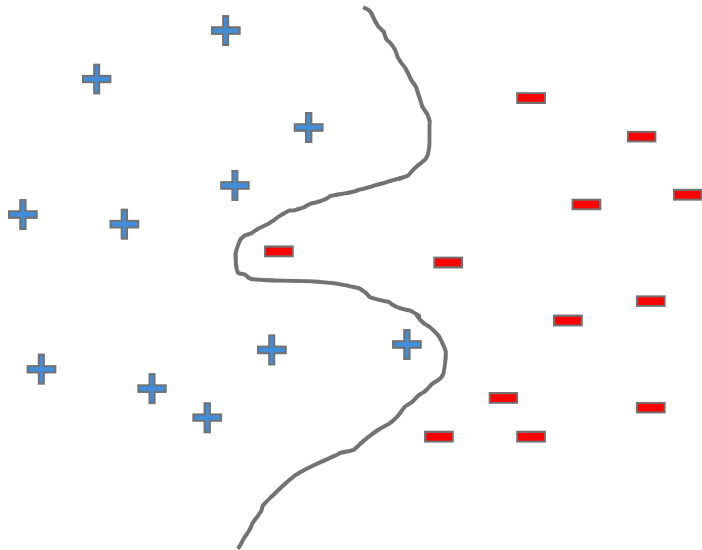
subject to $0 \leq \alpha_i \leq C$

$$\sum_i \alpha_i y_i = 0$$

So why solve the dual SVM?

- There are some quadratic programming algorithms that can solve the dual faster than the primal, specially in high dimensions $m \gg n$
- But, more importantly, the “**kernel trick**”!!!

What if data is not linearly separable?



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

For example polynomials

$$\Phi(\mathbf{x}) = (x_1^3, x_2^3, x_3^3, x_1^2x_2x_3, \dots)$$

Dot Product of Polynomials

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

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

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

$d=2$

$$\begin{aligned} \Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) &= \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix} \cdot \begin{bmatrix} z_1^2 \\ \sqrt{2}z_1z_2 \\ z_2^2 \end{bmatrix} = x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 x_2 z_1 z_2 \\ &= (x_1 z_1 + x_2 z_2)^2 \\ &= (\mathbf{x} \cdot \mathbf{z})^2 \end{aligned}$$

Dot Product of Polynomials

$$d=3$$

$$\begin{pmatrix} x_1^3 \\ \sqrt{3}x_1^2x_2 \\ \sqrt{3}x_1x_2^2 \\ x_2^3 \end{pmatrix} \cdot \begin{pmatrix} z_1^3 \\ \sqrt{3}z_1^2z_2 \\ \sqrt{3}z_1z_2^2 \\ z_2^3 \end{pmatrix}$$

$$\begin{aligned} & x_1^3z_1^3 + 3x_1^2x_2z_1^2z_2 + 3x_1z_1x_2^2z_2^2 + x_2^3z_2^3 \\ & = (x_1z_1 + z_2z_2)^3 \end{aligned}$$

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

Higher Order Polynomials

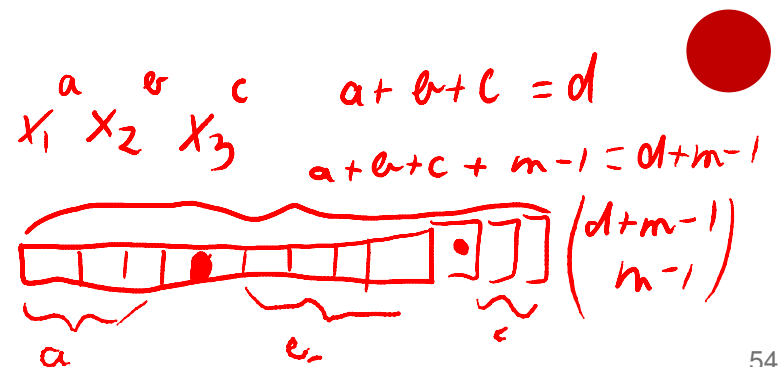
Feature space becomes really large very quickly!

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



Dual formulation only depends on dot-products, not on w !

$$\text{maximize}_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \underbrace{\mathbf{x}_i \cdot \mathbf{x}_j}$$
$$C \geq \alpha_i \geq 0$$



$$\text{maximize}_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \underbrace{K(\mathbf{x}_i, \mathbf{x}_j)}$$
$$K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$$
$$C \geq \alpha_i \geq 0$$

$\Phi(\mathbf{x})$ – High-dimensional feature space, but never need it explicitly as long as we can compute the dot product fast using some Kernel K

Finally: The Kernel Trick!

$$\text{maximize}_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

$$K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$$

$$\sum_i \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0$$

- Never represent features explicitly
 - Compute dot products in closed form
- Constant-time high-dimensional dot-products for many classes of features

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

$$b = y_k - \mathbf{w} \cdot \Phi(\mathbf{x}_k)$$

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

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 (polynomials of all orders – recall series expansion)
$$K(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u} - \mathbf{v}\|^2}{2\sigma^2}\right)$$
- Sigmoid
$$K(\mathbf{u}, \mathbf{v}) = \tanh(\eta\mathbf{u} \cdot \mathbf{v} + \nu)$$

Which functions can be used as kernels???

...and why are they called kernels???

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

What about classification time?

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

$$b = y_k - \mathbf{w} \cdot \Phi(\mathbf{x}_k)$$

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

- For a new input \mathbf{x} , if we need to represent $\Phi(\mathbf{x})$, we are in trouble!
- Recall classifier: $\text{sign}(\mathbf{w} \cdot \Phi(\mathbf{x}) + b)$
- Using kernels we are cool!

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

Kernels in Logistic Regression

$$P(Y = 1 | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{-(\mathbf{w} \cdot \Phi(\mathbf{x}) + b)}}$$

- Define weights in terms of features:

$$\mathbf{w} = \sum_i \alpha_i \Phi(\mathbf{x}_i)$$

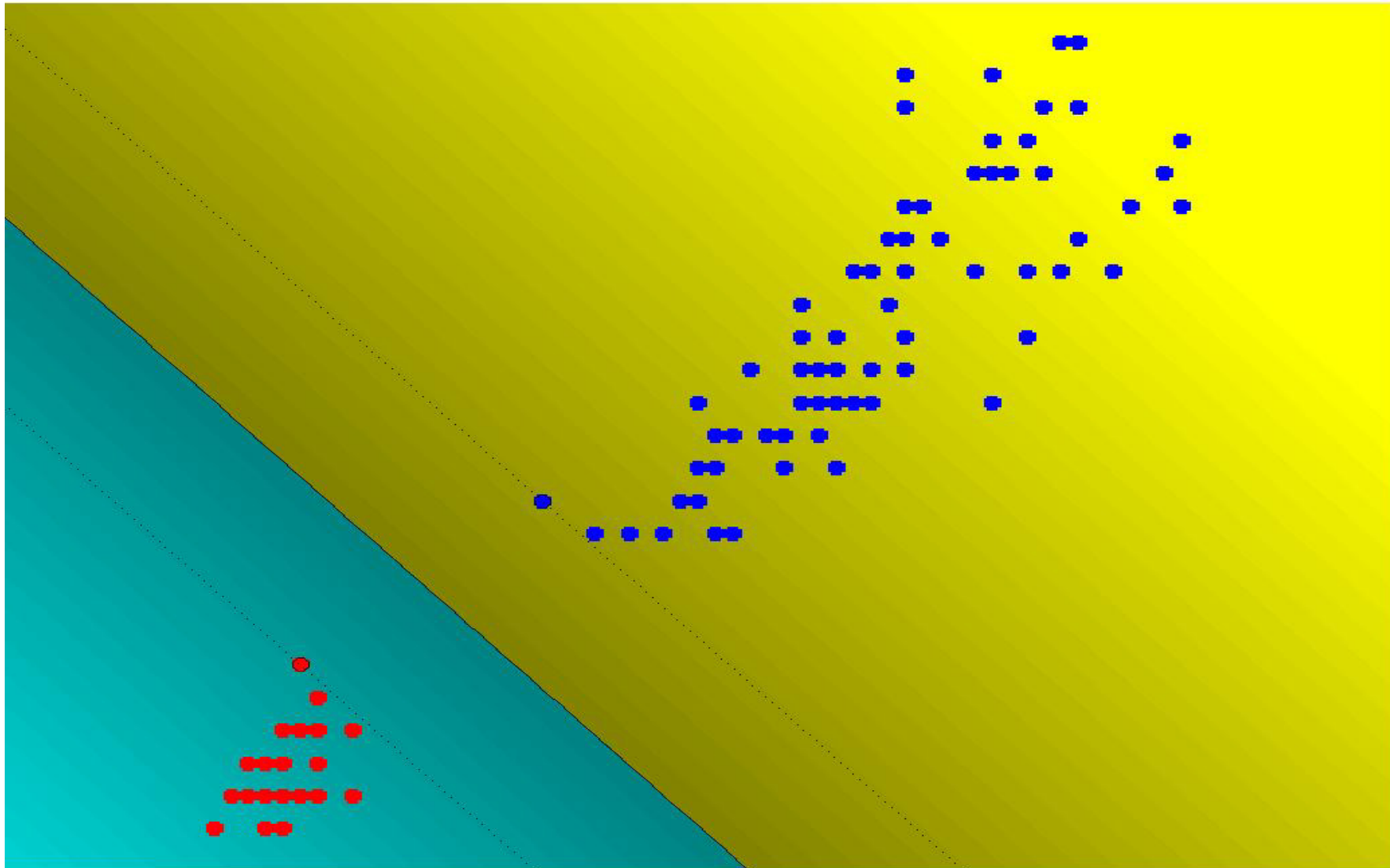
$$\begin{aligned} P(Y = 1 | \mathbf{x}, \mathbf{w}) &= \frac{1}{1 + e^{-(\sum_i \alpha_i \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}) + b)}} \\ &= \frac{1}{1 + e^{-(\sum_i \alpha_i K(\mathbf{x}, \mathbf{x}_i) + b)}} \end{aligned}$$

- Derive simple gradient descent rule on α_i

A few results

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

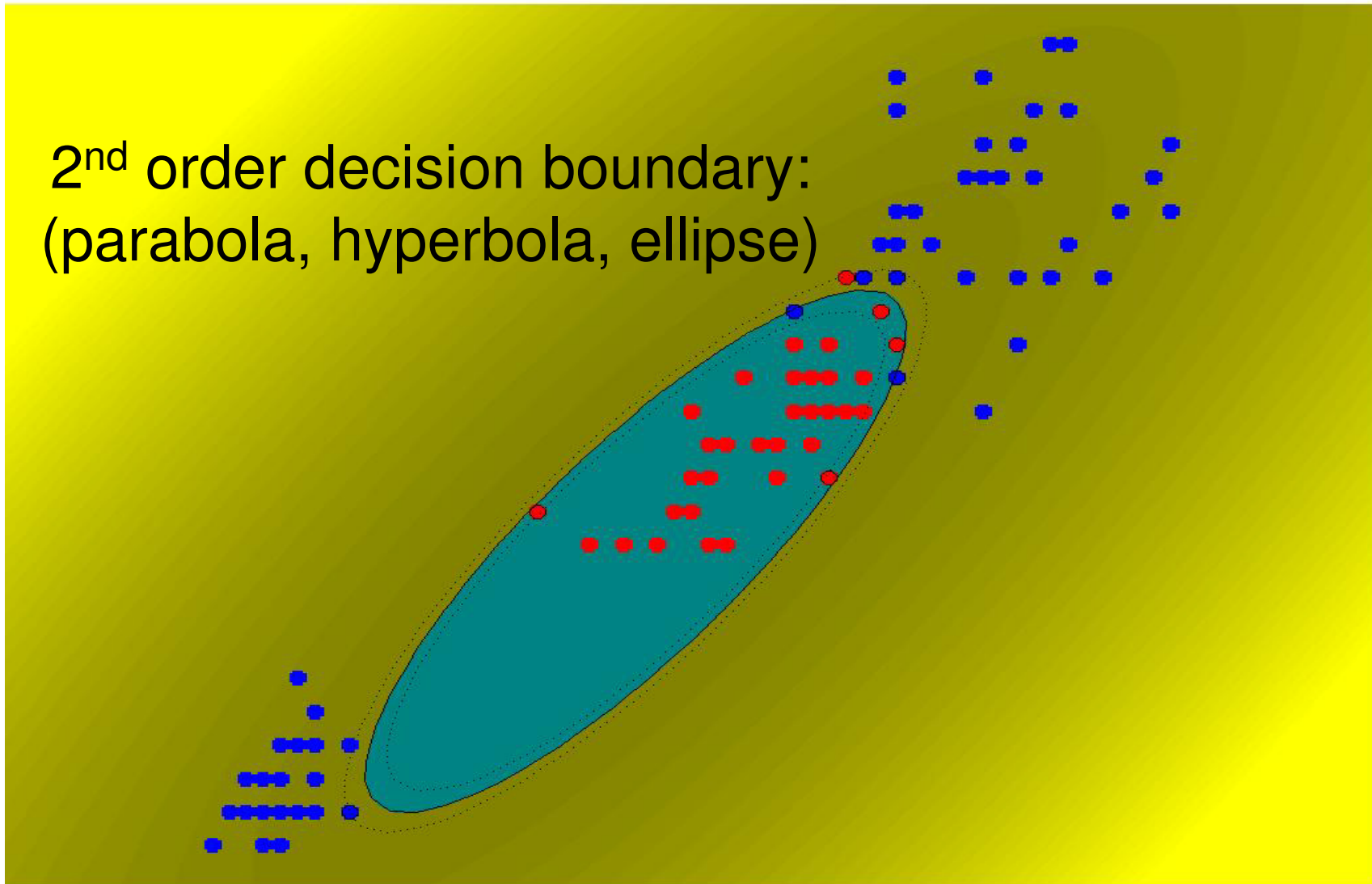
2

Separable

Bound

Inf

2nd order decision boundary:
(parabola, hyperbola, ellipse)



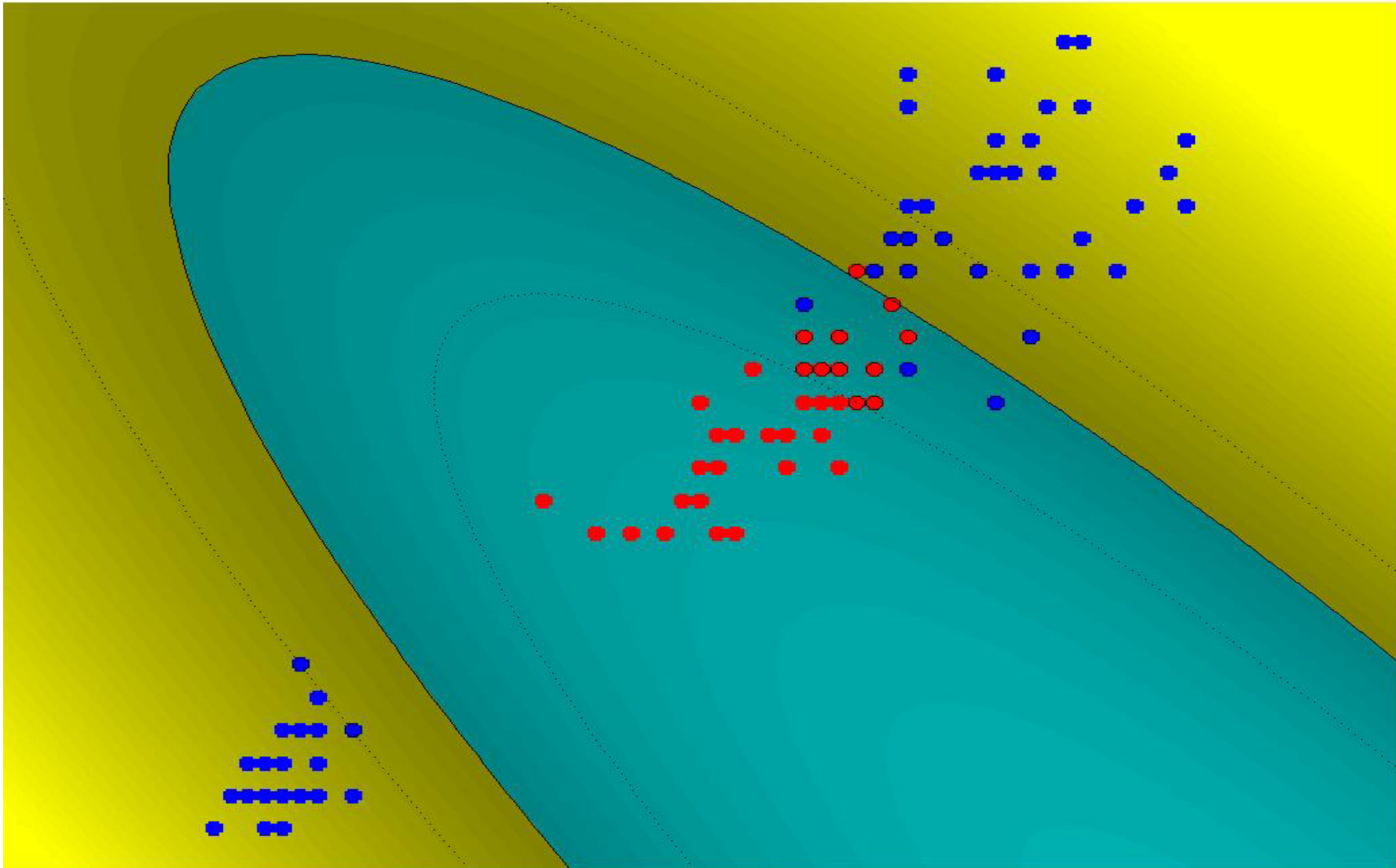
No. of Support Vectors: 12 (10.0%)

Polynomial

Degree 2

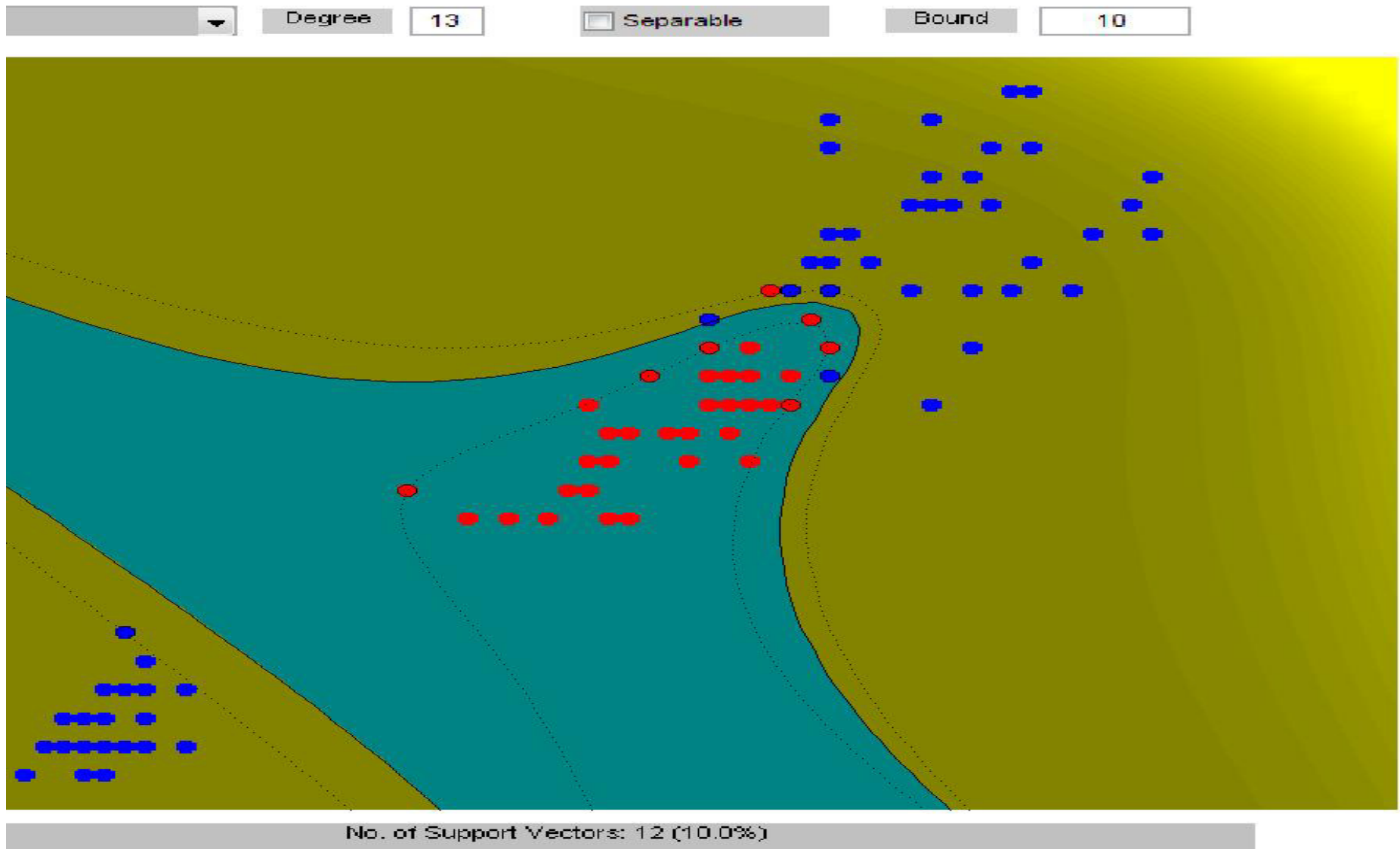
Separable

Bound 1



No. of Support Vectors: 30 (25.0%)

Results, Iris 1vs23, 13th order kernel



Gaussian RBF



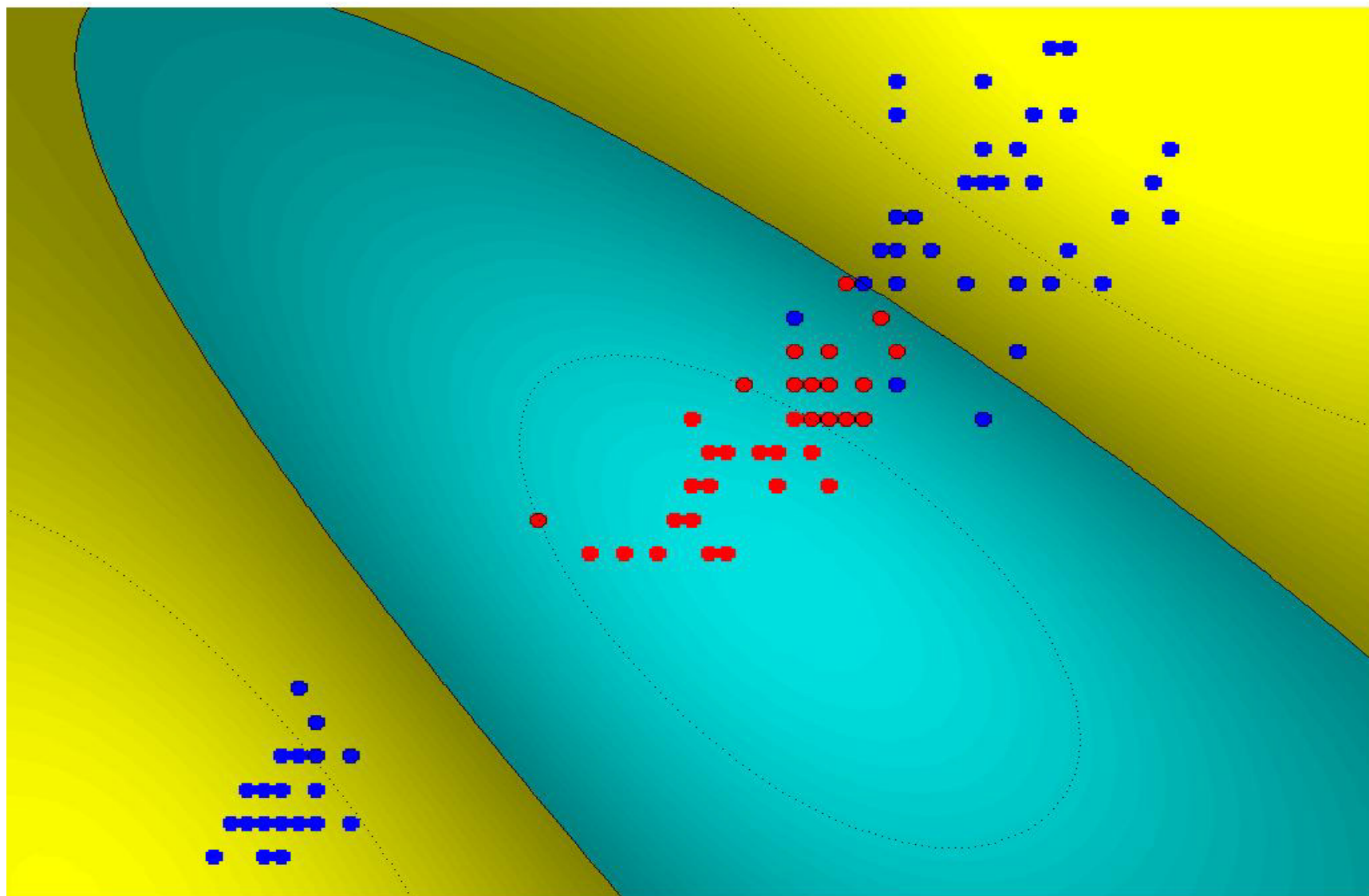
Sigma

1

Separable

Bound

1



No. of Support Vectors: 41 (34.2%)

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

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

Gaussian RBF

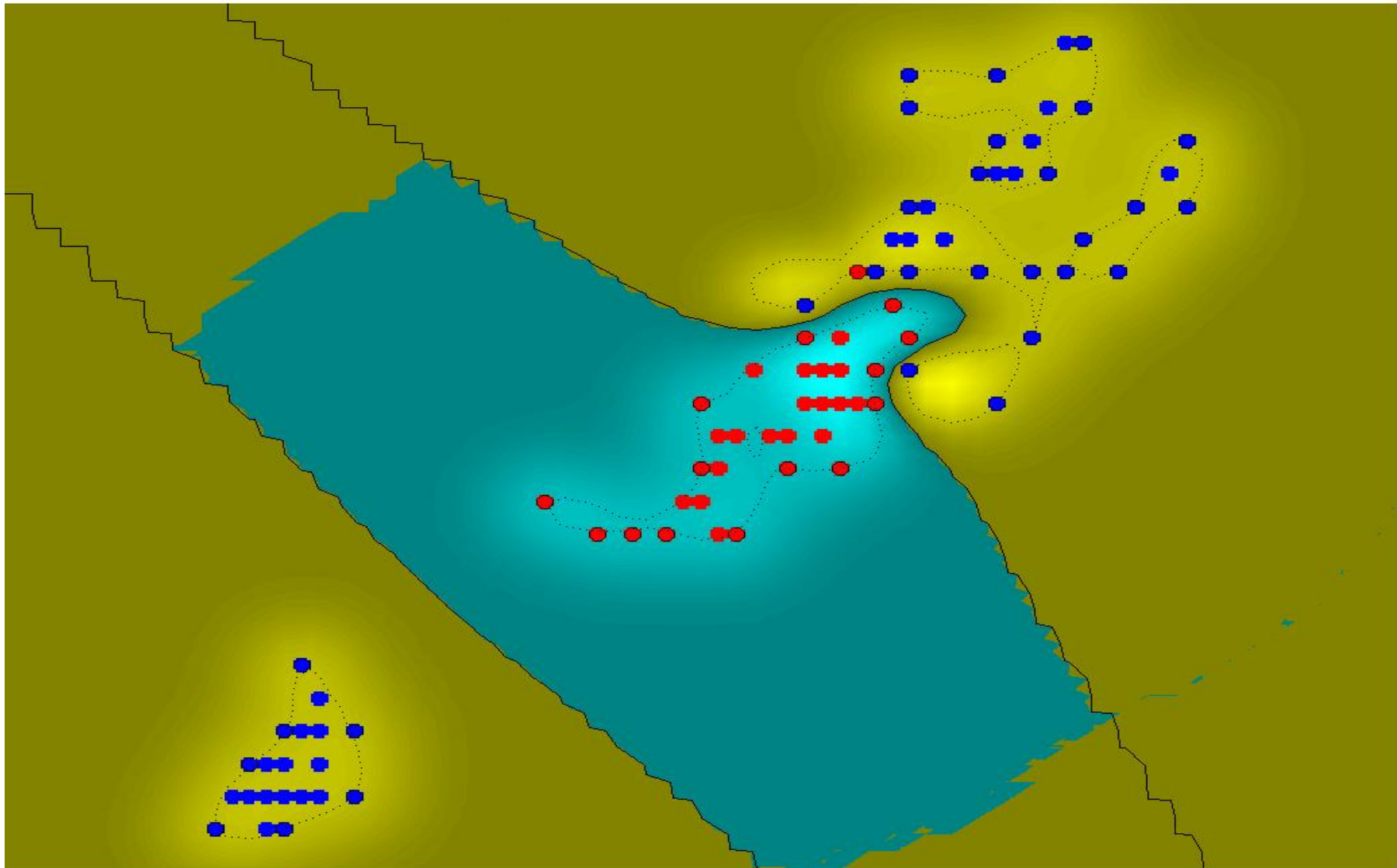
Sigma

.1

Separable

Bound

10



No. of Support Vectors: 55 (45.8%)

Chessboard dataset

Polynomial



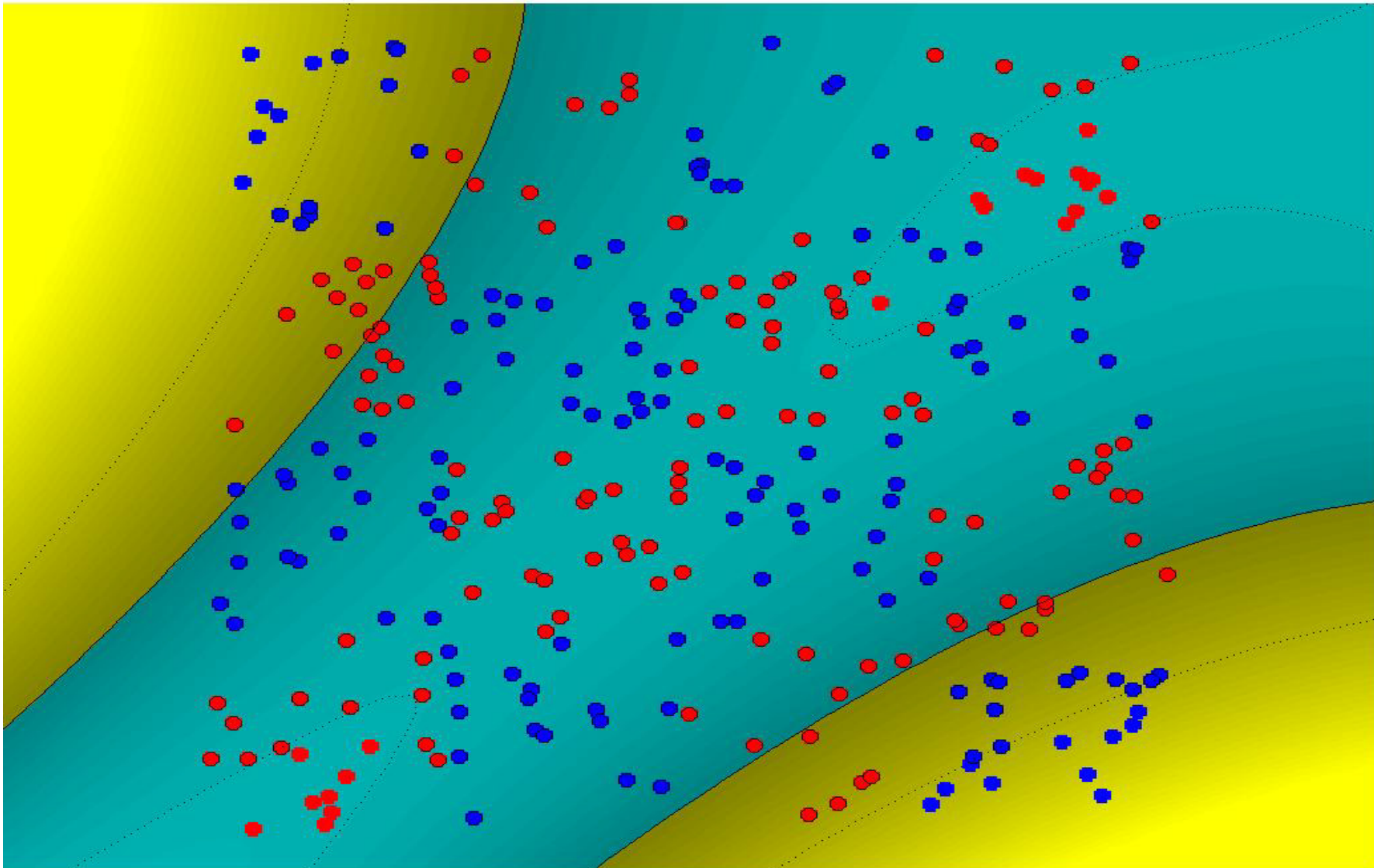
Degree

3

Separable

Bound

1



No. of Support Vectors: 263 (87.7%)

Polynomial



Degree

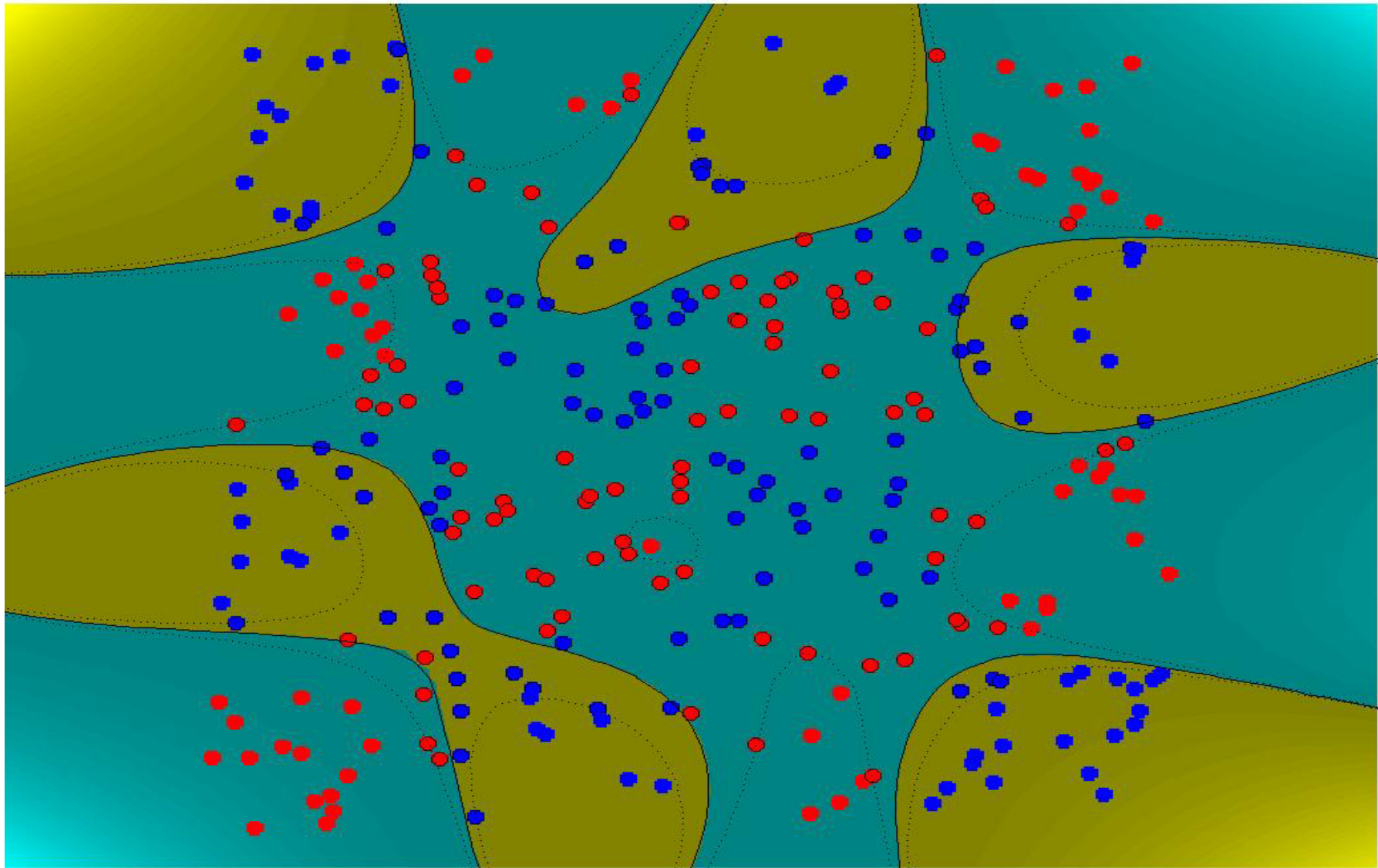
8



Separable

Bound

1



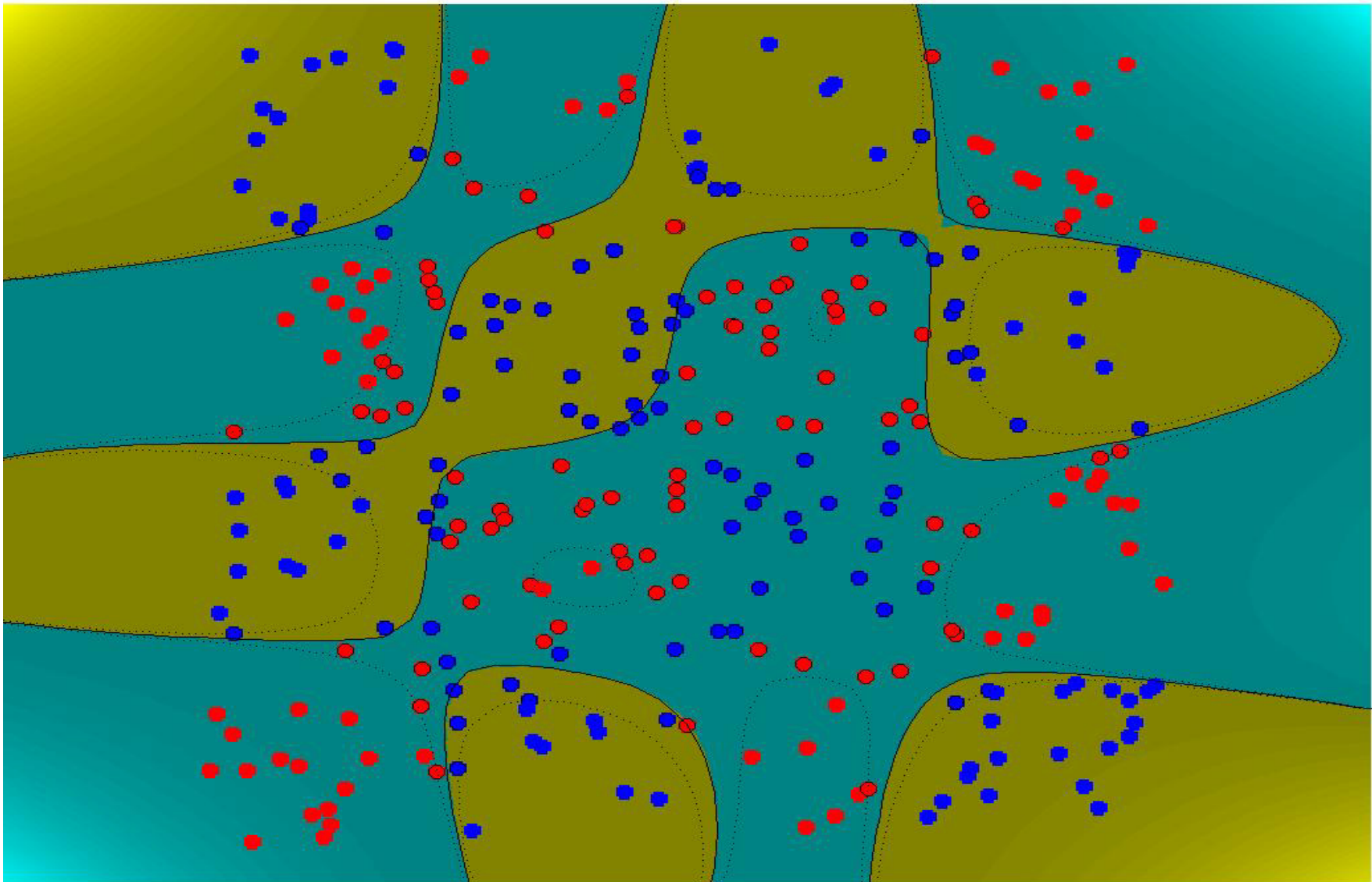
No. of Support Vectors: 183 (61.0%)

Polynomial

Degree 9

Separable

Bound 1



No. of Support Vectors: 164 (54.7%)

Polynomial

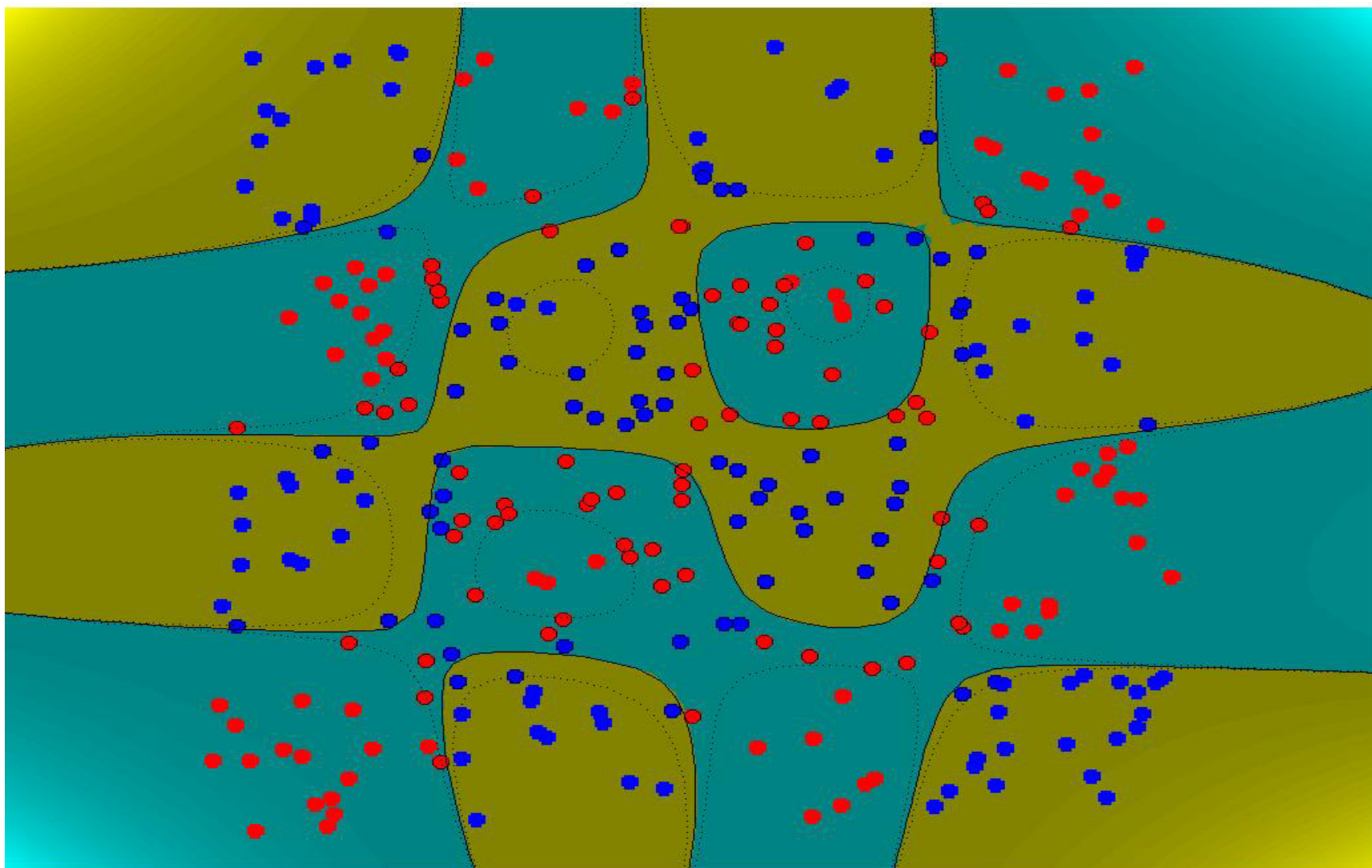
Degree

10

Separable

Bound

1



No. of Support Vectors: 147 (49.0%)

Polynomial



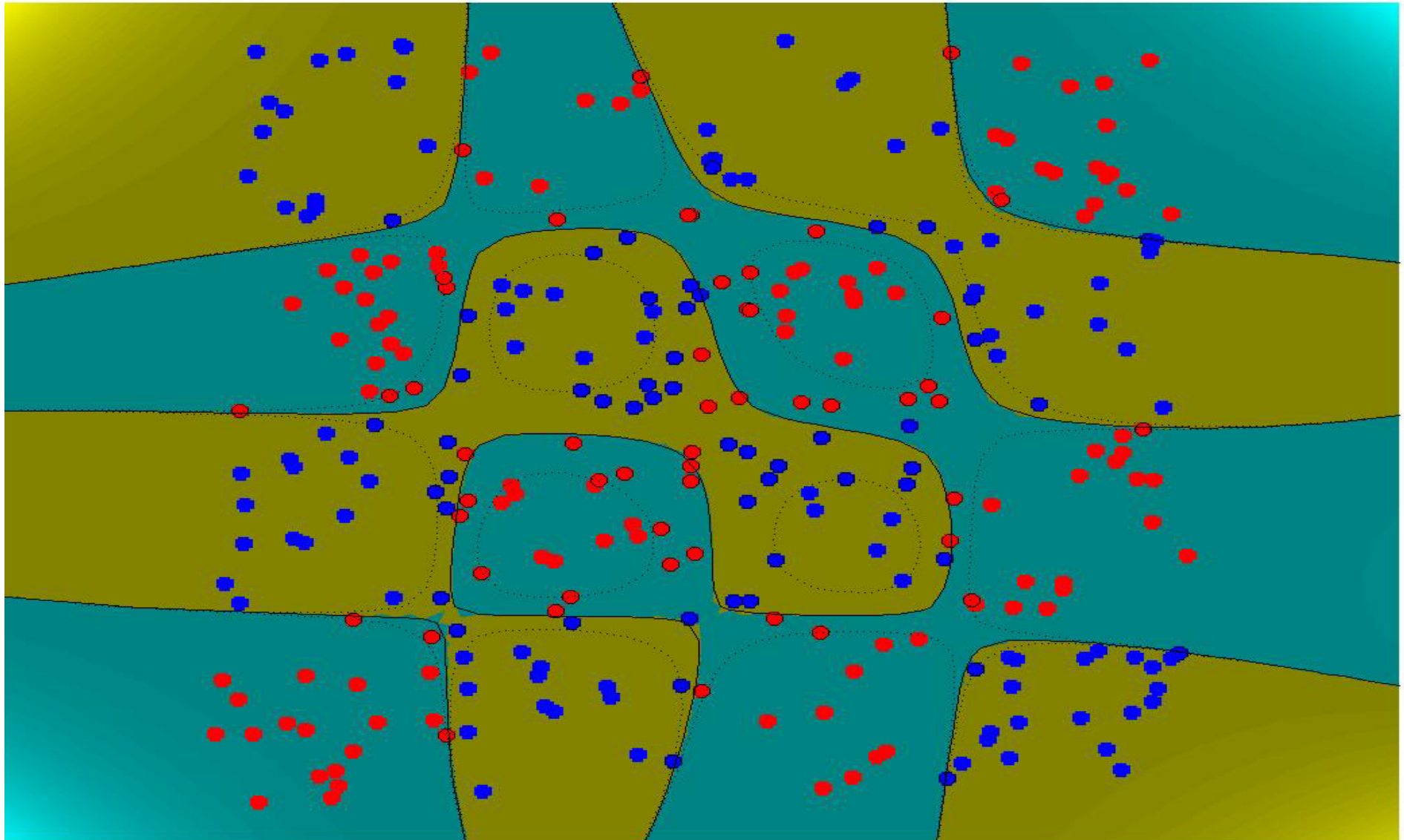
Degree

13

Separable

Bound

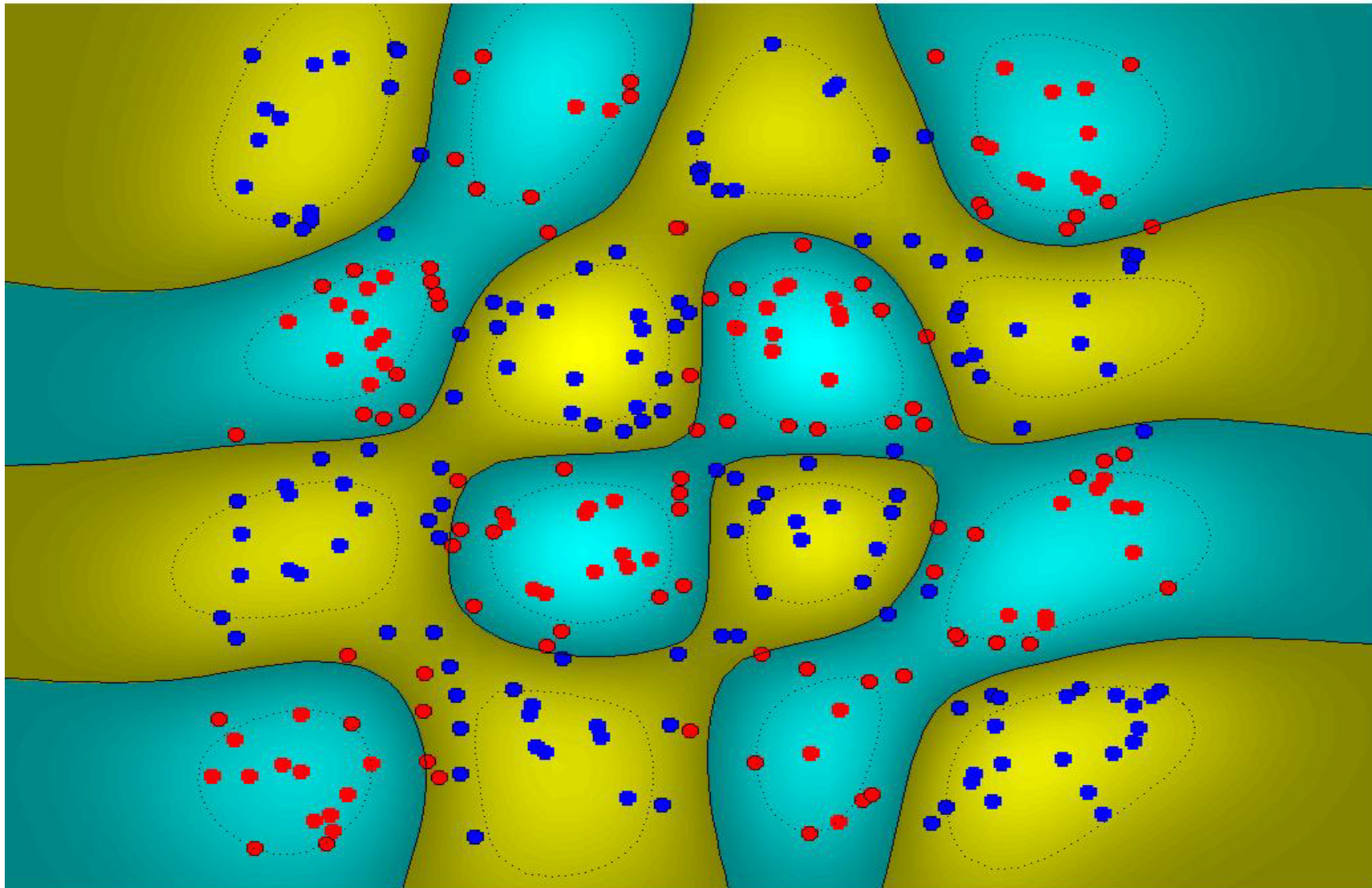
1



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%)