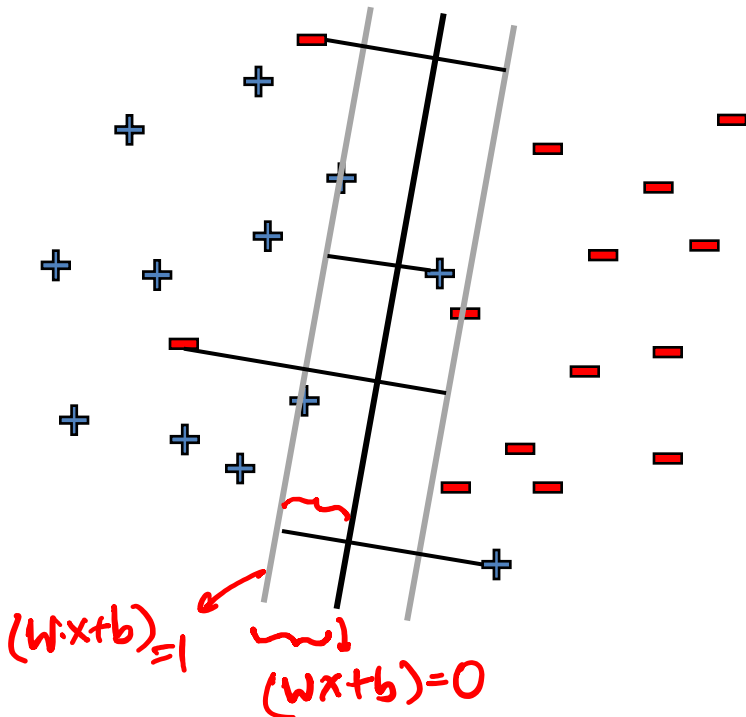


Soft margin SVM

Allow “error” in classification



$$\begin{aligned} \min_{w, b, \{\xi_j\}} \quad & \underline{w \cdot w} + C \sum_j \xi_j \\ \text{s.t.} \quad & \underline{(w \cdot x_j + b)} y_j \geq 1 - \xi_j \quad \forall j \\ & \xi_j \geq 0 \quad \forall j \end{aligned}$$

ξ_j - “slack” variables
= (>1 if x_j misclassified)

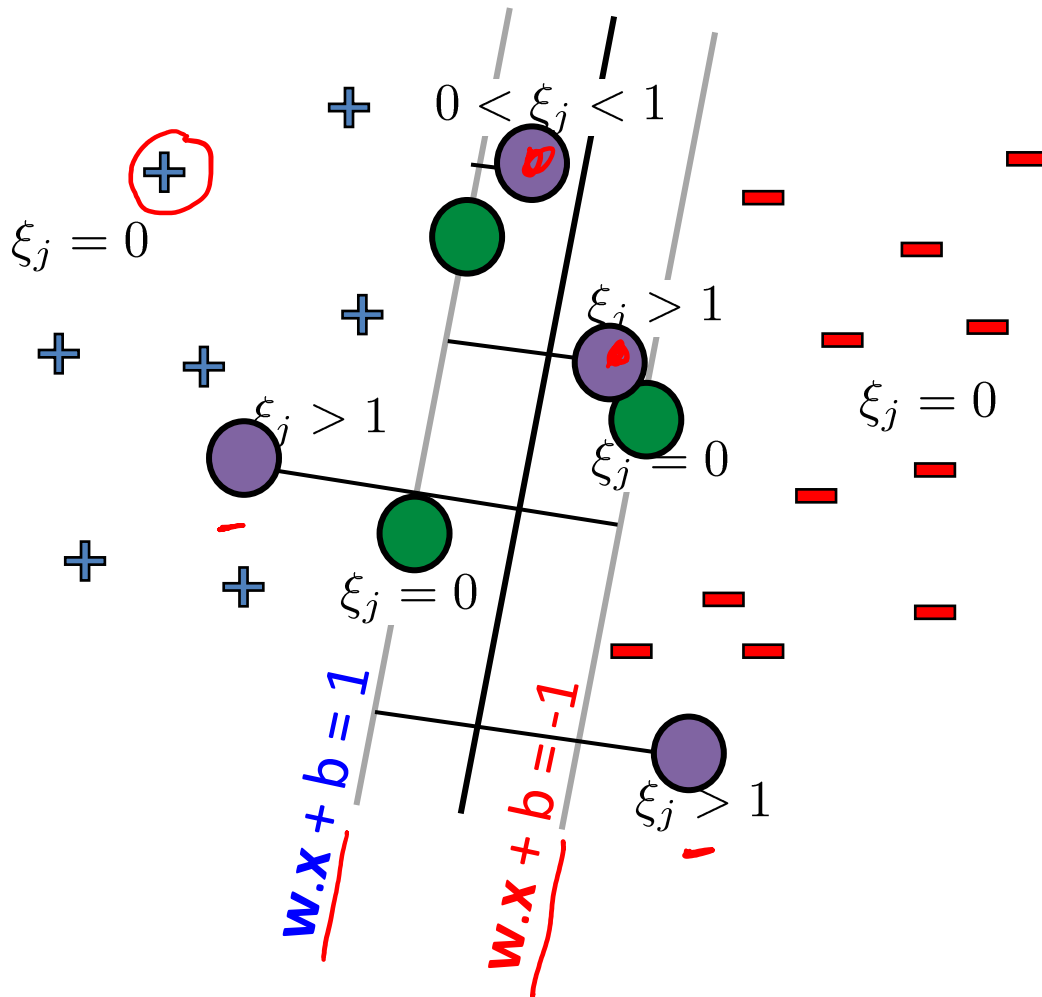
pay linear penalty if mistake

C - tradeoff parameter (C = ∞
recovers hard margin SVM)

Still QP 😊

$$(w \cdot x_j + b) y_j \geq 1 - \xi_j$$

Support Vectors



Margin support vectors

$\xi_j = 0$, $(w \cdot x_j + b) y_j = 1$ ←
 (don't contribute to objective but enforce constraints on solution)

Correctly classified but on margin

Non-margin support vectors

$\xi_j > 0$
 (contribute to both objective and constraints)

$1 > \xi_j > 0$ Correctly classified but inside margin

$\xi_j > 1$ Incorrectly classified 2

Support Vector Machines - Dual formulation and Kernel Trick

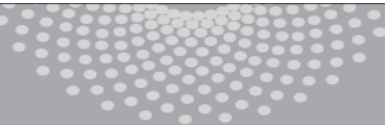
Aarti Singh & Geoff Gordon

Machine Learning 10-701

Mar 24, 2021



MACHINE LEARNING DEPARTMENT



Carnegie Mellon.
School of Computer Science

SVM – linearly separable case

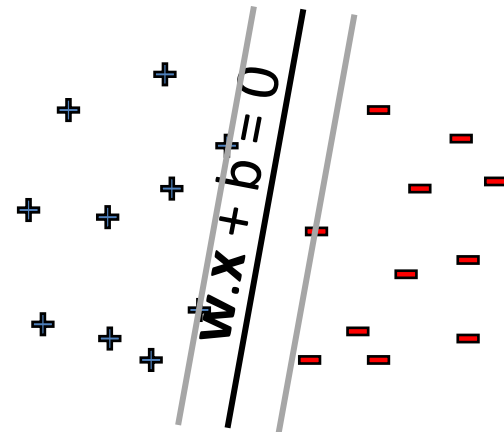
n training points

$(\mathbf{x}_1, \dots, \mathbf{x}_n)$

d features

\mathbf{x}_j is a d-dimensional vector

- Primal problem: minimize _{w, b} $\frac{1}{2} \mathbf{w} \cdot \mathbf{w}$
 $(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1, \forall j$



w - weights on features (d-dim problem)

- Convex quadratic program – quadratic objective, linear constraints
- But expensive to solve if d is very large
- Often solved in dual form (n-dim problem)

Dual SVM – linearly separable case

n training points, d features $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ where \mathbf{x}_i is a d-dimensional vector

- Primal problem: minimize_{w,b} $\frac{1}{2}\mathbf{w} \cdot \mathbf{w}$
 $(\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1, \forall j \leftarrow \alpha_j \geq 0$
 $\alpha_j [(\mathbf{w} \cdot \mathbf{x}_j + b) y_j - 1] \geq 0$
w - weights on features (d-dim problem)

- Dual problem (derivation):

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2}\mathbf{w} \cdot \mathbf{w} - \sum_{j=1}^n \alpha_j [(\mathbf{w} \cdot \mathbf{x}_j + b) y_j - 1]$$
$$\alpha_j \geq 0, \forall j$$

α - weights on training pts (n-dim problem)

Dual SVM – linearly separable case

$$\max_{\alpha_j \geq 0} \underline{\underline{d(\alpha)}}$$

- Dual problem (derivation):

$$\max_{\alpha} \min_{\underline{\underline{\mathbf{w}, b}}} L(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \sum_j \alpha_j \left[(\underline{\underline{\mathbf{w} \cdot \mathbf{x}_j}} + \underline{\underline{b}}) y_j - 1 \right]$$

$$\alpha_j \geq 0, \forall j$$

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_j \alpha_j \mathbf{x}_j y_j$$

$$\rightarrow \frac{\partial L}{\partial \mathbf{w}} = 0 \quad \Rightarrow \underline{\underline{\mathbf{w}}} = \sum_j \underline{\underline{\alpha_j y_j}} \mathbf{x}_j$$

$$\rightarrow \frac{\partial L}{\partial b} = 0 \quad \Rightarrow \sum_j \alpha_j y_j = 0$$
$$\frac{\partial L}{\partial b} = \sum_j \alpha_j y_j$$

Dual SVM – linearly separable case

- Dual problem:

$$\sum_j \alpha_j b y_j = b \sum_j \alpha_j y_j = 0$$

$$\max_{\alpha} \min_{\mathbf{w}, b} L(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w} \cdot \mathbf{w} - \sum_j \alpha_j [(\mathbf{w} \cdot \mathbf{x}_j + b) y_j - 1]$$

$$\alpha_j \geq 0, \forall j$$

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

$$\Rightarrow \sum_j \alpha_j y_j = 0$$

$$\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i \cdot x_j - \sum_j \alpha_j \left(\sum_i \alpha_i y_i x_i \cdot x_j \right) y_j + \sum_j \alpha_j$$

$$\sum_j \alpha_j - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i \cdot x_j =: d(\alpha)$$

Dual SVM – linearly separable case

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

$$\sum_i \alpha_i y_i = 0 \quad \leftarrow$$

$$\alpha_i \geq 0$$

Dual problem is also QP

Solution gives α_j s



$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i \quad \leftarrow$$

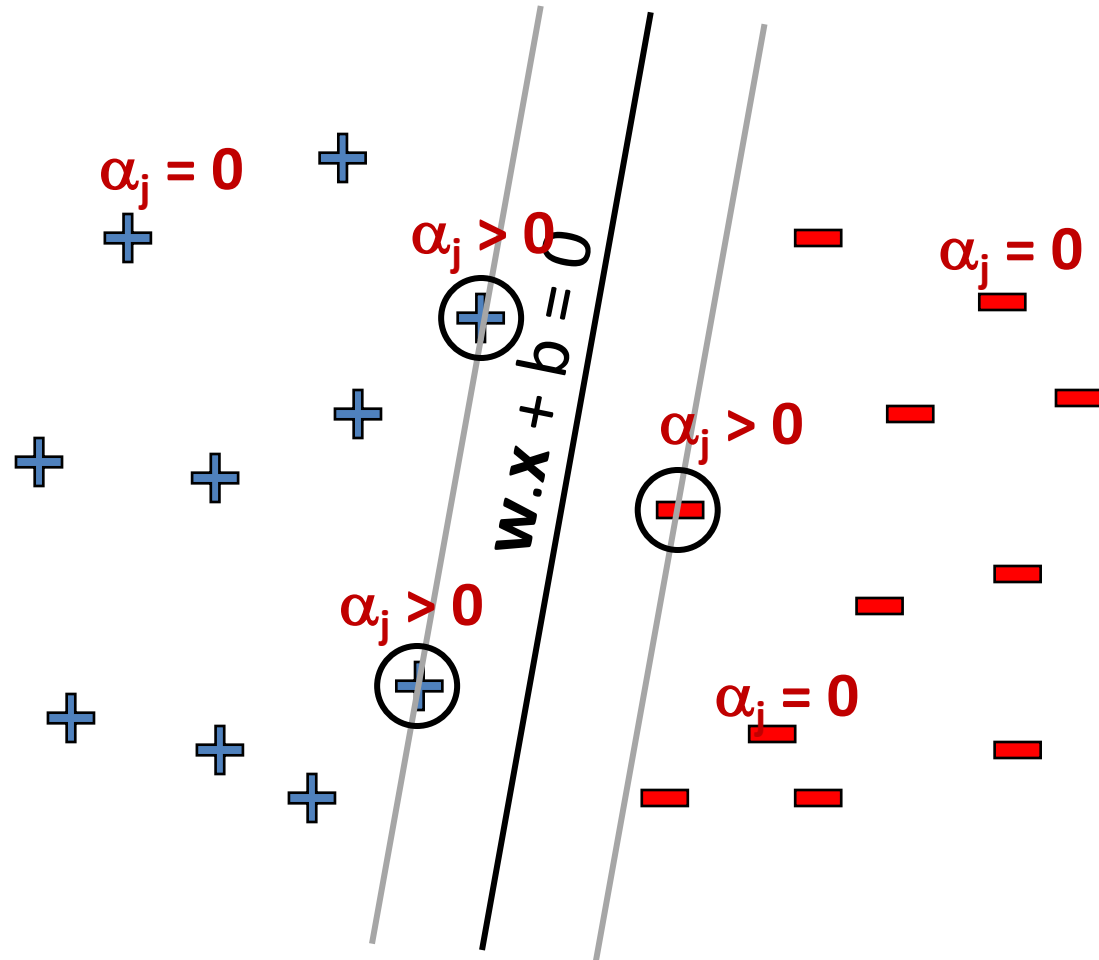
What about b?

Dual SVM: Sparsity of dual solution

KKT complementary slackness

$$\alpha_j [(w \cdot x_j + b) y_j - 1] = 0 \quad \leftarrow \quad \alpha_j \geq 0$$

$$w = \sum_j \alpha_j y_j x_j$$



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

$$(w \cdot x_j + b) y_j = 1$$

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

Dual SVM – linearly separable case

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

Dual problem is also QP

Solution gives α_j s \longrightarrow

Use any one of support vectors with $\alpha_k > 0$ to compute b since constraint is tight $(\mathbf{w} \cdot \mathbf{x}_k + b) y_k = 1$

$$\begin{aligned} \mathbf{w} &= \sum_i \alpha_i y_i \mathbf{x}_i \\ b &= y_k - \mathbf{w} \cdot \mathbf{x}_k \\ &\text{for any } k \text{ where } \alpha_k > 0 \end{aligned}$$

$$\begin{aligned} (\mathbf{w} \cdot \mathbf{x}_k + b) y_k &= 1 \\ \mathbf{w} \cdot \mathbf{x}_k + b &= y_k \end{aligned}$$

Dual SVM – non-separable case

$$L(\mathbf{w}, b, \xi, \alpha, \mu) = \frac{\mathbf{w} \cdot \mathbf{w}}{2} + C \sum_j \xi_j - \sum_j \alpha_j [(w \cdot x_j + b) y_j - 1 + \xi_j] - \sum_j \mu_j \xi_j$$

- Primal problem:

$$\begin{aligned} & \text{minimize}_{\mathbf{w}, b, \{\xi_j\}} \frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_j \xi_j \\ \rightarrow & (w \cdot x_j + b) y_j \geq 1 - \xi_j, \quad \forall j \\ & \xi_j \geq 0, \quad \forall j \end{aligned} \quad \begin{array}{l} \alpha_j \geq 0 \\ \mu_j \geq 0 \end{array}$$

Lagrange
Multipliers

- Dual problem:

$$\begin{aligned} & \max_{\alpha, \mu} \min_{\mathbf{w}, b, \{\xi_j\}} L(\mathbf{w}, b, \xi, \alpha, \mu) \\ & s.t. \alpha_j \geq 0 \quad \forall j \\ & \mu_j \geq 0 \quad \forall j \end{aligned} \quad \begin{array}{l} \frac{\partial L}{\partial \xi_j} = C - \alpha_j - \mu_j = 0 \\ \Rightarrow \alpha_j \leq C \end{array}$$

Dual SVM – non-separable case

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

$$\sum_i \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0 \quad \leftarrow$$

comes from $\frac{\partial L}{\partial \xi} = 0$

Intuition:

If $C \rightarrow \infty$, recover hard-margin SVM

Dual problem is also QP

Solution gives α_j



$$\underline{\mathbf{w}} = \sum_i \underline{\alpha}_i y_i \mathbf{x}_i$$

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

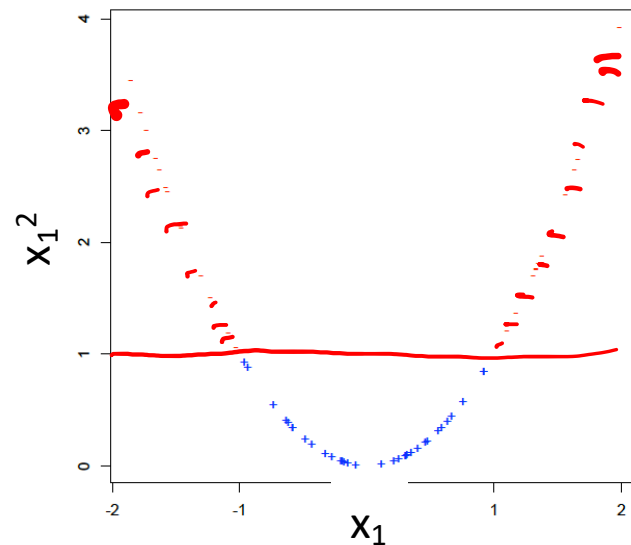
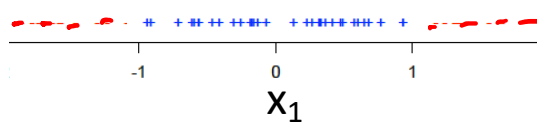
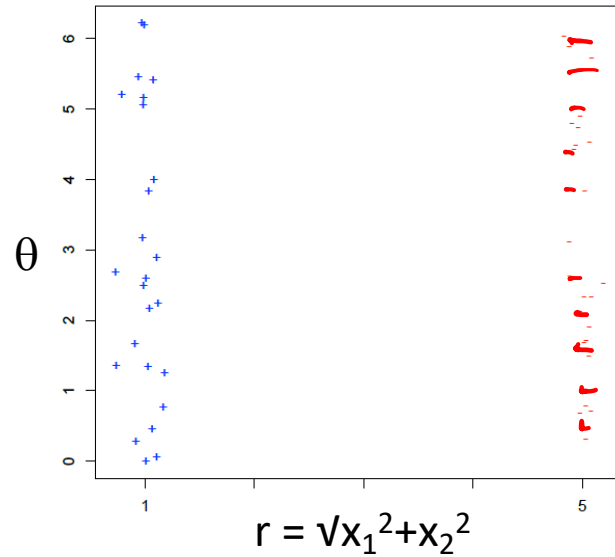
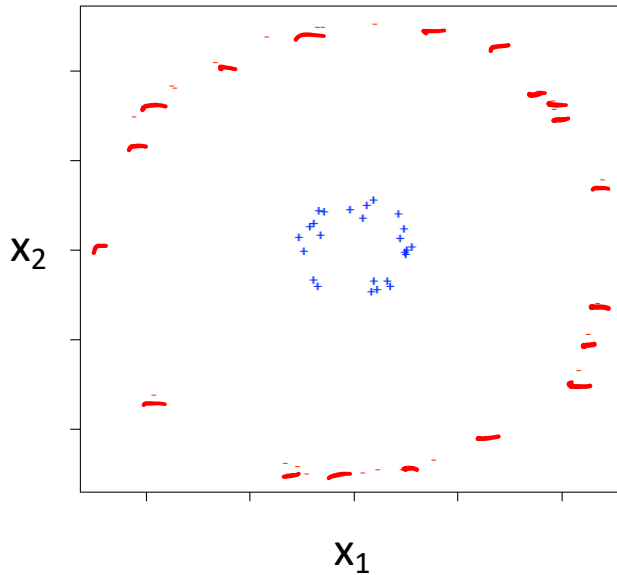
for any k where $C > \underline{\alpha}_k > 0$

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 $d \gg n$)
- But, more importantly, the “**kernel trick**”!!!

Separable using higher-order features

$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$



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

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



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

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

$$\begin{aligned} & \sum_i \alpha_i y_i = 0 \\ & C \geq \alpha_i \geq 0 \end{aligned}$$

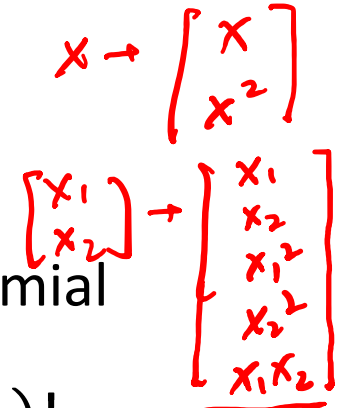
$\mathbf{x}_i \rightarrow \Phi(\mathbf{x}_i)$
 $\in \mathbb{R}^d \rightarrow \in \mathbb{R}^D$
 $D \gg d$

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

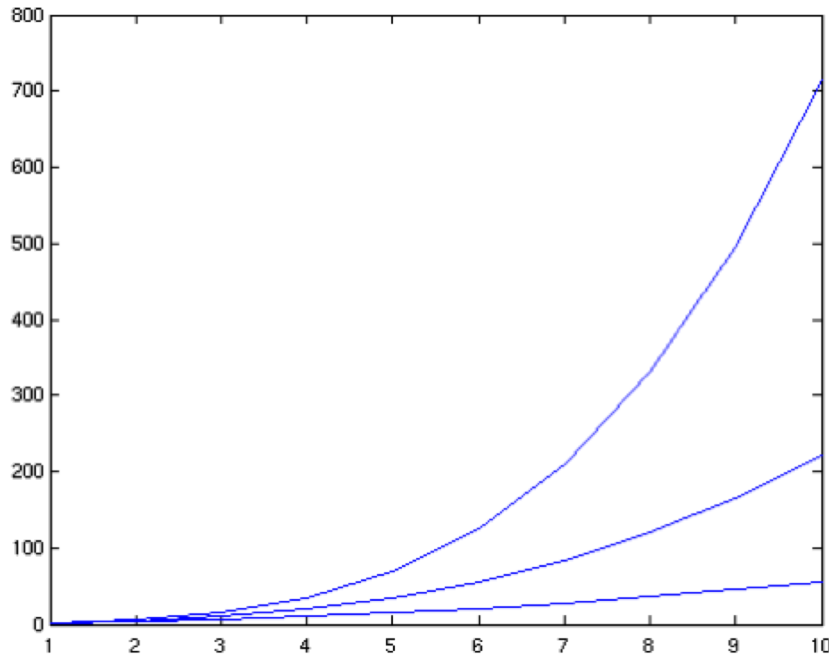
Polynomial features $\phi(\mathbf{x})$

m – input features

d – degree of polynomial



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



grows fast!

$d = 6, m = 100$

about 1.6 billion terms

Dot Product of Polynomial features

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

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

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

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

$\phi(\mathbf{x})$ $\phi(\mathbf{z})$
↑ ↑

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

The Kernel Trick!

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

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

$$\sum_i \alpha_i y_i = 0$$

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

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

Common Kernels

- Polynomials of degree d

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

- Polynomials of degree up to d

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

- Gaussian/Radial kernels (polynomials of all orders – recall series expansion of exp)

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

- Sigmoid

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

Mercer Kernels

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

Answer: Mercer kernels K

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

$$K_{ij} = \phi(x_i) \cdot \phi(x_j)$$

Ensures optimization is concave maximization

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?

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

$$\begin{aligned} \rightarrow \mathbf{w} &= \sum_i \alpha_i y_i \Phi(\mathbf{x}_i) \\ \rightarrow b &= y_k - \mathbf{w} \cdot \Phi(\mathbf{x}_k) \\ &\text{for any } k \text{ where } C > \alpha_k > 0 \end{aligned}$$

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

- Using kernels we are cool!

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

SVMs with Kernels

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

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

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

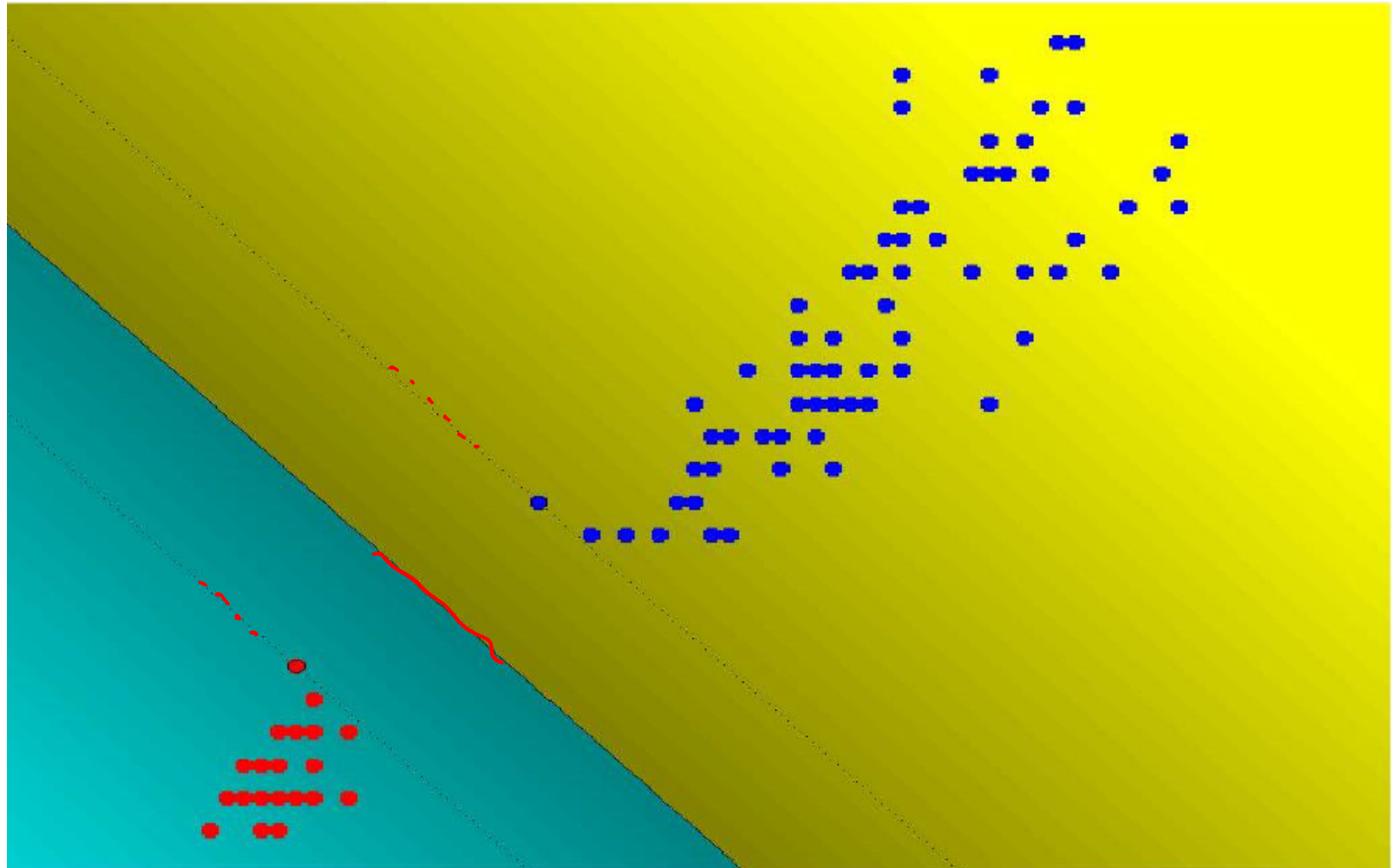
for any k where $C > \alpha_k > 0$

Classify as

$$\underline{\underline{\text{sign}(\mathbf{w} \cdot \Phi(\mathbf{x}) + b)}}$$

SVMs with Kernels

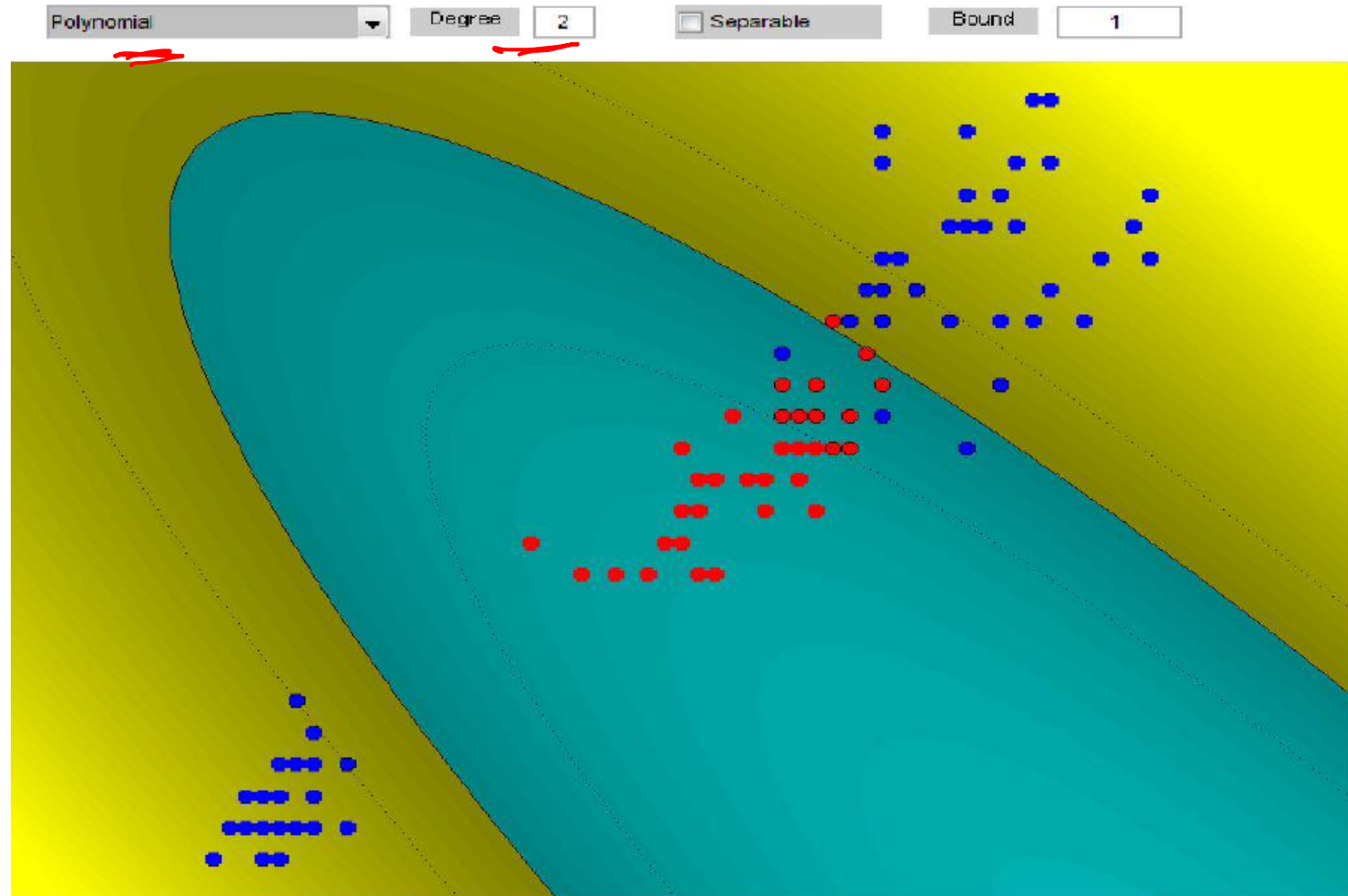
- Iris dataset, 2 vs 13, Linear Kernel



No. of Support Vectors: 2 (1.7%)

SVMs with Kernels

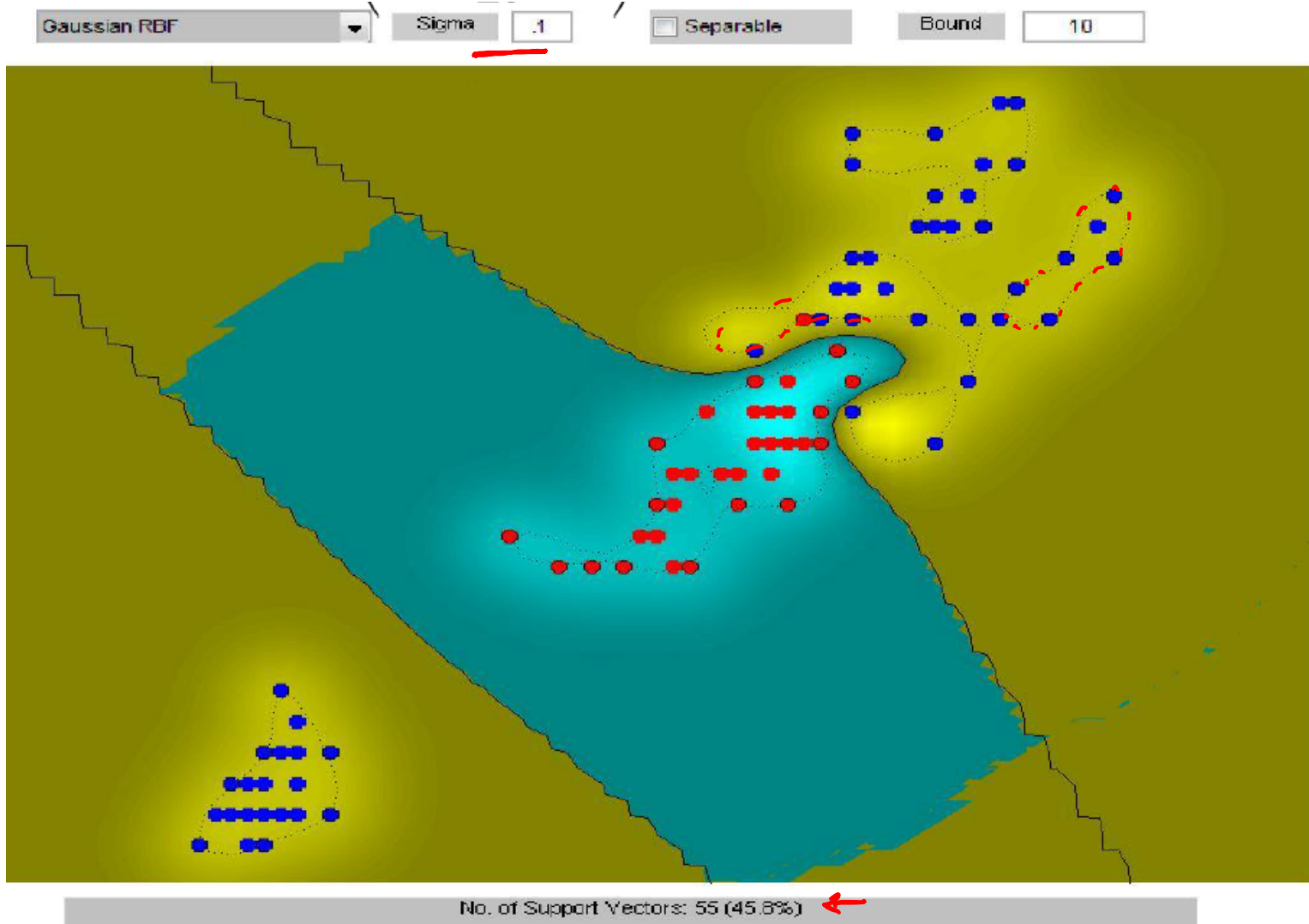
- Iris dataset, 1 vs 23, Polynomial Kernel degree 2



No. of Support Vectors: 30 (25.0%)

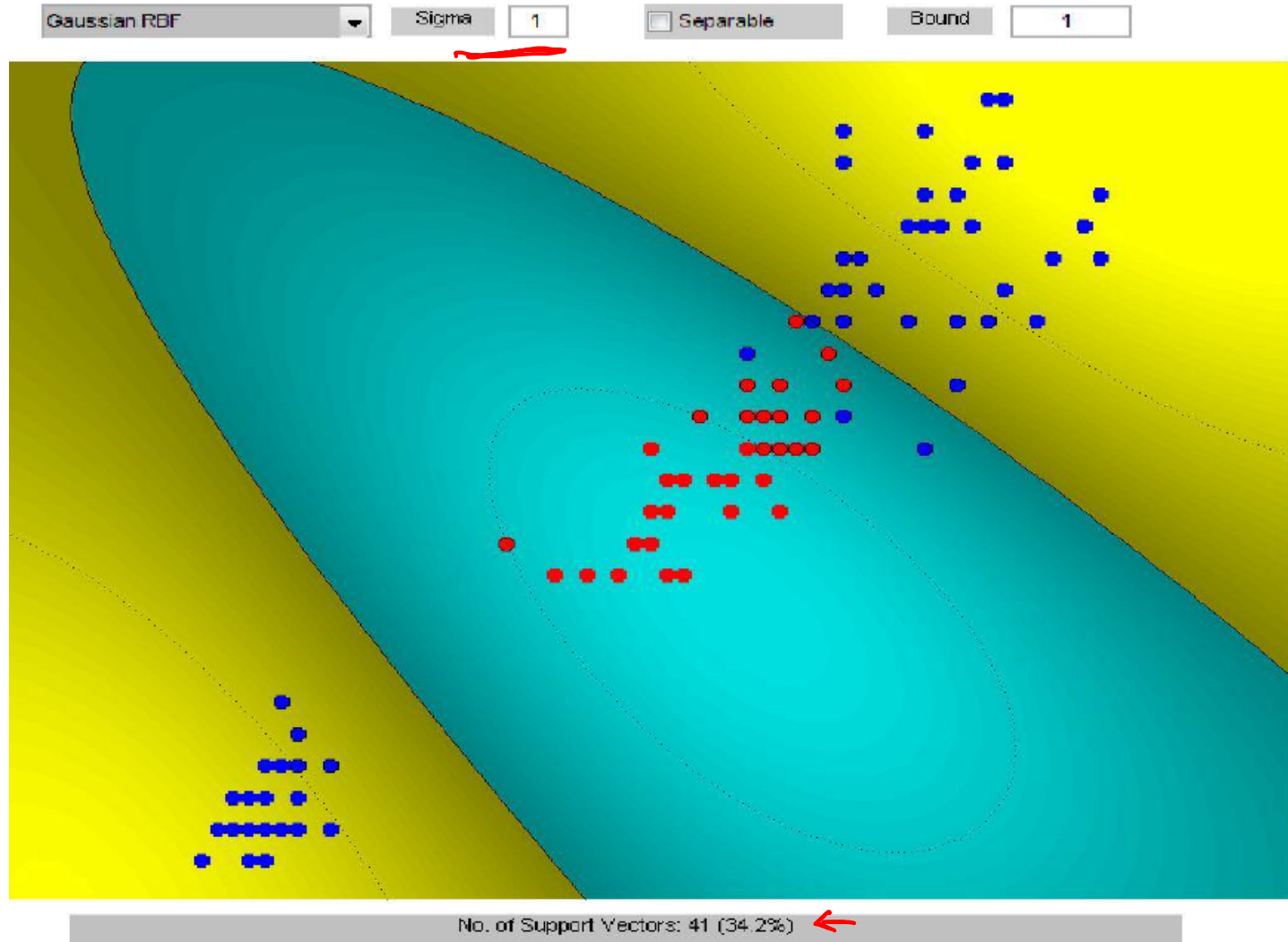
SVMs with Kernels

- Iris dataset, 1 vs 23, Gaussian RBF kernel



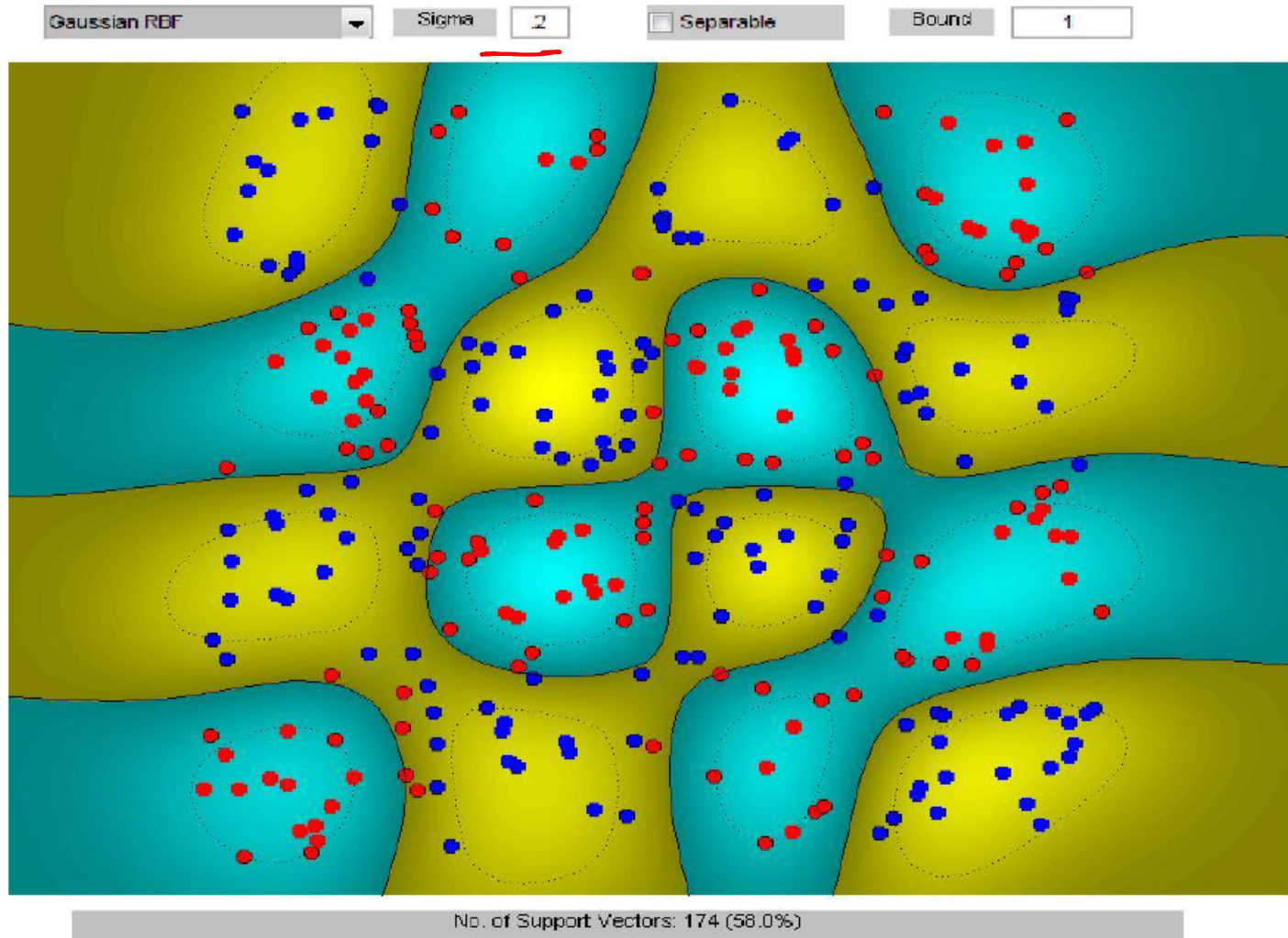
SVMs with Kernels

- Iris dataset, 1 vs 23, Gaussian RBF kernel



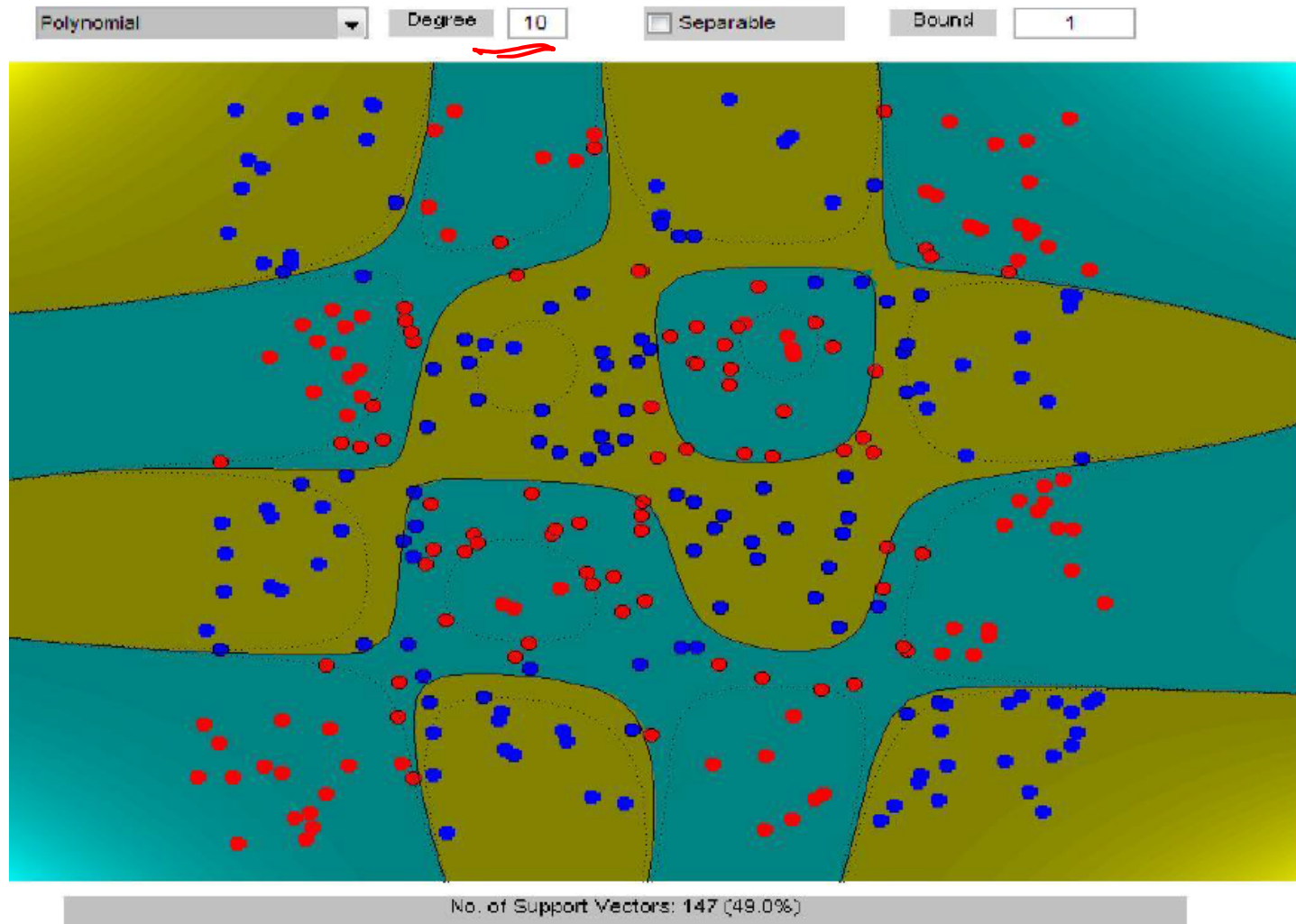
SVMs with Kernels

- Chessboard dataset, Gaussian RBF kernel

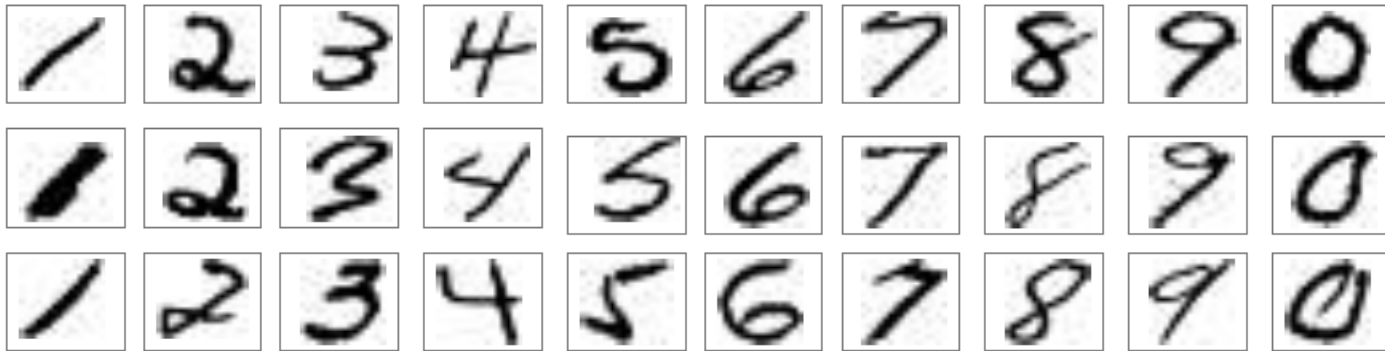


SVMs with Kernels

- Chessboard dataset, Polynomial kernel



USPS Handwritten digits



- 1000 training and 1000 test instances

Results:

SVM on raw images ~97% accuracy

SVMs vs. Logistic Regression

	SVMs	Logistic Regression
Loss function	Hinge loss	Log-loss

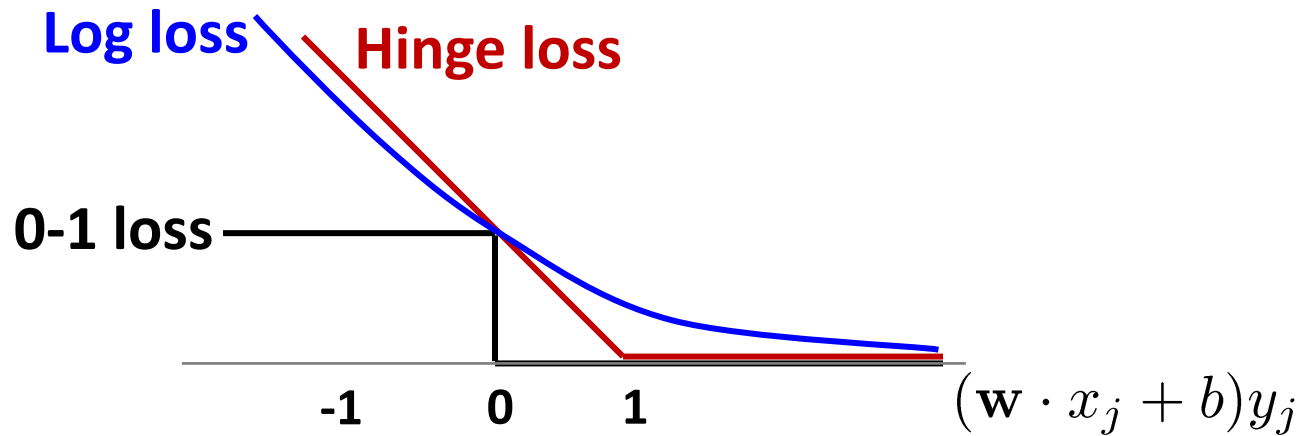
SVMs vs. Logistic Regression

SVM : **Hinge loss**

$$\text{loss}(f(x_j), y_j) = \max(0, 1 - (\mathbf{w} \cdot x_j + b)y_j)$$

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

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



SVMs vs. Logistic Regression

	SVMs	Logistic Regression
Loss function	Hinge loss	Log-loss
High dimensional features with kernels	Yes!	Yes!

Kernels in Logistic Regression

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

- Define weights in terms of features:

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

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

$$P(Y=1 | x, \alpha) = \frac{1}{1 + e^{-\left(\sum_i \alpha_i K(\mathbf{x}, \mathbf{x}_i) + b\right)}}$$

- Derive simple gradient descent rule on α_i

SVMs vs. Logistic Regression

	SVMs	Logistic Regression
Loss function	Hinge loss	Log-loss
High dimensional features with kernels	Yes!	Yes!
<u>Solution</u> sparse	Often yes!	Almost always no!
Semantics of output	“ <u>Margin</u> ”	Real <u>probabilities</u>