

Support Vector Machines - Dual formulation and Kernel Trick

Aarti Singh

Machine Learning 10-315

Mar 28, 2022



MACHINE LEARNING DEPARTMENT



SVM summary so far

n training points
d features

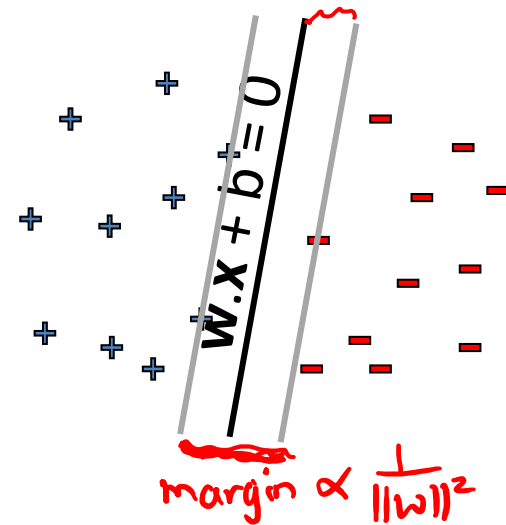
$(\mathbf{x}_1, \dots, \mathbf{x}_n)$

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

Hard-margin:

$$\text{minimize}_{\mathbf{w}, b} \quad \frac{1}{2} \mathbf{w} \cdot \mathbf{w} \quad \left\{ \begin{array}{l} (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1, \quad \forall j \end{array} \right.$$

linearly separable

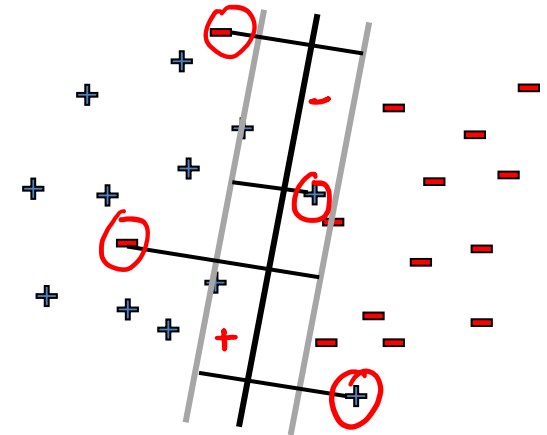


Soft-margin:

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



not linearly separable



SVM primal vs dual

n training points

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

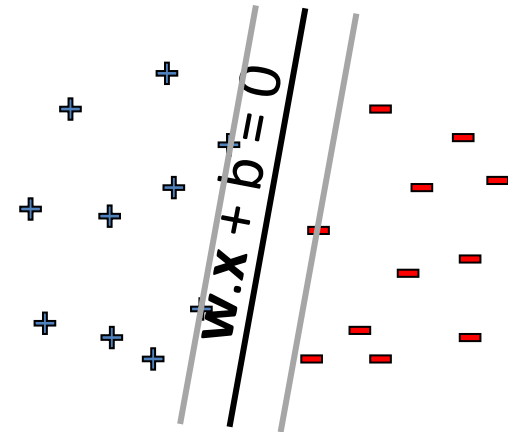
d features

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

Hard-margin:

Primal problem

$$\text{minimize}_{\mathbf{w}, b} \quad \frac{1}{2} \mathbf{w} \cdot \mathbf{w}$$
$$\left(\mathbf{w} \cdot \mathbf{x}_j + b \right) y_j \geq 1, \quad \forall j$$



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)

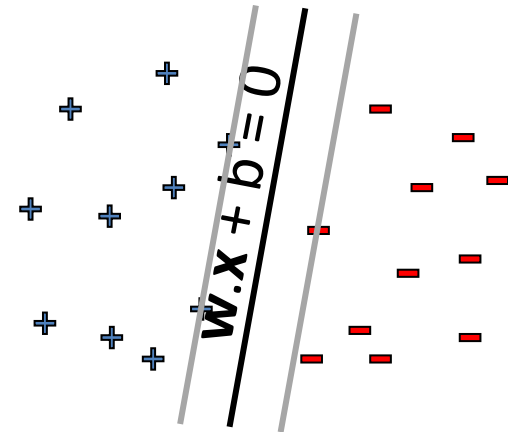
SVM primal vs dual

n training points

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

d features

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



Hard-margin:

Primal problem

$$\text{minimize}_{\mathbf{w}, b} \quad \frac{1}{2} \mathbf{w} \cdot \mathbf{w}$$

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

\mathbf{w} - weights on features (d-dim problem)

Dual problem

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

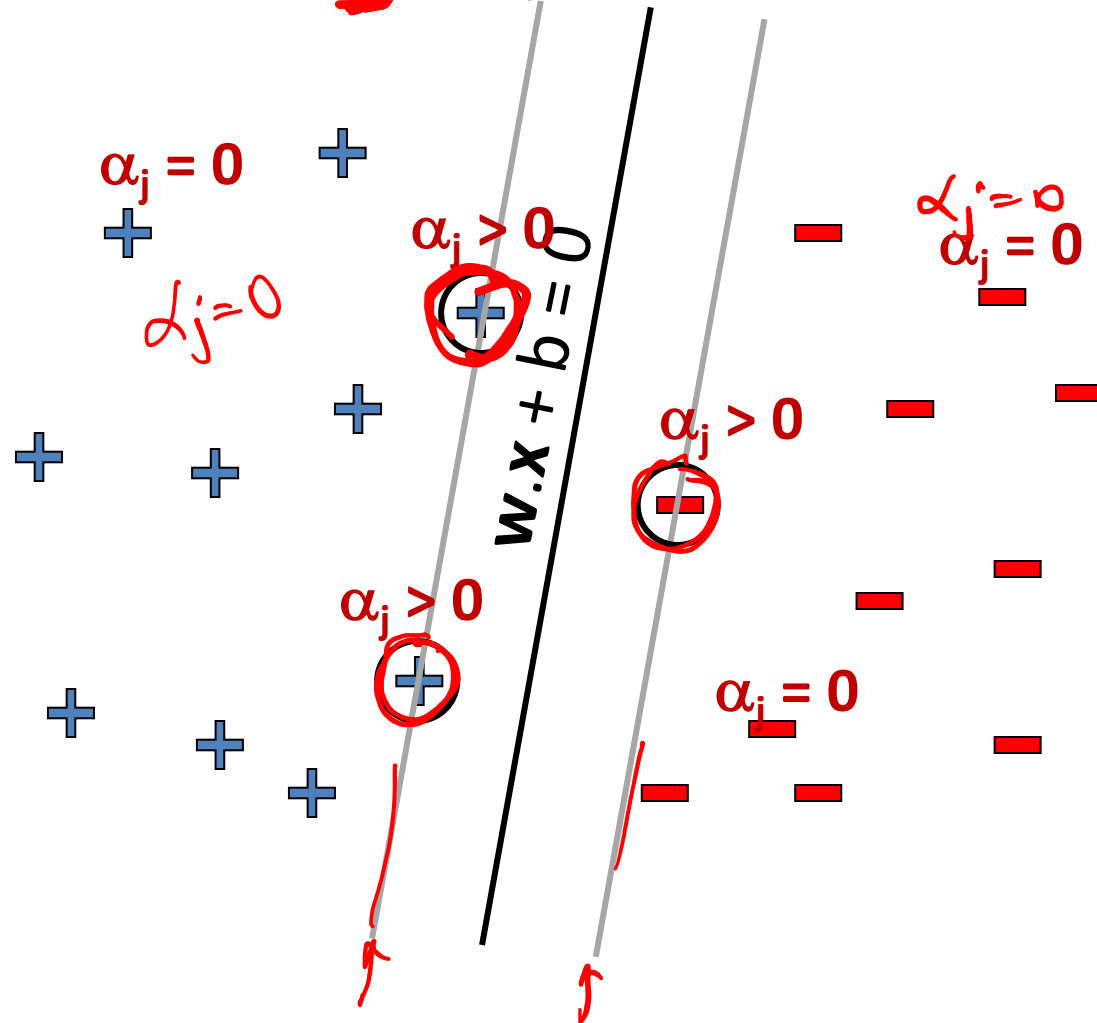
$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0$$

α - weights on data points (n-dim problem)

Dual SVM: Sparsity of dual solution

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



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

Complementary slackness implies

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

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

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

Dual SVM – linearly separable (aka hard margin) case

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

$$* \quad b = \frac{1}{y_k} - \mathbf{w} \cdot \mathbf{x}_k \quad \checkmark$$

Dual problem is also QP

Solution gives α_j s \longrightarrow

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

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

$\Rightarrow (x_k, y_k)$ – support vector

Dual SVM – non-separable case

soft-margin SVM

- Primal problem:

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

$$\begin{array}{|c|} \hline \alpha_j \\ \hline \mu_j \\ \hline \end{array}$$

**Lagrange
Multipliers**

- Dual problem:

$$\max_{\alpha, \mu} \min_{\mathbf{w}, b, \{\xi_j\}} L(\mathbf{w}, b, \xi, \alpha, \mu)$$

$$s.t. \alpha_j \geq 0 \quad \forall j$$

$$\mu_j \geq 0 \quad \forall j$$

$$\frac{\partial L}{\partial \mathbf{w}}$$

$$\frac{\partial L}{\partial b}$$

$$\frac{\partial L}{\partial \xi_j}$$

Dual SVM – non-separable case

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

$$\sum_i \alpha_i y_i = 0$$

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

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 s



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

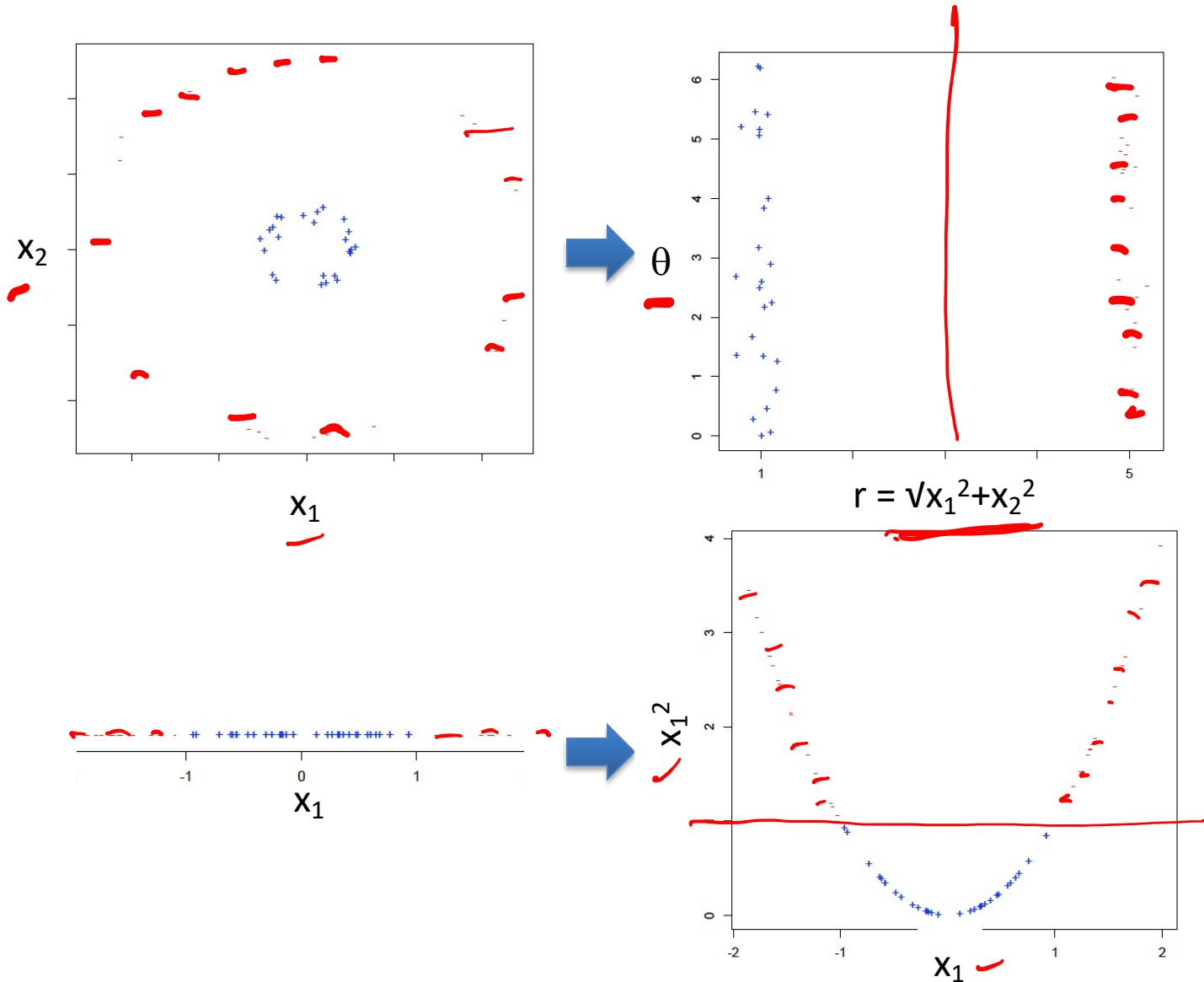
$$b = y_k - \mathbf{w} \cdot \mathbf{x}_k$$

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

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



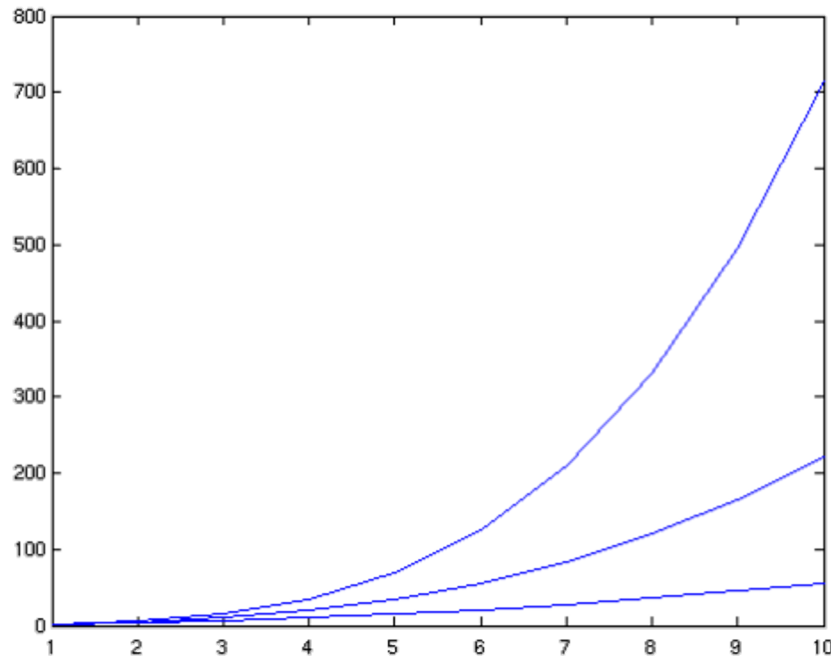
$$x_1 \dots x_m \quad x_1 x_1^2 \dots x_1^d \quad x_2 x_2^2 \dots x_2^d \quad x_1 x_2^{d-1}$$

Polynomial features $\phi(x)$

m – input features

d – degree of polynomial

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

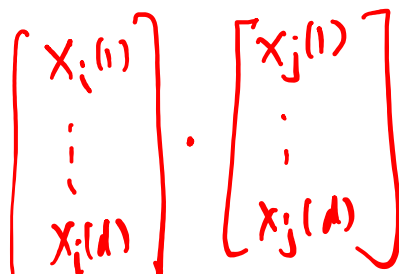


grows fast!

$d = 6, m = 100$

about 1.6 billion terms

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

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




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

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

Dot Product of Polynomial features

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

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

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

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

$d^2 - m^d$ \uparrow \uparrow

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

The Kernel Trick!

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

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

$$\sum_i \alpha_i y_i = 0$$

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

$$(x \cdot z)^d$$
$$\alpha^2$$

- 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

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

$d=2$

$$x_1^2 \quad x_2^2 \quad x_1 x_2$$

- Polynomials of degree up to d

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

$1 \quad x_1 \quad x_2 \quad x_1^2 \quad x_2^2 \quad x_1 x_2$

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

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

$$\equiv \phi(\mathbf{u}) \cdot \phi(\mathbf{v})$$

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


- Sigmoid

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

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 ✓ $K(u,v) = K(v,u) = \phi(v) \cdot \phi(u)$
- K is positive semi-definite, i.e. $\mathbf{x}^T K \mathbf{x} \geq 0$ for all \mathbf{x} 

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?

$$\mathbf{w} \cdot \phi(\mathbf{x}) = \sum_i \alpha_i y_i \underbrace{\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x})}_{\rightarrow K(\mathbf{x}_i, \mathbf{x})}$$

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

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

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

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

$$\phi(\mathbf{x}_i) = \begin{bmatrix} x_i^1 \\ x_i^2 \\ \vdots \\ x_i^d \end{bmatrix}$$

- 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:

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

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

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

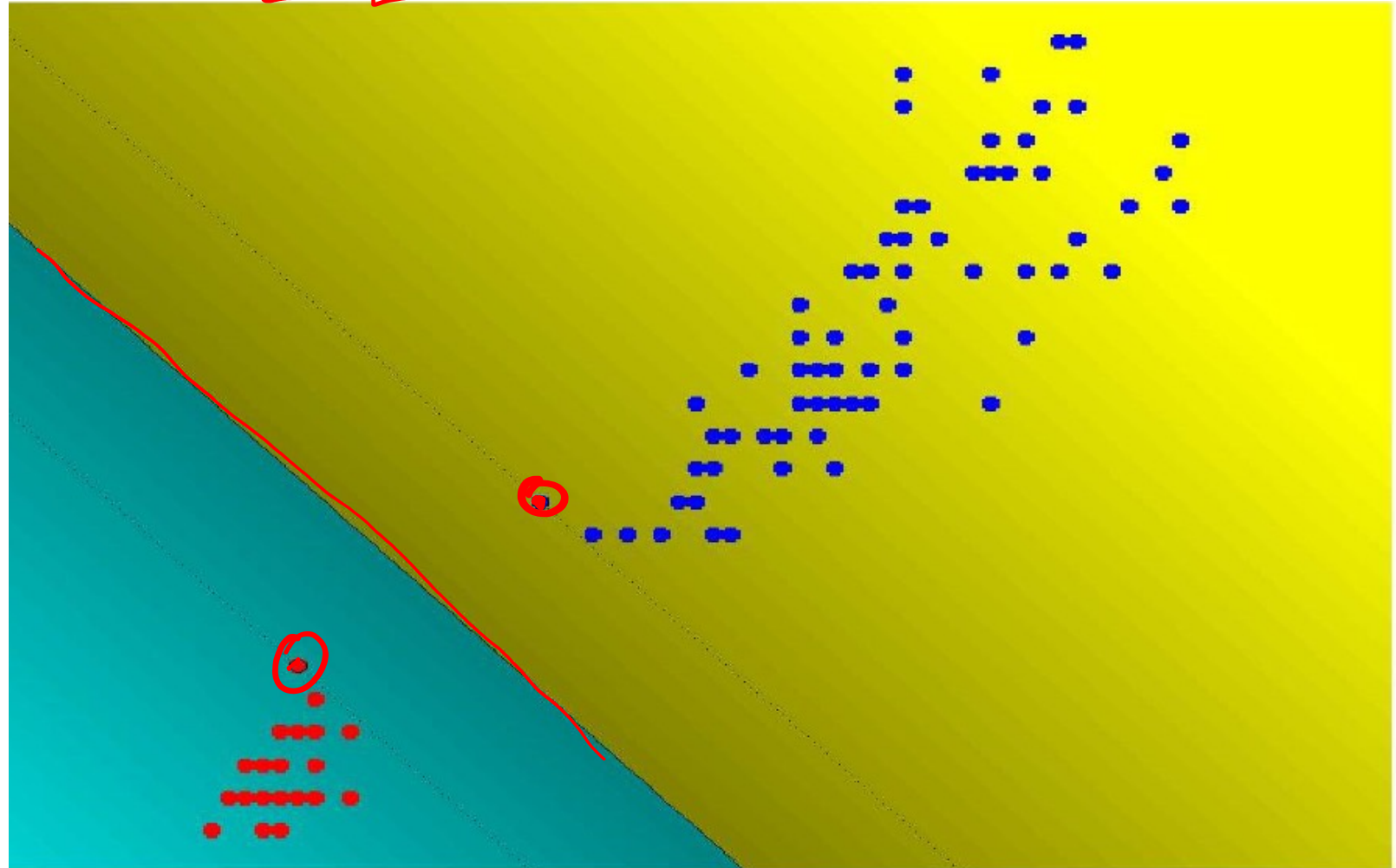
Classify as

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

SVMs with Kernels

- Iris dataset, 2 vs 13, Linear Kernel

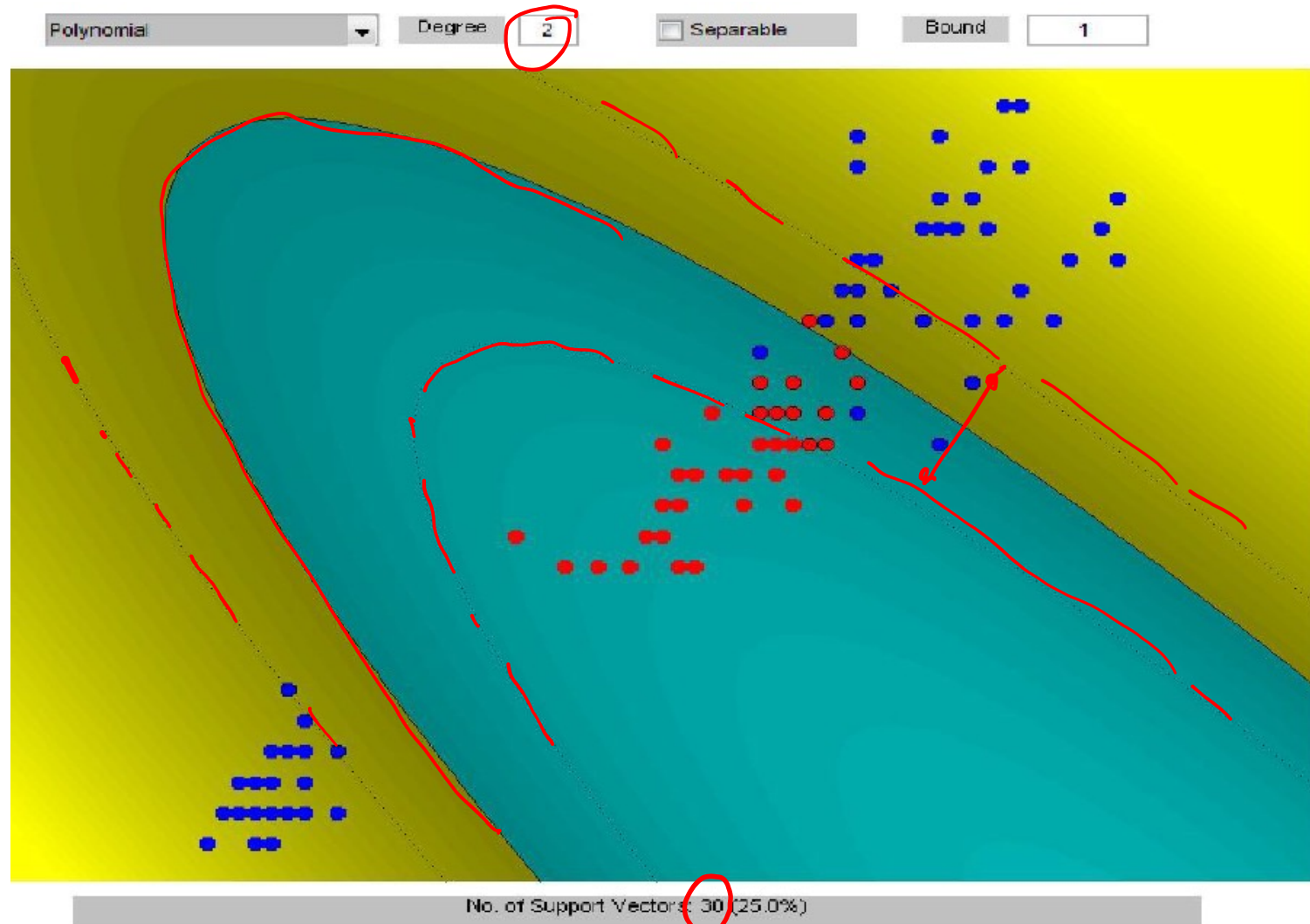
$$K(x_i, x_j) = x_i \cdot x_j$$



No. of Support Vectors (2 / 1.7%)

SVMs with Kernels

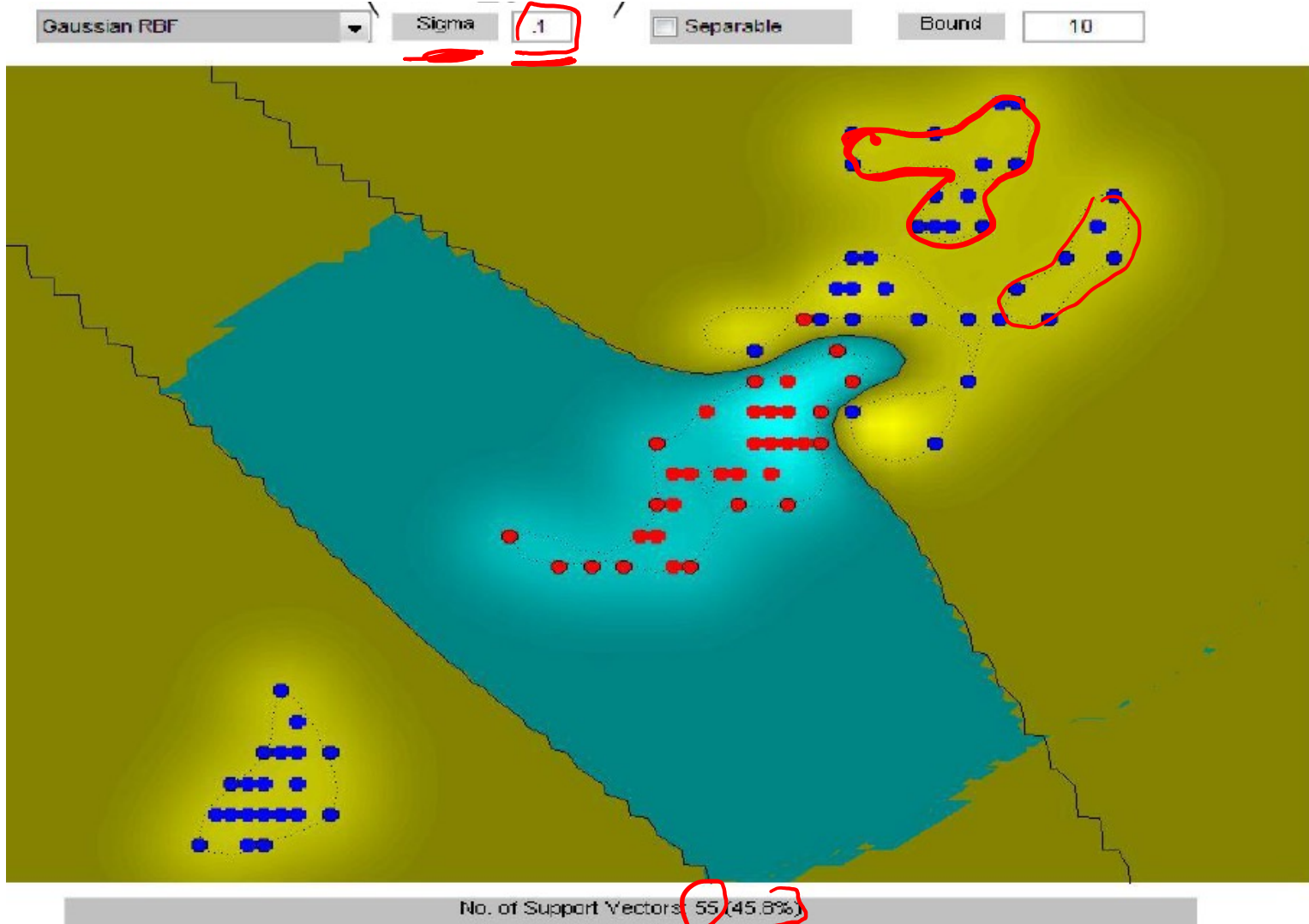
- Iris dataset, 1 vs 23, Polynomial Kernel degree 2



SVMs with Kernels

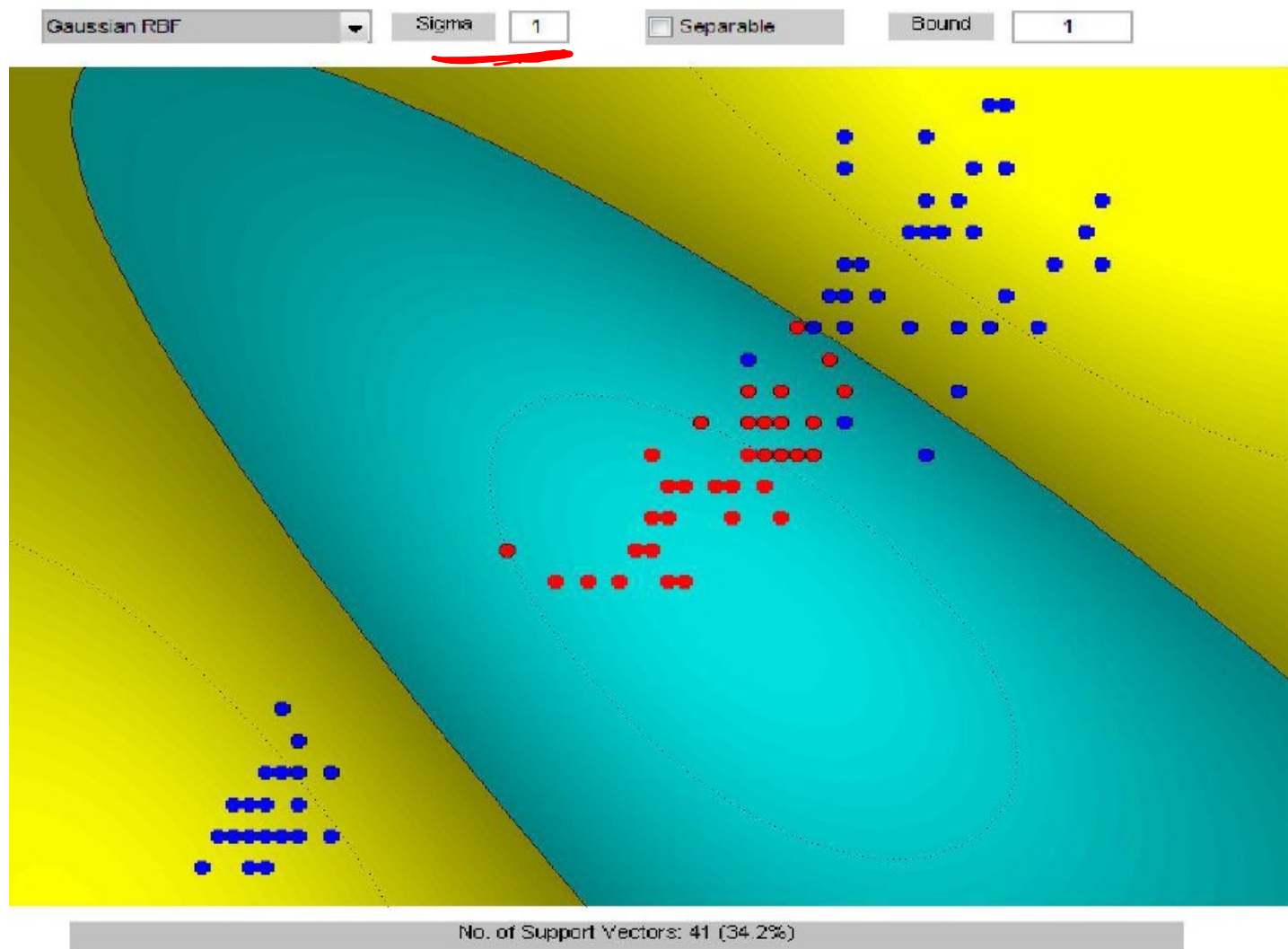
$$K(x, z) = e^{\left(\frac{-\|x - z\|^2}{\sigma^2} \right)}$$

- Iris dataset, 1 vs 23, Gaussian RBF kernel



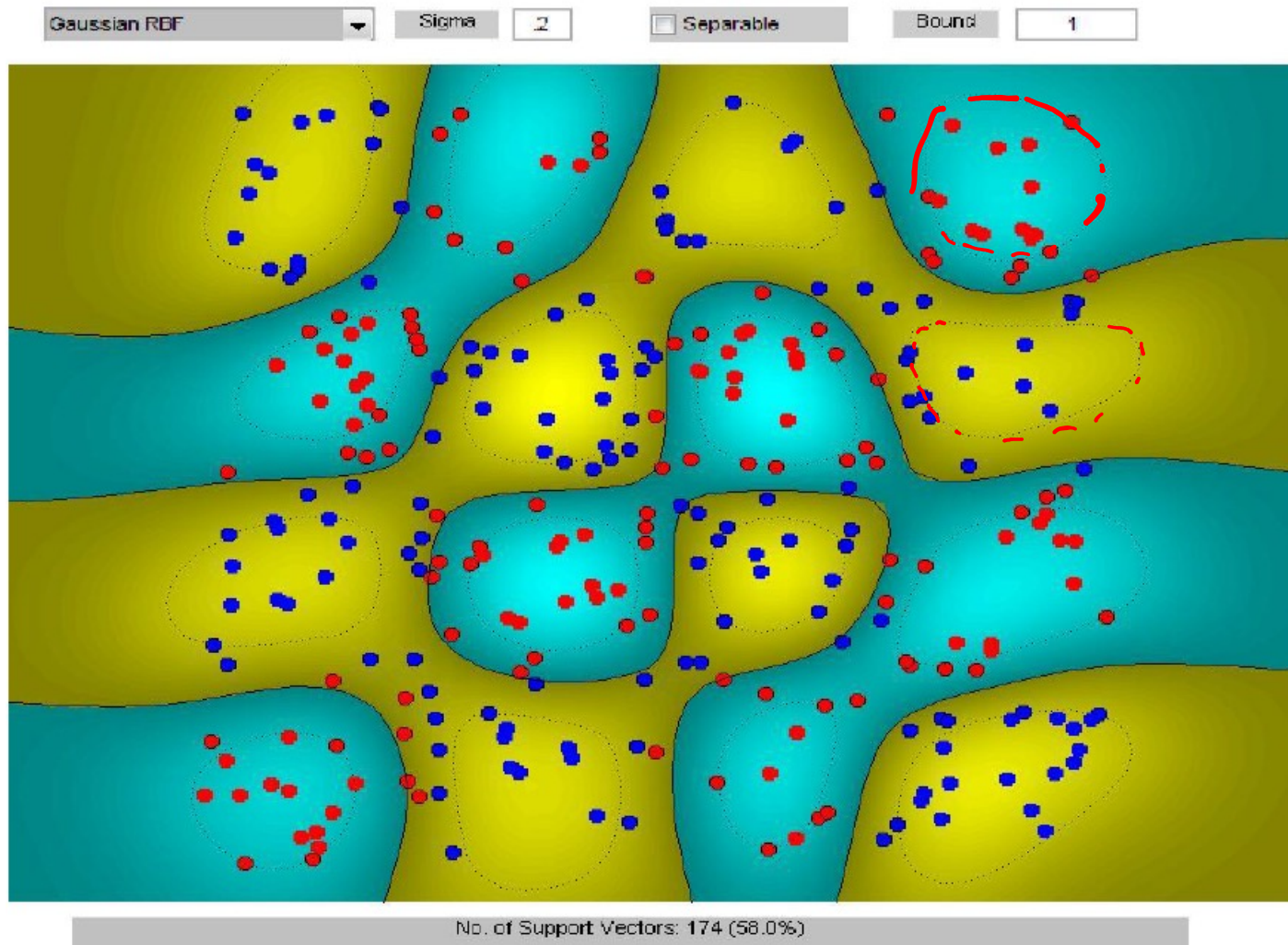
SVMs with Kernels

- Iris dataset, 1 vs 23, Gaussian RBF kernel



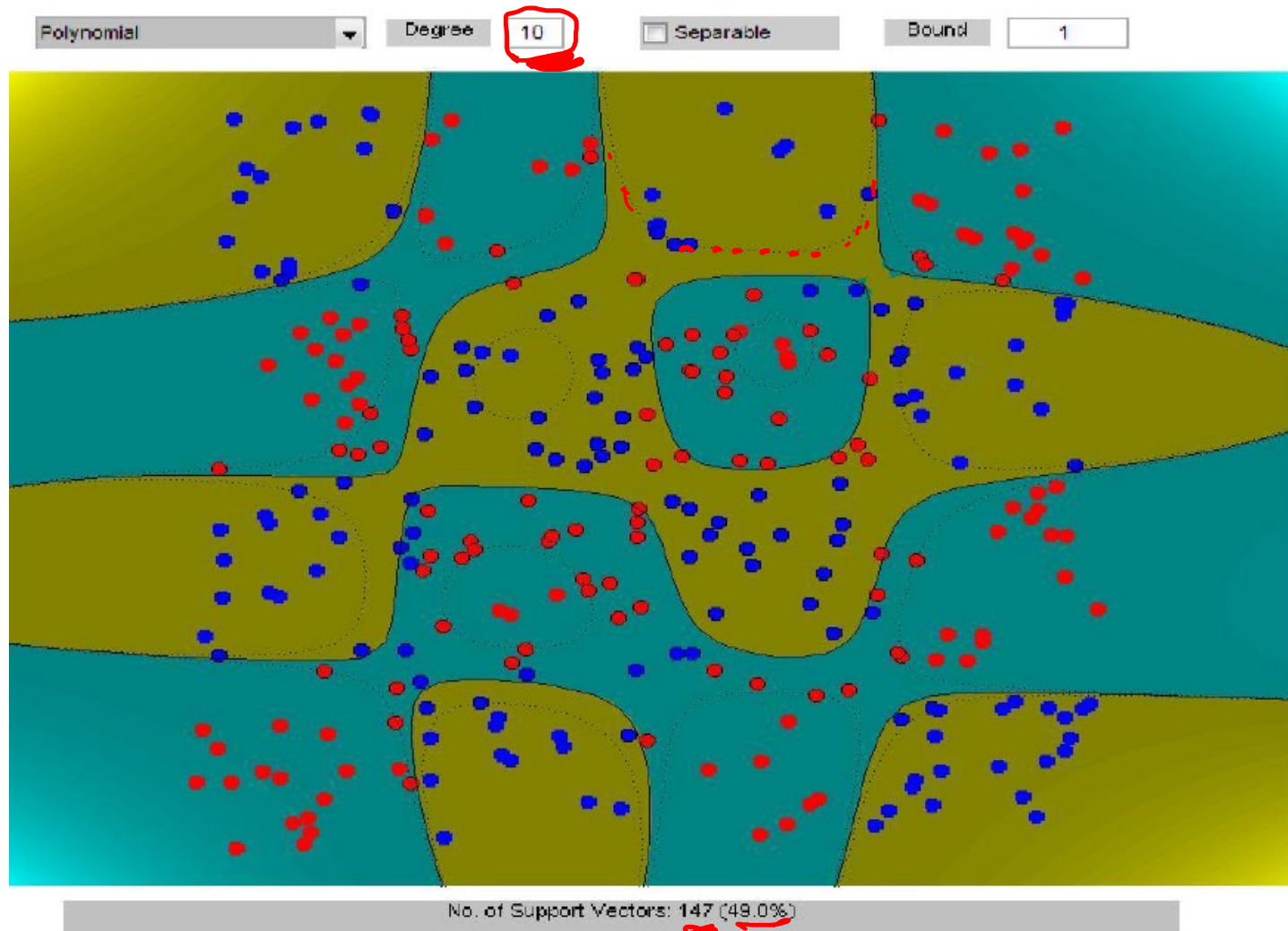
SVMs with Kernels

- Chessboard dataset, Gaussian RBF kernel

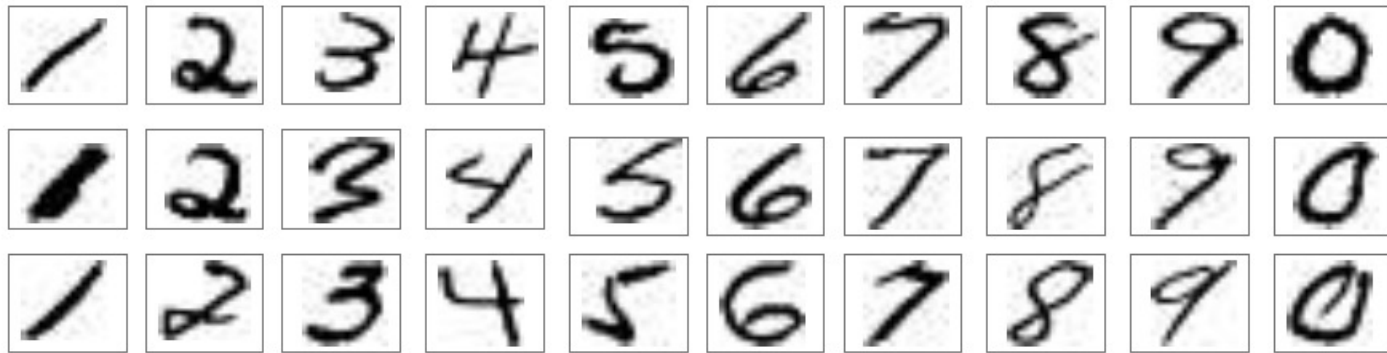


SVMs with Kernels

- Chessboard dataset, Polynomial kernel



USPS Handwritten digits



□ 1000 training and 1000 test instances

Results:

SVM on raw images ~97% accuracy

Kernels in Logistic Regression

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

- Define weights in terms of features:

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

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

- 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! ✓
Solution sparse	Often yes! ✓	Almost always no! ✓
Semantics of output	“Margin” ✓	Real probabilities ✓

Can we kernelize linear regression?

Linear (Ridge) regression

$$\min_{\beta} \sum_{i=1}^n (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2 \quad \hat{\beta} = (\underbrace{\mathbf{A}^T \mathbf{A}}_{p \times p} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y} \leftarrow$$

Recall

$$\mathbf{A} = \begin{bmatrix} \rightarrow X_1 \\ \vdots \\ \leftarrow X_n \end{bmatrix} = \begin{bmatrix} X_1^{(1)} & \dots & X_1^{(p)} \\ \vdots & \ddots & \vdots \\ X_n^{(1)} & \dots & X_n^{(p)} \end{bmatrix}$$

Hence $\mathbf{A}^T \mathbf{A}$ is a p x p matrix whose entries denote the (sample) correlation between the features

NOT inner products between the data points – the inner product matrix would be $\mathbf{A} \mathbf{A}^T$ which is n x n (also known as Gram matrix)

Using dual formulation, we can write the solution in terms of $\mathbf{A} \mathbf{A}^T$

Kernelized ridge regression

$$\hat{\beta} = (\underbrace{\mathbf{A}^T \mathbf{A}}_{\text{feat} \times \text{feat}} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$

$$(y_i - \hat{y}_i)^2$$

Using dual, can re-write solution as:

$$\hat{\beta} = \mathbf{A}^T (\underbrace{\mathbf{A} \mathbf{A}^T}_{n \times n} + \lambda \mathbf{I})^{-1} \mathbf{Y}$$

$$(AA^T)_{ij} = x_i \cdot x_j$$

$$\hat{y} = \underline{X} \hat{\beta} = \underline{X} \mathbf{A}^T (\underline{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{Y}$$

How does this help?

- Only need to invert $n \times n$ matrix (instead of $p \times p$ or $m \times m$)
- More importantly, kernel trick!

$\mathbf{A} \mathbf{A}^T$ involves only inner products between the training points
BUT still have an extra \mathbf{A}^T

$$\begin{aligned} \text{Recall the predicted label is } \hat{f}_n(X) &= \mathbf{X} \hat{\beta} \\ &= \mathbf{X} \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{Y} \end{aligned}$$

$\mathbf{X} \mathbf{A}^T$ contains inner products between test point \mathbf{X} and training points!

Kernelized ridge regression

$$\hat{\beta} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$

$$\hat{f}_n(X) = \underline{\mathbf{X} \hat{\beta}}$$

Using dual, can re-write solution as:

$$\hat{\beta} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{Y}$$

How does this help?

- Only need to invert $n \times n$ matrix (instead of $p \times p$ or $m \times m$)
- More importantly, kernel trick!

$$\underline{\hat{f}_n(X)} = \underline{\mathbf{K}_X (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{Y}} \quad \text{where} \quad \begin{aligned} \mathbf{K}_X(i) &= \phi(X) \cdot \phi(X_i) \\ \mathbf{K}(i, j) &= \phi(X_i) \cdot \phi(X_j) \end{aligned} \quad \leftarrow$$

Work with kernels, never need to write out the high-dim vectors

Ridge Regression with (implicit) nonlinear features $\phi(X)$!

$$f(X) = \phi(X) \beta$$

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
- Dual SVM formulation
 - Easier to solve when dimension high $d > n$
 - Kernel Trick