# Kernel Trick contd…

Aarti Singh

Machine Learning 10-315
Nov 2, 2020

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

*$\alpha_i$ – Lagrange = multiplier*

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

*n-dim problem*

$$\sum_i \alpha_i y_i = 0$$
$$C \geq \alpha_i \geq 0$$

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

$$\left[ K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j) \right]$$

*← high-dim features*

$$\sum_i \alpha_i y_i = 0$$
$$C \geq \alpha_i \geq 0$$

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

# Dot Product of Polynomials

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

$$\mathbf{x} = \left[ \begin{array}{c} x_1 \\ x_2 \end{array} \right] \qquad \mathbf{z} = \left[ \begin{array}{c} z_1 \\ z_2 \end{array} \right]$$

d=1 $\quad \Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) = \left[ \begin{array}{c} x_1 \\ x_2 \end{array} \right] \cdot \left[ \begin{array}{c} z_1 \\ z_2 \end{array} \right] = x_1 z_1 + x_2 z_2 = \mathbf{x} \cdot \mathbf{z}$

d=2 $\quad \Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) = \left[ \begin{array}{c} x_1^2 \\ \sqrt{2} x_1 x_2 \\ x_2^2 \end{array} \right] \cdot \left[ \begin{array}{c} z_1^2 \\ \sqrt{2} z_1 z_2 \\ z_2^2 \end{array} \right]$

$$= \; x_1^2 z_1^2 + x_2^2 z_2^2 + 2 x_1 x_2 z_1 z_2$$

$$= \; (x_1 z_1 + x_2 z_2)^2$$

$$= \; (\mathbf{x} \cdot \mathbf{z})^2$$

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

3

# 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 of exp)

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

- Sigmoid

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

Using kernels, cost of computing dot products depends on dimension of original features $x$, and NOT transformed features $\phi(x)$

# Mercer Kernels

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

$$K(x, x') = \phi(x) \cdot \phi(x')$$
$$= \phi(x') \cdot \phi(x)$$
$$= K(x', x)$$

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 $\mathbf{x}$

$$\bar{\Phi}_{D \times n} = \left[ \begin{array}{cccc} [\phi(x_1)] & [\phi(x_2)] & \cdots \end{array} \right]$$

Ensures optimization is concave maximization

$$x^T K x = x^T \underline{\phi(x) \cdot \phi(x)} x \geq 0 \quad \Rightarrow K_{ij} = \phi(x_i) \cdot \phi(x_j)$$

$$K = \Phi \cdot \Phi = \Phi^T \Phi$$

$$_{n \times n} K = \Phi^T \Phi$$
$$= _{n \times D} D \times n$$

# **Overfitting**

- Huge feature space with kernels, what about overfitting???
    - Maximizing margin leads to sparse set of support vectors     $\alpha_i = 0$
    - 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 **x**, if we need to represent $\Phi(\textbf{x})$, we are in trouble!
- Recall classifier: sign(**w**.$\Phi$(**x**)+b)

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

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

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

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

$$\mathbf{w} \cdot \Phi(x_k)$$

- 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 $\alpha_i$
- At classification time, compute:

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

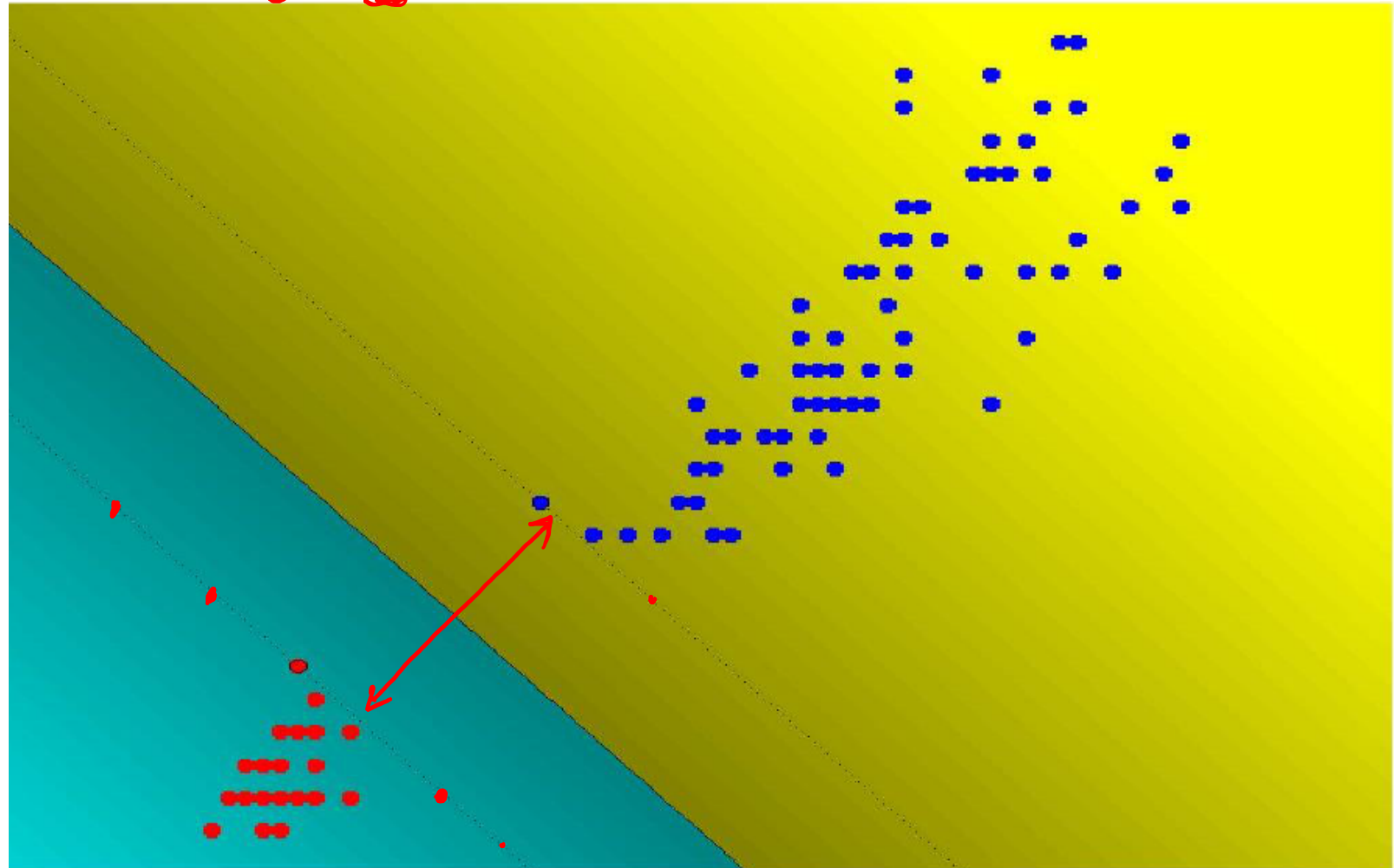$$b = y_k - \sum_i \alpha_i y_i K(\mathbf{x}_k, \mathbf{x}_i)$$

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

**Classify as** $\longrightarrow$ $sign\left(\mathbf{w} \cdot \Phi(\mathbf{x}) + b\right)$
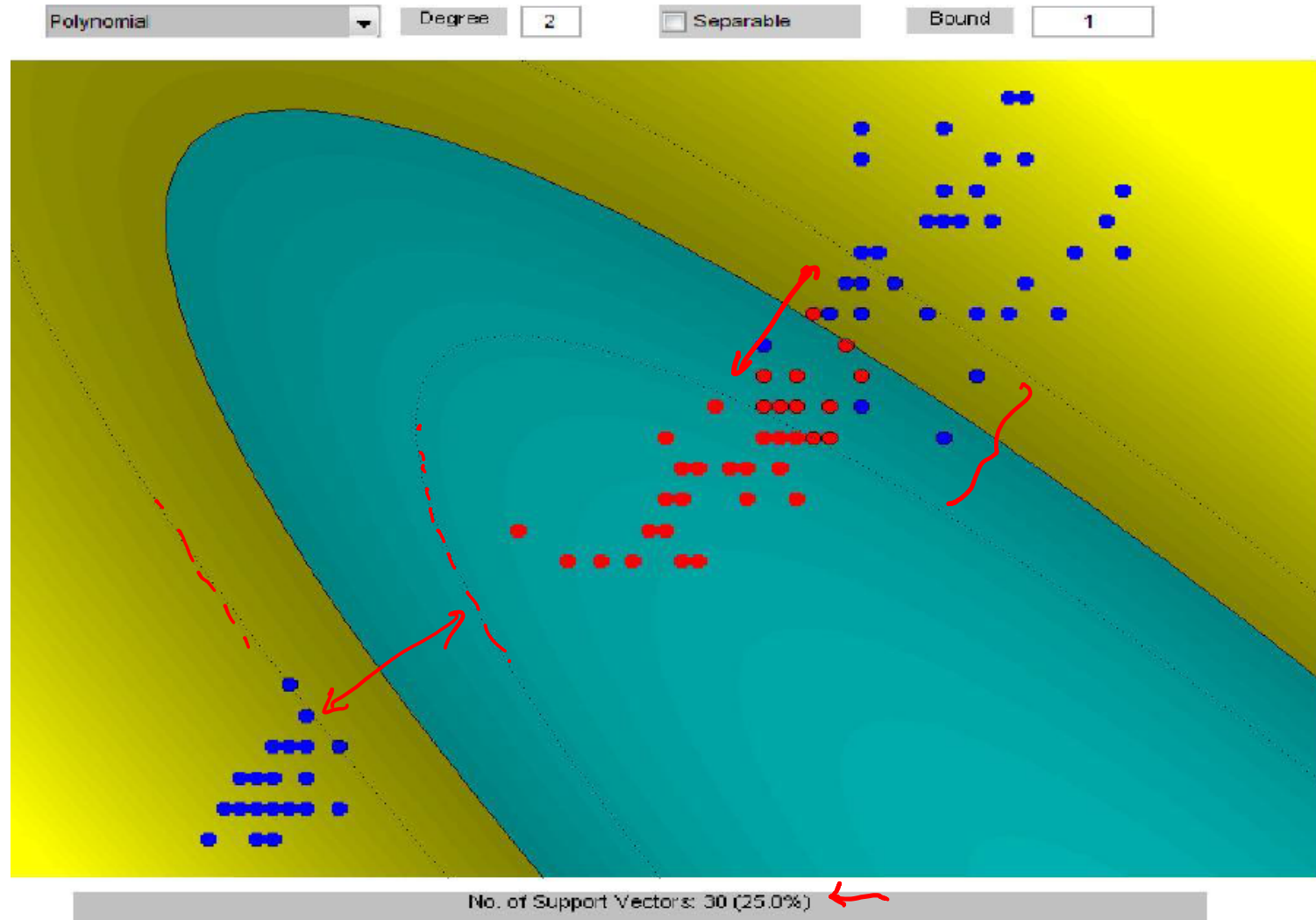
# SVMs with Kernels

- Iris dataset, 2 vs 13, Linear Kernel
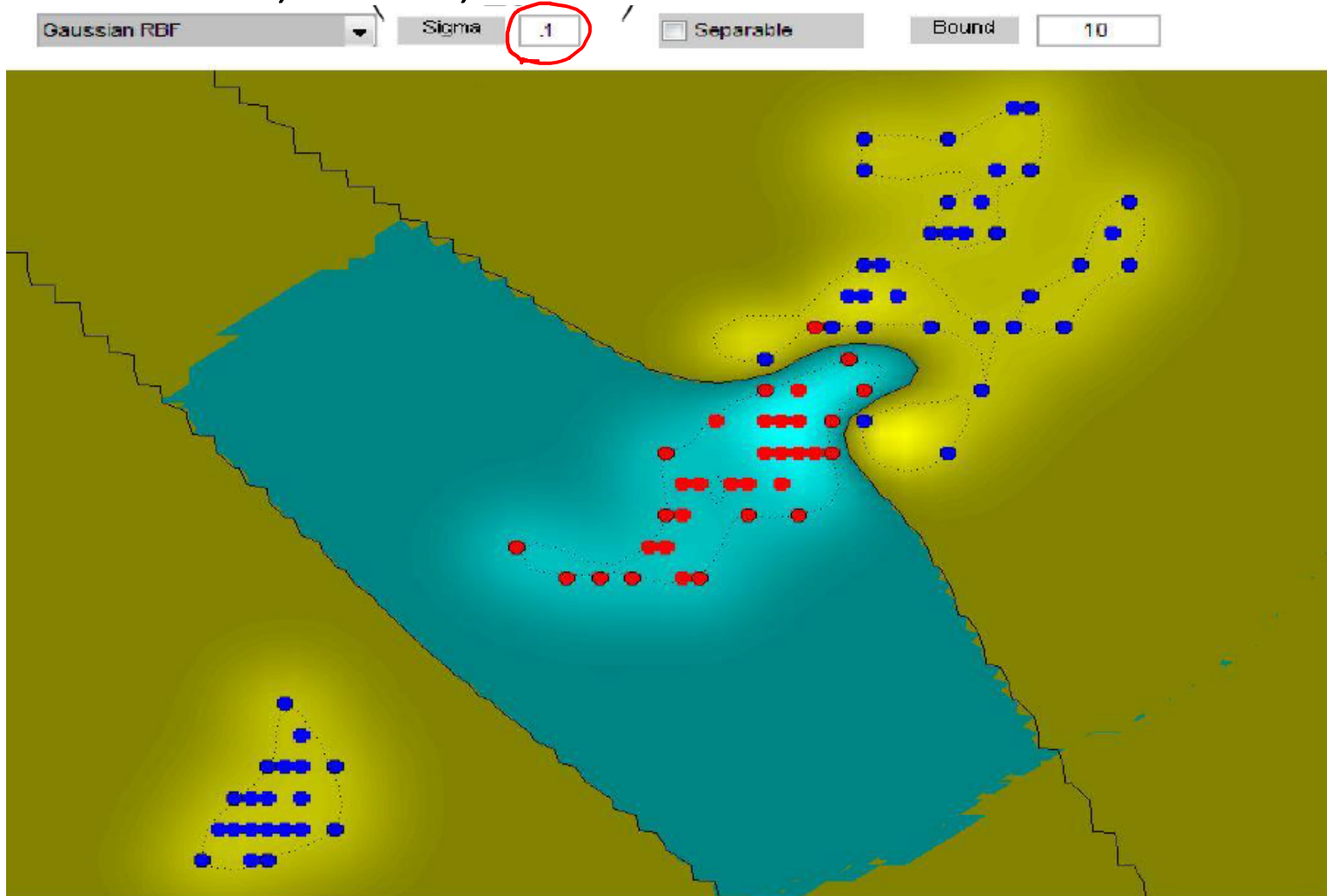


No. of Support Vectors: 2 ( 1.7%)

# SVMs with Kernels

- Iris dataset, 1 vs 23, Polynomial Kernel degree 2
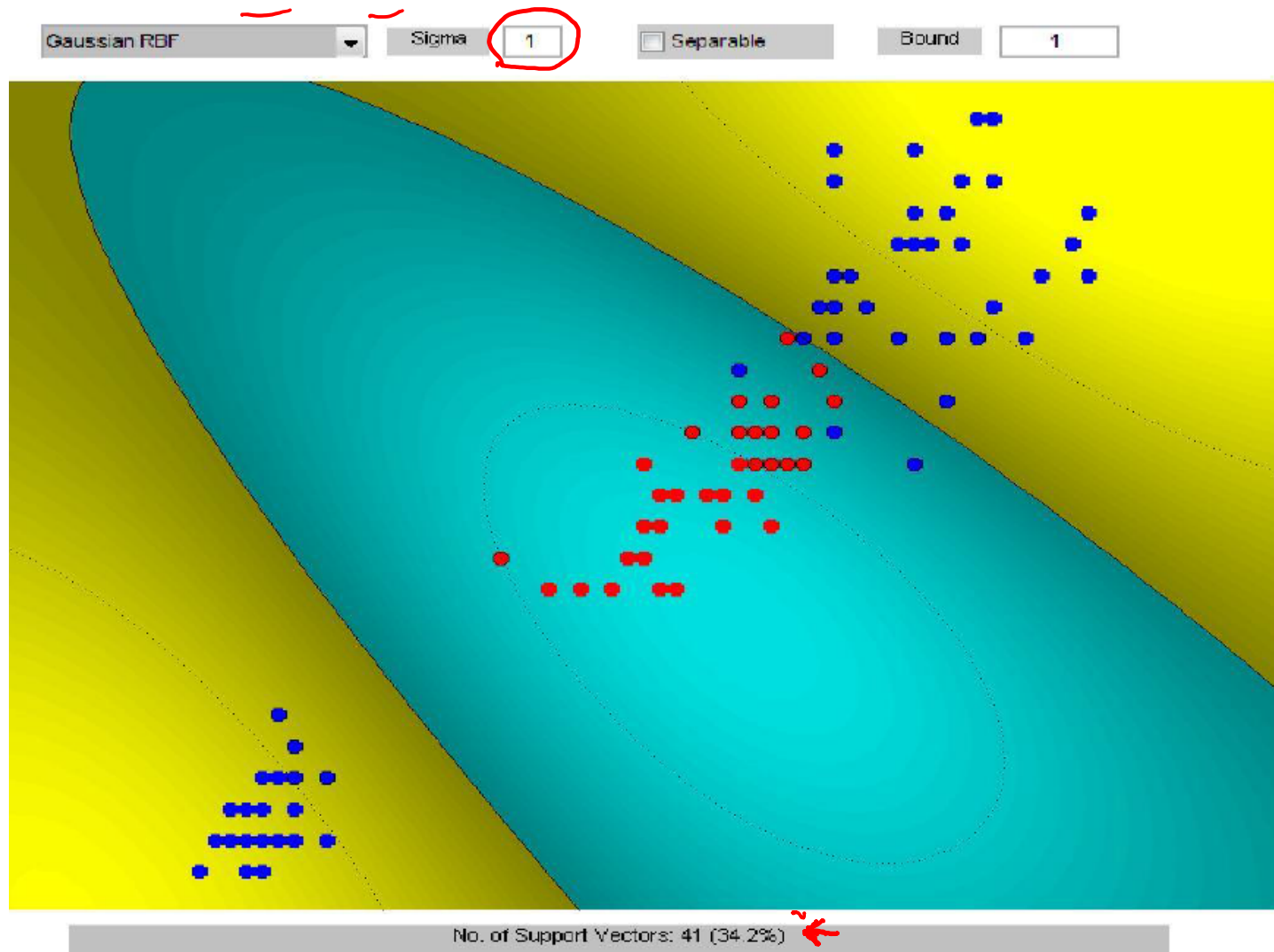
# SVMs with Kernels

- Iris dataset, 1 vs 23, Gaussian RBF kernel



No. of Support Vectors: 55 (45.8%)
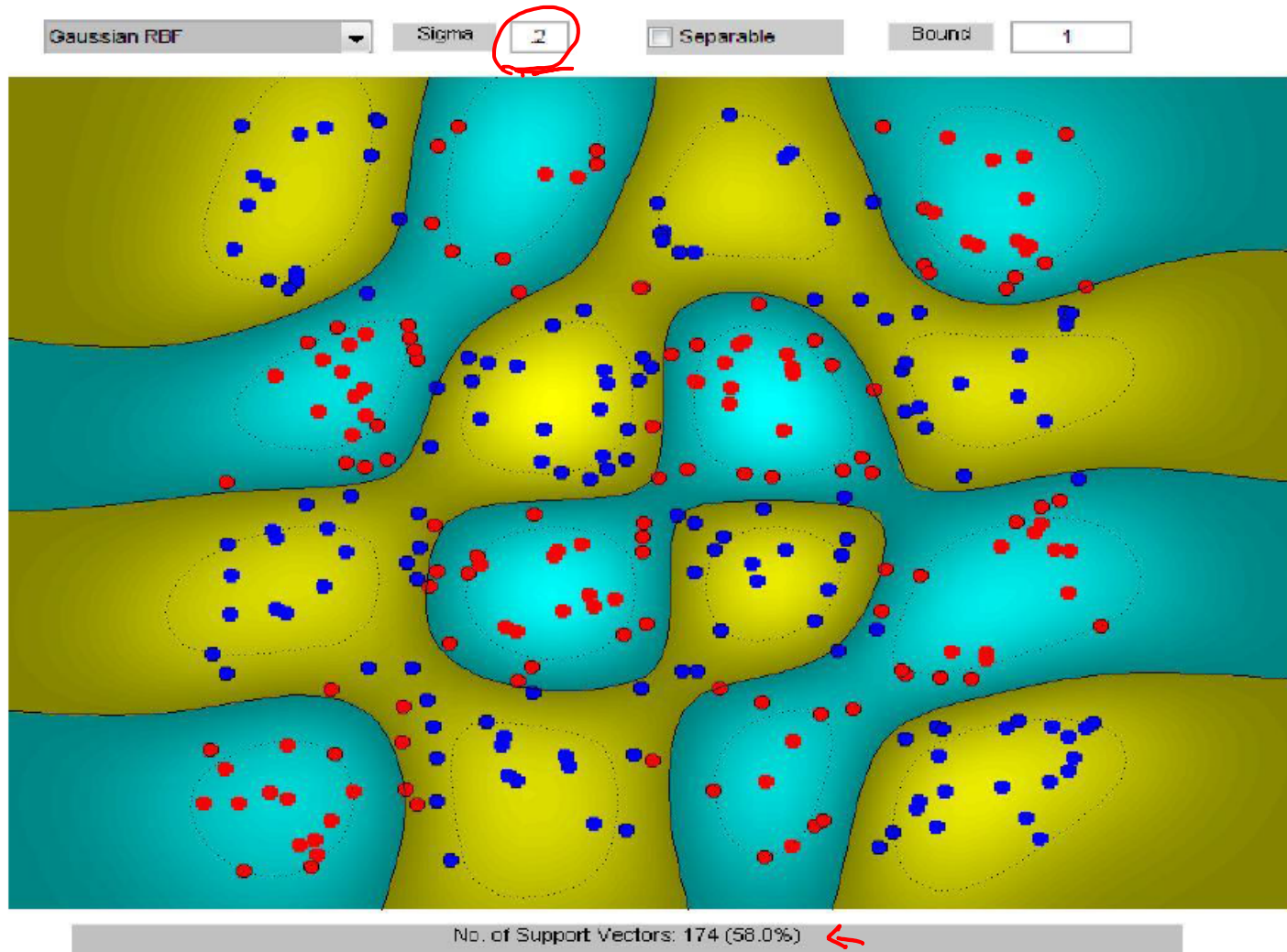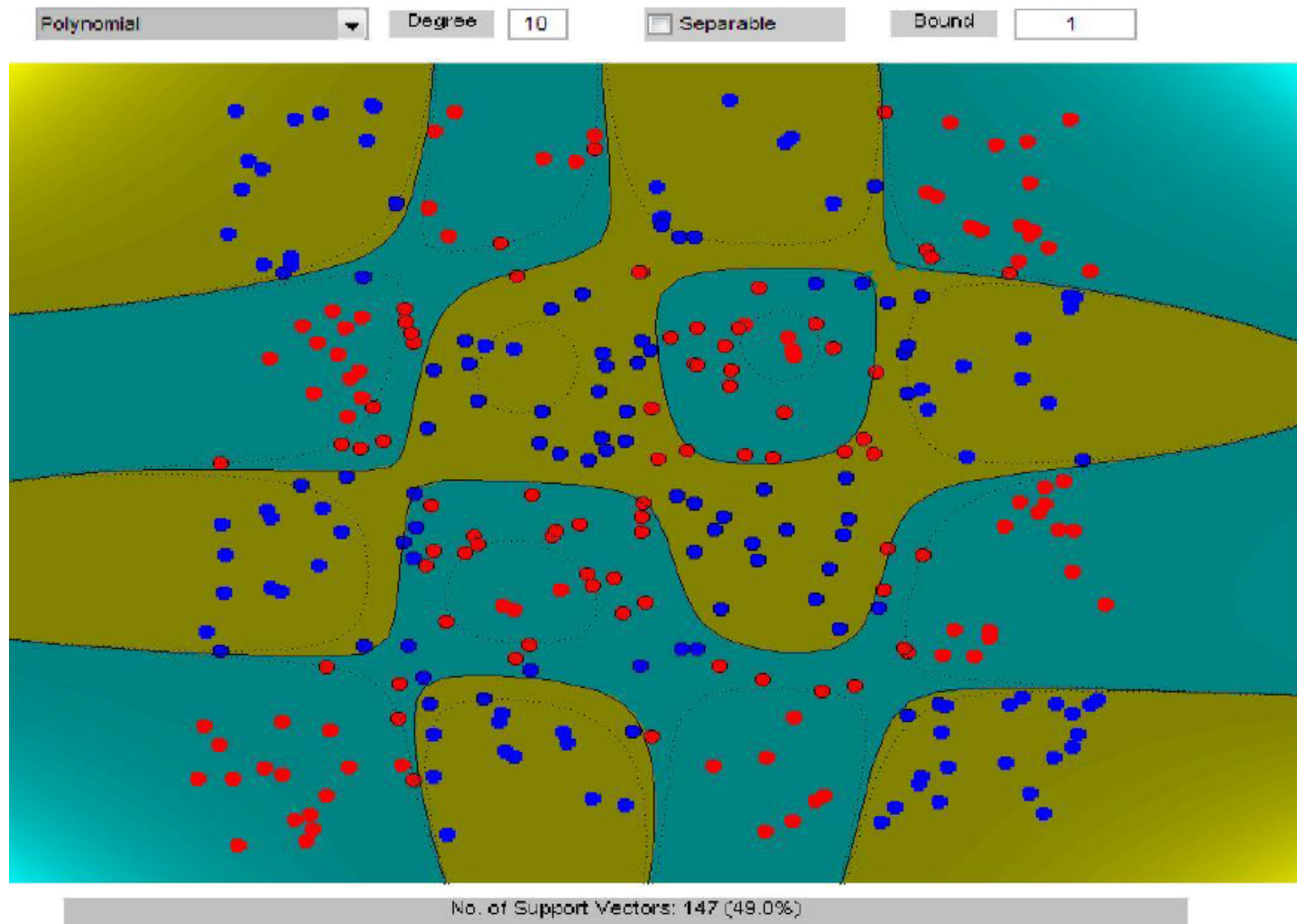
# SVMs with Kernels

- Iris dataset, 1 vs 23, Gaussian RBF kernel

# SVMs with Kernels

- Chessboard dataset, Gaussian RBF kernel

# SVMs with Kernels

- Chessboard dataset, Polynomial kernel



No. of Support Vectors: 147 (49.0%)

# 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 |

M(C) LE

**Log loss**   **Hinge loss**

**0-1 loss**

-1    0    1

$$(\mathbf{w} \cdot x_j + b)y_j$$

# SVMs vs. Logistic Regression

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

# Kernels in Logistic Regression

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

- Define weights in terms of features:

$$w = \sum_i \alpha_i y_i \phi(x_i)$$

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

$i \leftarrow$ data points

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

$$= \frac{1}{1 + e^{-(\sum_i \alpha_i K(\mathbf{x}, \mathbf{x}_i) + b)}}$$

- Derive simple gradient descent rule on $\alpha_i$

# SVMs vs. Logistic Regression

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

# 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! |
| | | |

# 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 |

21

# Can we kernelize linear regression?

# Linear (Ridge) regression

$$\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i\beta)^2 + \lambda\|\beta\|_2^2$$

$$x_i \cdot x_j = \sum_{k=1}^{p} x_i^{(k)} x_j^{(k)} \quad \overset{\phi(x_i) \cdot \phi(x_j)}{}$$

$$\widehat{\beta} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y} \qquad \hat{f}_n(X) = X\hat{\beta}$$

Recall

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

$$\sum_{k=1}^{n} x_k^{(i)} x_k^{(j)}$$

**A$^T$A** is a p x p matrix whose entries denote the (sample) correlation between the features

$$(A^TA)_{ij} = \begin{bmatrix} x_1^{(i)} & \cdots & x_n^{(i)} \end{bmatrix} \begin{bmatrix} x_1^{(j)} \\ \vdots \\ x_n^{(j)} \end{bmatrix}$$

NOT inner products between the data points – the inner product matrix would be **AA$^T$** which is n x n (also known as Gram matrix)

# Ridge regression (dual)

$$\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i\beta)^2 + \lambda\|\beta\|_2^2$$

$\phi(X_i)$

$$\hat{f}_n(X) = \sum_i \hat{\alpha}_i \Phi(X) \cdot \Phi(X_i)$$

$\Phi(X)\hat{\beta}$

$K(X_i, X_i)$

- Define weights in terms of features: $\beta = \sum_i \alpha_i \Phi(X_i)$

$$\min_{\boldsymbol{\alpha}} \sum_{i=1}^{n} (Y_i - \sum_j \alpha_j \Phi(X_i) \cdot \Phi(X_j))^2 + \lambda \sum_{ij} \alpha_i \alpha_j \Phi(X_i) \cdot \Phi(X_j)$$

$K(X_i, X_j)$

$K(X_i, X_j)$

$$(Y_{n\times 1} - \underset{n\times n}{K} \underset{n\times 1}{\alpha})^\top (Y - K\alpha)$$

$$\underset{1\times n}{\alpha^\top} \underset{n\times n}{K} \underset{n\times 1}{\alpha}$$

$$\min_{\boldsymbol{\alpha}} (\mathbf{Y} - \mathbf{K}\boldsymbol{\alpha})^\top (\mathbf{Y} - \mathbf{K}\boldsymbol{\alpha}) + \lambda \boldsymbol{\alpha}^\top \mathbf{K}\boldsymbol{\alpha}$$

$$Y^\top Y + \alpha^\top K^\top K \alpha - 2\alpha^\top K^\top Y + \lambda \alpha^\top K \alpha$$

$$2K^2\alpha - 2K^\top Y + 2\lambda K \alpha = 0$$

$$\hat{\alpha} = (K + \lambda I)^{-1} Y$$

$$K(K + \lambda I)\alpha = KY$$

# Kernel ridge regression

$$\hat{f}_n(X) = \sum_i \hat{\alpha}_i K(X, X_i) = \mathbf{K}_X \hat{\boldsymbol{\alpha}}$$

$1 \times n$   $n \times 1$

where $\quad \hat{\boldsymbol{\alpha}} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{Y}$

$K_X$ vector $1 \times n$

$$\mathbf{K}_X(i) = \Phi(X) \cdot \Phi(X_i)$$

$$\mathbf{K}(i, j) = \Phi(X_i) \cdot \Phi(X_j)$$

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

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

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

# Kernel ridge regression vs. (local) Kernel Regression

$$\hat{f}_n(X) = \sum_i \hat{\alpha}_i K(X, X_i)$$

$$\sum_i w_i (f(x_i) - y_i)^2$$

$$\sum_i w_i y_i = \sum_i \frac{K(x_i, x_i) y_i}{\sum_j K(x_j, x_j)}$$

## Kernel Ridge Regression

$$\hat{\boldsymbol{\alpha}} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{Y}$$

$a$ — training points    $0$    $K_{ij} = K(X_i, X_j)$

## (Local) Kernel Regression

$$\hat{\alpha}_i = \frac{Y_i}{\sum_i K(X, X_i)} = (\mathbf{1}^\top \mathbf{K}_X)^{-1} \mathbf{Y}$$

Weights depend on test point $X$

Global fit

Interpret as weighted Nonlinear features

Local fit

Interpret as weighted Least Squares

# What you need to know

- Maximizing margin
- Derivation of SVM formulation
- Slack variables and hinge loss
- Tackling multiple class
  - One against All
  - Multiclass SVMs
- Dual SVM formulation
  - Easier to solve when dimension high d > n
  - Kernel Trick
- Relationship between SVMs and logistic regression
- Kernelizing linear regression e.g. Kernel Ridge Regression