# Support Vector Machines - Dual formulation and Kernel Trick

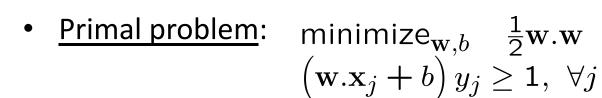
Aarti Singh

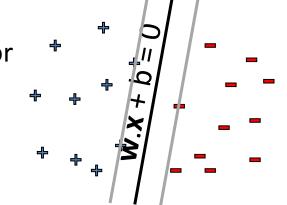
Machine Learning 10-315 Nov 1, 2021



# **SVM** – linearly separable case

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





#### 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**

maximize
$$_{\alpha}$$
  $\sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} \mathbf{x}_{i} \cdot \mathbf{x}_{j}$   $\sum_{i} \alpha_{i} y_{i} = 0$   $\alpha_{i} \geq 0$ 

Dual problem is also QP Solution gives  $\alpha_j s$ 

Use any one of support vectors with  $\alpha_k>0$  to compute b since constraint is tight  $(w.x_k + b)y_k = 1$ 

$$\mathbf{w} = \sum_{i} \alpha_i y_i \mathbf{x}_i$$

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

for any k where  $\alpha_k > 0$ 

# **Dual SVM – non-separable case**

Primal problem:

$$\begin{aligned} & \text{minimize}_{\mathbf{w},b,\{\xi_j\}} \frac{1}{2} \mathbf{w}.\mathbf{w} + C \sum_{j} \xi_j \\ & \left( \mathbf{w}.\mathbf{x}_j + b \right) y_j \geq 1 - \xi_j, \ \forall j \\ & \xi_j \geq 0, \ \forall j \end{aligned}$$

Lagrange Multipliers

• Dual problem:

$$\begin{aligned} \max_{\alpha,\mu} \min_{\mathbf{w},b,\{\xi_{\mathbf{j}}\}} L(\mathbf{w},b,\xi,\alpha,\mu) \\ s.t.\alpha_{j} &\geq \mathbf{0} \quad \forall j \\ \mu_{j} &\geq \mathbf{0} \quad \forall j \end{aligned}$$

# **Dual SVM – non-separable case**

$$\begin{aligned} \text{maximize}_{\alpha} \quad & \sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} \mathbf{x}_{i}. \mathbf{x}_{j} \\ & \sum_{i} \alpha_{i} y_{i} = \mathbf{0} \\ & C \geq \alpha_{i} \geq \mathbf{0} \end{aligned}$$
 
$$\text{comes from } \frac{\partial L}{\partial \xi} = \mathbf{0} \qquad \begin{aligned} & \underbrace{\text{Intuition:}}_{\text{If } C \rightarrow \infty \text{, recover hard-margin SVM}} \end{aligned}$$

Dual problem is also QP Solution gives  $\alpha_i$ s

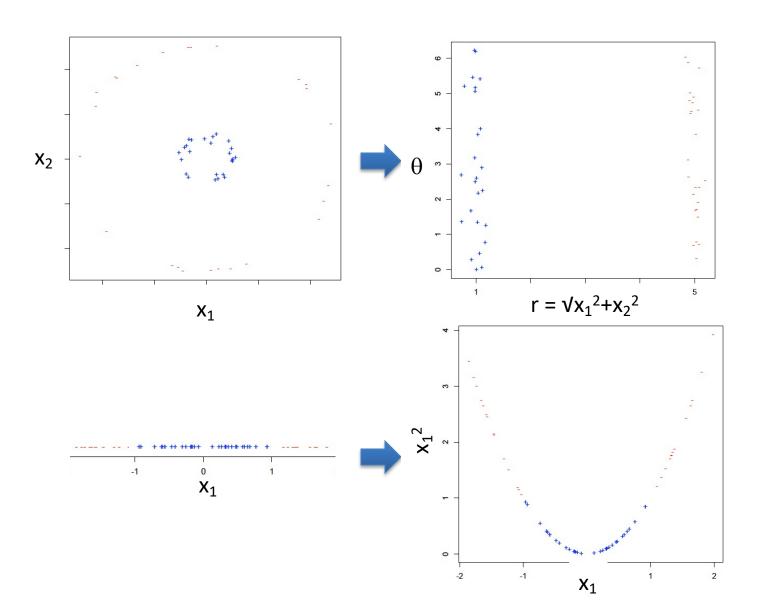
$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$
 
$$b = y_k - \mathbf{w}.\mathbf{x}_k$$
 for any  $k$  where  $C > \alpha_k > 0$ 

# 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>>n)

But, more importantly, the "kernel trick"!!!

# Separable using higher-order features



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

$$\begin{aligned} \text{maximize}_{\alpha} & \sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} \mathbf{x}_{i}. \mathbf{x}_{j} \\ & \sum_{i} \alpha_{i} y_{i} = 0 \\ & C \geq \alpha_{i} \geq 0 \end{aligned}$$

$$\mathbf{maximize}_{\alpha} & \sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} K(\mathbf{x}_{i}, \mathbf{x}_{j}) \\ & K(\mathbf{x}_{i}, \mathbf{x}_{j}) = \Phi(\mathbf{x}_{i}) \cdot \Phi(\mathbf{x}_{j}) \\ & \sum_{i} \alpha_{i} y_{i} = 0 \\ & C \geq \alpha_{i} \geq 0 \end{aligned}$$

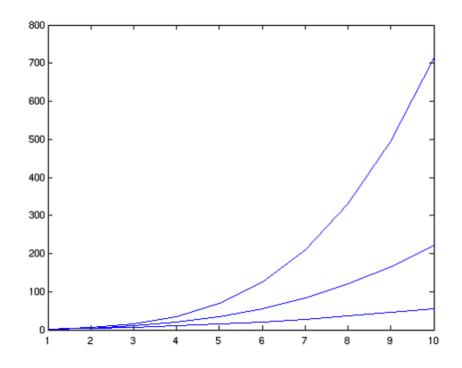
 $\Phi(\mathbf{x})$  – High-dimensional feature space, but never need it explicitly as long as we can compute the dot product fast using some Kernel K

# Polynomial features $\phi(x)$

m – input features

d – degree of polynomial

num. terms 
$$= \begin{pmatrix} d+m-1 \\ d \end{pmatrix} = \frac{(d+m-1)!}{d!(m-1)!} \sim m^d$$



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

## **Dot Product of Polynomial features**

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

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

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

$$d=2 \ \Phi(\mathbf{x}) \cdot \Phi(\mathbf{z}) = \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix} \cdot \begin{bmatrix} z_1^2 \\ \sqrt{2}z_1z_2 \\ z_2^2 \end{bmatrix} = x_1^2z_1^2 + x_2^2z_2^2 + 2x_1x_2z_1z_2$$
$$= (x_1z_1 + x_2z_2)^2$$
$$= (\mathbf{x} \cdot \mathbf{z})^2$$

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

#### The Kernel Trick!

maximize<sub>$$\alpha$$</sub>  $\sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} K(\mathbf{x}_{i}, \mathbf{x}_{j})$ 

$$K(\mathbf{x}_{i}, \mathbf{x}_{j}) = \Phi(\mathbf{x}_{i}) \cdot \Phi(\mathbf{x}_{j})$$

$$\sum_{i} \alpha_{i} y_{i} = 0$$

$$C \geq \alpha_{i} \geq 0$$

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

#### **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)

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

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  $\varphi(\mathbf{x})$ ?

Answer: **Mercer kernels** K

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

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

$$\mathbf{w} = \sum_i lpha_i y_i \Phi(\mathbf{x}_i)$$
  $b = y_k - \mathbf{w}.\Phi(\mathbf{x}_k)$  for any  $k$  where  $C > lpha_k > 0$ 

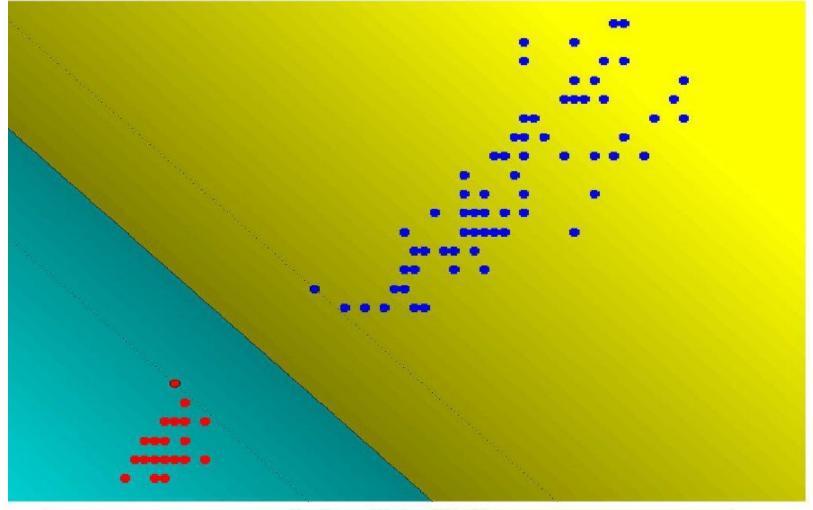
Using kernels we are cool!

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

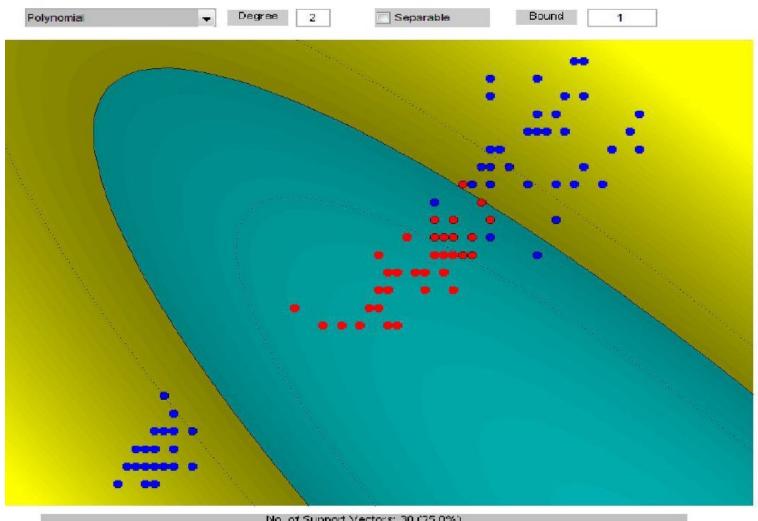
- Choose a set of features and kernel function
- Solve dual problem to obtain support vectors  $\alpha_i$
- At classification time, compute:

$$\begin{aligned} \mathbf{w} \cdot \Phi(\mathbf{x}) &= \sum_i \alpha_i y_i K(\mathbf{x}, \mathbf{x}_i) \\ b &= y_k - \sum_i \alpha_i y_i K(\mathbf{x}_k, \mathbf{x}_i) \\ \text{for any } k \text{ where } C > \alpha_k > 0 \end{aligned} \qquad \text{Classify as} \qquad sign\left(\mathbf{w} \cdot \Phi(\mathbf{x}) + b\right)$$

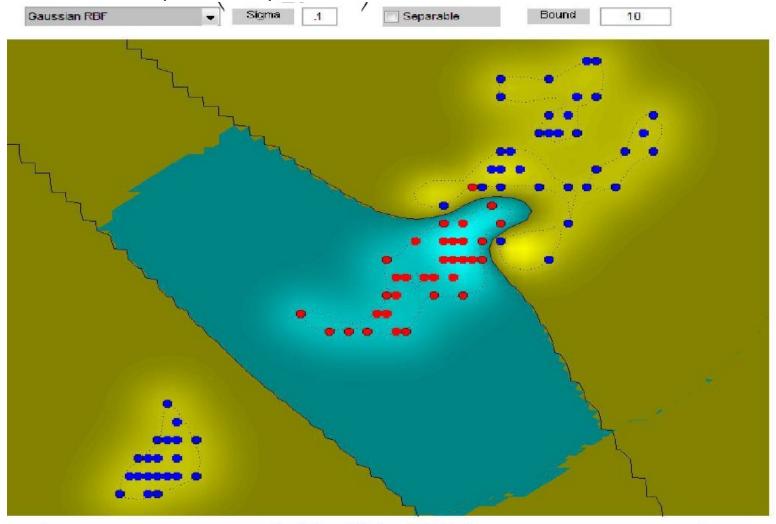
Iris dataset, 2 vs 13, Linear Kernel



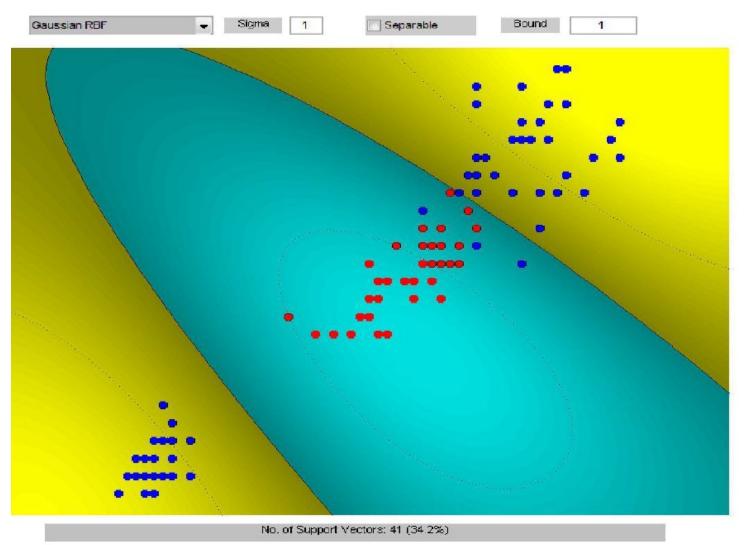
Iris dataset, 1 vs 23, Polynomial Kernel degree 2



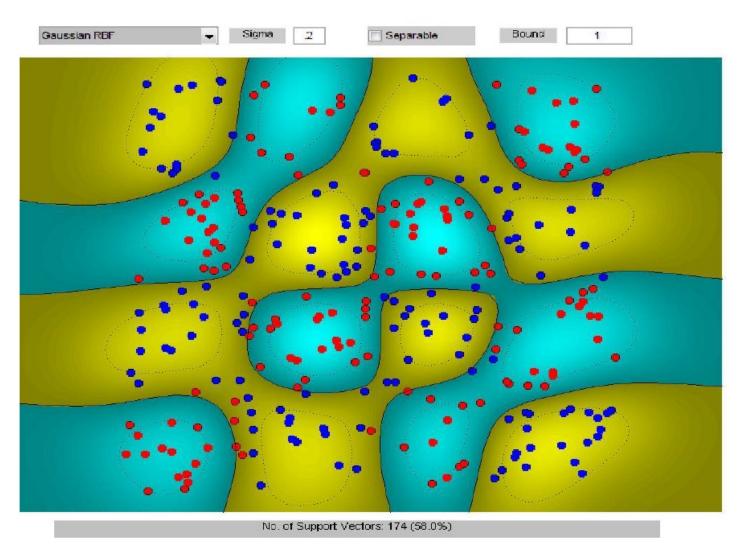
• Iris dataset, 1 vs 23, Gaussian RBF kernel



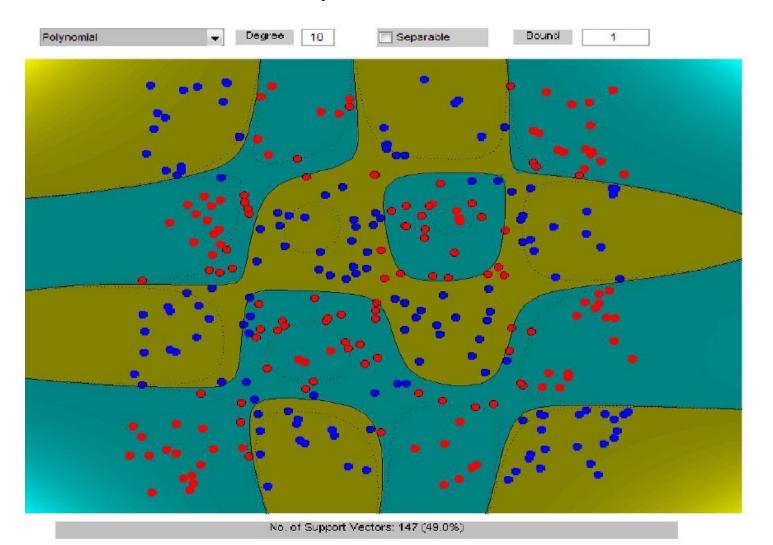
Iris dataset, 1 vs 23, Gaussian RBF kernel



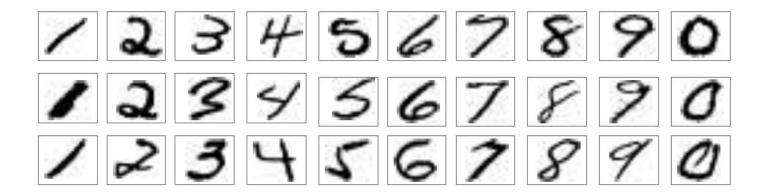
Chessboard dataset, Gaussian RBF kernel



Chessboard dataset, Polynomial kernel



# **USPS Handwritten digits**



■ 1000 training and 1000 test instances

#### Results:

**SVM** on raw images ~97% accuracy

# **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:

$$\mathbf{w} = \sum_{i} \alpha_{i} \Phi(\mathbf{x}_{i}) \, \mathbf{y}_{i}$$

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

$$= \frac{1}{1 + e^{-(\sum_{i} \alpha_{i} K(\mathbf{x}, \mathbf{x}_{i}) + b)}}$$

• Derive simple gradient descent rule on  $\alpha_{\rm i}$ 

# **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

$$\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2 \qquad \widehat{\beta} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$

Recall

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

Hence  $A^TA$  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  $\mathbf{A}\mathbf{A}^{\mathsf{T}}$  which is n x n (also known as Gram matrix)

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

# Kernelized ridge regression

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

Using dual, can re-write solution as:

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

How does this help?

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

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

Recall the predicted label is 
$$\widehat{f}_n(X) = \mathbf{X}\widehat{\beta}$$
 
$$= \mathbf{X}\mathbf{A}^T(\mathbf{A}\mathbf{A}^T + \lambda\mathbf{I})^{-1}\mathbf{Y}$$

**XA**<sup>T</sup> contains inner products between test point **X** and training points!

# Kernelized ridge regression

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

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

Using dual, can re-write solution as:

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

How does this help?

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

$$\widehat{f}_n(X) = \mathbf{K}_X(\mathbf{K} + \lambda \mathbf{I})^{-1}\mathbf{Y}$$
 where  $\mathbf{K}_X(i) = \phi(X) \cdot \phi(X_i)$   
 $\mathbf{K}(i,j) = \phi(X_i) \cdot \phi(X_j)$ 

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

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

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