Support Vector Machines - Dual formulation and Kernel Trick

Aarti Singh

Machine Learning 10-315 Nov 1, 2021



SVM – linearly separable case

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

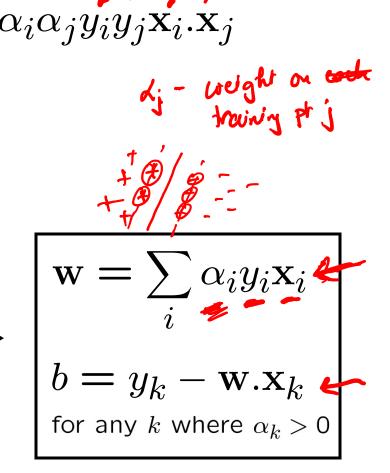
w - weights on features (d-dim problem)

- 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

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 $(w.x_k + b)y_k = 1$,



Dual SVM – non-separable case

Primal problem:

minimize
$$\mathbf{w}, b, \{\xi_j\} \frac{1}{2} \mathbf{w}. \mathbf{w} + C \sum_j \xi_j$$
 (w.x,+b) $y_j = 1 - \xi_j$, $\forall j \leftarrow \mathbf{n}$ α_j $\xi_j \geq 0$, $\forall j \leftarrow \mathbf{n}$ α_j

Dual problem:

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

$$s.t.\alpha_j \geq 0 \quad \forall j - \underbrace{\mathbf{w},\mathbf{w}}_{2} + \underbrace{\mathbf{c}}_{2},\underbrace{\mathbf{c}}_{3},\underbrace{\mathbf{c}}_{4},\underbrace{\mathbf{c}}_{2},\underbrace{\mathbf{c}}_{3},\underbrace{\mathbf{c}}_{3},\underbrace{\mathbf{c}}_{4},\underbrace{\mathbf{c}}_{3},\underbrace{\mathbf{c}$$

d(x,u)

Lagrange

Dual SVM – non-separable case

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

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

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

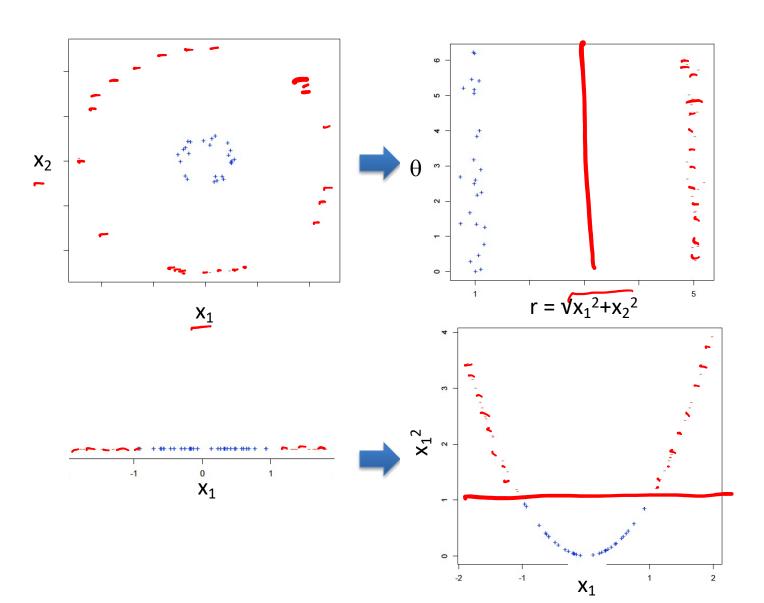
$$b = y_k - \mathbf{w}.\mathbf{x}_k \checkmark$$
 for any k where $C > \alpha_k > 0$

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!

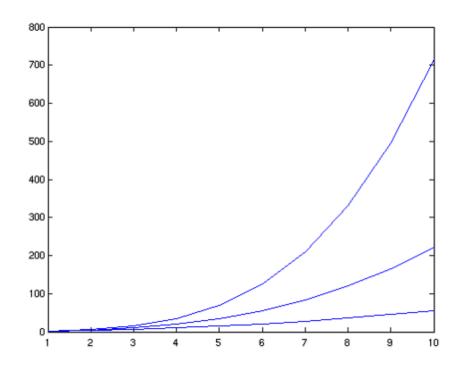
 $\Phi(\mathbf{x})$ – High-dimensional feature space, but never need it explicitly as long as we can compute the dot product fast using some Kernel K

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

$$\begin{pmatrix} \phi(x) \\ \chi_1 \\ \chi_2 \end{pmatrix} \rightarrow \begin{pmatrix} \chi_1 \\ \chi_2 \\ \chi_1 \\ \chi_1$$

Dot Product of Polynomial features

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

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

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

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

$$= (x_1z_1 + x_2z_2)^2 = (\mathbf{x} \cdot \mathbf{z})^2 \quad 3 \text{ multiplical in substitution}$$

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

The Kernel Trick!



maximize_{$$\alpha$$} $\sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} y_{i} y_{j} K(\mathbf{x}_{i}, \mathbf{x}_{j})$

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

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

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

- 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

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

Gaussian/Radial kernels (polynomials of all orders – recall e= 1+2+2+1-. series expansion of exp)

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

Sigmoid

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

d=2 [1 x1 x2 x x x2 x1 x2]T

Mercer Kernels

What functions are valid kernels that correspond to feature vectors $\varphi(\mathbf{x})$? $\mathsf{K}(\mathbf{x},\mathbf{v}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{v})$

Answer: **Mercer kernels** K

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

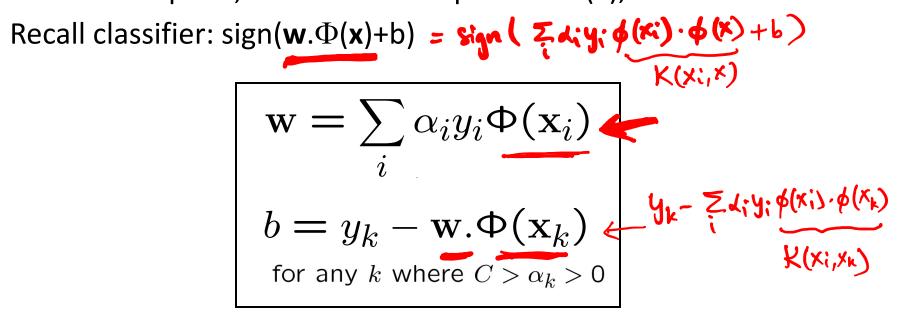
Ensures optimization is concave maximization

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!



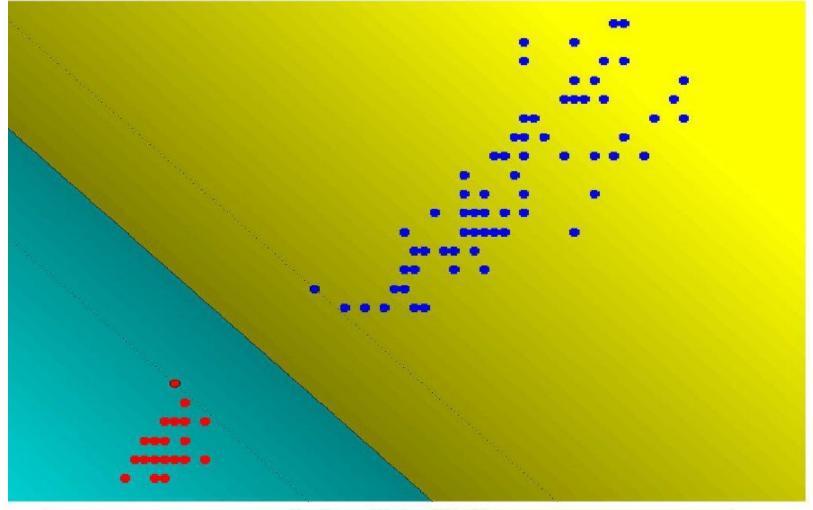
Using kernels we are cool!

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

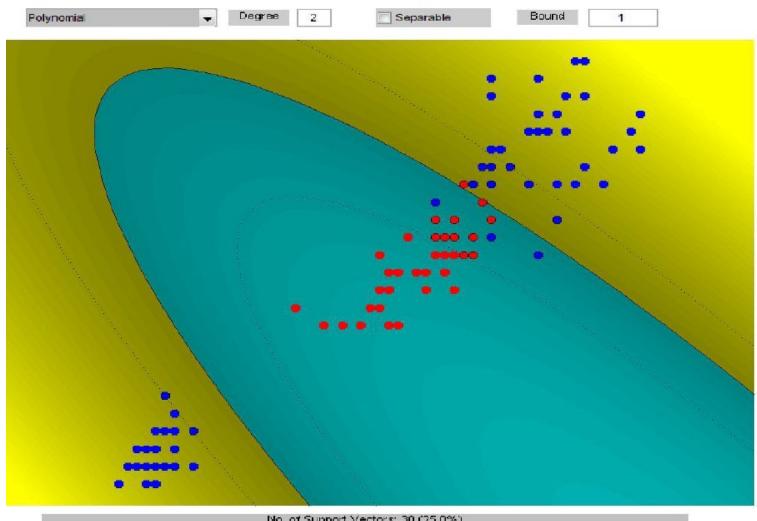
- Choose a set of features and kernel function
- Solve dual problem to obtain support vectors α_{i}
- At classification time, compute:

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

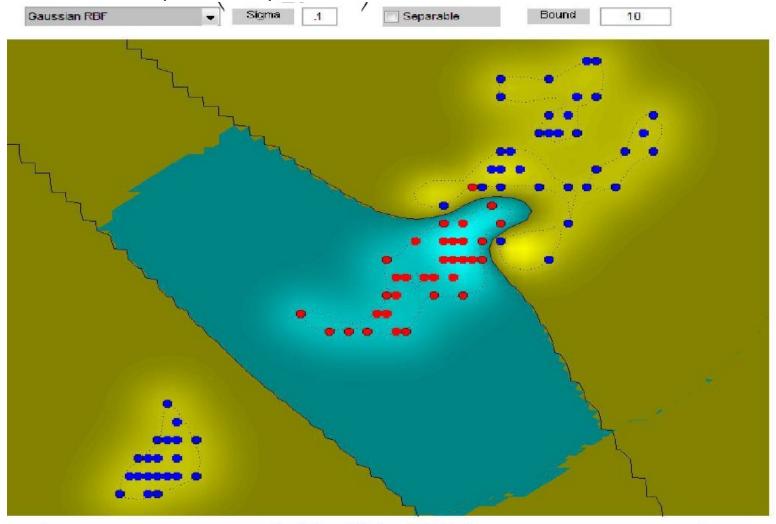
Iris dataset, 2 vs 13, Linear Kernel



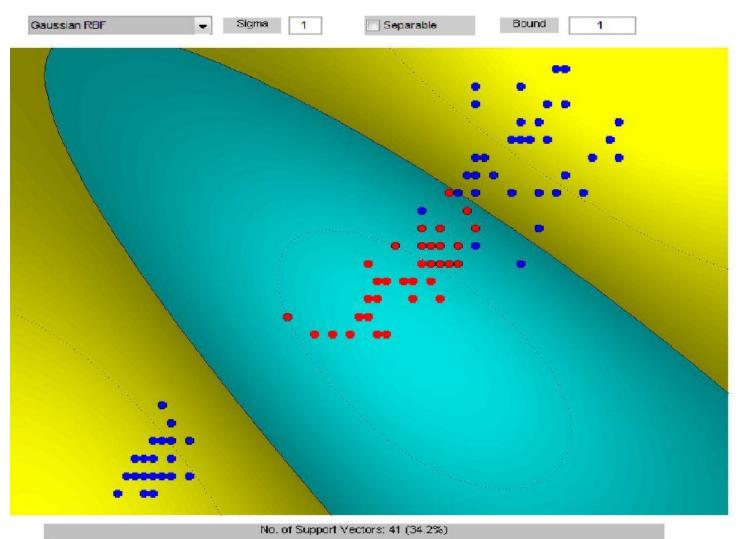
Iris dataset, 1 vs 23, Polynomial Kernel degree 2



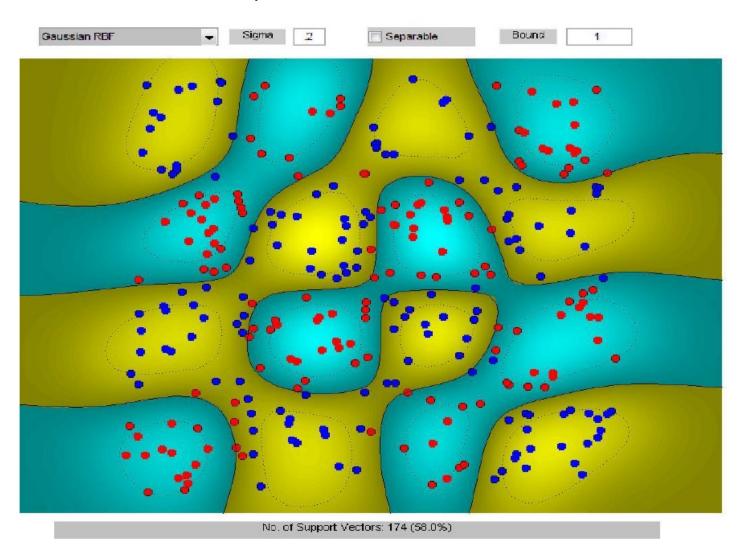
• Iris dataset, 1 vs 23, Gaussian RBF kernel



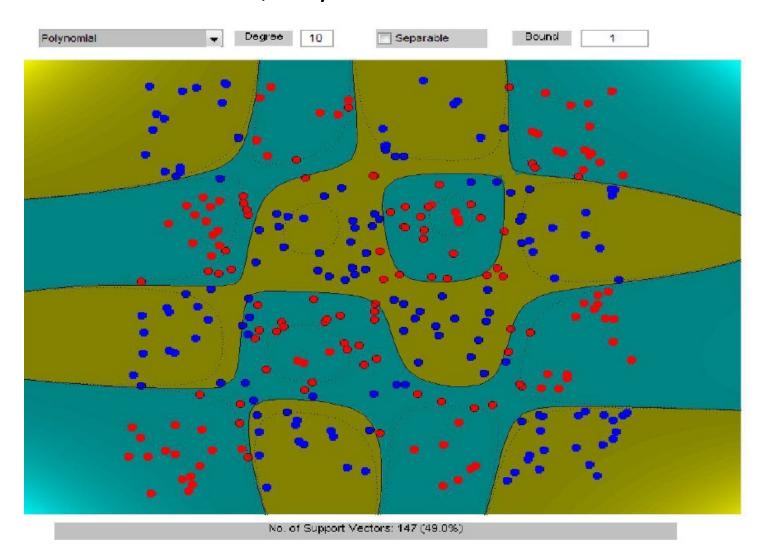
Iris dataset, 1 vs 23, Gaussian RBF kernel



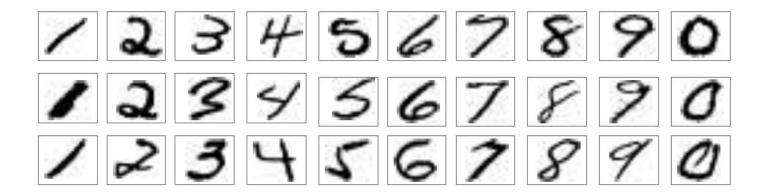
Chessboard dataset, Gaussian RBF kernel



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

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

$$\mathbf{Recall}$$

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

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 **AA**^T

Ridge regression

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

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

Similarity with SVMs

Primal problem:

$$\min_{\beta, z_i} \sum_{i=1}^n z_i^2 + \lambda \|\beta\|_2^2$$
s.t. $z_i = Y_i - X_i\beta$

SVM Primal problem:

$$\min_{w,\xi_i} C \sum_{i=1}^n \xi_i + \frac{1}{2} ||w||_2^2$$
s.t. $\xi_i = \max(1 - Y_i X_i w, 0)$

Lagrangian:

$$L(z_{i}, x, x) = \sum_{i=1}^{n} z_{i}^{2} + \lambda \|\beta\|^{2} + \sum_{i=1}^{n} \alpha_{i}(z_{i} - Y_{i} + X_{i}\beta)$$

 α_i – Lagrange parameter, one per training point

Ridge regression (dual)

$$\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2 \qquad \widehat{\beta} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$
 Dual problem:
$$\max_{\alpha} \min_{\beta, z_i} \sum_{i=1}^{n} z_i^2 + \lambda \|\beta\|^2 + \sum_{i=1}^{n} \alpha_i (z_i - Y_i + X_i \beta)$$

 $\alpha = {\alpha_i}$ for i = 1,..., n

Taking derivatives of Lagrangian wrt β and z_i we get:

$$\beta = -\frac{1}{2\lambda} \mathbf{A}^{\top} \alpha \qquad z_i = -\frac{\alpha_i}{2}$$
Dual problem:
$$\max_{\alpha} -\frac{\alpha^{\top} \alpha}{4} - \frac{1}{4\lambda} \alpha^{\top} \mathbf{A} \mathbf{A}^{\top} \alpha - \alpha^{\top} \mathbf{Y} \qquad \boldsymbol{\leftarrow}$$

n-dimensional optimization problem

Ridge regression (dual)

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

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

Dual problem:

$$\max \ -\frac{\alpha^{\top} \alpha}{4} - \frac{1}{4} \alpha^{\top} \mathbf{A} \mathbf{A}^{\top} \alpha - \alpha^{\top} \mathbf{Y}$$

$$\max_{\alpha} \ -\frac{\alpha^{\top} \alpha}{4} - \frac{1}{4\lambda} \alpha^{\top} \mathbf{A} \mathbf{A}^{\top} \alpha - \alpha^{\top} \mathbf{Y} \qquad \Rightarrow \widehat{\alpha} = -\left(\frac{\mathbf{A} \mathbf{A}^{\top}}{\lambda} + \mathbf{I}\right)^{-1} 2 \mathbf{Y}$$

can get back
$$\hat{\beta} = -\frac{1}{2\lambda} \mathbf{A}^{\top} \hat{\alpha} = \mathbf{A}^{\top} (\mathbf{A} \mathbf{A}^{\top} + \lambda \mathbf{I})^{-1} \mathbf{Y}$$



Weight of each training point (but typically not sparse)

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:

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

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}^{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^T 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!

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

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