

Kernel Machines

Stéphane Canu, Gilles Gasso
{stephane.canu, gilles.gasso}@insa-rouen.fr

DM_2 2016

February 1, 2017

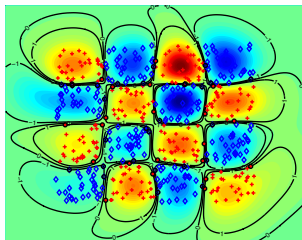
Plan

1 Non sparse kernel machines

- Interpolation problem
- From Regression to classification: kernel logistic regression

2 Sparse kernel machines: SVM

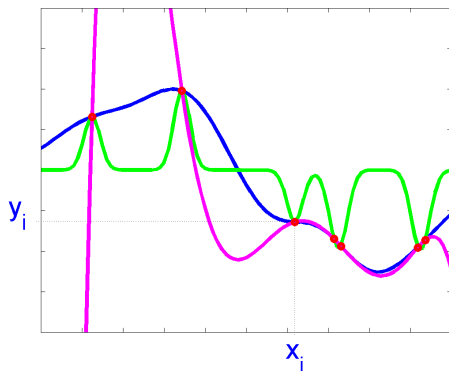
- Sparsity of kernel SVM for classification
 - SVM: variations on a theme
- Sparse kernel machines for regression



Splines interpolation

Interpolation and regression problems

Find out a function $f \in \mathcal{H}$ such that $f(\mathbf{x}_i) = y_i, \quad i = 1, \dots, n$



It is an ill posed problem

Interpolation splines: minimum norm interpolation

$$\begin{cases} \min_{f \in \mathcal{H}} \frac{1}{2} \|f\|_{\mathcal{H}}^2 \\ \text{such that } f(\mathbf{x}_i) = y_i, \quad i = 1, \dots, n \end{cases}$$

Remark: we assume \mathcal{H} is a Hilbert space induced by reproducing kernel k

The lagrangian ($\alpha_i \in \mathbb{R}$ are Lagrange multipliers)

$$L(f, \alpha) = \frac{1}{2} \|f\|^2 - \sum_{i=1}^n \alpha_i (f(\mathbf{x}_i) - y_i)$$

Interpolation splines: minimum norm interpolation

$$\begin{cases} \min_{f \in \mathcal{H}} & \frac{1}{2} \|f\|_{\mathcal{H}}^2 \\ \text{such that} & f(\mathbf{x}_i) = y_i, \quad i = 1, \dots, n \end{cases}$$

Remark: we assume \mathcal{H} is a Hilbert space induced by reproducing kernel k

The lagrangian ($\alpha_i \in \mathbb{R}$ are Lagrange multipliers)

$$L(f, \alpha) = \frac{1}{2} \|f\|^2 - \sum_{i=1}^n \alpha_i (f(\mathbf{x}_i) - y_i)$$

optimality for f : $\nabla_f L(f, \alpha) = 0 \quad \Leftrightarrow \quad f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$

Interpolation splines: minimum norm interpolation

$$\begin{cases} \min_{f \in \mathcal{H}} & \frac{1}{2} \|f\|_{\mathcal{H}}^2 \\ \text{such that} & f(\mathbf{x}_i) = y_i, \quad i = 1, \dots, n \end{cases}$$

Remark: we assume \mathcal{H} is a Hilbert space induced by reproducing kernel k

The lagrangian ($\alpha_i \in \mathbb{R}$ are Lagrange multipliers)

$$L(f, \alpha) = \frac{1}{2} \|f\|^2 - \sum_{i=1}^n \alpha_i (f(\mathbf{x}_i) - y_i)$$

optimality for f : $\nabla_f L(f, \alpha) = 0 \Leftrightarrow f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}_i, \mathbf{x})$

dual formulation (remove f from the lagrangian):

$$Q(\alpha) = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^n \alpha_i y_i \quad \text{solution: } \max_{\alpha \in \mathbb{R}^n} Q(\alpha)$$

$$K\alpha = y$$

Representer theorem

Theorem (Representer theorem)

Let \mathcal{H} be a RKHS with kernel $k(s, t)$. Let ℓ be a function from \mathcal{X} to \mathbb{R} (loss function) and Φ a non decreasing function from \mathbb{R}^+ to \mathbb{R} . If there exists a function f^* minimizing:

$$f^* = \underset{f \in \mathcal{H}}{\operatorname{argmin}} \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i)) + \Phi(\|f\|_{\mathcal{H}}^2)$$

then there exists a vector $\alpha \in \mathbb{R}^n$ such that:

$$f^*(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i)$$

Elements of a proof

- ① $\mathcal{H}_s = \text{span}\{k(\cdot, \mathbf{x}_1), \dots, k(\cdot, \mathbf{x}_i), \dots, k(\cdot, \mathbf{x}_n)\}$
- ② orthogonal decomposition: $\mathcal{H} = \mathcal{H}_s \oplus \mathcal{H}_\perp \Rightarrow \forall f \in \mathcal{H}; f = f_s + f_\perp$
- ③ pointwise evaluation decomposition

$$\begin{aligned} f(\mathbf{x}_i) &= f_s(\mathbf{x}_i) + f_\perp(\mathbf{x}_i) \\ &= \langle f_s(\cdot), k(\cdot, \mathbf{x}_i) \rangle_{\mathcal{H}} + \underbrace{\langle f_\perp(\cdot), k(\cdot, \mathbf{x}_i) \rangle_{\mathcal{H}}}_{=0} \\ &= f_s(\mathbf{x}_i) \end{aligned}$$

- ④ norm decomposition
- ⑤ decompose the global cost

$$\|f\|_{\mathcal{H}}^2 = \|f_s\|_{\mathcal{H}}^2 + \underbrace{\|f_\perp\|_{\mathcal{H}}^2}_{\geq 0} \geq \|f_s\|_{\mathcal{H}}^2$$

$$\begin{aligned} \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i)) + \Phi(\|f\|_{\mathcal{H}}^2) &= \sum_{i=1}^n \ell(y_i, f_s(\mathbf{x}_i)) + \Phi(\|f_s\|_{\mathcal{H}}^2 + \|f_\perp\|_{\mathcal{H}}^2) \\ &\geq \sum_{i=1}^n \ell(y_i, f_s(\mathbf{x}_i)) + \Phi(\|f_s\|_{\mathcal{H}}^2) \end{aligned}$$

⑥

$$\underset{f \in \mathcal{H}}{\operatorname{argmin}} = \underset{f \in \mathcal{H}_s}{\operatorname{argmin}}$$

Smoothing splines

introducing the error (the slack) $\xi = f(x_i) - y_i$

$$(S) \quad \left\{ \begin{array}{ll} \min_{f \in \mathcal{H}} & \frac{1}{2} \|f\|_{\mathcal{H}}^2 + \frac{1}{2\lambda} \sum_{i=1}^n \xi_i^2 \\ \text{such that} & f(x_i) = y_i + \xi_i, \quad i = 1, n \end{array} \right.$$

one equivalent formulation

$$(S') \quad \min_{f \in \mathcal{H}} \quad \frac{1}{2} \sum_{i=1}^n (f(x_i) - y_i)^2 + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2$$

using the representer theorem

$$(S') \quad \min_{\alpha \in \mathbb{R}^n} \quad \frac{1}{2} \|K\alpha - \mathbf{y}\|^2 + \frac{\lambda}{2} \alpha^\top K\alpha$$

solution: $(S) \Leftrightarrow (S') \quad \Leftrightarrow \quad \alpha = (K + \lambda I)^{-1} \mathbf{y}$

Remark: this is different from ridge regression problem

$$\min_{\alpha \in \mathbb{R}^n} \quad \frac{1}{2} \|K\alpha - \mathbf{y}\|^2 + \frac{\lambda}{2} \alpha^\top \alpha \quad \text{with} \quad \text{solution} \quad \alpha = (K^\top K + \lambda I)^{-1} K^\top \mathbf{y}$$

Kernel logistic regression

inspiration: the Bayes rule

$$D(\mathbf{x}) = \text{sign}(f(\mathbf{x}) + \alpha_0) \implies \log \left(\frac{\mathbb{P}(Y=1|\mathbf{x})}{\mathbb{P}(Y=-1|\mathbf{x})} \right) = f(\mathbf{x}) + \alpha_0$$

probabilities:

$$\mathbb{P}(Y = 1|\mathbf{x}) = \frac{\exp^{f(\mathbf{x})+\alpha_0}}{1 + \exp^{f(\mathbf{x})+\alpha_0}} \quad \mathbb{P}(Y = -1|\mathbf{x}) = \frac{1}{1 + \exp^{f(\mathbf{x})+\alpha_0}}$$

Rademacher distribution

$$\mathcal{L}(x_i, y_i, f, \alpha_0) = \mathbb{P}(Y = 1|\mathbf{x}_i)^{\frac{y_i+1}{2}} (1 - \mathbb{P}(Y = 1|\mathbf{x}_i))^{\frac{1-y_i}{2}}$$

penalized log-likelihood

$$\begin{aligned} J(f, \alpha_0) &= - \sum_{i=1}^n \log(\mathcal{L}(x_i, y_i, f, \alpha_0)) + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2 \\ &= \sum_{i=1}^n \log \left(1 + \exp^{-y_i(f(\mathbf{x}_i) + \alpha_0)} \right) + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2 \end{aligned}$$

Kernel logistic regression (2)

$$(\mathcal{R}) \quad \begin{cases} \min_{f \in \mathcal{H}} & \frac{1}{2} \|f\|_{\mathcal{H}}^2 + \frac{1}{\lambda} \sum_{i=1}^n \log(1 + \exp^{-\xi_i}) \\ \text{with} & \xi_i = y_i (f(\mathbf{x}_i) + \alpha_0), \quad i = 1, n \end{cases}$$

Representer theorem leads to

$$J(\alpha, \alpha_0) = \mathbb{1}^\top \log \left(\mathbb{1} + \exp^{\text{diag}(\mathbf{y}) K \alpha + \alpha_0 \mathbf{y}} \right) + \frac{\lambda}{2} \alpha^\top K \alpha$$

gradient vector and Hessian matrix:

$$\nabla_{\alpha} J(\alpha, \alpha_0) = K(\mathbf{y} - (2\mathbf{p} - \mathbb{1})) + \lambda K \alpha$$

$$H_{\alpha} J(\alpha, \alpha_0) = K \text{diag}(\mathbf{p}(\mathbb{1} - \mathbf{p})) K + \lambda K$$

solve the problem using Newton iterations

$$\alpha^{\text{new}} = \alpha^{\text{old}} + (K \text{diag}(\mathbf{p}(\mathbb{1} - \mathbf{p})) K + \lambda K)^{-1} K(\mathbf{y} - (2\mathbf{p} - \mathbb{1}) + \lambda \alpha)$$

Let's summarize

- Kernel machines: pros
 - ▶ Universality
 - ▶ from \mathcal{H} to \mathbb{R}^n using the representer theorem
 - ▶ no (explicit) curse of dimensionality
- splines $\mathcal{O}(n^3)$ (can be reduced to $\mathcal{O}(n^2)$)
- logistic regression $\mathcal{O}(kn^3)$ (can be reduced to $\mathcal{O}(kn^2)$)
- no scalability!

sparsity comes to the rescue!

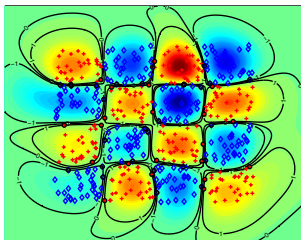
Roadmap

1 Non sparse kernel machines

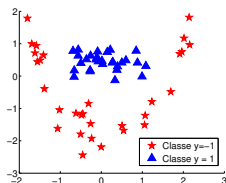
- Interpolation problem
- From Regression to classification: kernel logistic regression

2 Sparse kernel machines: SVM

- Sparsity of kernel SVM for classification
 - SVM: variations on a theme
- Sparse kernel machines for regression



SVM in a RKHS: the separable case (no noise)



- dataset

$$\mathcal{D} = \{(\mathbf{x}_i, y_i) \in \mathcal{X} \times \{-1, 1\}\}_{i=1}^n$$

- problem : learn a non linear SVM

$$\begin{cases} \min_{f,b} & \frac{1}{2} \|f\|_{\mathcal{H}}^2 \\ \text{with} & y_i (f(\mathbf{x}_i) + b) \geq 1 \end{cases}$$

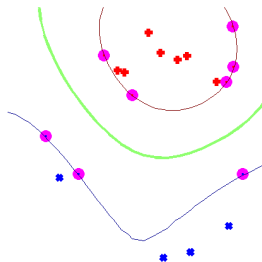
3 ways to represent function f

$$\underbrace{f(\mathbf{x})}_{\text{in the RKHS } \mathcal{H}} = \underbrace{\sum_{j=1}^d w_j \phi_j(\mathbf{x})}_{d \text{ features}} = \underbrace{\sum_{i=1}^n \alpha_i y_i k(\mathbf{x}, \mathbf{x}_i)}_{n \text{ data points}}$$

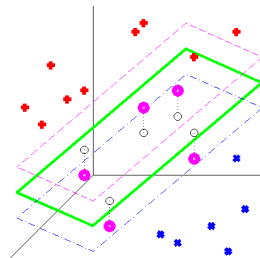
$$\begin{cases} \min_{\mathbf{w}, b} & \frac{1}{2} \|\mathbf{w}\|_{\mathbb{R}^d}^2 = \frac{1}{2} \mathbf{w}^\top \mathbf{w} \\ \text{with} & y_i (\mathbf{w}^\top \phi(\mathbf{x}_i) + b) \geq 1 \end{cases} \Leftrightarrow \begin{cases} \min_{\alpha, b} & \frac{1}{2} \alpha^\top K \alpha \\ \text{with} & y_i (\alpha^\top K(:, i) + b) \geq 1 \end{cases}$$

using relevant features...

a data point becomes a function $\mathbf{x} \rightarrow k(\mathbf{x}, \bullet)$



input space representation: \mathbf{x}



feature space: $k(\mathbf{x}, \cdot)$

Representer theorem for SVM

$$\begin{cases} \min_{f,b} & \frac{1}{2} \|f\|_{\mathcal{H}}^2 \\ \text{with} & y_i(f(\mathbf{x}_i) + b) \geq 1 \end{cases}$$

Lagrangian

$$L(f, b, \alpha) = \frac{1}{2} \|f\|_{\mathcal{H}}^2 - \sum_{i=1}^n \alpha_i (y_i(f(\mathbf{x}_i) + b) - 1) \quad \alpha \geq 0$$

optimality conditions: (i) $\nabla_f L(f, b, \alpha) = 0 \Leftrightarrow f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i k(\mathbf{x}_i, \mathbf{x})$
(ii) $\nabla_b L(f, b, \alpha) = 0 \Leftrightarrow \sum_{i=1}^n \alpha_i y_i = 0$

eliminate f from L :
$$\begin{cases} \|f\|_{\mathcal{H}}^2 = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) \\ \sum_{i=1}^n \alpha_i y_i f(\mathbf{x}_i) = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) \end{cases}$$

$$Q(b, \alpha) = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^n \alpha_i (y_i b - 1)$$

Dual formulation for SVM

the intermediate function

$$Q(b, \alpha) = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) - b \left(\sum_{i=1}^n \alpha_i y_i \right) + \sum_{i=1}^n \alpha_i$$

$$\max_{\alpha} \min_b Q(b, \alpha)$$

b can be seen as the Lagrange multiplier of the following (balanced) constraint $\sum_{i=1}^n \alpha_i y_i = 0$ which is also the optimality KKT condition on b

Dual formulation

$$\left\{ \begin{array}{ll} \max_{\alpha \in \mathbf{R}^n} & -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^n \alpha_i \\ \text{such that} & \sum_{i=1}^n \alpha_i y_i = 0 \\ \text{and} & 0 \leq \alpha_i, \quad i = 1, n \end{array} \right.$$

SVM dual formulation

Dual formulation

$$\left\{ \begin{array}{l} \max_{\alpha \in \mathbb{R}^n} \quad -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^n \alpha_i \\ \text{with} \quad \sum_{i=1}^n \alpha_i y_i = 0 \quad \text{and} \quad 0 \leq \alpha_i, \quad i = 1, n \end{array} \right.$$

The dual formulation gives a quadratic program (QP)

$$\left\{ \begin{array}{l} \min_{\alpha \in \mathbb{R}^n} \quad \frac{1}{2} \alpha^\top G \alpha - \mathbb{1}^\top \alpha \\ \text{with} \quad \alpha^\top \mathbf{y} = 0 \quad \text{and} \quad 0 \leq \alpha \end{array} \right.$$

with $G_{ij} = y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$

case of linear kernel: $f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i (\mathbf{x}^\top \mathbf{x}_i) = \sum_{j=1}^d \beta_j x_j$

remark: when d is small wrt. n solving the primal may be interesting.

the general case: C-SVM

Primal formulation

$$(\mathcal{P}) \begin{cases} \min_{f \in \mathcal{H}, b, \xi \in \mathbb{R}^n} & \frac{1}{2} \|f\|^2 + C \sum_{i=1}^n \xi_i \\ \text{such that} & y_i (f(\mathbf{x}_i) + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, n \end{cases}$$

C is the *regularization* parameter (to be tuned)

Dual formulation

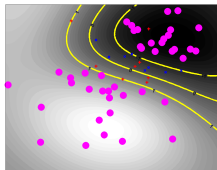
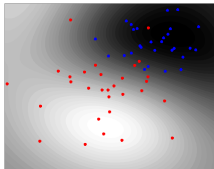
$$\begin{cases} \max_{\alpha \in \mathbb{R}^n} & -\frac{1}{2} \alpha^\top G \alpha + \alpha^\top \mathbb{I} \\ \text{such that} & \alpha^\top \mathbf{y} = 0 \text{ and } 0 \leq \alpha_i \leq C \quad i = 1, n \end{cases}$$

remark: regularization path is the set of solutions $\alpha(C)$ when C varies

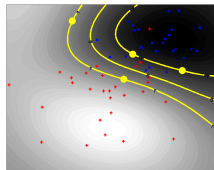
Data groups: illustration

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i)$$

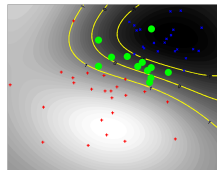
$$D(\mathbf{x}) = \text{sign}(f(\mathbf{x}) + b)$$



useless data
well classified
 $\alpha = 0$



important data
support
 $0 < \alpha < C$



suspicious data
 $\alpha = C$

the regularization path: is the set of solutions $\alpha(C)$ when C varies

The importance of being support

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i k(\mathbf{x}_i, \mathbf{x})$$

data point	α	constraint value	set
\mathbf{x}_i <i>useless</i>	$\alpha_i = 0$	$y_i(f(\mathbf{x}_i) + b) > 1$	l_0
\mathbf{x}_i <i>support</i>	$0 < \alpha_i < C$	$y_i(f(\mathbf{x}_i) + b) = 1$	l_α
\mathbf{x}_i <i>suspicious</i>	$\alpha_i = C$	$y_i(f(\mathbf{x}_i) + b) < 1$	l_C

Table: When a data point is « support » it lies exactly on the margin.

here lies the efficiency of the algorithm (and its complexity)!

sparsity: $\alpha_i = 0$

The active set method for SVM (1)

$$\left\{ \begin{array}{ll} \min_{\alpha \in \mathbb{R}^n} & \frac{1}{2} \alpha^\top G \alpha - \alpha^\top \mathbb{I} \\ \text{such that} & \alpha^\top \mathbf{y} = 0 \quad i = 1, n \\ \text{and} & 0 \leq \alpha_i \quad i = 1, n \end{array} \right. \quad \left\{ \begin{array}{l} G\alpha - \mathbb{I} - \beta + b\mathbf{y} = 0 \\ \alpha^\top \mathbf{y} = 0 \\ 0 \leq \alpha_i \quad i = 1, n \\ 0 \leq \beta_i \quad i = 1, n \\ \alpha_i \beta_i = 0 \quad i = 1, n \end{array} \right.$$

G_a	G_i^\top
G_i	G_0

$$\begin{array}{c} \alpha_a \\ 0 \end{array} - \begin{array}{c} 1 \\ 1 \end{array} - \begin{array}{c} 0 \\ \beta_0 \end{array} + b \begin{array}{c} y_a \\ y_0 \end{array} = \begin{array}{c} 0 \\ 0 \end{array}$$

$G \quad \alpha - \mathbb{I} - \beta + b \mathbf{y} = 0$

- (1) $G_a \alpha_a - \mathbb{I}_a + b y_a = 0$
- (2) $G_i \alpha_a - \mathbb{I}_0 - \beta_0 + b y_0 = 0$

- ① solve (1) (find α together with b)
- ② if $\alpha_j < 0$ move it from l_α to l_0
goto 1
- ③ else solve (2)
if $\beta_j < 0$ move it from l_0 to l_α
goto 1

The active set method for SVM (2)

Function $(\alpha, b, l_\alpha) \leftarrow \text{Solve_QP_Active_Set}(G, \mathbf{y})$

```
% Solve  $\min_{\alpha} \quad 1/2 \alpha^\top G \alpha - \mathbf{1}^\top \alpha$   
%           s.t.  $0 \leq \alpha$  and  $\mathbf{y}^\top \alpha = 0$ 
```

$(l_\alpha, l_0, \alpha) \leftarrow \text{initialization}$

while The_optimal_is_not_reached **do**

$(\alpha, b) \leftarrow \text{solve} \begin{cases} G_a \alpha_a - \mathbf{1}_a + b \mathbf{y}_a \\ \mathbf{y}_a^\top \alpha_a \end{cases} = 0$

if $\exists i \in l_\alpha$ such that $\alpha_i < 0$ **then**

$\alpha \leftarrow \text{projection}(\alpha_a, \alpha)$

move i from l_α to l_0

else if $\exists j \in l_0$ such that $\beta_j < 0$ **then**

use $\beta_0 = \mathbf{y}_0(K_i \alpha_a + b \mathbf{1}_0) - \mathbf{1}_0$

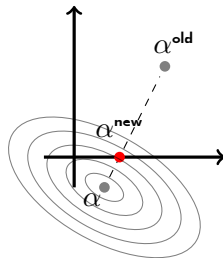
move j from l_0 to l_α

else

The_optimal_is_not_reached \leftarrow FALSE

end if

end while



Projection step of the active constraints algorithm

```
d = alpha - alphaold;  
alpha = alpha + t * d;
```

Caching Strategy

Save space and computing time by computing only the needed parts of kernel matrix G

Two more ways to get SVM

Using the hinge loss

$$\min_{f \in \mathcal{H}, b \in \mathbb{R}} \sum_{i=1}^n \max(0, 1 - y_i(f(\mathbf{x}_i) + b)) + \frac{1}{2C} \|f\|^2$$

Minimizing the distance between the convex hulls

$$\left\{ \begin{array}{l} \min_{\alpha} \quad \|u - v\|_{\mathcal{H}}^2 \\ \text{with} \quad u(\mathbf{x}) = \sum_{\{i|y_i=1\}} \alpha_i k(\mathbf{x}_i, \mathbf{x}), \quad v(\mathbf{x}) = \sum_{\{i|y_i=-1\}} \alpha_i k(\mathbf{x}_i, \mathbf{x}) \\ \text{and} \quad \sum_{\{i|y_i=1\}} \alpha_i = 1, \quad \sum_{\{i|y_i=-1\}} \alpha_i = 1, \quad 0 \leq \alpha_i \quad i = 1, n \end{array} \right.$$

$$f(\mathbf{x}) = \frac{2}{\|u - v\|_{\mathcal{H}}^2} (u(\mathbf{x}) - v(\mathbf{x})) \quad \text{and} \quad b = \frac{\|u\|_{\mathcal{H}}^2 - \|v\|_{\mathcal{H}}^2}{\|u - v\|_{\mathcal{H}}^2}$$

SVM with non symmetric costs

problem in the primal

$$\left\{ \begin{array}{ll} \min_{f \in \mathcal{H}, b, \xi \in \mathbb{R}^n} & \frac{1}{2} \|f\|_{\mathcal{H}}^2 + C^+ \sum_{\{i|y_i=1\}} \xi_i + C^- \sum_{\{i|y_i=-1\}} \xi_i \\ \text{with} & y_i(f(\mathbf{x}_i) + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, n \end{array} \right.$$

dual formulation

$$\left\{ \begin{array}{ll} \max_{\alpha \in \mathbb{R}^n} & -\frac{1}{2} \alpha^\top G \alpha + \alpha^\top \mathbb{I} \\ \text{with} & \alpha^\top \mathbf{y} = 0 \\ \text{and} & 0 \leq \alpha_i \leq C^+ \text{ or } C^- \quad i = 1, n \end{array} \right.$$

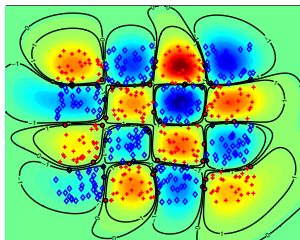
Roadmap

1 Non sparse kernel machines

- Interpolation problem
- From Regression to classification: kernel logistic regression

2 Sparse kernel machines: SVM

- Sparsity of kernel SVM for classification
 - SVM: variations on a theme
- Sparse kernel machines for regression



K-Lasso (Kernel Basis pursuit)

The Kernel Lasso : non-linear regression with sparsity penalization

$$(\mathcal{S}_1) \quad \left\{ \min_{\alpha \in \mathbb{R}^n} \quad \frac{1}{2} \|K\alpha - \mathbf{y}\|^2 + \lambda \sum_{i=1}^n |\alpha_i| \right.$$

- Typical parametric quadratic program (pQP)
- The ℓ_1 norm $\|\alpha\|_1 = \sum_{i=1}^n |\alpha_i|$ induces sparsity of the solution α (i.e. some $\alpha_i = 0$)

The dual:

$$(\mathcal{D}_1) \quad \left\{ \begin{array}{l} \min_{\alpha} \quad \frac{1}{2} \|K\alpha\|^2 \\ \text{such that} \quad K^\top (K\alpha - \mathbf{y}) \leq t \end{array} \right.$$

- require to compute $K^\top K$!

Obtained regression function

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i)$$

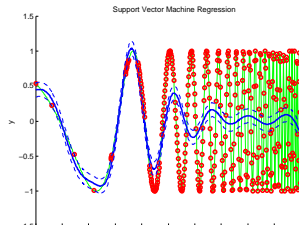
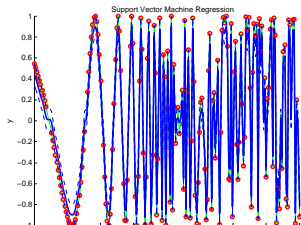
Support vector regression (SVR)

Regression with absolute value error $\begin{cases} \min_{f \in \mathcal{H}} & \frac{1}{2} \|f\|_{\mathcal{H}}^2 \\ \text{s. t.} & |f(\mathbf{x}_i) - y_i| \leq t, \quad i = 1, n \end{cases}$

The support vector regression introduce slack variables

$$(SVR) \quad \begin{cases} \min_{f \in \mathcal{H}} & \frac{1}{2} \|f\|_{\mathcal{H}}^2 + C \sum |\xi_i| \\ \text{such that} & |f(\mathbf{x}_i) - y_i| \leq t + \xi_i \quad 0 \leq \xi_i \quad i = 1, n \end{cases}$$

- a typical **multi** parametric quadratic program (mpQP)



SVM reduction (reduced set method)

- objective: compile the model

- $f(x) = \sum_{i=1}^{n_s} \alpha_i k(\mathbf{x}_i, \mathbf{x}), n_s \ll n, \quad n_s \text{ too big}$

- compiled model as the solution of: $g(\mathbf{x}) = \sum_{i=1}^{n_c} \beta_i k(\mathbf{z}_i, \mathbf{x}), n_c \ll n_s$

- β, \mathbf{z}_i and c are tuned by minimizing:

$$\min_{\beta, \mathbf{z}_i} \|g - f\|_H^2$$

where

$$\min_{\beta, \mathbf{z}_i} \|g - f\|_H^2 = \alpha^\top K_x \alpha + \beta^\top K_z \beta - 2\alpha^\top K_{xz} \beta$$

some authors advice $0,03 \leq \frac{n_c}{n_s} \leq 0,1$

- solve it by using use (stochastic) gradient (its a RBF problem)

logistic regression and the import vector machine

- Logistic regression is NON sparse
- kernalize it using the *dictionary* strategy
- Algorithm:
 - ▶ find the solution of the KLR using only a subset S of the data
 - ▶ build S iteratively using active constraint approach
- this trick brings sparsity
- it estimates probability
- it can naturally be generalized to the multiclass case

- efficient when uses:
 - ▶ a few import vectors
 - ▶ component-wise update procedure

- extention using L_1 KLR

Historical perspective on kernel machines

statistics

1960 Parzen, Nadaraya Watson

1970 Splines

1980 Kernels: Silverman, Hardle...

1990 sparsity: Donoho (pursuit),
Tibshirani (Lasso)...

Statistical learning

1985 Neural networks:

- ▶ non linear - universal
- ▶ structural complexity
- ▶ non convex optimization

1992 Vapnik et. al.

- ▶ theory - regularization - consistency
- ▶ convexity - Linearity
- ▶ **Kernel** - universality
- ▶ **sparsity**
- ▶ results: MNIST

what's new since 1995

- Applications

- ▶ kernlisation $w^\top \mathbf{x} \rightarrow \langle f, k(\mathbf{x}, \cdot) \rangle_{\mathcal{H}}$
- ▶ kernel engineering
- ▶ sturtured outputs
- ▶ applications: image, text, signal, bio-info...

- Optimization

- ▶ dual: mloss.org
- ▶ approximation
- ▶ primal

- Statistics

- ▶ proofs and bounds
- ▶ model selection
 - ★ span bound
 - ★ multikernel: tuning (k and σ)

challenges

- the size effect
 - ▶ ready to use: automatization
 - ▶ adaptative: on line context aware
 - ▶ beyond kenrels
- Automatic and adaptive model selection
 - ▶ variable selection
 - ▶ kernel tuning (k et σ)
 - ▶ hyper-parameters: C , duality gap, λ
- \mathbb{P} change
- Theory
 - ▶ non positive kernels
 - ▶ a more general representer theorem

biblio: kernel-machines.org

- John Shawe-Taylor and Nello Cristianini Kernel Methods for Pattern Analysis, Cambridge University Press, 2004
- Bernhard Schölkopf and Alex Smola. Learning with Kernels. MIT Press, Cambridge, MA, 2002.
- Trevor Hastie, Robert Tibshirani and Jerome Friedman, The Elements of Statistical Learning: Data Mining, Inference, and Prediction, springer, 2001
- Léon Bottou, Olivier Chapelle, Dennis DeCoste and Jason Weston Large-Scale Kernel Machines (Neural Information Processing, MIT press 2007
- Olivier Chapelle, Bernhard Scholkopf and Alexander Zien, Semi-supervised Learning, MIT press 2006
- Vladimir Vapnik. Estimation of Dependences Based on Empirical Data. Springer Verlag, 2006, 2nd edition.
- Vladimir Vapnik. The Nature of Statistical Learning Theory. Springer, 1995.
- Grace Wahba. Spline Models for Observational Data. SIAM CBMS-NSF Regional Conference Series in Applied Mathematics vol. 59, Philadelphia, 1990
- Alain Berline and Christine Thomas-Agnan, Reproducing Kernel Hilbert Spaces in Probability and Statistics, Kluwer Academic Publishers, 2003
- Marc Attéia et Jean Gaches , Approximation Hilbertienne - Splines, Ondelettes, Fractales, PUG, 1999