



**HAL**  
open science

## SVM and kernel machines

Stéphane Canu, Gaëlle Loosli, Alain Rakotomamonjy

► **To cite this version:**

Stéphane Canu, Gaëlle Loosli, Alain Rakotomamonjy. SVM and kernel machines. École thématique. SVM and kernel machines, ECI 2011, Buenos Aires, 2011, pp.100. cel-00643485

**HAL Id: cel-00643485**

**<https://cel.hal.science/cel-00643485>**

Submitted on 22 Nov 2011

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Recent advances in kernel machines

Julio 2011

Escuela de Ciencias Informáticas 2011

Laboratoire d'Informatique, de Traitement de l'Information et des Systèmes

Stéphane Canu, Gaëlle Loosli & Alain Rakotomamonjy

[stephane.canu@litislab.eu](mailto:stephane.canu@litislab.eu)

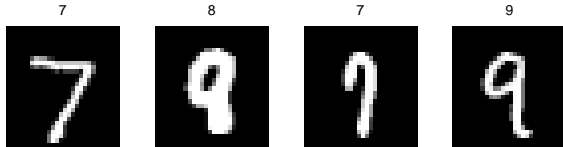
 litis

# Optical character recognition

## Example (The MNIST database)

- ▶ MNIST<sup>a</sup>, data = « image-label »
- ▶  $n = 60,000$ ;  $d = 700$ ; classes = 10
- ▶ Kernel error rate = 0.56 %,
- ▶ Best error rate = 0.4 % .

<sup>a</sup><http://yann.lecun.com/exdb/mnist/index.html>





# Historical perspective on kernel machines

## statistics

- 1960** Parzen, Nadaraya Watson
- 1970** Splines
- 1980** Kernels: Silverman, Hardsle...
- 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

# Notations

- ▶ inputs  $\mathbf{x} \in \mathcal{X} = \mathbb{R}^d$ ,  $d$  features
- ▶ outputs  $y$
- ▶ training set  $(\mathbf{x}_i, y_i)$ ,  $i = 1, n$ ,  $n$  exemples
- ▶ test set  $(\mathbf{x}_j, y_j)$ ,  $j = 1, \ell$ ,  $\ell$  exemples
- ▶ kernel  $k(\mathbf{x}_i, \mathbf{x}_j) \quad \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$
- ▶ RKHS  $\mathcal{H}$  (set of hypothesis associated with **positive** kernel  $k$ )
- ▶ RKKS  $\mathcal{K}$  (set of hypothesis associated with kernel  $k$ ) – Krein

## Definition (Kernel machines)

$$\mathcal{A}((\mathbf{x}_i, y_i)_{i=1,n})(\mathbf{x}) = \psi\left(\sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i) + \sum_{j=1}^p \beta_j q_j(\mathbf{x})\right)$$

$\alpha$  et  $\beta$ : parameters to be estimated

# Road map

## 1 Introduction

## 2 Tuning the kernel: MKL

- The multiple kernel problem
- SimpleMKL: the multiple kernel solution

## 3 Non positive kernel

- NON Positive kernels
- Functional estimation in a RKKS
- Non positive SVM

## 4 Distribution shift

- Distribution shift: the problem
- Density ratio estimation principle

## 5 Conclusion



# Multiple Kernel

The model

$$f(x) = \sum_{i=1}^n \alpha_i k(x, x_i) + b,$$

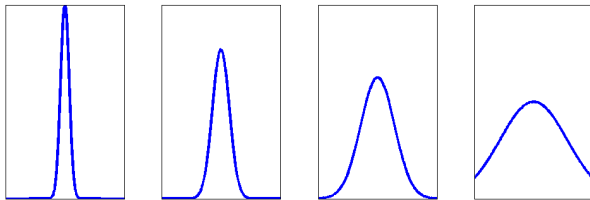
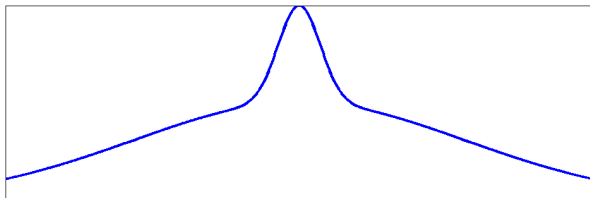
Given  $M$  kernel functions  $K_1, \dots, K_M$  that are potentially well suited for a given problem, find a positive linear combination of these kernels such that the resulting kernel  $k$  is “optimal”

$$k(\mathbf{x}, \mathbf{x}') = \sum_{m=1}^M d_m k_m(\mathbf{x}, \mathbf{x}'), \text{ with } d_m \geq 0, \sum_m d_m = 1$$

Need to learn together the kernel coefficients  $d_m$  and the SVR parameters  $\alpha_i, b$ .



# Multiple Kernel: illustration

 $k_1$  $k_2$  $k_3$  $k_4$ 

$$k = m_1 k_1 + m_2 k_2 + m_3 k_3 + m_4 k_4$$

$$m_2 = m_3 = 0$$

# Multiple Kernel functional Learning

The problem (for given  $C$  and  $t$ )

$$\min_{\{f_m\}, b, \xi, d} \frac{1}{2} \sum_m \frac{1}{d_m} \|f_m\|_{\mathcal{H}_m}^2 + C \sum_i \xi_i$$

$$\text{s.t.} \quad \left| \sum_m f_m(x_i) + b - y_i \right| \leq t + \xi_i \quad \forall i, \xi_i \geq 0 \quad \forall i$$

$$\sum_m d_m = 1, \quad d_m \geq 0 \quad \forall m,$$

## regularization formulation

$$\min_{\{f_m\}, b, d} \frac{1}{2} \sum_m \frac{1}{d_m} \|f_m\|_{\mathcal{H}_m}^2 + C \sum_i \max\left(\left| \sum_m f_m(x_i) + b - y_i \right| - t, 0\right)$$

$$\sum_m d_m = 1, \quad d_m \geq 0 \quad \forall m,$$

Equivalently

$$\min_{\{f_m\}, b, \xi, d} \sum_i \max\left(\left| \sum_m f_m(x_i) + b - y_i \right| - t, 0\right) + \frac{1}{2C} \sum_m \frac{1}{d_m} \|f_m\|_{\mathcal{H}_m}^2 + \mu \sum_m |d_m|$$



# Multiple Kernel Algorithm

Use a Reduced Gradient Algorithm<sup>1</sup>

$$\begin{aligned} &\min_{d \in \mathbb{R}^M} J(d) \\ &\text{s.t.} \quad \sum_m d_m = 1, \quad d_m \geq 0 \quad \forall m, \end{aligned}$$

## SimpleMKL algorithm

set  $d_m = \frac{1}{M}$  for  $m = 1, \dots, M$

**while** stopping criterion not met **do**

  compute  $J(d)$  using an QP solver with  $K = \sum_m d_m K_m$

  compute  $\frac{\partial J}{\partial d_m}$ , Hessian and descent direction  $D$

$\gamma \leftarrow$  compute optimal stepsize

$d \leftarrow d + \gamma D$

**end while**

→ Recent improvement reported using the Hessian

<sup>1</sup>Rakotomamonjy et al. JMLR 08

# Complexity

## For each iteration:

- ▶ SVM training:  $O(nn_{sv} + n_{sv}^3)$ .
- ▶ Inverting  $K_{sv,sv}$  is  $O(n_{sv}^3)$ , but might already be available as a by-product of the SVM training.
- ▶ Computing  $H$ :  $O(Mn_{sv}^2)$
- ▶ Finding  $d$ :  $O(M^3)$ .

The number of iterations is usually less than 10.

→ When  $M < n_{sv}$ , computing  $d$  is not more expensive than QP.



# Conclusion on multiple kernel (MKL)

- ▶ MKL: Kernel tuning, variable selection. . .
  - ▶ extension to classification and one class SVM
- ▶ SVM KM: an efficient Matlab toolbox (available at MLOSS)<sup>2</sup>
- ▶ Multiple Kernels for Image Classification: Software and Experiments on Caltech-101<sup>3</sup>
- ▶ new trend: Multi kernel and Multi task

---

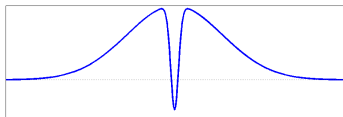
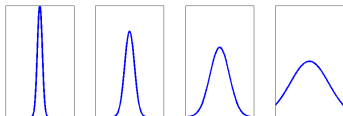
<sup>2</sup><http://mloss.org/software/view/33/>

<sup>3</sup><http://www.robots.ox.ac.uk/~vgg/software/MKL/>





# Learning with non positive kernel: why?



$$k = m_1 k_1 + m_2 k_2 + m_3 k_3 + m_4 k_4 \quad m_2 = m_3 = 0 \text{ and } m_1 < 0$$

- ▶ multiple non positive kernels

$$k(\mathbf{x}, \mathbf{x}') = \sum_{m=1}^M d_m k_m(\mathbf{x}, \mathbf{x}')$$

without  $d_m \geq 0$ ,

- ▶ Biological non positive kernels
- ▶ Positive radial kernels are Localized
- ▶  $\tanh(\mathbf{w}^\top \mathbf{x})$

# NON Positive kernels

## Definition of the associated pre Krein space

▶  $\mathcal{K}_0 = \{f \in \mathbb{R}^{\mathcal{X}} \mid f(x) = \sum_{i=1}^n \alpha_i k(x, x_i), \alpha_i \in \mathbb{R}, x_i \in \mathcal{X}\}$

▶ inner product on  $\mathcal{K}_0$  :

$$f(x) = \sum_{i=1}^n \alpha_i k(x, x_i), \quad g(x) = \sum_{i=1}^m \beta_i k(x, \tilde{x}_i)$$

$$\langle f(\cdot), g(\cdot) \rangle_{\mathcal{K}_0} = \sum_{i=1}^n \sum_{j=1}^m \alpha_i \beta_j k(x_i, \tilde{x}_j)$$

## the continuity of the evaluation functional

▶  $A_x f = f(x) = \langle f(\cdot), k(x, \cdot) \rangle_{\mathcal{K}_0}$       evaluation functional

▶  $\langle k(x, \cdot), k(y, \cdot) \rangle_{\mathcal{K}_0} = k(x, y)$       reproducing property

no more norm:  $\langle f, f \rangle_{\mathcal{K}_0}$  is NOT always positive

# Reproducing Kernel Krein Spaces (RKKS)

## Fundamental hypothesis

- ▶ There exist two positive kernels  $k_+$  and  $k_-$  such that

$$k(x, y) = k_+(x, y) - k_-(x, y)$$

## the RKKS space has to be complete

- ▶ define a topology using the positive kernels  $k_+$  and  $k_-$
- ▶  $\mathcal{K} = \widehat{\mathcal{K}}_0$
- ▶ remark about unicity

## Theorem: 3 equivalent statements

- ▶  $\mathcal{K}$  is a RKKS with kernel  $k$
- ▶  $k(x, y) = k_+(x, y) - k_-(x, y)$
- ▶  $k(x, y)$  is dominated by a positive kernel

# Examples

- ▶ Minkowski space time

$$\langle (x, y, z, t), (\tilde{x}, \tilde{y}, \tilde{z}, \tilde{t}) \rangle_{\mathcal{K}} = x\tilde{x} + y\tilde{y} + z\tilde{z} - t\tilde{t}$$

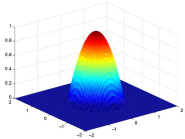
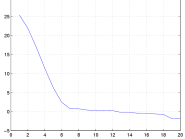
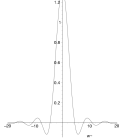
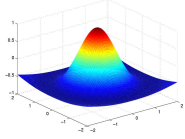
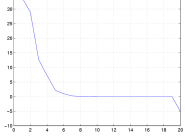
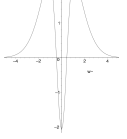
- ▶ Difference of two gaussians

$$k(s, t) := \alpha \exp^{-\frac{\|s-t\|^2}{b}} - \beta \exp^{-\frac{\|s-t\|^2}{c}}$$

- ▶ “Wavelets” kind: assume  $\mathcal{H} = V \oplus W$

$$\mathcal{K} = V \ominus W$$

# Examples

Kernel	2D kernel	Eigenvalues	Fourier
Epanechnikov kernel  $\left(1 - \frac{\ s-t\ ^2}{\sigma}\right)^p,$ for $\frac{\ s-t\ ^2}{\sigma} \leq 1$			
Gaussian Combination  $\exp\left(\frac{-\ s-t\ ^2}{\sigma_1}\right)$ $+ \exp\left(\frac{-\ s-t\ ^2}{\sigma_2}\right)$ $- \exp\left(\frac{-\ s-t\ ^2}{\sigma_3}\right)$			

Examples of indefinite kernels. Column 2 shows the 2D surface of the kernel with respect to the origin, column 3 shows plots of the 20 eigenvalues with largest magnitude of uniformly spaced data from the interval  $[-2, 2]$ , column 4 shows plots of the Fourier spectra.

# An other view on splines

## Interpolation is an ill posed problem

Let  $\mathcal{H}$  be a RKHS: Minimize  $\|f\|_{\mathcal{H}}^2$  such that  $f(\mathbf{x}_i) = y_i, i = 1, n$

In a Krein space

- ▶ NO more norm: how to regularize?
- ▶ project 0 on the set of constrains:  $K\alpha = \mathbf{y}$

## approximation is an ill posed problem

Let  $\mathcal{H}$  be a RKHS: Minimize  $\|f\|_{\mathcal{H}}^2 + \frac{C}{2} \sum \xi_i^2$  such that

$f(\mathbf{x}_i) - y_i = \xi_i, i = 1, n$  In a Krein space

- ▶ NO more norm: how to regularize?
- ▶ compute a **path** between 0 and the interpolating solution

## Non positive SVM: related work

- ▶ considering that the indefinite kernel is a perturbation of a true Mercer kernel.
- ▶ finding a stationary point, which is not unique but each of those performs correct separation. Moreover, it is shown that the problem is then cannot be seen as a margin maximization although a notion of margin can be defined.
- ▶ Krein space instead of a Hilbert space.
- ▶ Applying this to SVM requires to interpret this stabilization setting.
- ▶ a (unconstraint) quadratic program in a Krein space has a unique solution (if the involved matrix is non singular) which is in general a stationary point.

## SVM in a Krein space

$$\left\{ \begin{array}{l} \min_{f, \alpha_0} \quad \frac{1}{2} \|f\|_{\mathcal{H}}^2 \\ \text{with} \quad y_i(f(\mathbf{x}_i) + \alpha_0) \geq 1 \end{array} \right. \Leftrightarrow \left\{ \begin{array}{l} \text{stab}_{f, \alpha_0} \quad \frac{1}{2} \langle f, f \rangle_{\mathcal{K}} \\ \text{with} \quad y_i(f(\mathbf{x}_i) + \alpha_0) \geq 1 \end{array} \right.$$

The representer theorem holds:  $\langle f, f \rangle_{\mathcal{K}} = \alpha^\top K \alpha$

solve the problem using normal residuals (ie. solving  $Ax = b$  via  $A^\top Ax = A^\top b$ ).

$$\begin{aligned} \alpha^\top G &= \mathbf{1} - \lambda \mathbf{y}^\top - \mu^\top + \eta^\top \\ \alpha^\top GG^\top &= (\mathbf{1} - \lambda \mathbf{y}^\top - \mu^\top + \eta^\top) G^\top \end{aligned}$$

This can be seen as least squares. All the other conditions remain identical.



# SVM and KKT conditions of optimality

$$\text{KKT conditions for SVM} \quad \left\{ \begin{array}{l} \min_{\alpha} \quad \frac{1}{2} \alpha^{\top} G \alpha - \alpha^{\top} \mathbf{1} \\ \text{subject to} \quad \alpha^{\top} \mathbf{y} = 0 \\ \text{and} \quad 0 \leq \alpha_i \leq C \quad \forall i \in [1..n] \end{array} \right.$$

The stationarity condition is as follows:

$$-\alpha^{\top} G + \mathbf{1} - \lambda \mathbf{y}^{\top} - \mu^{\top} + \eta^{\top} = \mathbf{0}$$

The primal admissibility is given by

$$\begin{array}{l} -\alpha^{\top} \mathbf{y} = 0 \\ \alpha_i \leq C \quad \forall i \in [1..n] \end{array}$$

The dual admissibility is given by

$$\alpha_i \geq 0 \quad \forall i \in [1..n]$$

The dual admissibility is given by

$$\mu_i \geq 0 \quad \forall i \in [1..n]$$

The dual admissibility is given by

$$\eta_i \geq 0 \quad \forall i \in [1..n]$$

The complementary conditions are

$$\begin{array}{l} -\alpha_i \mu_i = 0 \quad \forall i \in [1..n] \\ (\alpha_i - C) \eta_i = 0 \quad \forall i \in [1..n] \end{array}$$

## Point of view 2 : a stabilization problem

Stabilizing  $\mathcal{J}$  is equivalent to minimizing  $\mathcal{M}$ :

$$\mathcal{J} = \frac{1}{2} \alpha^\top G \alpha - \alpha^\top \mathbf{1} \qquad \mathcal{M}(\alpha) = \langle \alpha^\top G - \mathbf{1}^\top, \alpha^\top G - \mathbf{1}^\top \rangle$$

This provides  $\left\{ \begin{array}{l} \min_{\alpha} \quad \langle \alpha^\top G - \mathbf{1}^\top, \alpha^\top G - \mathbf{1}^\top \rangle \\ \text{with} \quad \alpha^\top \mathbf{y} = 0 \\ \text{and} \quad 0 \leq \alpha_i \leq C \quad \forall i \in [1..n] \end{array} \right.$

### KKT conditions

The stationarity condition is as follows:

$$(\alpha^\top G - \mathbf{1}^\top + \lambda \mathbf{y}^\top - \mu^\top + \eta^\top) G^\top = \mathbf{0}$$

The primal admissibility is given by

$$-\alpha^\top \mathbf{y} = 0$$

$$\alpha_i \leq C \quad \forall i \in [1..n]$$

$$\alpha_i \geq 0 \quad \forall i \in [1..n]$$

The dual admissibility is given by

$$\mu_i \geq 0 \quad \forall i \in [1..n]$$

$$\eta_i \geq 0 \quad \forall i \in [1..n]$$

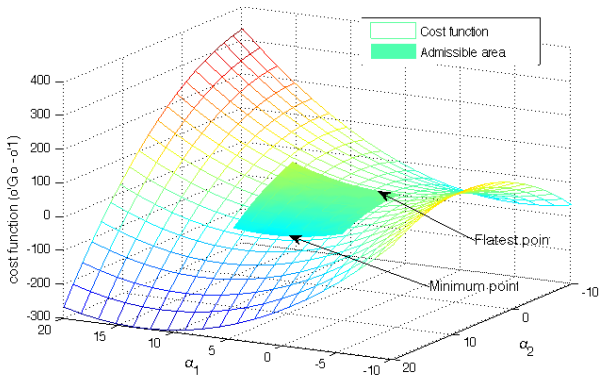
The complementary conditions are

$$-\alpha_i \mu_i = 0 \quad \forall i \in [1..n]$$

$$(\alpha_i - C) \eta_i = 0 \quad \forall i \in [1..n]$$

## Point of view 3 : The projection

Chasing the the most stable point, ie. the admissible point minimizing the gradient of the cost function (which is  $\alpha^T G - \mathbf{1}^T$ ).



**Figure:** SVM cost function with sigmoid kernel, illustrated for 2 support vectors. The plain area shows the admissible solutions.

## The solver

The proposed algorithm is derived from active set approach for SVM, The sets of points are defined according to the complementarity conditions (see table 2).

**Table:** Definition of groups for active set depending on the dual variable values

Group	$\alpha$	$\eta$	$\mu$
$l_0$	0	0	$> 0$
$l_C$	$C$	$> 0$	0
$l_w$	$0 < \alpha < C$	0	0

By default, all training points are in the non support vector set  $l_0$  except for a couple with opposite labels which is in  $l_w$ . Any other initial situation based on warm-start or a priori does not change the algorithm.

## Relaxing constraints in $l_0$ or $l_C$

If the current solution is admissible, we check the stationarity conditions for  $l_0$  and  $l_C$ . The most violating point is transferred from its group to  $l_w$ .

### The NPSVM algorithm

- 1: Initialize (one random point for each class in  $l_w$ , all others in  $l_0$ )
- 2: **while** solution is not optimal **do**
- 3:   solve linear system
- 4:   **if** primal admissibility is not satisfied **then**
- 5:     project solution in the admissible domain : remove a support vector from  $l_w$  (to  $l_0$  or  $l_C$ )
- 6:   **else if** stationarity condition is not satisfied **then**
- 7:     add new support vector to  $l_w$  (from  $l_0$  or  $l_C$ )
- 8:   **end if**
- 9: **end while**

# Experimental results

- ▶ sigmoid kernel (tanh) :  $k(x_i, x_j) = \tanh(\text{scale} \times \langle x_i, x_j \rangle + \text{bias})$
- ▶ the epanechnikov kernel:  $k(x_i, x_j) = \max(0, 1 - \gamma \langle x_i, x_j \rangle)$

## Validation protocol

- ▶ split randomly the dataset, 2/3 for cross validation, 1/3 for test.
- ▶ perform 10 fold cross validation on the validation set  
( $C \in [0.01, 0.1, 1, 10, 100, 1000]$ ,  
 $\sigma \in [0.1, 0.5, 1, 5, 10, 15, 25, 50, 100, 250, 500] * \sqrt{n}$  for rbf kernel,  
 $\text{scale} = [\text{pow}2(-5 : 1.5 : 2), -\text{pow}2(-5 : 1.5 : 2)]$  and  
 $\text{bias} = [\text{pow}2(-5 : 1.5 : 2), -\text{pow}2(-5 : 1.5 : 2)]$  for tanh kernel).
- ▶ train the svm on the full validation set with the parameters providing the best average performance during cross validation.
- ▶ test on the separate test set.

## Experimental results

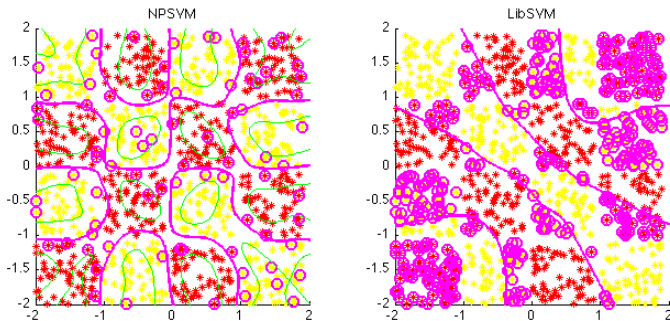
Solver	Kernel	Checkers	Checkers 10	Clown
C-SVM	rbf	87.1 % (51 sv)	79.5 % (151 sv)	99.93% (17 sv)
NPSVM	rbf	88.9 % (167 sv)	81.9 % (170 sv)	99.93% (53 sv)
Constraint-NPSVM	rbf	87.1 % (107 sv)	81.6 % (141 sv)	99.97% (54 sv)
NPSVM	tanh	74.9 % (35 sv)	71.8 % (41 sv)	99.97% (81 sv)
Constraint-NPSVM	tanh	86.6 % (114 sv)	81.6 % (145 sv)	99.87% (27 sv)
NPSVM	epanech	86.8 % (133 sv)	81.4 % (118 sv)	99.93% (32 sv)
Constraint-NPSVM	epanech	83.0 % (119 sv)	77.7 % (133 sv)	99.63% (56 sv)

**Table:** Results on synthetic dataset. Dataset sizes : 200 training points, 3000 testing points. Checkers 10 is a checker dataset with 10% of overlapping between classes. Clowns is also known as apple/banana.

**Table:** Results on some UCI dataset.

Solver	kernel	Heart	Sonar	Breast
C-SVM	rbf	82.22% (23.2 sv)	84.78% (90.3 sv)	97.47% (53.8 sv)
NPSVM	rbf	<b>83.44%</b> (35.9 sv)	<b>86.09%</b> (94.9 sv)	97.37% (56.6 sv)
NPSVM	tanh	<b>82.44%</b> (14.8 sv)	84.06% (70.1 sv)	<b>97.76 %</b> (116 sv)

# Comparison with libSVM

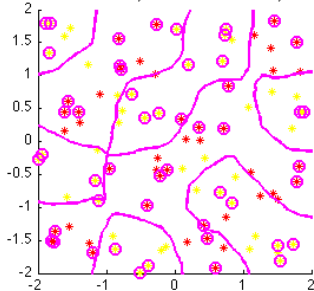


**Figure:** Results on checkers with NPSVM on the left and libSVM on the right, for an identical sigmoid kernel (scale = 2, bias = -2). Circles are support vectors.

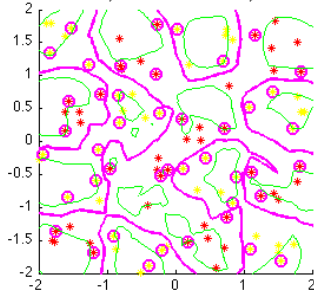


# Comparison to IndefiniteSVM

IndefiniteSVM, Perf = 83.8904 , 51 SV



NPSVM, Perf = 91.7447 , 44 SV



**Figure:** Results on checkers with IndefiniteSVM on the left and NPSVM on the right, for an identical epanech kernel. Circles are support vectors.

# Discussion

SVM with non positive kernels is possible

- ▶ representer theorem
- ▶ sparse solution
- ▶ efficient solver: NPSVM

SVM with non positive kernel is useful

- ▶ to be prove

# Road map

## 1 Introduction

## 2 Tuning the kernel: MKL

- The multiple kernel problem
- SimpleMKL: the multiple kernel solution

## 3 Non positive kernel

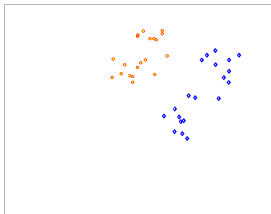
- NON Positive kernels
- Functional estimation in a RKKS
- Non positive SVM

## 4 Distribution shift

- Distribution shift: the problem
- Density ratio estimation principle

## 5 Conclusion

# Distribution shift: the problem

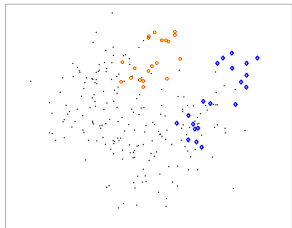


Training data

$$(\mathbf{x}_i^A, y_i^A), i = 1, n$$

i.i.d. from

$$\mathbb{P}_A(\mathbf{x}, y) = \mathbb{P}(y/\mathbf{x})\mathbb{P}_A(\mathbf{x})$$



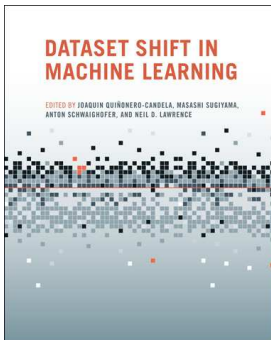
Test data

$$(\mathbf{x}_j^T, y_j^T), j = 1, \ell$$

i.i.d. from

$$\mathbb{P}_T(\mathbf{x}, y) = \mathbb{P}(y/\mathbf{x})\mathbb{P}_T(\mathbf{x})$$

# Distribution shift: references



- ▶ Masashi Sugiyama (Tokyo Institute of Technology) Density ratio estimation methods: Tutorial in ACML2009<sup>a</sup>
- ▶ Arthur Gretton (Max Planck Institute for Biological Cybernetics) Covariate Shift by Kernel Mean Matching Workshop at NIPS'09<sup>b</sup>
- ▶ Mahesan Niranjan (University of Southampton): application to intrusion detection

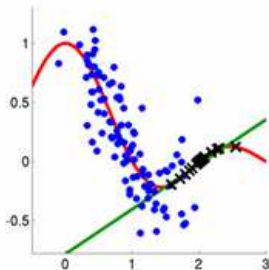
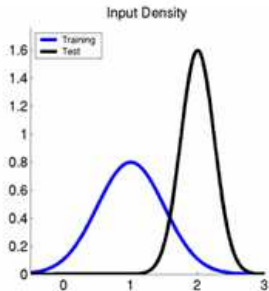
---

<sup>a</sup> [http://lamda.nju.edu.cn/conf/acml09/files/invited\\_sugi.pdf](http://lamda.nju.edu.cn/conf/acml09/files/invited_sugi.pdf)

<sup>b</sup> [http://videlectures.net/nipsworkshops09\\_gretton\\_cskm/](http://videlectures.net/nipsworkshops09_gretton_cskm/)

# Covariate Shift

Training and test input follow different distributions, but functional relation remains unchanged.



Goal: Estimate test output from  $\{(x_i, y_i)\}_{i=1}^n$

# Density ratio estimation principle

$$\begin{aligned}\min_{f \in \mathcal{H}} \mathbb{E}(J(X, Y)) &= \min_{f \in \mathcal{H}} \int_{\mathbf{x}} \int_y J(\mathbf{x}, y) \mathbb{P}_T(\mathbf{x}, y) \, d\mathbf{x} dy \\ &= \min_{f \in \mathcal{H}} \int_{\mathbf{x}} \int_y J(\mathbf{x}, y) \mathbb{P}(y|\mathbf{x}) \mathbb{P}_T(\mathbf{x}) \, d\mathbf{x} dy && \text{factorize} \\ &= \min_{f \in \mathcal{H}} \int_{\mathbf{x}} \left( \int_y J(\mathbf{x}, y) \mathbb{P}(y|\mathbf{x}) \, dy \right) \mathbb{P}_T(\mathbf{x}) \, d\mathbf{x} && \text{reorganize} \\ &= \min_{f \in \mathcal{H}} \int_{\mathbf{x}} \left( \int_y J(\mathbf{x}, y) \mathbb{P}(y|\mathbf{x}) \, dy \right) \mathbb{P}_T(\mathbf{x}) \frac{\mathbb{P}_A(\mathbf{x})}{\mathbb{P}_A(\mathbf{x})} \, d\mathbf{x} && \mathbb{P}_A(\mathbf{x}) \neq 0 \\ &= \min_{f \in \mathcal{H}} \int_{\mathbf{x}} \left( \int_y J(\mathbf{x}, y) \mathbb{P}(y|\mathbf{x}) \, dy \right) \mathbb{P}_T(\mathbf{x}) \frac{\mathbb{P}_A(\mathbf{x})}{\mathbb{P}_A(\mathbf{x})} \, d\mathbf{x} \\ &= \min_{f \in \mathcal{H}} \int_{\mathbf{x}} \left( \int_y J(\mathbf{x}, y) \mathbb{P}(y|\mathbf{x}) \, dy \right) w(\mathbf{x}) \mathbb{P}_A(\mathbf{x}) \, d\mathbf{x} && w(\mathbf{x}) = \frac{\mathbb{P}_T(\mathbf{x})}{\mathbb{P}_A(\mathbf{x})}\end{aligned}$$

## Importance weighting

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n J(\mathbf{x}_i, y_i) w(\mathbf{x}_i) \quad w(\mathbf{x}_i) = \frac{\mathbb{P}_T(\mathbf{x}_i)}{\mathbb{P}_A(\mathbf{x}_i)}$$

# Density ratio estimation principle

## the algorithm

1. estimate  $w(\mathbf{x}_i) = \frac{P_T(\mathbf{x}_i)}{P_A(\mathbf{x}_i)}$
2. solve a weighted version of our favorite learning algorithm

$$(\mathcal{PW}) \begin{cases} \min_{f \in \mathcal{H}, \alpha_0, \xi \in \mathbb{R}^n} & \frac{1}{2} \|f\|^2 + \frac{C}{p} \sum_{i=1}^n w(\mathbf{x}_i) \xi_i^p \\ \text{with} & y_i (f(\mathbf{x}_i) + \alpha_0) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad i = 1, n \end{cases}$$

$p = 1$ : L1 weighted SVM

$p = 2$ : L2 weighted SVM

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

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

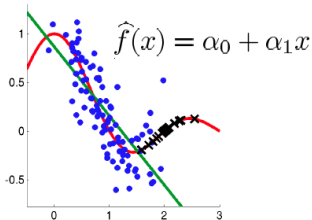


# Adaptation Using Density Ratios

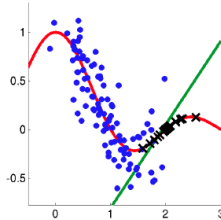
Shimodaira (JSPI2000), Sugiyama & Müller (ICANN2005, Stat&Deci2005)

- Ordinary least-squares is **not consistent**.
- Density-ratio weighted least-squares is **consistent**.

$$\min_{\alpha} \left[ \sum_{i=1}^n (\hat{f}(x_i) - y_i)^2 \right]$$



$$\min_{\alpha} \left[ \sum_{i=1}^n \frac{p_{\text{test}}(x_i)}{p_{\text{train}}(x_i)} (\hat{f}(x_i) - y_i)^2 \right]$$



From M. Sugiyama Tutorial in ACML2009

# Estimating the weights

$$w(\mathbf{x}_i) = \frac{\mathbb{P}_T(\mathbf{x}_i)}{\mathbb{P}_A(\mathbf{x}_i)}$$

- ▶ estimating the distributions –BAD–

1. Parzen (or other) estimate on test data  $\hat{\mathbb{P}}_T(\mathbf{x})$
2. Parzen (or other) estimate on training data  $\hat{\mathbb{P}}_A(\mathbf{x})$
3.  $\hat{w}(\mathbf{x}) = \frac{\hat{\mathbb{P}}_T(\mathbf{x})}{\hat{\mathbb{P}}_A(\mathbf{x})}$

- ▶ Direct estimation:  $w(\mathbf{x}) = \frac{\mathbb{P}_T(\mathbf{x})}{\mathbb{P}_A(\mathbf{x})} \iff w(\mathbf{x}) \mathbb{P}_A(\mathbf{x}) = \mathbb{P}_T(\mathbf{x})$ 
  - ▶ Kullback-Leibler Importance Estimation Procedure

$$\min_{\hat{w}} KL(\mathbb{P}_T(\mathbf{x}) \parallel \hat{w}(\mathbf{x}) \mathbb{P}_A(\mathbf{x}))$$

- ▶ Least-Squares Importance Fitting

$$\min_{\hat{w}} \mathbb{E}_A \left( \hat{w}(\mathbf{x}) - \frac{\mathbb{P}_T(\mathbf{x}_i)}{\mathbb{P}_A(\mathbf{x}_i)} \right)^2$$

# Estimating the weights

$$w(\mathbf{x}_i) = \frac{\mathbb{P}_T(\mathbf{x}_i)}{\mathbb{P}_A(\mathbf{x}_i)}$$

- ▶ estimating the distributions **-BAD-**

1. Parzen (or other) estimate on test data  $\hat{\mathbb{P}}_T(\mathbf{x})$
2. Parzen (or other) estimate on training data  $\hat{\mathbb{P}}_A(\mathbf{x})$
3.  $\hat{w}(\mathbf{x}) = \frac{\hat{\mathbb{P}}_T(\mathbf{x})}{\hat{\mathbb{P}}_A(\mathbf{x})}$

- ▶ Direct estimation:  $w(\mathbf{x}) = \frac{\mathbb{P}_T(\mathbf{x})}{\mathbb{P}_A(\mathbf{x})} \iff w(\mathbf{x}) \mathbb{P}_A(\mathbf{x}) = \mathbb{P}_T(\mathbf{x})$ 
  - ▶ Kullback-Leibler Importance Estimation Procedure

$$\min_{\hat{w}} KL(\mathbb{P}_T(\mathbf{x}) \parallel \hat{w}(\mathbf{x}) \mathbb{P}_A(\mathbf{x}))$$

- ▶ Least-Squares Importance Fitting

$$\min_{\hat{w}} \mathbb{E}_A \left( \hat{w}(\mathbf{x}) - \frac{\mathbb{P}_T(\mathbf{x}_i)}{\mathbb{P}_A(\mathbf{x}_i)} \right)^2$$

# Estimating the weights: empirical distribution

- ▶ Empirical distributions

$$\hat{\mathbb{P}}_T(\mathbf{x}) = \frac{1}{\ell} \sum_{j=1}^{\ell} \delta_{\mathbf{x}_j}$$

$$\hat{\mathbb{P}}_A(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i}$$

- ▶ test data (T)

$$\min_{\hat{w}} \int_{\mathbf{x}} \varphi(\hat{w}(\mathbf{x})) \mathbb{P}_T(\mathbf{x}) d\mathbf{x} \quad \longrightarrow \quad \min_{\hat{w}} \frac{1}{\ell} \sum_{j=1}^{\ell} \varphi(\hat{w}(\mathbf{x}_j))$$

- ▶ training data (A)

$$\min_{\hat{w}} \int_{\mathbf{x}} \varphi(\hat{w}(\mathbf{x})) \mathbb{P}_A(\mathbf{x}) d\mathbf{x} \quad \longrightarrow \quad \min_{\hat{w}} \frac{1}{n} \sum_{i=1}^n \varphi(\hat{w}(\mathbf{x}_i))$$

# Kullback-Leibler Importance Estimation

$$\begin{cases} \min_{\hat{w}} & KL(\mathbb{P}_T(\mathbf{x}) \parallel \hat{w}(\mathbf{x}) \mathbb{P}_A(\mathbf{x})) \\ \text{with} & \int_{\mathbf{x}} \hat{w}(\mathbf{x}) \mathbb{P}_A(\mathbf{x}) d\mathbf{x} = 1 \\ \text{and} & 0 \leq \hat{w}(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{X} \end{cases}$$

$$\begin{aligned} KL(\mathbb{P}_T(\mathbf{x}) \parallel \hat{w}(\mathbf{x}) \mathbb{P}_A(\mathbf{x})) &= \int_{\mathbf{x}} \mathbb{P}_T(\mathbf{x}) \log \frac{\mathbb{P}_T(\mathbf{x})}{\hat{w}(\mathbf{x}) \mathbb{P}_A(\mathbf{x})} d\mathbf{x} \\ &= \underbrace{\int_{\mathbf{x}} \mathbb{P}_T(\mathbf{x}) \log \frac{\mathbb{P}_T(\mathbf{x})}{\mathbb{P}_A(\mathbf{x})} d\mathbf{x}}_{\text{constant}} - \int_{\mathbf{x}} \mathbb{P}_T(\mathbf{x}) \log \hat{w}(\mathbf{x}) d\mathbf{x} \end{aligned}$$

$$\begin{cases} \min_{\hat{w}} & - \int_{\mathbf{x}} \mathbb{P}_T(\mathbf{x}) \log \hat{w}(\mathbf{x}) d\mathbf{x} \\ \text{with} & \int_{\mathbf{x}} \hat{w}(\mathbf{x}) \mathbb{P}_A(\mathbf{x}) d\mathbf{x} = 1 \\ \text{and} & 0 \leq \hat{w}(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{X} \end{cases} \quad \begin{cases} \min_{\hat{w}} & - \sum_{j=1}^{\ell} \log \hat{w}(\mathbf{x}_j) \\ \text{with} & \sum_{i=1}^n \hat{w}(\mathbf{x}_i) = n \\ \text{and} & 0 \leq \hat{w}(\mathbf{x}_i) \quad i = 1, n \end{cases}$$

# Kullback-Leibler Importance Estimation

Use a kernel representation

$$\hat{w}(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}, \mathbf{x}_i)$$

$$\left\{ \begin{array}{l} \min_{\hat{w} \in \mathcal{H}} - \sum_{j=1}^{\ell} \log \hat{w}(\mathbf{x}_j) \\ \text{with } \sum_{i=1}^n \hat{w}(\mathbf{x}_i) = n \\ \text{and } 0 \leq \hat{w}(\mathbf{x}_i) \quad i = 1, n \end{array} \right. \quad \left\{ \begin{array}{l} \min_{\alpha \in \mathbb{R}^n} - \sum_{j=1}^{\ell} \log K_T \alpha \\ \text{with } \mathbf{e}^T K_A \alpha = n \\ \text{and } 0 \leq \alpha_j \quad i = 1, n \end{array} \right.$$

Convex (non linear, non quadratic) problem with a sparse solution

# Least-Squares Importance Fitting

$$\begin{cases} \min_{\hat{w}} \mathbb{E}_A \left( \hat{w}(\mathbf{x}) - \frac{\mathbb{P}_T(\mathbf{x}_i)}{\mathbb{P}_A(\mathbf{x}_i)} \right)^2 \\ \text{with } \int_{\mathbf{x}} \hat{w}(\mathbf{x}) \mathbb{P}_A(\mathbf{x}) d\mathbf{x} = 1 \\ \text{and } 0 \leq \hat{w}(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{X} \end{cases}$$

$$\begin{aligned} \mathbb{E}_A \left( \hat{w}(\mathbf{x}) - \frac{\mathbb{P}_T(\mathbf{x}_i)}{\mathbb{P}_A(\mathbf{x}_i)} \right)^2 &= \int_{\mathbf{x}} \left( \hat{w}(\mathbf{x}) - \frac{\mathbb{P}_T(\mathbf{x}_i)}{\mathbb{P}_A(\mathbf{x}_i)} \right)^2 \mathbb{P}_A(\mathbf{x}) d\mathbf{x} \\ &= \int_{\mathbf{x}} \hat{w}(\mathbf{x})^2 \mathbb{P}_A(\mathbf{x}) d\mathbf{x} - 2 \int \hat{w}(\mathbf{x}) \mathbb{P}_T(\mathbf{x}) d\mathbf{x} + \text{constant} \end{aligned}$$

$$\begin{cases} \min_{\hat{w}} \frac{1}{2} \int_{\mathbf{x}} \hat{w}(\mathbf{x})^2 \mathbb{P}_A(\mathbf{x}) - \int \hat{w}(\mathbf{x}) \mathbb{P}_T(\mathbf{x}) \\ \text{with } \int_{\mathbf{x}} \hat{w}(\mathbf{x}) \mathbb{P}_A(\mathbf{x}) d\mathbf{x} = 1 \\ \text{and } 0 \leq \hat{w}(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{X} \end{cases} \quad \begin{cases} \min_{\hat{w}} \frac{1}{2} \sum_{i=1}^n \hat{w}(\mathbf{x}_i)^2 - \sum_{j=1}^{\ell} \hat{w}(\mathbf{x}_j) \\ \text{with } \sum_{i=1}^n \hat{w}(\mathbf{x}_i) = n \\ \text{and } 0 \leq \hat{w}(\mathbf{x}_i) \quad i = 1, n \end{cases}$$

# Least-Squares Importance Fitting

Use a feature space representation in a RKHS

$$\hat{w}(\mathbf{x}) = \sum_{k=1}^{\kappa} \alpha_k \phi_k(\mathbf{x})$$

$$\left\{ \begin{array}{l} \min_{\hat{w}} \quad \frac{1}{2} \sum_{i=1}^n \hat{w}(\mathbf{x}_i)^2 - \sum_{j=1}^{\ell} \hat{w}(\mathbf{x}_j) \\ \text{with} \quad \sum_{i=1}^n \hat{w}(\mathbf{x}_i) = n \\ \text{and} \quad 0 \leq \hat{w}(\mathbf{x}_i) \quad i = 1, n \end{array} \right. \quad \left\{ \begin{array}{l} \min_{\hat{w}} \quad \frac{1}{2} \alpha K \alpha - \mathbf{e}^T \Phi_T \alpha \\ \text{with} \quad \mathbf{e}^T \Phi_A \alpha = n \\ \text{and} \quad 0 \leq \alpha_i \quad i = 1, n \end{array} \right.$$

Convex quadratic program with a sparse solution



# Conclusion

- ▶ use weights
- ▶ compute weights at a reasonable cost
- ▶ solve weighted SVM
- ▶ provides a full RKHS embedding with representer theorem

# Road map

## 1 Introduction

## 2 Tuning the kernel: MKL

- The multiple kernel problem
- SimpleMKL: the multiple kernel solution

## 3 Non positive kernel

- NON Positive kernels
- Functional estimation in a RKKS
- Non positive SVM

## 4 Distribution shift

- Distribution shift: the problem
- Density ratio estimation principle

## 5 Conclusion

# what's new since 1995

## ▶ Applications

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

## ▶ Optimization

- ▶ dual: [mloss.org](http://mloss.org)
- ▶ regularization path
- ▶ approximation
- ▶ primal

## ▶ Statistic

- ▶ proofs and bounds
- ▶ model selection
  - ▶ span bound
  - ▶ multikernel: tuning ( $k$  and  $\sigma$ )

# challenges: towards tough learning

- ▶ the size effect
  - ▶ ready to use: automatization
  - ▶ adaptative: on line context aware
  - ▶ beyond kenrels: deep learning
- ▶ Automatic and adaptive model selection
  - ▶ variable selection
  - ▶ kernel tuning: coarse-to-fine
  - ▶ hyperparametres:  $C$ , duality gap,  $\lambda$
- ▶  $\mathbb{P}$  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 Berlinet and Christine Thomas-Agnan, Reproducing Kernel Hilbert Spaces in Probability and Statistics, Kluwer Academic Publishers, 2003
- ▶ Marc Atteia et Jean Gaches , Approximation Hilbertienne - Splines, Ondelettes, Fractales, PUG, 1999