

# Sparsifying LS-SVM Models via $L_0$ -Norm Minimization

Jorge López<sup>1</sup>   Kris De Brabanter<sup>2</sup>   Johan A.K. Suykens<sup>2</sup>

<sup>1</sup>Dpto. de Ingeniería Informática and Instituto de Ingeniería del Conocimiento  
Universidad Autónoma de Madrid

<sup>2</sup>Departement Electrotechniek (ESAT)  
Katholieke Universiteit Leuven

Tuesday, September 28th 2010

# Contents

- 1 Introduction
- 2 Sparse LS-SVMs
- 3 Experiments
- 4 Discussion

# LS-SVMs primal

## Characteristics

- Simplification of standard SVMs via equality constraints.
- Common formulation for classification ( $\vec{y} \in \{+1, -1\}^N$ ) and regression ( $\vec{y} \in \mathbb{R}^N$ ).
- $\phi$  feature map of given Mercer kernel  $k(\vec{x}_i, \vec{x}_j) = \phi(\vec{x}_i) \cdot \phi(\vec{x}_j)$ .
- Lagrangian coefficients  $\alpha_i$  introduced to dualize.

## Primal

$$\begin{aligned} \min_{\vec{w}, b, \xi} \quad & \frac{1}{2} \|\vec{w}\|^2 + \frac{C}{2} \sum_{i=1}^N \xi_i^2 \\ \text{s.t.} \quad & \vec{w} \cdot \phi(\vec{x}_i) + b = y_i - \xi_i, \quad \forall i = 1, \dots, N \end{aligned} \quad (1)$$

## Lagrangian

$$\mathcal{L}(\vec{w}, b, \vec{\xi}, \vec{\alpha}) = \frac{1}{2} \|\vec{w}\|^2 + \frac{C}{2} \sum_{i=1}^N \xi_i^2 - \sum_{i=1}^N \alpha_i [\vec{w} \cdot \phi(\vec{x}_i) + b - y_i + \xi_i]$$

# LS-SVMs dual

## Characteristics

- Dual obtained by setting Lagrangian's derivatives to 0.
- Reduces to a KKT system of equations with  $\tilde{K}_{ij} = k(\vec{x}_i, \vec{x}_j) + \delta_{ij}/C$ .

## Derivatives

$$\frac{\partial L}{\partial \vec{w}} = \vec{w} - \sum_{i=1}^N \alpha_i \phi(\vec{x}_i) = 0 \quad \Rightarrow \quad \vec{w} = \sum_{i=1}^N \alpha_i \phi(\vec{x}_i),$$

$$\frac{\partial L}{\partial b} = -\sum_{i=1}^N \alpha_i = 0 \quad \Rightarrow \quad \sum_{i=1}^N \alpha_i = 0,$$

$$\frac{\partial L}{\partial \xi} = C\xi_i - \alpha_i = 0 \quad \Rightarrow \quad \alpha_i = C\xi_i.$$

## Dual Formulation

$$\left[ \begin{array}{c|c} 0 & \vec{1}^T \\ \hline \vec{1} & \tilde{K} \end{array} \right] \left[ \begin{array}{c} b \\ \vec{\alpha} \end{array} \right] = \left[ \begin{array}{c} 0 \\ \vec{y} \end{array} \right] \quad (2)$$

## Some observations

- Final decision function of test point  $\vec{x}$  given by  $f(\vec{x}) = \text{sign}(\vec{w} \cdot \phi(\vec{x}) + b)$  (classification) or  $f(\vec{x}) = \vec{w} \cdot \phi(\vec{x}) + b$  (regression).
- $\phi(\cdot)$  usually unknown, so change  $\vec{w} \cdot \phi(\vec{x})$  by  $\sum_i \alpha_i \phi(\vec{x}_i) \cdot \phi(\vec{x}) = \sum_i \alpha_i k(\vec{x}_i, \vec{x})$ .
- Thus, test time proportional to number of points with  $\alpha_i \neq 0$  (SVs).
- But  $\alpha_i = C\xi_i$ , so  $\alpha_i = 0$  only when  $\xi_i = 0$ .
- That means when  $\vec{x}_i$  lies exactly on its support hyperplane (classification) or when  $y_i$  is exactly the output estimation for  $\vec{x}_i$  (regression).
- Very unlikely to happen, so **in practice all patterns are SVs**.
- Is there any way to reduce the number of SVs without degrading the model?

## Previous approaches (1)

### Pruning after training, then retraining

- Patterns with smallest  $|\alpha_i|$  <sup>1</sup> (Performance can drop quickly).
- Correctly classified patterns and furthest from boundary <sup>2</sup> (Performance can drop quickly).
- Basing on the SMO algorithm while solving dual <sup>3</sup> (Only homogeneous LS-SVM).
- Single pattern that introduces smallest error when omitted <sup>4</sup> (Very costly computationally).

---

<sup>1</sup>J.A.K. Suykens, L. Lukas, J. Vandewalle. Sparse Approximation using Least Squares Support Vector Machines. Proceedings of the IEEE International Symposium on Circuits and Systems (ISCASS'2000), pp. 757–760, 2000.

<sup>2</sup>Y. Li, C. Lin, W. Zhang. Improved Sparse Least-Squares Support Vector Machine Classifiers, Neurocomputing 69 (13–15), pp. 1655–1658, 2006.

<sup>3</sup>X. Zeng, X.W. Chen. SMO-based Pruning Methods for Sparse Least Squares Support Vector Machines. IEEE Transactions on Neural Networks 16 (6), pp. 1541–1546, 2005.

<sup>4</sup>B.J. De Kruif, T.J.A. De Vries. Pruning Error Minimization in Least-Squares Support Vector Machines. IEEE Transactions on Neural Networks 14 (3), pp. 696–702, 2003.

## Previous approaches (2)

### Enforcing before training

- Hierarchical model with  $L_1$ -norm minimization <sup>5</sup>(Resulting problem difficult to solve).
- Searching for linearly independent subset of patterns <sup>6</sup>(Not always good results).
- Fixing the size of the final model <sup>7</sup>(Feature map only approximated).

---

<sup>5</sup>K. Pelckmans, J.A.K. Suykens, B. De Moor. Building Sparse Representations and Structure Determination on LS-SVM Substrates. *Neurocomputing* 64, pp. 137–159, 2005.

<sup>6</sup>J. Vallyon, G. Horvath. A Sparse Least Squares Support Vector Machine Classifier. *Proceedings of the IEEE International Joint Conference on Neural Networks (IJCNN'04)*, pp. 543–548, 2004.

<sup>7</sup>K. De Brabanter, J. De Brabanter, J.A.K. Suykens, B. De Moor. Optimal Fixed-Size Kernel Models for Large Data Sets. *Computational Statistics and Data Analysis* 54(6), pp. 1484–1504, 2010.

## $L_0$ -norm

### Basic facts

- Defined as limit of  $p$ -norms when  $p \rightarrow 0$ :

$$\|\vec{v}\|_0 = \lim_{p \rightarrow 0} \left( \sum_i |v_i|^p \right)^{\frac{1}{p}}$$

- Alternatively,  $\|\vec{v}\|_0 = |\{v_i : v_i \neq 0\}|$ .
- It counts number of non-zero elements, so minimizing it implies sparseness.
- In our context, we can think of minimizing  $\|\vec{w}\|_0$  or  $\|\vec{\alpha}\|_0$ .
- The former corresponds to using minimal number of features for prediction, the latter to using minimal number of patterns.
- Problems: **nonconvex, NP-hard to optimize**<sup>8</sup>.
- Solutions: approximations<sup>9</sup>, iterative procedures.

---

<sup>8</sup>E. Amaldi, V. Kann. On the Approximability of Minimizing Nonzero Variables or Unsatisfied Relations in Linear Systems. Theoretical Computer Science 209 (1–2), pp. 237–260, 1998.

<sup>9</sup>J. Weston, A. Elisseeff, B. Schölkopf, M. Tipping. Use of the Zero Norm with Linear Models and Kernel Methods. Journal of Machine Learning Research 3, pp. 1439–1461, 2003.



## An iterative approach for LS-SVMs (1)

### Primal problem

$$\begin{aligned} \min_{\vec{v}^t, b^t, \vec{\xi}^t} \quad & \frac{1}{2} \sum_{i=1}^N \lambda_i^t (v_i^t)^2 + \frac{C}{2} \sum_{i=1}^N (\xi_i^t)^2 \\ \text{s.t.} \quad & \sum_j v_j^t k(\vec{x}_i, \vec{x}_j) + b^t = y_i - \xi_i^t, \quad \forall i = 1, \dots, N. \end{aligned} \quad (3)$$

### Observations

- Implicit vector  $\vec{w}^t = \sum_j v_j^t \phi(\vec{x}_j)$  still underlying in constraints.
- Regularization not on  $\|\vec{w}^t\|^2$ , but on the weighted  $L_2$ -norm of coefficients  $v_j^t$ .
- These  $v_j^t$  are no longer Lagrangian coefficients.
- Weights of regularization given by  $\lambda_i^t$ .
- Solution of (3) gives  $\vec{v}^{t+1}, b^{t+1}$ .
- Under probabilistic framework, can be shown that  $\lim_{t \rightarrow \infty} \vec{v}^t = \vec{v}^*$  and that  $\lim_{t \rightarrow \infty} \sum_i \lambda_i^t (v_i^t)^2 = \|\vec{v}^*\|_0$ , provided we update  $\lambda_i^{t+1} = \frac{1}{(v_i^{t+1})^2}^{10}$ .

<sup>10</sup>K. Huang, D. Zheng, J. Sun, Y. Hotta, K. Fujimoto, S. Naoi. Sparse Learning for Support Vector Classification.

## An iterative approach for LS-SVMs (2)

### Algorithm (IS-LSSVM)

- 1 Compute the LS-SVM solution  $\vec{\alpha}$  for given sample, kernel and  $C$  using (2).
- 2 Set  $t \leftarrow 0$ ,  $\vec{\lambda}^0 \leftarrow \alpha$  and  $\vec{v}^0 \leftarrow \vec{\infty}$ .
- 3 Solve problem (3) to give  $\vec{v}^{t+1}$  and  $b^{t+1}$ .
- 4 Update  $\lambda_i^{t+1} \leftarrow \frac{1}{(v_i^{t+1})^2}, i = 1, \dots, N$ .
- 5 Set  $t \leftarrow t + 1$  and go back to 3 until convergence.
- 6 Return model  $(v^t, b^t)$ .

### Remarks

- (3) can be solved either in primal or in dual (more on this later).
- Convergence criterion based on similarity between  $\vec{v}^t$  and  $\vec{v}^{t+1}$ , when  $\frac{1}{N} \|\vec{v}^t - \vec{v}^{t+1}\|^2 \leq \epsilon$ .
- Final result dependent on choice of  $\vec{\lambda}^0$ : only local optimum attained.

## Solution of the problem

### In the primal

- Writing (3) in matrix notation and substituting  $\xi^t$  yields

$$\min_{\vec{v}^t, b^t} \frac{1}{2} (\vec{v}^t)^T \text{diag}(\vec{\lambda}^t) \vec{v}^t + \frac{C}{2} (\vec{y} - K \vec{v}^t - b^t \vec{1})^T (\vec{y} - K \vec{v}^t - b^t \vec{1})$$

- Differentiating w.r.t.  $\vec{v}^t, b^t$  and equalling to 0 produces system

$$\left[ \begin{array}{c|c} N & \vec{1}^T K \\ \hline K^T \vec{1} & K^T K + \frac{\text{diag}(\vec{\lambda}^t)}{C} \end{array} \right] \begin{bmatrix} b^t \\ \vec{v}^t \end{bmatrix} = \begin{bmatrix} \vec{y}^T \vec{1} \\ K^T \vec{y} \end{bmatrix} \quad (4)$$

### In the dual

- Differentiating Lagrangian and equalling to 0 produces another system

$$\left[ \begin{array}{c|c} 0 & \vec{1}^T \\ \hline \vec{1} & K^T \text{diag}(\vec{\lambda}^t)^{-1} K + \frac{1}{C} \end{array} \right] \begin{bmatrix} b^t \\ \vec{\beta}^t \end{bmatrix} = \begin{bmatrix} 0 \\ \vec{y} \end{bmatrix} \quad (5)$$

- Same as (2), but with kernel switched to

$$\hat{k}(X_i, X_j) = \sum_m \frac{k(X_i, X_m) k(X_j, X_m)}{\lambda_m} + \frac{\delta_{ij}}{C}$$

- Primal variable recovered with  $\vec{v}^t = \text{diag}(\vec{\lambda}^t)^{-1} K \vec{\beta}^t$ .

## Experimental framework

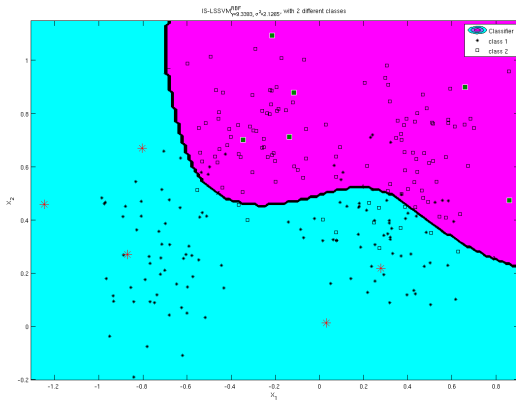
### Datasets

- Well-known datasets for classification (*Ripley*, *Fourclass*) and regression (*Motorcycle*, *Fossil*).
- All bidimensional for plotting purposes.

### Setting

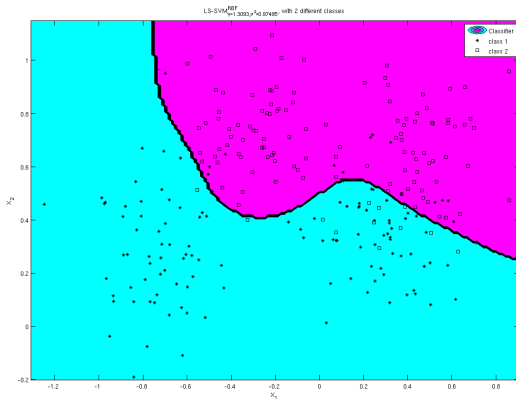
- RBF kernel used for all experiments.
- Stopping criterion with  $\epsilon = 10^{-4}$ .
- Hyperparameters  $C$  and  $\sigma$  for both ISLS-SVM and LS-SVM tuned with CSA (broad search) + simplex (fine search).
- CSA uses 10-fold CV as score function.
- Implemented in LS-SVM Matlab Toolbox (future version 2.0).

# Ripley dataset: ISLS-SVM



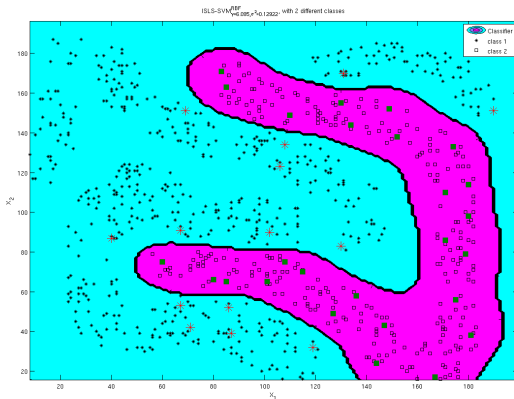
- 11 SVs out of 250 patterns.
- Smooth decision border.

# Ripley dataset: LS-SVM



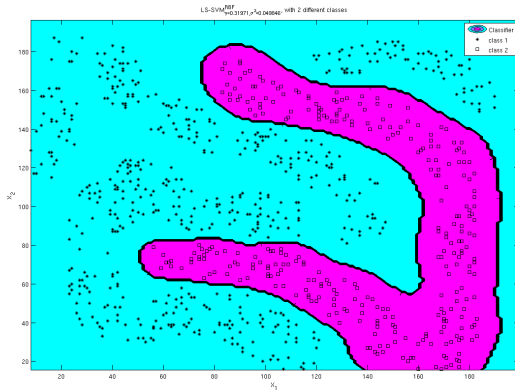
- All 250 patterns are SVs.
- Similar decision function, a bit fitter to the data.

## Fourclass dataset: ISLS-SVM



- 40 SVs out of 862 patterns, spread close to decision border.
- Easy problem (no error), but highly nonlinear.

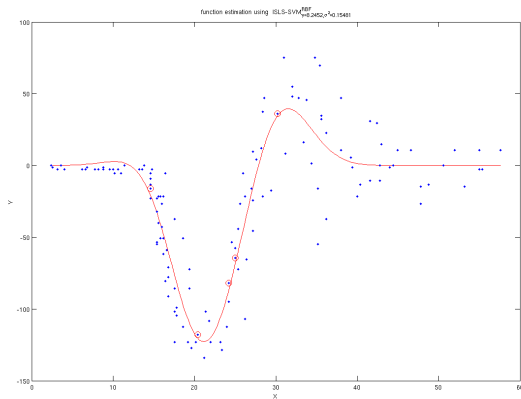
## Fourclass dataset: LS-SVM



- All 862 patterns are SVs.
- Also perfect performance and very similar decision function.

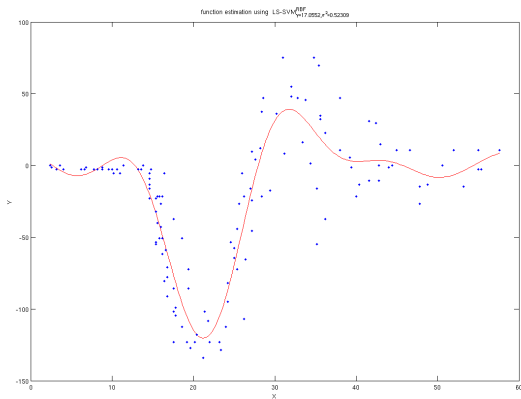


## Motorcycle dataset: ISLS-SVM



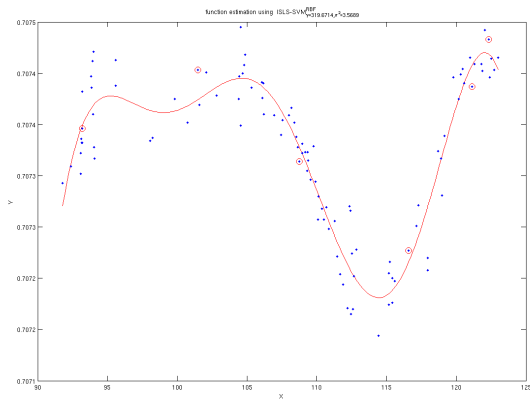
- 5 SVs out of 133 patterns.
- Difficult problem and highly nonlinear.

## Motorcycle dataset: LS-SVM



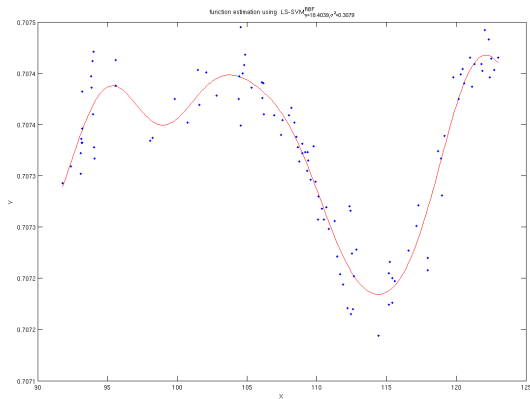
- All 133 patterns are SVs.
- Fitter to data at end, more oscillating at beginning.

## Fossil dataset: ISLS-SVM



- 6 SVs out of 106 patterns.
- Highly nonlinear problem.

## Fossil dataset: LS-SVM



- All 106 patterns are SVs.
- More oscillating at beginning.

# Summary

## Pros

- Straightforward algorithm to find sparse LS-SVM models.
- Small number of iterations (10-20 for most cases).
- Final models an order of magnitude (10-30x) sparser than LS-SVMs.
- Performance only degrades slightly.
- Applicable to other formulations e.g. Multiclass LS-SVMs.

## Cons

- Only **local minimum attained**, basing on initial weights given by LS-SVM.
- **Assumed model**  $W = \sum_i v_i \phi(X_i)$ .
- Computationally **costly** (each iteration is  $\mathcal{O}(N^3)$ ).
- Occasional **numerical instabilities**.

## (Possible) solutions (I)

### Numerical instabilities

- Solving (4) instead of (5).
- System looks more robust since  $\lambda$  and  $C$  are in the same term.

### Local minimum

- Apparently no solution, since finding global minimum is NP-hard.

### Computational cost

- $\mathcal{O}(N^3)$  cost cause every system solved with operator  $\setminus$ .
- Too costly for medium or large-scale problems.
- Possible to reduce to  $\mathcal{O}(N^2)$  using the SMO algorithm in the dual <sup>11</sup>.
- This also allows for caching kernel matrix, so large-scale can be issued.

---

<sup>11</sup>J. López, J.A.K. Suykens. First and Second Order SMO Algorithms for Large Scale LS-SVM Training. Internal Report 09-179, ESAT-SISTA, K.U. Leuven, 2009.

## (Possible) solutions (II)

### Assumption on $\vec{w}$

- Is it possible to remove assumption on  $\vec{w}$  in (3)?

$$\min_{\vec{w}^t, b^t, \vec{\xi}^t} \frac{1}{2} \sum_{i=1}^{n_F} \lambda_i^t (w_i^t)^2 + \frac{C}{2} \sum_{i=1}^N (\xi_i^t)^2$$

$$\text{s.t.} \quad \vec{w}^t \cdot \phi(\vec{x}_i) + b^t = y_i - \xi_i^t, \quad \forall i = 1, \dots, N.$$

- Solving in the primal we get

$$\left[ \begin{array}{c|c} N & \vec{1}^T \Phi \\ \hline \Phi^T \vec{1} & \Phi^T \Phi + \frac{\text{diag}(\vec{\lambda}^t)}{C} \end{array} \right] \begin{bmatrix} b^t \\ \vec{w}^t \end{bmatrix} = \begin{bmatrix} \vec{y}^T \vec{1} \\ \Phi^T \vec{y} \end{bmatrix}$$

- Same as (4), but with matrix  $\Phi^T = (\phi(X_1) \dots \phi(X_N))$ .
- Intractable cause  $\Phi$  usually unknown, unless we have an estimate  $\hat{\Phi}$ .
- Solving in the dual we get

$$\left[ \begin{array}{c|c} 0 & \vec{1}^T \\ \hline \vec{1} & \Phi \text{diag}(\vec{\lambda}^t)^{-1} \Phi^T + \frac{1}{C} \end{array} \right] \begin{bmatrix} b^t \\ \vec{\beta}^t \end{bmatrix} = \begin{bmatrix} 0 \\ \vec{y} \end{bmatrix}$$

- Same as (5), but with matrix  $\Phi^T \text{diag}(\vec{\lambda}^t)^{-1} \Phi$  instead of  $K$ .
- Intractable cause  $\Phi \Phi^T = K$ , but  $\Phi \text{diag}(\vec{\lambda}^t)^{-1} \Phi^T$  usually unknown, unless we have an estimate  $\hat{\Phi}$ .

# Farewell

Thank you for your attention!