

Decomposition Algorithms for Training Large-scale Semiparametric Support Vector Machines

Sangkyun Lee and Stephen J. Wright

Department of Computer Sciences, University of Wisconsin-Madison, USA



Abstract: We describe a method for solving large-scale semiparametric SVMs for regression problems. Most of the approaches proposed to date for large-scale SVMs cannot accommodate the multiple equality constraints that appear in semiparametric problems. Our approach uses a decomposition framework, with a primal-dual algorithm to find an approximate saddle point for the min-max formulation of each subproblem. We compare our method with algorithms previously proposed for semiparametric SVMs, and show that it scales well as the number of training examples grows.

SEMPARAMETRIC SVR (SP-SVR)

SVR with ϵ -insensitive loss:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{C}{M} \sum_{i=1}^M \ell_\epsilon(h; \mathbf{x}_i, \mathbf{y}_i),$$

$$\ell_\epsilon(h; \mathbf{x}_i, \mathbf{y}_i) := \max\{|\mathbf{y}_i - h(\mathbf{x}_i)| - \epsilon, 0\}.$$

- Dataset $\mathcal{D} = \{\mathbf{x}_i, \mathbf{y}_i\}, i = 1, \dots, M,$
- $\mathbf{x}_i \in \mathbb{R}^n, \mathbf{y}_i \in \mathbb{R}.$

- Kernel $\kappa(\mathbf{x}_i, \mathbf{x}_j) := \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle.$

Standard (nonparametric) SVR:

$$h(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b, \quad b \in \mathbb{R}.$$

Semiparametric SVR:

$$h(\mathbf{x}) = \underbrace{\langle \mathbf{w}, \phi(\mathbf{x}) \rangle}_{\text{Nonparametric}} + \sum_{j=1}^K \underbrace{\beta_j \psi_j(\mathbf{x})}_{\text{Parametric}}.$$

PRIMAL FORMULATION

$$\min_{\mathbf{w}, \beta, \xi, \xi^*} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{M} \sum_{i=1}^M (\xi_i + \xi_i^*) \quad (1a)$$

$$\text{s.t.} \quad \mathbf{y}_i - \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle - \sum_{j=1}^K \beta_j \psi_j(\mathbf{x}_i) \leq \epsilon + \xi_i, \quad i = 1, \dots, M \quad (1b)$$

$$-\left[\mathbf{y}_i - \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle - \sum_{j=1}^K \beta_j \psi_j(\mathbf{x}_i) \right] \leq \epsilon + \xi_i^*, \quad i = 1, \dots, M \quad (1c)$$

$$\xi \geq \mathbf{0}, \xi^* \geq \mathbf{0}. \quad (1d)$$

DUAL FORMULATION

$$\min_{\mathbf{z}} F(\mathbf{z}) := \frac{1}{2} \mathbf{z}^T \mathbf{Q} \mathbf{z} + \mathbf{p}^T \mathbf{z} \quad (2)$$

$$\text{s.t.} \quad \mathbf{A} \mathbf{z} = \mathbf{0},$$

$$\mathbf{0} \leq \mathbf{z} \leq (C/M) \mathbf{1},$$

- $n := 2M,$
- $\mathbf{z}, \mathbf{p}, \mathbf{1} = \{1, \dots, 1\}^T \in \mathbb{R}^n,$
- $\mathbf{Q} \in \mathbb{R}^{n \times n}$ p.s.d.,
- $\mathbf{A} \in \mathbb{R}^{K \times n}, K \geq 1.$

$\mathbf{z} = \begin{bmatrix} \alpha \\ \alpha^* \end{bmatrix} \in \mathbb{R}^{2M}$ for the dual vectors of (1b) and (1c),

$\mathbf{p} = [\epsilon - \mathbf{y}_1, \dots, \epsilon - \mathbf{y}_M, \epsilon + \mathbf{y}_1, \dots, \epsilon + \mathbf{y}_M]^T \in \mathbb{R}^{2M},$

$\mathbf{Q}_{ij} = \begin{cases} \mathbf{y}_i \mathbf{y}_j^T \kappa(\mathbf{x}_i, \mathbf{x}_j) & \text{if } 1 \leq i, j \leq M, \text{ or } M+1 \leq i, j \leq 2M \\ -\mathbf{y}_i \mathbf{y}_j^T \kappa(\mathbf{x}_i, \mathbf{x}_j) & \text{otherwise} \end{cases}$

$\mathbf{A} = \begin{bmatrix} \psi_1(\mathbf{x}_1) & \dots & \psi_1(\mathbf{x}_M) & -\psi_1(\mathbf{x}_1) & \dots & -\psi_1(\mathbf{x}_M) \\ \psi_2(\mathbf{x}_1) & \dots & \psi_2(\mathbf{x}_M) & -\psi_2(\mathbf{x}_1) & \dots & -\psi_2(\mathbf{x}_M) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \psi_K(\mathbf{x}_1) & \dots & \psi_K(\mathbf{x}_M) & -\psi_K(\mathbf{x}_1) & \dots & -\psi_K(\mathbf{x}_M) \end{bmatrix} \in \mathbb{R}^{K \times 2M}.$

MOST SOLVERS NEED $K = 1$

SVM^{light} [3], GPDT [7, 8]

- Working-set selection is for $K = 1.$
- Subproblem solver requires $K = 1.$

SMO [6], LIBSVM [1]

- Working-set selection assumes $K = 1.$
- The analytic solutions of the subproblems are available only when $K = 1.$

DECOMPOSITION ALGORITHM

- In each **outer** iteration, split variables \mathbf{z} into
 - Basic variables $\mathbf{z}_B, B \subset \{1, 2, \dots, n\}.$
 - Nonbasic variables $\mathbf{z}_N, N = \{1, 2, \dots, n\} \setminus B.$
 - B is our **working set**, of which the size $n_B \ll n.$
- Fix $\mathbf{z}_N,$ change $\mathbf{z}_B.$
- Given $\mathbf{z}^k = (\mathbf{z}_B^k, \mathbf{z}_N^k),$ solve the subproblem

$$\min_{\mathbf{z}_B} f(\mathbf{z}_B) := \frac{1}{2} \mathbf{z}_B^T \mathbf{Q}_{BB} \mathbf{z}_B + (\mathbf{Q}_{BN} \mathbf{z}_N^k + \mathbf{p}_B)^T \mathbf{z}_B$$

$$\text{s.t.} \quad \mathbf{A}_B \mathbf{z}_B = -\mathbf{A}_N \mathbf{z}_N^k + \mathbf{b},$$

$$0 \leq \mathbf{z}_B \leq (C/M) \mathbf{1}.$$

- $\mathbf{z}^{k+1} \leftarrow (\mathbf{z}_B^{k+1}, \mathbf{z}_N^k).$

I. WORKING-SET SELECTION

- $n_B:$ working set size.
- $n_c:$ max. number of “fresh” indices, $\ll n_B.$

Consider Lagrangian relaxation \mathcal{L} of (2),

$$\mathcal{L}(\mathbf{z}; \boldsymbol{\eta}) = F(\mathbf{z}) + \boldsymbol{\eta}^T \mathbf{A} \mathbf{z}.$$

Given $(\mathbf{z}^k, \boldsymbol{\eta}^k),$ find a solution \mathbf{d} of

$$\min_{\mathbf{d}} \left(\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}^k; \boldsymbol{\eta}^k) \right)^T \mathbf{d}$$

$$\text{s.t.} \quad \begin{aligned} 0 \leq \mathbf{d}_i \leq 1 & \quad \text{if } \mathbf{z}_i^{k+1} = 0, \\ -1 \leq \mathbf{d}_i \leq 0 & \quad \text{if } \mathbf{z}_i^{k+1} = C/M, \\ -1 \leq \mathbf{d}_i \leq 1 & \quad \text{if } \mathbf{z}_i^{k+1} \in (0, C/M), \\ \#\{\mathbf{d}_i | \mathbf{d}_i \neq 0\} & \leq n_c. \end{aligned} \quad (3)$$

- Solved efficiently in $\mathcal{O}(n \log n).$
- Convergence proofs in [5, 10].

II. PRIMAL-DUAL SUBPROBLEM SOLVER (PDSG)

Consider

$$\max_{\boldsymbol{\eta}} \min_{\mathbf{z}_B \in \Omega} \tilde{\mathcal{L}}(\mathbf{z}_B, \boldsymbol{\eta}),$$

where

$$\tilde{\mathcal{L}}(\mathbf{z}_B, \boldsymbol{\eta}) := f(\mathbf{z}_B) + \boldsymbol{\eta}^T (\mathbf{A}_B \mathbf{z}_B + \mathbf{A}_N \mathbf{z}_N^k),$$

$$\Omega = \{\mathbf{z}_B \in \mathbb{R}^{n_B} | \mathbf{0} \leq \mathbf{z}_B \leq (C/M) \mathbf{1}\}.$$

In each **inner** iteration, update primal and dual variables by,

$$\begin{cases} \mathbf{z}_B^{\ell+1} \leftarrow \mathbf{z}_B^\ell + s(\mathbf{z}_B^\ell, \boldsymbol{\eta}^\ell) \\ \boldsymbol{\eta}^{\ell+1} \leftarrow \boldsymbol{\eta}^\ell + t(\mathbf{z}_B^{\ell+1}, \boldsymbol{\eta}^\ell), \end{cases}$$

- $s(\cdot, \cdot):$ two-metric GP [2] followed by line-search, on a sub-workingset of size 2.
- $t(\cdot, \cdot): \nabla_{\boldsymbol{\eta}} \tilde{\mathcal{L}},$ scaled by dual Hessian diagonal [4], on a sub-workingset of size 2.

III. UPDATE

- Update primal-dual iterate pair $(\mathbf{z}^{k+1}, \boldsymbol{\eta}^{k+1}).$
 - $\mathbf{z}^{k+1} \leftarrow (\mathbf{z}_B^{k+1}, \mathbf{z}_N^k).$
 - $\boldsymbol{\eta}^{k+1}$ is provided by the subproblem solver.
- Incremental **full gradient** $\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}; \boldsymbol{\eta})$ update,

$$\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}^{k+1}, \boldsymbol{\eta}^{k+1})$$

$$= \nabla F(\mathbf{z}^k) + \begin{bmatrix} \mathbf{Q}_{BB} \\ \mathbf{Q}_{NB} \end{bmatrix} (\mathbf{z}_B^{k+1} - \mathbf{z}_B^k) + (\boldsymbol{\eta}^{k+1})^T \mathbf{A}.$$

EXPERIMENT: TOY PROBLEM

Modified Mexican hat function [4, 9]:

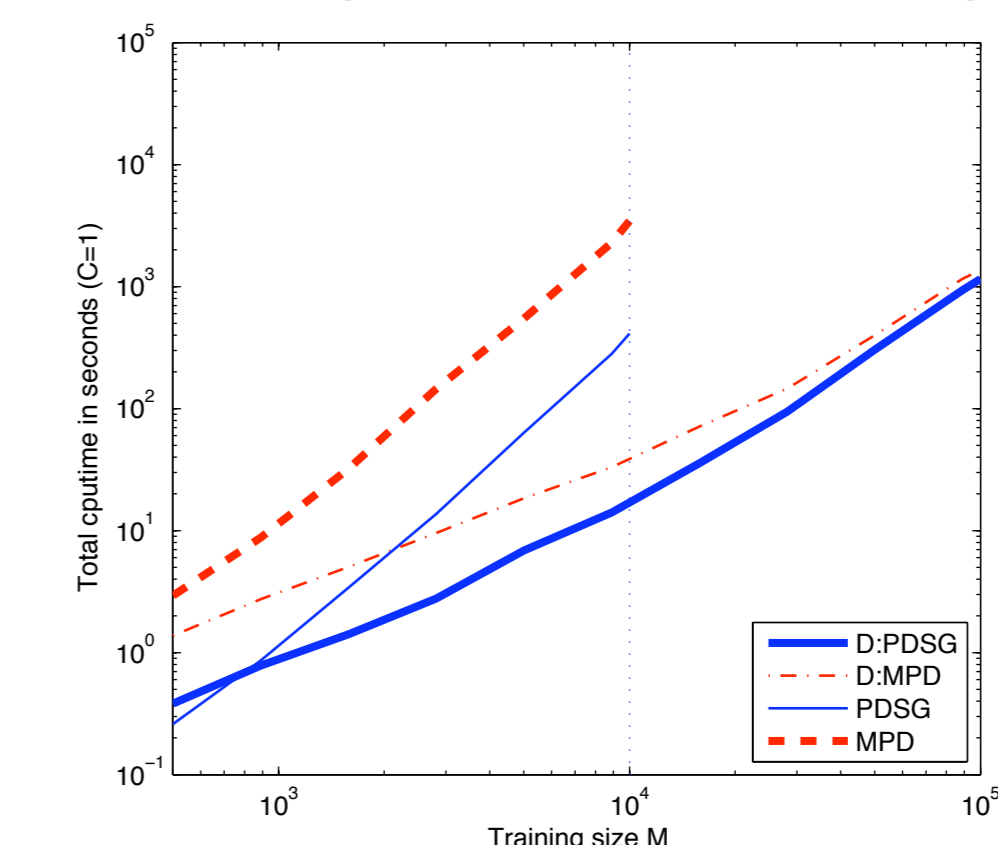
$$\omega(x) = \sin(x) + \text{sinc}(2\pi(x - 5)).$$

- $\mathbf{y}_i = \omega(x_i) + \zeta_i; x_i \sim \text{Uniform}[0, 10], \zeta_i \sim \mathcal{N}(0, 0.2^2).$
- $\psi_1(x) = \sin(x), \psi_2(x) = \text{sinc}(2\pi(x - 5)).$
- Gaussian kernel $\kappa(x, y) = \exp(-0.25|x - y|^2).$
- $\epsilon = 0.05, n_B = 500, n_c = 100.$

Compare to the current best solver MPD [4].

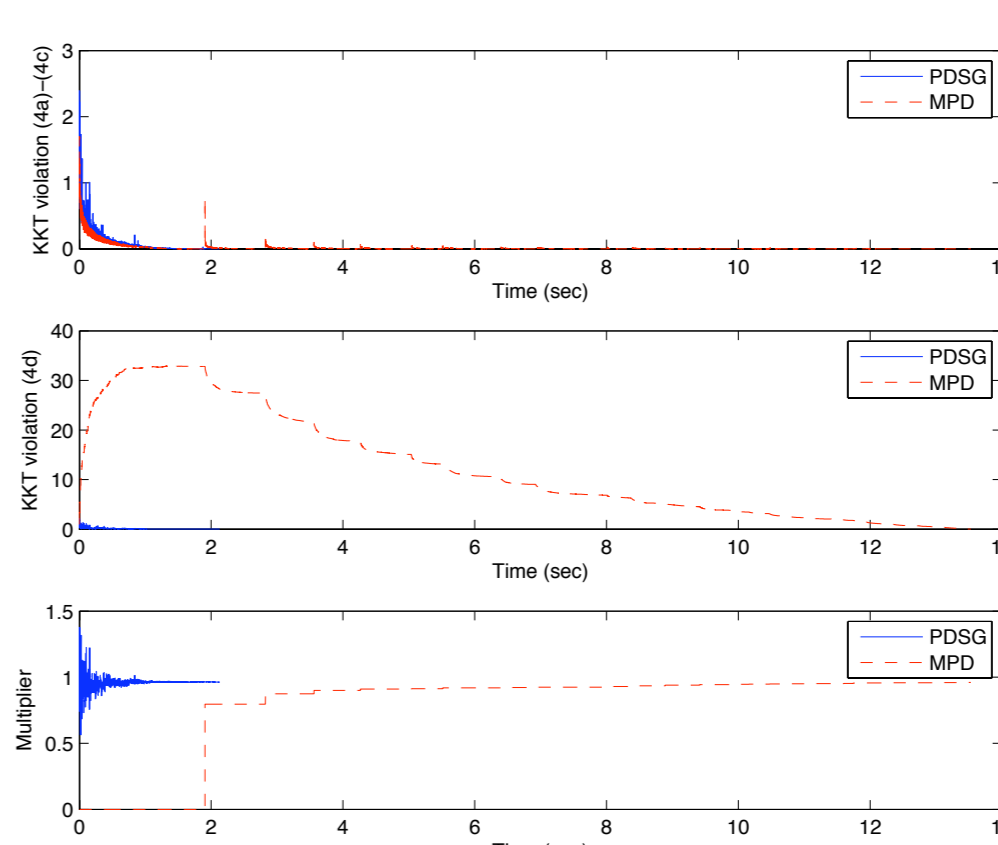
- No decomposition.
- Primal-dual method (method of multipliers).
- Gradient projection (primal), scaled gradient ascent (dual) on (sub-)workingset of size 1.

Scaling w.r.t. Training Size



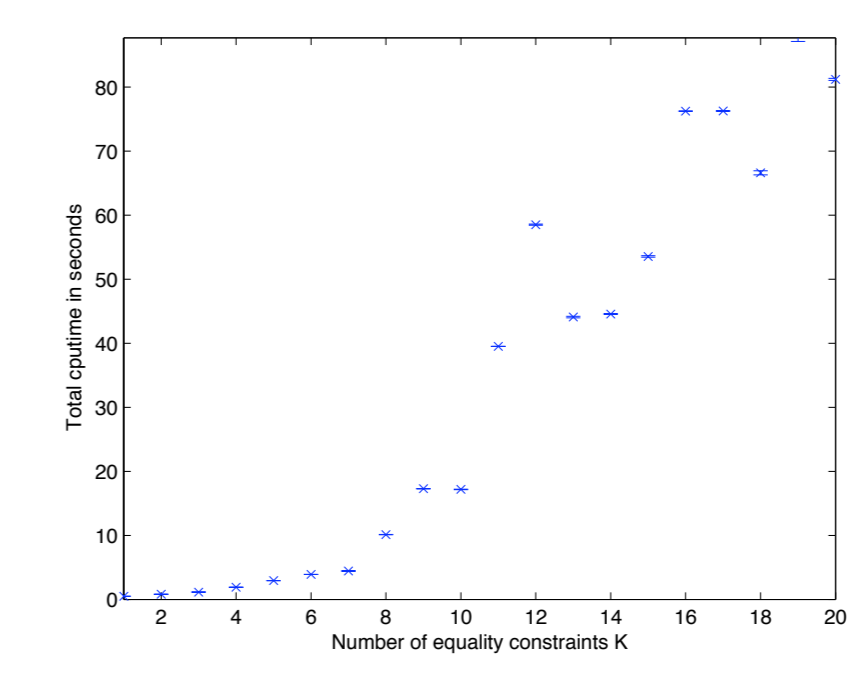
- PDSG vs. MPD (stand-alone).
- D:PDSG vs. D:MPD (in decomposition).
- D:MPD catches up D:PDSG when $M \uparrow:$ the full gradient update step becomes dominant as $M \uparrow.$

Convergence ($M = 1000$)



- PDSG vs. MPD (stand-alone).
- (T) max. dual feasibility violation.
- (M) max. primal equality constraints violation.
- (B) the first parametric component coeff.

Scaling of D:PDSG w.r.t. K ($M = 1000$)



- Time complexity of D:PDSG is $\mathcal{O}(uKn_B),$ u is the no. of outer iterations.
- Solver time seems to increase linearly with K for $K \geq 6.$

$$\psi_j(x) = \begin{cases} \cos(j\pi x), & j = 0, 2, 4, \dots \\ \sin(j\pi x), & j = 1, 3, 5, \dots \end{cases}$$

EXPERIMENT: REAL DATASET

Air pollution effect on respiratory illness in Milan, Italy, '80–'89 [11].

- 3652 daily records of env. conditions and the no. of respiratory deaths. Features are rescaled to $[0, 1].$
- Semiparametric model to predict the no. of deaths,

$$h_{sp}(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle + \beta_1(\mathbf{x}_{temp}) + \beta_2(\mathbf{x}_{SO_2}) + \beta_3(\mathbf{x}_{temp})^2 + \beta_4(\mathbf{x}_{SO_2})^2 + \beta_5,$$

- The fitted model is,

$$h_{sp}(\mathbf{x}) = \langle \mathbf{w}^*, \phi(\mathbf{x}) \rangle + 0.22(\mathbf{x}_{temp} - 0.47)^2 + 0.26(\mathbf{x}_{SO_2}) - 0.07(\mathbf{x}_{SO_2})^2 + 0.12.$$

- Deaths are lower in the middle of the **temp.** range.
- Linear increase of death rate with **SO₂ level.**

SUMMARY

- Introduced an efficient algorithm for solving a generalized dual formulation of SVMs, which has multiple equality constraints.
- Easily extends to other SVMs, e.g. ν -SVMs, semiparametric extensions of SVMs, etc.
- Future work: theoretical analysis, fine-tuning working-set selection, stochastic approximation alternative method.

Decomposition Algorithms for Training Large-scale Semiparametric Support Vector Machines

Sangkyun Lee and Stephen J. Wright

Department of Computer Sciences, University of Wisconsin-Madison, USA



THE UNIVERSITY
of
WISCONSIN
MADISON

- [1] Chih-Chung Chang and Chih-Jen Lin. *LIBSVM: a library for support vector machines*, April 2009. version 2.89, <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [2] E. M. Gafni and D. P. Bertsekas. Two-metric projection methods for constrained optimization. *SIAM Journal on Control and Optimization*, 22:936–964, 1984.
- [3] Thorsten Joachims. Making large-scale support vector machine learning practical. In B. Schölkopf, C. Burges, and A. Smola, editors, *Advances in Kernel Methods - Support Vector Learning*, chapter 11, pages 169–184. MIT Press, Cambridge, MA, 1999.
- [4] Wolf Kienzle and Bernhard Schölkopf. Training support vector machines with multiple equality constraints. In *Machine Learning: ECML 2005*, volume 16, October 2005.
- [5] C. J. Lin. Linear convergence of a decomposition method for support vector machines. Technical report, Department of Computer Science and Information Engineering, National Taiwan University, 2001.
- [6] John C. Platt. Fast training of support vector machines using sequential minimal optimization. In B. Schölkopf, C. Burges, and A. Smola, editors, *Advances in Kernel Methods - Support Vector Learning*, chapter 12, pages 185–208. MIT Press, Cambridge, MA, 1999.
- [7] T. Serafini, G. Zanghirati, and L. Zanni. Gradient projection methods for large quadratic programs and applications in training support vector machines. *Optimization Methods and Software*, 20(2–3):353–378, 2004.
- [8] T. Serafini and L. Zanni. On the working set selection in gradient projection-based decomposition techniques for support vector machines. *Optimization Methods and Software*, 20:583–596, 2005.
- [9] Alex J. Smola, Thilo T. Frieß, and Bernhard Schölkopf. Semiparametric support vector and linear programming machines. In *Advances in Neural Information Processing Systems 11*, pages 585–591, Cambridge, MA, USA, 1999. MIT Press.
- [10] P. Tseng and S. Yun. A coordinate gradient descent method for linearly constrained smooth optimization and support vector machines training. Published online in *Computational Optimization and Applications*, October 2008.
- [11] M A Vigotti, G Rossi, L Bisanti, A Zanobetti, and J Schwartz. Short term effects of urban air pollution on respiratory health in Milan, Italy, 1980-89. *Journal of Epidemiology Community Health*, 50:s71–75, 1996.