

Decomposition and Stochastic Subgradient Algorithms for Support Vector Machines

Sangkyun Lee and Stephen J. Wright

University of Wisconsin-Madison

ISMP 2009



Support Vector Machines

- Support Vector Machines (SVMs) are popular in many areas.

Support Vector Machines

- Support Vector Machines (SVMs) are popular in many areas.
 - Decision making, Machine learning, Statistics.
 - Bio-informatics, Neuroscience, Geophysics ...

Support Vector Machines

- Support Vector Machines (SVMs) are popular in many areas.
 - Decision making, Machine learning, Statistics.
 - Bio-informatics, Neuroscience, Geophysics ...
- For classification, regression and many other tasks.

Support Vector Machines

- Support Vector Machines (SVMs) are popular in many areas.
 - Decision making, Machine learning, Statistics.
 - Bio-informatics, Neuroscience, Geophysics ...
- For classification, regression and many other tasks.
- Result in two different types of convex programs,

“Primal” {
Number of variables = length of an input vector.
Obj. consists of a quadratic term and a piecewise linear function.
Costly obj. function evaluation with many input points.

“Dual” {
Number of variables = number of input points.
QP with dense and ill-conditioned Hessian.
A single equality constraint and bound constraints.

Support Vector Machines

- Support Vector Machines (SVMs) are popular in many areas.
 - Decision making, Machine learning, Statistics.
 - Bio-informatics, Neuroscience, Geophysics ...
- For classification, regression and many other tasks.
- Result in two different types of convex programs,

“Primal” {
Number of variables = length of an input vector.
Obj. consists of a quadratic term and a piecewise linear function.
Costly obj. function evaluation with many input points.

“Dual” {
Number of variables = number of input points.
QP with dense and ill-conditioned Hessian.
A single equality constraint and bound constraints.

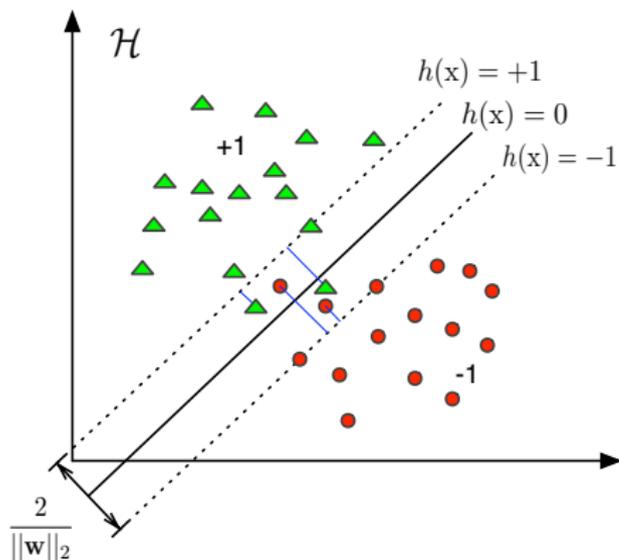


SVMs for Classification (SVC)

■ $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^M$ i.i.d. $\sim P(X, Y)$,

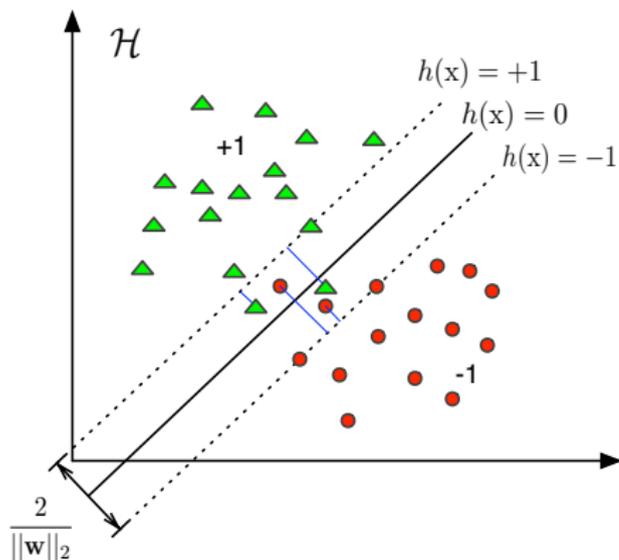
■ $\mathbf{x}_i \in \mathbb{R}^N$.

■ $\mathbf{y}_i \in \{-1, +1\}$.

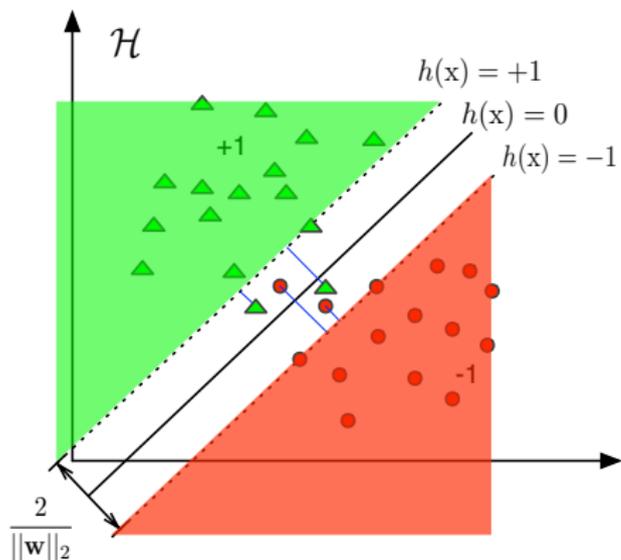


SVMs for Classification (SVC)

- $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^M$ *i.i.d.* $\sim P(X, Y)$,
 - $\mathbf{x}_i \in \mathbb{R}^N$.
 - $\mathbf{y}_i \in \{-1, +1\}$.
- $\phi : \mathbb{R}^N \rightarrow \mathcal{H}$.

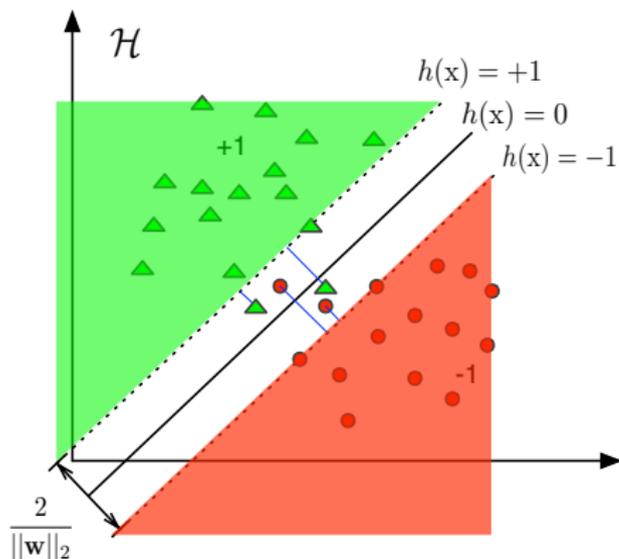


SVMs for Classification (SVC)



- $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^M$ *i.i.d.* $\sim P(X, Y)$,
 - $\mathbf{x}_i \in \mathbb{R}^N$.
 - $\mathbf{y}_i \in \{-1, +1\}$.
- $\phi : \mathbb{R}^N \rightarrow \mathcal{H}$.
- Find a classifier
$$h(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b,$$
 - $h(\mathbf{x}_i) \geq +1$ for $\mathbf{y}_i = +1$,
 - $h(\mathbf{x}_i) \leq -1$ for $\mathbf{y}_i = -1$,
 - Maximizing the “margin” $2/\|\mathbf{w}\|_2$.

SVMs for Classification (SVC)

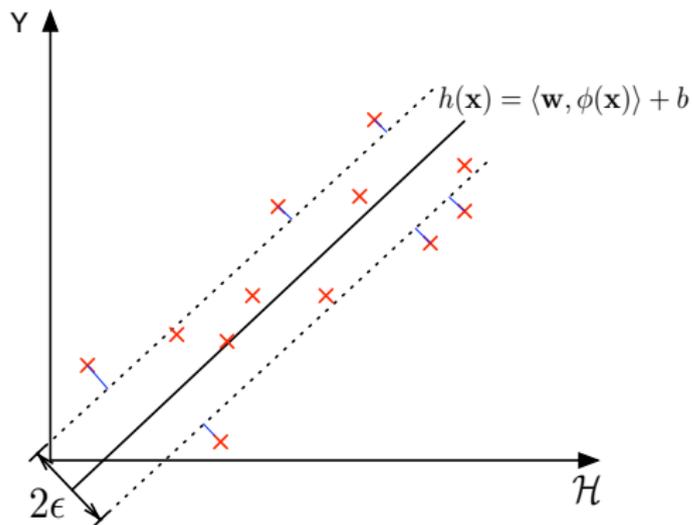


- $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^M$ i.i.d. $\sim P(X, Y)$,
 - $\mathbf{x}_i \in \mathbb{R}^N$.
 - $\mathbf{y}_i \in \{-1, +1\}$.
- $\phi : \mathbb{R}^N \rightarrow \mathcal{H}$.
- Find a classifier $h(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b$,
 - $h(\mathbf{x}_i) \geq +1$ for $\mathbf{y}_i = +1$,
 - $h(\mathbf{x}_i) \leq -1$ for $\mathbf{y}_i = -1$,
 - Maximizing the “margin” $2/\|\mathbf{w}\|_2$.

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

Hinge loss: $\ell_{\mathcal{H}}(h; \mathbf{x}_i, \mathbf{y}_i) := \max\{1 - \mathbf{y}_i h(\mathbf{x}_i), 0\}$.

SVMs for Regression (SVR)

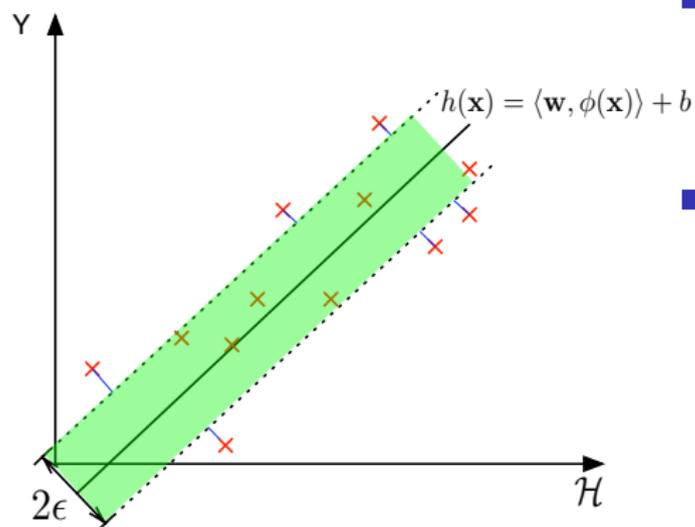


■ $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^M$ i.i.d $\sim P(X, Y)$,

■ $\mathbf{x}_i \in \mathbb{R}^N$.

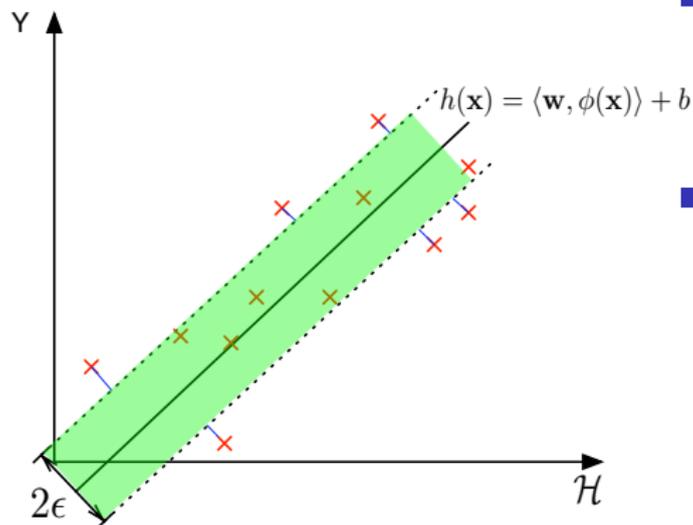
■ $\mathbf{y}_i \in \mathbb{R}$.

SVMs for Regression (SVR)



- $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^M$ i.i.d $\sim P(X, Y)$,
 - $\mathbf{x}_i \in \mathbb{R}^N$.
 - $\mathbf{y}_i \in \mathbb{R}$.
- Find a regression function,
 $h(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b$
 - Minimizing prediction error.
 - Capture data points in an ϵ -radius hyper-tube surrounding $h(\mathbf{x})$.

SVMs for Regression (SVR)



■ $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^M$ i.i.d $\sim P(X, Y)$,

■ $\mathbf{x}_i \in \mathbb{R}^N$.

■ $\mathbf{y}_i \in \mathbb{R}$.

■ Find a regression function,
 $h(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b$

■ Minimizing prediction error.

■ Capture data points in an ϵ -radius hyper-tube surrounding $h(\mathbf{x})$.

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

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

SVM Formulations of Interest

Primal

$$\min_{\mathbf{w}, b} \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + R_{\text{emp}}(h; \mathbf{x}, \mathbf{y}),$$

where

$$R_{\text{emp}} = \begin{cases} \frac{1}{M} \sum_{i=1}^M \ell_H(h; \mathbf{x}_i, \mathbf{y}_i), & (\text{SVC}) \\ \frac{1}{M} \sum_{i=1}^M \ell_\epsilon(h; \mathbf{x}_i, \mathbf{y}_i), & (\text{SVR}) \end{cases}$$

and $\lambda = 1/C$. The objective function is convex but non-smooth.

SVM Formulations of Interest

Primal

$$\min_{\mathbf{w}, b} \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + R_{\text{emp}}(h; \mathbf{x}, \mathbf{y}),$$

where

$$R_{\text{emp}} = \begin{cases} \frac{1}{M} \sum_{i=1}^M \ell_H(h; \mathbf{x}_i, \mathbf{y}_i), & (\text{SVC}) \\ \frac{1}{M} \sum_{i=1}^M \ell_\epsilon(h; \mathbf{x}_i, \mathbf{y}_i), & (\text{SVR}) \end{cases}$$

and $\lambda = 1/C$. The objective function is convex but non-smooth.

Dual

$$\begin{aligned} \min_{\mathbf{z}} \quad & \frac{1}{2} \mathbf{z}^T \mathbf{Q} \mathbf{z} + \mathbf{p}^T \mathbf{z} \\ \text{s.t.} \quad & \mathbf{c}^T \mathbf{z} = d \\ & \ell \leq \mathbf{z} \leq \mathbf{u}, \end{aligned} \quad (1)$$

- \mathbf{Q} is a p.s.d. $n \times n$ matrix, usually dense and ill-conditioned.

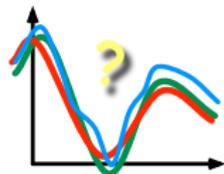
- $n = M$ (SVC) or $n = 2M$ (SVR)

- Determined by \mathbf{y} and **kernel function**

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

- $\mathbf{z}, \mathbf{p}, \mathbf{c}, \ell, \mathbf{u} \in \mathbb{R}^n$, and $d \in \mathbb{R}$.

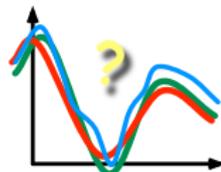
Semiparametric SVM



- Standard (nonparametric) SVR: use a linear model

$$h(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b ,$$

Semiparametric SVM



- Standard (nonparametric) SVR: use a linear model

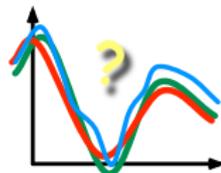
$$h(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b ,$$

- Semiparametric SVR [SFS99]: use an extended linear model

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

where $\psi_j(\cdot)$'s are user-defined (basis) functions.

Semiparametric SVM



- Standard (nonparametric) SVR: use a linear model

$$h(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b ,$$

- Semiparametric SVR [SFS99]: use an extended linear model

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

where $\psi_j(\cdot)$'s are user-defined (basis) functions.

Benefits of semiparametric models

- No explicit modeling is necessary (nonparametric).
- Embedding of prior knowledge / model interpretation (parametric).

Primal Formulation

The “primal” SVR formulation is,

$$\min_{\mathbf{w}, b} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{M} \sum_{i=1}^M \ell_{\epsilon}(\tilde{\mathbf{h}}; \mathbf{x}_i, \mathbf{y}_i), \quad \ell_{\epsilon}(\tilde{\mathbf{h}}; \mathbf{x}_i, \mathbf{y}_i) := \max\{|\mathbf{y}_i - \tilde{\mathbf{h}}(\mathbf{x}_i)| - \epsilon, 0\}.$$

Introducing slack variables ξ_i and ξ_i^* to represent the deviations from the ϵ -tube, we obtain

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

$$\text{s.t.} \quad \mathbf{y}_i - \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle - \sum_{j=1}^K \beta_j \psi_j(\mathbf{x}_i) \leq \epsilon + \xi_i \quad \text{for } i = 1, \dots, M \quad (2b)$$

$$- \left[\mathbf{y}_i - \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle - \sum_{j=1}^K \beta_j \psi_j(\mathbf{x}_i) \right] \leq \epsilon + \xi_i^* \quad \text{for } i = 1, \dots, M \quad (2c)$$

$$\xi \geq \mathbf{0}, \xi^* \geq \mathbf{0} .$$



Dual Formulation

$$\min_{\mathbf{z}} F(\mathbf{z}) := \frac{1}{2} \mathbf{z}^T \mathbf{Q} \mathbf{z} + \mathbf{p}^T \mathbf{z} \quad \text{s.t.} \quad \mathbf{A} \mathbf{z} = \mathbf{0}, \quad \mathbf{0} \leq \mathbf{z} \leq \frac{C}{M} \mathbf{1}, \quad (3)$$

where $\mathbf{z}, \mathbf{p} \in \mathbb{R}^{2M}$, $\mathbf{Q} \in \mathbb{R}^{2M \times 2M}$ p.s.d., and $\mathbf{A} \in \mathbb{R}^{K \times 2M}$.

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

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

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

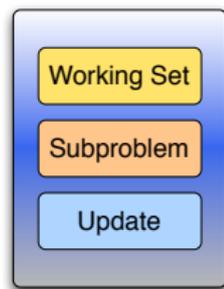
$$\mathbf{A} = \begin{bmatrix} \psi_1(\mathbf{x}_1) & \cdots & \psi_1(\mathbf{x}_M) & -\psi_1(\mathbf{x}_1) & \cdots & -\psi_1(\mathbf{x}_M) \\ \psi_2(\mathbf{x}_1) & \cdots & \psi_2(\mathbf{x}_M) & -\psi_2(\mathbf{x}_1) & \cdots & -\psi_2(\mathbf{x}_M) \\ \vdots & \ddots & \vdots & \vdots & \ddots &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

This is a generalization of the standard SVM dual problem. $n := 2M$.



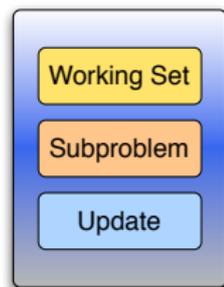
Decomposition Framework [LW09]

- In each outer iteration, we split variables \mathbf{z} into
 - Basic variables $\mathbf{z}_{\mathcal{B}}$, $\mathcal{B} \subset \{1, 2, \dots, n\}$.
 - Nonbasic variables $\mathbf{z}_{\mathcal{N}}$, $\mathcal{N} = \{1, 2, \dots, n\} \setminus \mathcal{B}$.



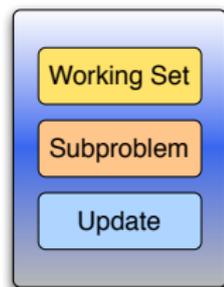
Decomposition Framework [LW09]

- In each outer iteration, we split variables \mathbf{z} into
 - Basic variables $\mathbf{z}_{\mathcal{B}}$, $\mathcal{B} \subset \{1, 2, \dots, n\}$.
 - Nonbasic variables $\mathbf{z}_{\mathcal{N}}$, $\mathcal{N} = \{1, 2, \dots, n\} \setminus \mathcal{B}$.
 - \mathcal{B} is our **working set**, of which the size $n_{\mathcal{B}} \ll n$.



Decomposition Framework [LW09]

- In each outer iteration, we split variables \mathbf{z} into
 - Basic variables $\mathbf{z}_{\mathcal{B}}$, $\mathcal{B} \subset \{1, 2, \dots, n\}$.
 - Nonbasic variables $\mathbf{z}_{\mathcal{N}}$, $\mathcal{N} = \{1, 2, \dots, n\} \setminus \mathcal{B}$.
 - \mathcal{B} is our **working set**, of which the size $n_{\mathcal{B}} \ll n$.
- Fix $\mathbf{z}_{\mathcal{N}}$, change $\mathbf{z}_{\mathcal{B}}$.



Decomposition Framework [LW09]

- In each outer iteration, we 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 $\mathbf{z}^k = (\mathbf{z}_B^k, \mathbf{z}_N^k)$, we solve the subproblem to get \mathbf{z}_B^{k+1} .

Working Set

Subproblem

Update

Subproblem

$$\begin{aligned} \min_{\mathbf{z}_B} \quad & f(\mathbf{z}_B) := \frac{1}{2} \mathbf{z}_B^T \mathbf{Q}_{BB} \mathbf{z}_B + (\mathbf{Q}_{BN} \mathbf{z}_N^k + \mathbf{p}_B)^T \mathbf{z}_B & (4) \\ \text{s.t.} \quad & \mathbf{A}_B \mathbf{z}_B = -\mathbf{A}_N \mathbf{z}_N^k + \mathbf{b}, \quad 0 \leq \mathbf{z}_B \leq \frac{C}{M} \mathbf{1}. \end{aligned}$$



Decomposition Framework [LW09]

- In each outer iteration, we 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 $\mathbf{z}^k = (\mathbf{z}_B^k, \mathbf{z}_N^k)$, we solve the subproblem to get \mathbf{z}_B^{k+1} .

Working Set

Subproblem

Update

Subproblem

$$\begin{aligned} \min_{\mathbf{z}_B} \quad & f(\mathbf{z}_B) := \frac{1}{2} \mathbf{z}_B^T \mathbf{Q}_{BB} \mathbf{z}_B + (\mathbf{Q}_{BN} \mathbf{z}_N^k + \mathbf{p}_B)^T \mathbf{z}_B \\ \text{s.t.} \quad & \mathbf{A}_B \mathbf{z}_B = -\mathbf{A}_N \mathbf{z}_N^k + \mathbf{b}, \quad 0 \leq \mathbf{z}_B \leq \frac{C}{M} \mathbf{1}. \end{aligned} \quad (4)$$

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



Choosing \mathcal{B} : Working Set Selection



- Inspired by the approach of [Joa99], later improved by [SZ05].
 - $n_{\mathcal{B}}$: working set size.
 - n_c : max. number of “fresh” indices. $n_c \ll n_{\mathcal{B}}$.

Choosing \mathcal{B} : Working Set Selection



- Inspired by the approach of [Joa99], later improved by [SZ05].
 - $n_{\mathcal{B}}$: working set size.
 - n_c : max. number of “fresh” indices. $n_c \ll n_{\mathcal{B}}$.

Consider Lagrangian relaxation \mathcal{L} of the dual formulation (3),

$$\mathcal{L}(\mathbf{z}; \boldsymbol{\eta}) = F(\mathbf{z}) + \boldsymbol{\eta}^T \mathbf{A} \mathbf{z} . \quad (5)$$

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

$$\begin{aligned} \min_{\mathbf{d}} \quad & (\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}^k; \boldsymbol{\eta}^k))^T \mathbf{d} \\ \text{s.t.} \quad & 0 \leq \mathbf{d}_i \leq 1 && \text{if } \mathbf{z}_i^{k+1} = 0, \\ & -1 \leq \mathbf{d}_i \leq 0 && \text{if } \mathbf{z}_i^{k+1} = C/M, \\ & -1 \leq \mathbf{d}_i \leq 1 && \text{if } \mathbf{z}_i^{k+1} \in (0, C/M), \\ & \#\{\mathbf{d}_i | \mathbf{d}_i \neq 0\} \leq n_c. \end{aligned} \quad (6)$$

Choosing \mathcal{B} : Working Set Selection



- Inspired by the approach of [Joa99], later improved by [SZ05].
 - $n_{\mathcal{B}}$: working set size.
 - n_c : max. number of “fresh” indices. $n_c \ll n_{\mathcal{B}}$.

Consider Lagrangian relaxation \mathcal{L} of the dual formulation (3),

$$\mathcal{L}(\mathbf{z}; \boldsymbol{\eta}) = F(\mathbf{z}) + \boldsymbol{\eta}^T \mathbf{A} \mathbf{z} . \quad (5)$$

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

$$\begin{aligned} \min_{\mathbf{d}} \quad & (\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}^k; \boldsymbol{\eta}^k))^T \mathbf{d} \\ \text{s.t.} \quad & 0 \leq \mathbf{d}_i \leq 1 && \text{if } \mathbf{z}_i^{k+1} = 0, \\ & -1 \leq \mathbf{d}_i \leq 0 && \text{if } \mathbf{z}_i^{k+1} = C/M, \\ & -1 \leq \mathbf{d}_i \leq 1 && \text{if } \mathbf{z}_i^{k+1} \in (0, C/M), \\ & \#\{\mathbf{d}_i | \mathbf{d}_i \neq 0\} \leq n_c. \end{aligned} \quad (6)$$

- Solved efficiently, $\mathcal{O}(n \log n)$.
- Convergence of decomposition + workingset selection [Lin01, TY08].



Subproblem: Primal-dual Solver (PDSG)



- We consider the following formulation of (4):

$$\max_{\eta} \min_{\mathbf{z}_B \in \Omega} \tilde{\mathcal{L}}(\mathbf{z}_B, \eta) , \quad (7)$$

where

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

$$\Omega = \{ \mathbf{z}_B \in \mathbb{R}^{n_B} \mid \mathbf{0} \leq \mathbf{z}_B \leq \frac{C}{M} \mathbf{1} \} .$$

Subproblem: Primal-dual Solver (PDSG)



- We consider the following formulation of (4):

$$\max_{\eta} \min_{\mathbf{z}_B \in \Omega} \tilde{\mathcal{L}}(\mathbf{z}_B, \eta) , \quad (7)$$

where

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

$$\Omega = \{ \mathbf{z}_B \in \mathbb{R}^{n_B} \mid \mathbf{0} \leq \mathbf{z}_B \leq \frac{C}{M} \mathbf{1} \} .$$

In each “inner” iteration, update primal and dual variables by,

$$\begin{cases} \mathbf{z}_B^{\ell+1} \leftarrow \mathbf{z}_B^{\ell} + s(\mathbf{z}_B^{\ell}, \eta^{\ell}) \\ \eta^{\ell+1} \leftarrow \eta^{\ell} + t(\mathbf{z}_B^{\ell+1}, \eta^{\ell}) , \end{cases} \quad (8)$$

- Primal step $s(\cdot, \cdot)$ is chosen by two-metric GP [GB84] followed by line-search, on a sub-workingset of size 2.
- Dual step $t(\cdot, \cdot)$ is a direction $\nabla_{\eta} \tilde{\mathcal{L}}$, scaled by dual Hessian diagonal [KS05], on a sub-workingset of size 2.



Update



- Update primal-dual iterate pair $(\mathbf{z}^{k+1}, \eta^{k+1})$.
 - $\mathbf{z}^{k+1} \leftarrow (\mathbf{z}_B^{k+1}, \mathbf{z}_N^k)$.
 - η^{k+1} is provided by the subproblem solver.

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.
- “Full gradient” $\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}; \boldsymbol{\eta})$ has to be updated.
 - To check KKT conditions violation.
 - For the next working set selection.

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.
- “Full gradient” $\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}; \boldsymbol{\eta})$ has to be updated.
 - To check KKT conditions violation.
 - For the next working set selection.

Update incrementally,

$$\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}^{k+1}, \boldsymbol{\eta}^{k+1}) = \nabla F(\mathbf{z}^{k+1}) + (\boldsymbol{\eta}^{k+1})^T \mathbf{A} \quad (9)$$

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

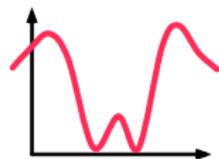


(10)

Experiments

- A toy test problem: modified Mexican hat function [SFS99, KS05]:

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

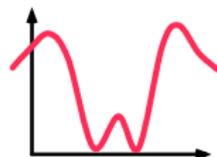


Experiments

- A toy test problem: modified Mexican hat function [SFS99, KS05]:

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

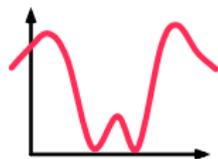
- Sample y_i 's from ω at uniform random points x_i 's in $[0, 10]$ with additive noise $\zeta_i \sim \mathcal{N}(0, 0.2^2)$: $y_i = \omega(x_i) + \zeta_i$.



Experiments

- A toy test problem: modified Mexican hat function [SFS99, KS05]:

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

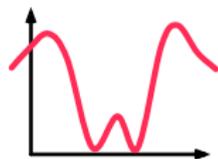


- Sample y_i 's from ω at uniform random points x_i 's in $[0, 10]$ with additive noise $\zeta_i \sim \mathcal{N}(0, 0.2^2)$: $y_i = \omega(x_i) + \zeta_i$.
- Experiment settings
 - Parametric components: $\psi_1(x) = \sin(x)$, $\psi_2(x) = \text{sinc}(2\pi(x - 5))$.
 - Gaussian kernel $\kappa(x, y) = \exp(-\gamma\|x - y\|^2)$ with $\gamma = 0.25$.
 - Loss function parameter $\epsilon = 0.05$.
 - $n_B = 500$, $n_c = n_B/5$.

Experiments

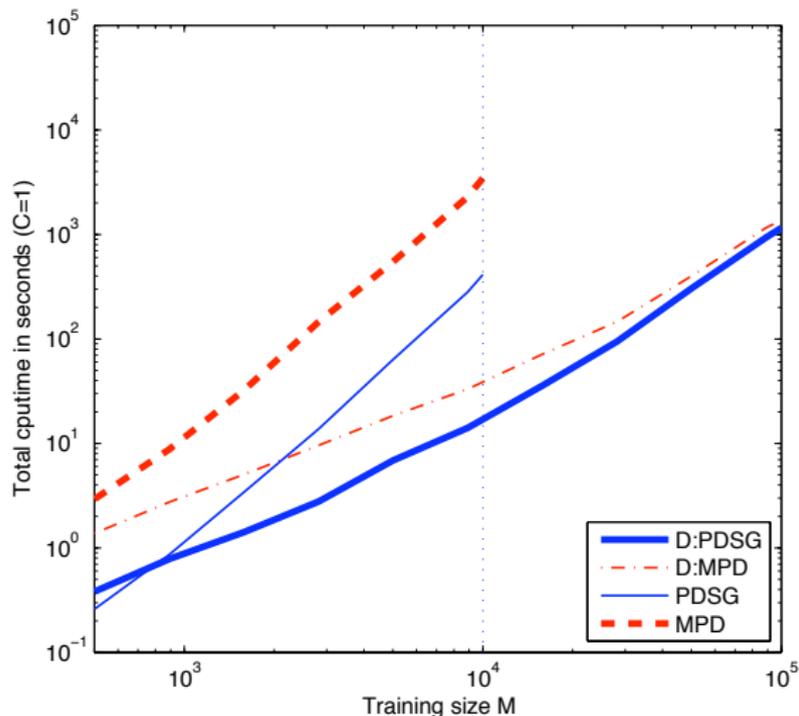
- A toy test problem: modified Mexican hat function [SFS99, KS05]:

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



- Sample y_i 's from ω at uniform random points x_i 's in $[0, 10]$ with additive noise $\zeta_i \sim \mathcal{N}(0, 0.2^2)$: $y_i = \omega(x_i) + \zeta_i$.
- Experiment settings
 - Parametric components: $\psi_1(x) = \sin(x)$, $\psi_2(x) = \text{sinc}(2\pi(x - 5))$.
 - Gaussian kernel $\kappa(x, y) = \exp(-\gamma\|x - y\|^2)$ with $\gamma = 0.25$.
 - Loss function parameter $\epsilon = 0.05$.
 - $n_B = 500$, $n_c = n_B/5$.
- Compare to the current best solver, MPD [KS05].
 - Handles the problem as a whole. Working set size is 1.
 - Primal-dual method, based on the method of multipliers.
 - Primal: gradient projection, dual: scaled gradient ascent.

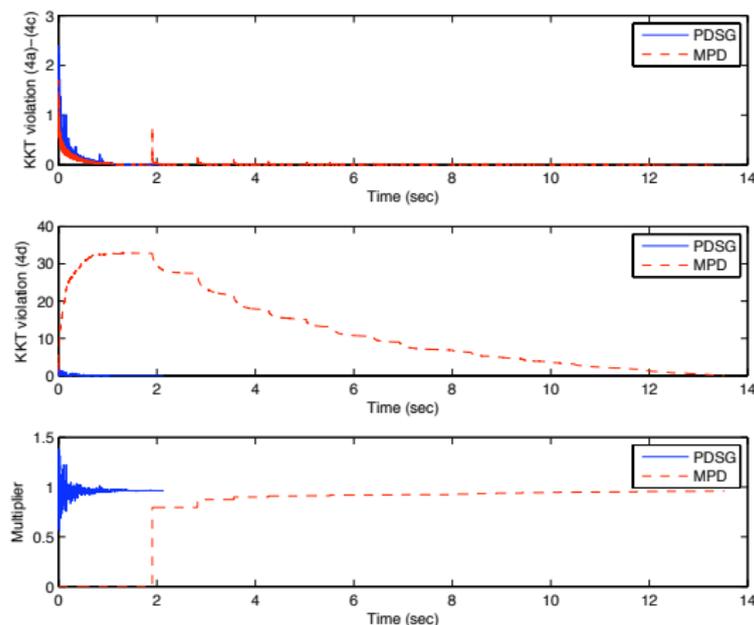
Scaling w.r.t. Training Size



► Varying K

- PDSG vs. MPD (stand-alone).
- D:PDSG vs. D:MPD (in decomposition).
- $C = 1$.
- D : MPD catches up D : PDSG when $M \uparrow$: the full gradient update step becomes dominant as M grows.

Convergence Behavior



- PDSG vs. MPD (stand-alone).
- $M = 1000$.
- PDSG: 2 sec.
- MPD: 14 sec.
- (Top) max. violation of the dual feasibility conditions.
- (Middle) max. violation of the primal equality constraints.
- (Bottom) convergence of the coefficient of the first parametric basis function.

Stochastic Subgradient Methods for SVMs



- Recent ML research on solving the primal formulation¹,

$$\min_{\mathbf{w}, b} f(\mathbf{w}, \mathcal{D}) = \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} + \frac{1}{M} \sum_{i=1}^M \ell_H(\mathbf{w}; \mathbf{x}_i, \mathbf{y}_i). \quad (11)$$

- A large dataset $\mathcal{D} := \{(\mathbf{x}_i, \mathbf{y}_i) : i = 1, \dots, M\}$.

^{1,*}; $\mathbf{w} \leftarrow (\mathbf{w}, b)$

Stochastic Subgradient Methods for SVMs



- Recent ML research on solving the primal formulation¹,

$$\min_{\mathbf{w}, b} f(\mathbf{w}, \mathcal{D}) = \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} + \frac{1}{M} \sum_{i=1}^M \ell_H(\mathbf{w}; \mathbf{x}_i, \mathbf{y}_i). \quad (11)$$

- A large dataset $\mathcal{D} := \{(\mathbf{x}_i, \mathbf{y}_i) : i = 1, \dots, M\}$.
- The objective function is strongly convex*.

¹,*; $\mathbf{w} \leftarrow (\mathbf{w}, b)$

Stochastic Subgradient Methods for SVMs



- Recent ML research on solving the primal formulation¹,

$$\min_{\mathbf{w}, b} f(\mathbf{w}, \mathcal{D}) = \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} + \frac{1}{M} \sum_{i=1}^M \ell_H(\mathbf{w}; \mathbf{x}_i, \mathbf{y}_i). \quad (11)$$

- A large dataset $\mathcal{D} := \{(\mathbf{x}_i, \mathbf{y}_i) : i = 1, \dots, M\}$.
 - The objective function is strongly convex*.
- These has a connection to stochastic approximation methods that have developed in the past 50 years and still active.

^{1,*}; $\mathbf{w} \leftarrow (\mathbf{w}, b)$

Stochastic Subgradient Methods for SVMs



- Recent ML research on solving the primal formulation¹,

$$\min_{\mathbf{w}, b} f(\mathbf{w}, \mathcal{D}) = \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} + \frac{1}{M} \sum_{i=1}^M \ell_H(\mathbf{w}; \mathbf{x}_i, \mathbf{y}_i). \quad (11)$$

- A large dataset $\mathcal{D} := \{(\mathbf{x}_i, \mathbf{y}_i) : i = 1, \dots, M\}$.
- The objective function is strongly convex*.
- These has a connection to stochastic approximation methods that have developed in the past 50 years and still active.
- New issues arise when applied to machine learning problems.

^{1,*}; $\mathbf{w} \leftarrow (\mathbf{w}, b)$

Large-scale Linear SVM Training [Bot, SSSS07]

Given \mathcal{D} , consider the subgradient of an approximate objective function $\tilde{f}(\mathbf{w}; \mathcal{D}_t)$ of $f(\mathbf{w}; \mathcal{D})$ in (11) for a sample dataset $\mathcal{D}_t \subseteq \mathcal{D}$:

$$\tilde{f}(\mathbf{w}; \mathcal{D}_t) := \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} + \frac{1}{|\mathcal{D}_t|} \sum_{(\mathbf{x}, y) \in \mathcal{D}_t} \ell_H(\mathbf{w}; (\mathbf{x}, y))$$

$$g(\mathbf{w}_t; \mathcal{D}_t) := \lambda \mathbf{w}_t - \frac{1}{|\mathcal{D}_t|} \sum_{(\mathbf{x}, y) \in \mathcal{D}_t^+} y \mathbf{x} \in \partial \tilde{f}(\mathbf{w}; \mathcal{D}_t) ,$$

where $\mathcal{D}_t^+ := \{(\mathbf{x}, y) \in \mathcal{D}_t : 1 - y(\mathbf{w}^T \mathbf{x}) > 0\}$.

Large-scale Linear SVM Training [Bot, SSSS07]

Given \mathcal{D} , consider the subgradient of an approximate objective function $\tilde{f}(\mathbf{w}; \mathcal{D}_t)$ of $f(\mathbf{w}; \mathcal{D})$ in (11) for a sample dataset $\mathcal{D}_t \subseteq \mathcal{D}$:

$$\tilde{f}(\mathbf{w}; \mathcal{D}_t) := \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} + \frac{1}{|\mathcal{D}_t|} \sum_{(\mathbf{x}, y) \in \mathcal{D}_t} \ell_H(\mathbf{w}; (\mathbf{x}, y))$$
$$g(\mathbf{w}_t; \mathcal{D}_t) := \lambda \mathbf{w}_t - \frac{1}{|\mathcal{D}_t|} \sum_{(\mathbf{x}, y) \in \mathcal{D}_t^+} y \mathbf{x} \in \partial \tilde{f}(\mathbf{w}; \mathcal{D}_t) ,$$

where $\mathcal{D}_t^+ := \{(\mathbf{x}, y) \in \mathcal{D}_t : 1 - y(\mathbf{w}^T \mathbf{x}) > 0\}$.

Update the iterate \mathbf{w} by

$$\mathbf{w}_{t+1} = \mathbf{P}_{\mathcal{W}}(\mathbf{w}_t - \eta_t g(\mathbf{w}_t; \mathcal{D}_t)) . \quad (12)$$

where

$$\eta_t = \frac{1}{\lambda t}, \quad \mathcal{W} := \{\mathbf{w} : \|\mathbf{w}\|_2 \leq \frac{1}{\sqrt{\lambda}}\}, \quad |\mathcal{D}_t| = 1.$$

Stochastic Approximation (SA)

Classical SA methods

- Choice of $\eta_t = \mathcal{O}(1/t)$ has a history back to [RM51, KW52, Chu54, Sac58].
- Require the objective function to be strongly convex.
 - SVM objective function $f(\cdot)$ is strongly convex with modulus λ .
- Highly sensitive to the scaling of η_t [NJLS09].
- Asymptotic convergence of $\mathcal{O}(1/t)$ in expectation.

Stochastic Approximation (SA)

Classical SA methods

- Choice of $\eta_t = \mathcal{O}(1/t)$ has a history back to [RM51, KW52, Chu54, Sac58].
- Require the objective function to be strongly convex.
 - SVM objective function $f(\cdot)$ is strongly convex with modulus λ .
- Highly sensitive to the scaling of η_t [NJLS09].
- Asymptotic convergence of $\mathcal{O}(1/t)$ in expectation.

Robust SA methods

- Choice of $\eta_t = \mathcal{O}(1/\sqrt{t})$ suggested in [NY83].
- Useful when the objective is convex but not strongly convex, or the curvature is not known.
 - $\lambda \approx 0$ for some choices of C ($\lambda = 1/C$).
- Asymptotic convergence of $\mathcal{O}(1/\sqrt{t})$ in expectation.
- Similar analysis in online learning [Zin03].

Stochastic Approximation (SA)

Classical SA methods

- Choice of $\eta_t = \mathcal{O}(1/t)$ has a history back to [RM51, KW52, Chu54, Sac58].
- Require the objective function to be strongly convex.
 - SVM objective function $f(\cdot)$ is strongly convex with modulus λ .
- Highly sensitive to the scaling of η_t [NJLS09].
- Asymptotic convergence of $\mathcal{O}(1/t)$ in expectation.

Robust SA methods

- Choice of $\eta_t = \mathcal{O}(1/\sqrt{t})$ suggested in [NY83].
- Useful when the objective is convex but not strongly convex, or the curvature is not known.
 - $\lambda \approx 0$ for some choices of C ($\lambda = 1/C$).
- Asymptotic convergence of $\mathcal{O}(1/\sqrt{t})$ in expectation.
- Similar analysis in online learning [Zin03].

Both requires a bound on $\mathbb{E}(\|g(\mathbf{w}; \mathcal{D})\|^2)$.

A Stopping Criterion?

- SA algorithms require the number of iterations T to run.
- An efficient stopping criterion is important,
 - Slow convergence of SA methods.
 - Data sets are large.

A Stopping Criterion?

- SA algorithms require the number of iterations T to run.
- An efficient stopping criterion is important,
 - Slow convergence of SA methods.
 - Data sets are large.

Elements of statistical learning theory,

- Unknown $P(X, Y)$, and a dataset $D = \{x_i, y_i\}_{i=1}^M$ *i.i.d.* $\sim P(X, Y)$.

A Stopping Criterion?

- SA algorithms require the number of iterations T to run.
- An efficient stopping criterion is important,
 - Slow convergence of SA methods.
 - Data sets are large.

Elements of statistical learning theory,

- Unknown $P(X, Y)$, and a dataset $D = \{x_i, y_i\}_{i=1}^M$ *i.i.d.* $\sim P(X, Y)$.
- Hypothesis $f \in \mathcal{F}$, \mathcal{F} is a chosen family of hypotheses.

A Stopping Criterion?

- SA algorithms require the number of iterations T to run.
- An efficient stopping criterion is important,
 - Slow convergence of SA methods.
 - Data sets are large.

Elements of statistical learning theory,

- Unknown $P(X, Y)$, and a dataset $D = \{x_i, y_i\}_{i=1}^M$ *i.i.d.* $\sim P(X, Y)$.
- Hypothesis $f \in \mathcal{F}$, \mathcal{F} is a chosen family of hypotheses.
- Loss function $\ell(f(X), Y)$.

A Stopping Criterion?

- SA algorithms require the number of iterations T to run.
- An efficient stopping criterion is important,
 - Slow convergence of SA methods.
 - Data sets are large.

Elements of statistical learning theory,

- Unknown $P(X, Y)$, and a dataset $D = \{x_i, y_i\}_{i=1}^M$ *i.i.d.* $\sim P(X, Y)$.
- Hypothesis $f \in \mathcal{F}$, \mathcal{F} is a chosen family of hypotheses.
- Loss function $\ell(f(X), Y)$.
- Risk $R(f) := \mathbb{E}(\ell(f(X), Y)) = \int \ell(f(X), Y) dP(X, Y)$.

A Stopping Criterion?

- SA algorithms require the number of iterations T to run.
- An efficient stopping criterion is important,
 - Slow convergence of SA methods.
 - Data sets are large.

Elements of statistical learning theory,

- Unknown $P(X, Y)$, and a dataset $D = \{x_i, y_i\}_{i=1}^M$ *i.i.d.* $\sim P(X, Y)$.
- Hypothesis $f \in \mathcal{F}$, \mathcal{F} is a chosen family of hypotheses.
- Loss function $\ell(f(X), Y)$.
- Risk $R(f) := \mathbb{E}(\ell(f(X), Y)) = \int \ell(f(X), Y) dP(X, Y)$.
- Empirical Risk $R_{emp}(f) := \frac{1}{M} \sum_{i=1}^M \ell(f(x_i), y_i)$.

A Stopping Criterion?

- SA algorithms require the number of iterations T to run.
- An efficient stopping criterion is important,
 - Slow convergence of SA methods.
 - Data sets are large.

Elements of statistical learning theory,

- Unknown $P(X, Y)$, and a dataset $D = \{x_i, y_i\}_{i=1}^M$ *i.i.d.* $\sim P(X, Y)$.
- Hypothesis $f \in \mathcal{F}$, \mathcal{F} is a chosen family of hypotheses.
- Loss function $\ell(f(X), Y)$.
- Risk $R(f) := \mathbb{E}(\ell(f(X), Y)) = \int \ell(f(X), Y) dP(X, Y)$.
- Empirical Risk $R_{emp}(f) := \frac{1}{M} \sum_{i=1}^M \ell(f(x_i), y_i)$.
- $R^* := \inf_f R(f)$.

A Stopping Criterion?

- SA algorithms require the number of iterations T to run.
- An efficient stopping criterion is important,
 - Slow convergence of SA methods.
 - Data sets are large.

Elements of statistical learning theory,

- Unknown $P(X, Y)$, and a dataset $D = \{x_i, y_i\}_{i=1}^M$ i.i.d. $\sim P(X, Y)$.
- Hypothesis $f \in \mathcal{F}$, \mathcal{F} is a chosen family of hypotheses.
- Loss function $\ell(f(X), Y)$.
- Risk $R(f) := \mathbb{E}(\ell(f(X), Y)) = \int \ell(f(X), Y) dP(X, Y)$.
- Empirical Risk $R_{emp}(f) := \frac{1}{M} \sum_{i=1}^M \ell(f(x_i), y_i)$.
- $R^* := \inf_f R(f)$.
- Error decomposition,

$$\underbrace{\inf_{f \in \mathcal{F}} R_{emp}(f) - R^*}_{\text{generalization error}} = \underbrace{\left(\inf_{f \in \mathcal{F}} R(f) - R^* \right)}_{\text{approximation error}} + \underbrace{\left(\inf_{f \in \mathcal{F}} R_{emp}(f) - \inf_{f \in \mathcal{F}} R(f) \right)}_{\text{estimation error}}.$$



A Stopping Criterion?

[SSS08] suggested a new error decomposition

$$(\text{gen. err}) = (\text{approx. err}) + (\text{est. err}) + (\text{optimization err}) .$$

A Stopping Criterion?

[SSS08] suggested a new error decomposition

$$(\text{gen. err}) = (\text{approx. err}) + (\text{est. err}) + (\text{optimization err}) .$$

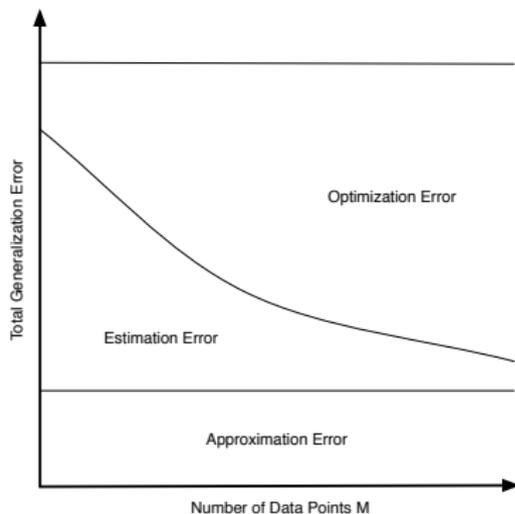
- Approx. error doesn't change for fixed \mathcal{F} .
- As $M \rightarrow \infty$, $(\text{est. err}) \rightarrow 0$ if f is consistent.

A Stopping Criterion?

[SSS08] suggested a new error decomposition

$$(\text{gen. err}) = (\text{approx. err}) + (\text{est. err}) + (\text{optimization err}) .$$

- Approx. error doesn't change for fixed \mathcal{F} .
- As $M \rightarrow \infty$, (est. err) $\rightarrow 0$ if f is consistent.
- Allow larger opt. err to achieve the same level of gen. err with large M .



[Due to N. Srebro at MLSS09]

Conclusions

Decomposition Algorithm

- Can solve other SVMs, ν -SVM, semiparametric SlapSVM, etc.
- Proofs are on the way.

SA Algorithms

- More work is needed.
- SA methods are inherently serial, each iterate is an instantiation.
 - Reduce the variation of the final iterate distribution, possibly by running several SA algorithms in parallel.
- Nonlinear $\phi(\mathbf{x})$ (other than $\phi(\mathbf{x}) = \mathbf{x}$).
 - Initial work by [JY09].
- Explicit consideration of the intercept b .

Thank you.

Optimality Condition of the Dual Formulation

Lagrangian function \mathcal{L} of (3) and its gradient w.r.t. \mathbf{z} :

$$\mathcal{L}(\mathbf{z}; \boldsymbol{\eta}) = F(\mathbf{z}) + \boldsymbol{\eta}^T \mathbf{A} \mathbf{z} . \quad (13)$$

$$\nabla_{\mathbf{z}} \mathcal{L}(\mathbf{z}; \boldsymbol{\eta}) = \mathbf{Q} \mathbf{z} + \mathbf{p} + \mathbf{A}^T \boldsymbol{\eta} . \quad (14)$$

From Karush-Kuhn-Tucker (KKT) first-order optimality conditions,

$$\left(\mathbf{Q} \mathbf{z} + \mathbf{p} + \mathbf{A}^T \boldsymbol{\eta} \right)_i \geq 0 \quad \text{if } \mathbf{z}_i = 0 \quad (15a)$$

$$\left(\mathbf{Q} \mathbf{z} + \mathbf{p} + \mathbf{A}^T \boldsymbol{\eta} \right)_i \leq 0 \quad \text{if } \mathbf{z}_i = C \quad (15b)$$

$$\left(\mathbf{Q} \mathbf{z} + \mathbf{p} + \mathbf{A}^T \boldsymbol{\eta} \right)_i = 0 \quad \text{if } \mathbf{z}_i \in (0, C/M) \quad (15c)$$

$$\mathbf{A} \mathbf{z} = \mathbf{b} \quad (15d)$$

$$\mathbf{0} \leq \mathbf{z} \leq (C/M) \mathbf{1} . \quad (15e)$$

which is necessary and sufficient. [Return](#)

Decomposition Framework

Algorithm 1 Decomposition Framework

1. **Initialization.** Choose an initial \mathbf{z}^1 (3) (possibly infeasible), initial guess of $\boldsymbol{\eta}^1$, positive integers $n_B \geq K$ and $0 < n_c < n_B$, and tolD . Choose an initial working set \mathcal{B} . $k \leftarrow 1$.

2. **Subproblem.** Solve the subproblem (4) for the current working set \mathcal{B} , to obtain \mathbf{z}_B^{k+1} and $\boldsymbol{\eta}^{k+1}$. Set $\mathbf{z}^{k+1} = (\mathbf{z}_B^{k+1}, \mathbf{z}_N^k)$.

3. **Gradient Update.**

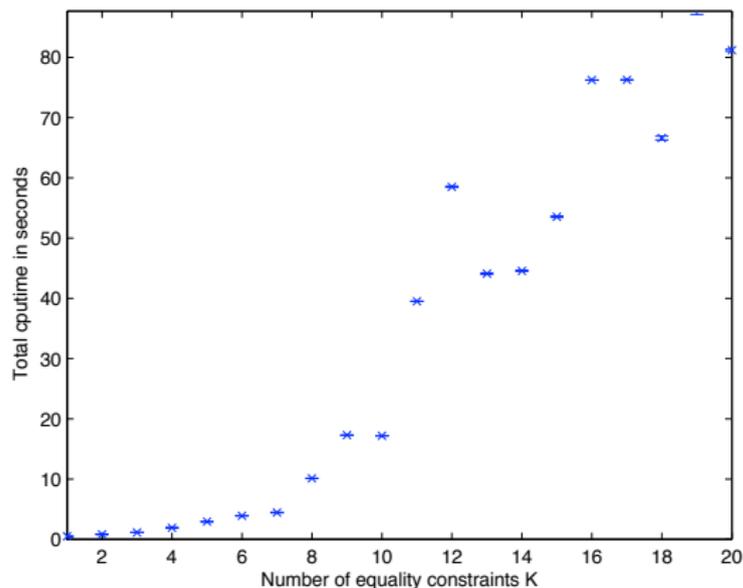
$$\nabla F(\mathbf{z}^{k+1}) + (\boldsymbol{\eta}^{k+1})^T \mathbf{A} = \nabla F(\mathbf{z}^k) + \begin{bmatrix} \mathbf{Q}_{\mathcal{B}\mathcal{B}} \\ \mathbf{Q}_{\mathcal{N}\mathcal{B}} \end{bmatrix} (\mathbf{z}_B^{k+1} - \mathbf{z}_B^k) + (\boldsymbol{\eta}^{k+1})^T \mathbf{A} .$$

4. **Convergence Check.** If the maximal violation of the KKT conditions falls below tolD , terminate with the primal-dual solution $(\mathbf{z}^{k+1}, \boldsymbol{\eta}^{k+1})$.

5. **Working Set Update.** Find a new working set \mathcal{B} by solving (6).

6. Set $k \leftarrow k + 1$ and go to step 2.

Scaling of D:PDSG w.r.t K



← Return

- Total solution time of D:PDSG with increasing number of parametric components K .
- $M = 1000$.
- Time complexity of D:PDSG is $\mathcal{O}(uKn_B)$, u is the number of outer iterations.
- Solver time appears to increase linearly with K for $K \geq 6$.

$$\psi_j(x) =$$

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



Reference I

- [Bot] Léon Bottou, <http://leon.bottou.org/projects/sgd>.
- [Chu54] K.L Chung, *On a stochastic approximation method*, The Annals of Mathematical Statistics **25** (1954), no. 3, 463–483.
- [GB84] E. M. Gafni and D. P. Bertsekas, *Two-metric projection methods for constrained optimization*, SIAM Journal on Control and Optimization **22** (1984), 936–964.
- [Joa99] Thorsten Joachims, *Making large-scale support vector machine learning practical*, Advances in Kernel Methods - Support Vector Learning (B. Schölkopf, C. Burges, and A. Smola, eds.), MIT Press, Cambridge, MA, 1999, pp. 169–184.
- [JY09] T. Joachims and Chun-Nam John Yu, *Sparse kernel svms via cutting-plane training*, Machine Learning (2009), European Conference on Machine Learning (ECML).
- [KS05] Wolf Kienzle and Bernhard Schölkopf, *Training support vector machines with multiple equality constraints*, Machine Learning: ECML 2005, vol. 16, October 2005.
- [KW52] J. Kiefer and J. Wolfowitz, *Stochastic estimation of the maximum of a regression function*, Annals of Mathematical Statistics **23** (1952), no. 3, 462–466.
- [Lin01] C. J. Lin, *Linear convergence of a decomposition method for support vector machines*, Tech. report, Department of Computer Science and Information Engineering, National Taiwan University, 2001.
- [LW09] Sangkyun Lee and Stephen J. Wright, *Decomposition algorithms for training large-scale semiparametric support vector machines*, Lecture Notes in Artificial Intelligence: ECML, 2009.
- [NJLS09] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro, *Robust stochastic approximation approach to stochastic programming*, SIAM Journal on Optimization **19** (2009), no. 4, 1574–1609.
- [NY83] A. Nemirovski and D. Yudin, *Problem complexity and method efficiency in optimization*, Wiley-Intersci. Ser. Discrete Math., 1983.
- [RM51] Herbert Robbins and Sutton Monro, *A stochastic approximation method*, Annals of Mathematical Statistics **22** (1951), no. 3, 400–407.

Reference II

- [Sac58] Jerome Sacks, *Asymptotic distribution of stochastic approximation procedures*, The Annals of Mathematical Statistics **29** (1958), no. 2, 373–405.
- [SFS99] Alex J. Smola, Thilo T. Frieß, and Bernhard Schölkopf, *Semiparametric support vector and linear programming machines*, Advances in Neural Information Processing Systems 11 (Cambridge, MA, USA), MIT Press, 1999, pp. 585–591.
- [SSS08] S. Shalev-Shwartz and N. Srebro, *Svm optimization: inverse dependence on training set size*, ICML '08: Proceedings of the 25th Annual International Conference on Machine Learning (2008), 928–935.
- [SSSS07] Shai Shalev-Shwartz, Yoram Singer, and Nathan Srebro, *Pegasos: Primal estimated sub-gradient solver for svm*, ICML '07: Proceedings of the 24th international conference on Machine learning (New York, NY, USA), ACM, 2007, pp. 807–814.
- [SZ05] 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** (2005), 583–596.
- [TY08] 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.
- [Zin03] Martin Zinkevich, *Online convex programming and generalized infinitesimal gradient ascent*, ICML '03: Proceedings of the twentieth international conference on Machine learning, 2003, pp. 928–936.