

# Lecture 12: Online solvers for Convex Risk Minimization

Vojtěch Franc

May 14, 2015

## 10.A: Sequential Dual Ascent

- ◆ Generic framework
- ◆ Steepest feasible ascent variable selection strategy
- ◆ Near maximal gain variable selection strategy

## 10.B: Stochastic Gradient Descent

**XEP33SML – Structured Model Learning, Summer 2015**

## 10.A: Optimization zoo

**Batch methods** like the cutting plane algorithms (previous lecture, 9.A, 9.B)

- Updates: expensive after processing all training examples.
- + Stopping condition: certificate of optimality.
- + Tunable memory requirements.

**Primal on-line methods** like the Stochastic Gradient Descent (this lecture, 10.B)

- + Very simple algorithm.
- + Updates: cheap after processing a single example.
- + Small memory requirements.
- Stopping condition: missing.
- Step-size: requires sensitive setting.

**Sequential Dual Ascent** (this lecture, 10.A)

- + Updates: cheap after processing a single example.
- + Stopping condition: certificate of optimality.
- + Step-size: optimal can be computed analytically.
- Memory requirements.

## 10.A: SDA - formulation of the primal problem

- ◆ Consider learning of a linear classifier  $h(\mathbf{x}; \mathbf{w}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{w}, \Psi(\mathbf{x}, \mathbf{y}) \rangle$  from a training set  $\{(\mathbf{x}^1, \mathbf{y}^1), \dots, (\mathbf{x}^m, \mathbf{y}^m)\} \in (\mathcal{X} \times \mathcal{Y})^m$  by solving the SO-SVM problem

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left[ \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \max_{\mathbf{y} \in \mathcal{Y}} (\ell(\mathbf{y}^i, \mathbf{y}) + \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) - \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle) \right]$$

where  $\lambda > 0$  is a reg. constant,  $\Psi: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^n$  is a feature map and  $\ell: \mathcal{Y}^2 \rightarrow \mathbb{R}_+$  is a loss such that  $\ell(\mathbf{y}, \mathbf{y}') = 0$  iff  $\mathbf{y} = \mathbf{y}'$ .

- ◆ We will use a compact formulation of the same problem:

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} P(\mathbf{w}) \quad \text{where} \quad P(\mathbf{w}) = \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i \in \mathcal{I}} \max_{\mathbf{y} \in \mathcal{Y}} (\ell_i(\mathbf{y}) + \langle \mathbf{w}, \Psi_i(\mathbf{y}) \rangle)$$

where  $\mathcal{I} = \{1, \dots, m\}$ ,  $\Psi_i(\mathbf{y}) = \Psi(\mathbf{x}^i, \mathbf{y}) - \Psi(\mathbf{x}^i, \mathbf{y}^i)$ ,  $\ell_i(\mathbf{y}) = \ell(\mathbf{y}^i, \mathbf{y})$  and  $\lambda$ .

- ◆ The formulations for other proxies like the slack-rescaling or "PosLearn" are similar.

## 10.A: SDA - dual problem

- ◆ The SO-SVM dual problem reads

$$\boldsymbol{\alpha}^* = \underset{\boldsymbol{\alpha} \in \mathbb{R}^d}{\operatorname{argmin}} D(\boldsymbol{\alpha}) \quad \text{where} \quad D(\boldsymbol{\alpha}) = \langle \mathbf{b}, \boldsymbol{\alpha} \rangle - \frac{1}{2} \|\mathbf{A}\boldsymbol{\alpha}\|^2$$

where

$\boldsymbol{\alpha} = (\boldsymbol{\alpha}_1; \dots; \boldsymbol{\alpha}_m)^d$  is a vector of  $d = m|\mathcal{Y}|$  dual variables where  
 $\boldsymbol{\alpha}_i = (\alpha_i(\mathbf{y}) \mid \mathbf{y} \in \mathcal{Y}) \in \mathbb{R}^{|\mathcal{Y}|}$ .

$\mathbf{b} = (\mathbf{b}_1; \dots; \mathbf{b}_m) \in \mathbb{R}^d$  is a vector of losses where  $\mathbf{b}_i = (\ell_i(\mathbf{y}) \mid \mathbf{y} \in \mathcal{Y}) \in \mathbb{R}^{|\mathcal{Y}|}$ .

$\mathbf{A} = (\mathbf{A}_1, \dots, \mathbf{A}_m) \in \mathbb{R}^{n \times d}$  is a matrix of features where  
 $\mathbf{A}_i = (\Psi_i(\mathbf{y})/\sqrt{\lambda} \mid \mathbf{y} \in \mathcal{Y}) \in \mathbb{R}^{n \times |\mathcal{Y}|}$ .

$\mathcal{A} = \{\boldsymbol{\alpha} \in \mathbb{R}^d \mid \boldsymbol{\alpha} \geq \mathbf{0} \wedge \sum_{\mathbf{y} \in \mathcal{Y}} \alpha_i(\mathbf{y}) = \frac{1}{m}, i \in \mathcal{I}\}$  is feasibility set.

- ◆ The optimal primal variables can be computed from the dual variables by

$$\mathbf{w}^* = -\frac{1}{\sqrt{\lambda}} \mathbf{A}\boldsymbol{\alpha}^*$$

- ◆ Thanks to the weak duality we have  $P(\mathbf{w}) \geq D(\boldsymbol{\alpha}), \forall \boldsymbol{\alpha} \in \mathcal{A}$ , and thanks to the strong duality  $P(\mathbf{w}^*) = D(\boldsymbol{\alpha}^*)$ .

## 10.A: SDA - generic algorithm

- ◆ Solving the SO-SVM dual can be converted into a series of simpler (reduced) QP tasks which can be solved analytically.
- ◆ The generic SDA: starting from  $\alpha^0 \in \mathcal{A}$  it recursively solves a reduced problem

$$\alpha^{t+1} = \operatorname{argmax}_{\alpha \in \mathcal{A}^t} D(\alpha)$$

where  $\mathcal{A}^t \subset \mathcal{A}$  is a line between the current solution  $\alpha^t$  and a point  $\beta^t$  selected from  $\mathcal{A}$ , i.e.

$$\mathcal{A}^t = \left\{ \alpha \in \mathbb{R}^d \mid \alpha = (1 - \tau)\alpha^t + \tau\beta^t, \tau \in [0, 1] \right\}$$

- ◆ All iterates of the SDA algorithm are feasible, i.e.  $\alpha^t \in \mathcal{A}, \forall t$ .
- ◆ The instances of the SDA algorithm differ in the strategy used to construct the point  $\beta^t$ :
  - In order to query the oracle for the  $i$ -th example we will modify only those dual variables associated with  $i$ -th example:  $\beta_j^t(y) = \alpha_j^t(y), \forall j \in \mathcal{I} \setminus \{i\}, \forall y \in \mathcal{Y}$ .
  - **Steepest feasible ascent:** maximize the derivative  $D'_{\mathcal{A}^t}(0)$  where  $D_{\mathcal{A}^t}(\tau) = D(\alpha^t(1 - \tau) + \tau\beta^t)$
  - **Near maximal gain:** maximize the improvement  $\delta^t = D(\alpha^{t+1}) - D(\alpha^t)$

## 10.A: SDA - solving the reduced problem analytically

- ◆ The dual objective  $D(\boldsymbol{\alpha}) = \langle \mathbf{b}, \boldsymbol{\alpha} \rangle - \frac{1}{2} \|\mathbf{A}\boldsymbol{\alpha}\|^2$  restricted to the line segment  $\mathcal{A}^t$  reads

$$D_{\mathcal{A}^t}(\tau) = D((1 - \tau)\boldsymbol{\alpha}^t + \tau\boldsymbol{\beta}^t) = D(\boldsymbol{\alpha}^t) + \tau \langle \boldsymbol{\beta}^t - \boldsymbol{\alpha}^t, \mathbf{b} - \mathbf{A}^T \mathbf{A} \boldsymbol{\alpha}^t \rangle - \frac{\tau^2}{2} \|\mathbf{A}(\boldsymbol{\beta}^t - \boldsymbol{\alpha}^t)\|^2$$

and its derivative is

$$D'_{\mathcal{A}^t}(\tau) = \langle \boldsymbol{\beta}^t - \boldsymbol{\alpha}^t, \mathbf{b} - \mathbf{A}^T \mathbf{A} \boldsymbol{\alpha}^t \rangle - \tau \|\mathbf{A}(\boldsymbol{\beta}^t - \boldsymbol{\alpha}^t)\|^2$$

- ◆ The reduced problem  $\boldsymbol{\alpha}^{t+1} = \operatorname{argmax}_{\boldsymbol{\alpha} \in \mathcal{A}^t} D(\boldsymbol{\alpha})$  can be solved by computing

$$\tau^t = \operatorname{argmax}_{\tau \in [0,1]} D(\boldsymbol{\alpha}^t(1 - \tau) + \boldsymbol{\beta}^t)$$

and setting  $\boldsymbol{\alpha}^{t+1} = \boldsymbol{\alpha}^t(1 - \tau^t) + \boldsymbol{\beta}^t \tau^t$ .

- ◆ Provided  $D'_{\mathcal{A}^t}(0) > 0$ , the optimal  $\tau^t$  is either 1 or a solution of  $D'_{\mathcal{A}^t}(\tau) = 0$ , i.e.

$$\tau^t = \min \left\{ 1, \frac{\langle \boldsymbol{\beta}^t - \boldsymbol{\alpha}^t, \mathbf{b} - \mathbf{A}^T \mathbf{A} \boldsymbol{\alpha}^t \rangle}{\|\mathbf{A}(\boldsymbol{\beta}^t - \boldsymbol{\alpha}^t)\|^2} \right\}$$

# 10.A: SDA - The stopping condition inspired variable selection strategy

- ◆ Let  $\mathbf{w}^t = -\frac{1}{\sqrt{\lambda}}\mathbf{A}\boldsymbol{\alpha}^t$  be a primal solution. Then the duality gap reads:

$$\begin{aligned}
 G(\mathbf{w}^t, \boldsymbol{\alpha}^t) &= P(\mathbf{w}^t) - D(\boldsymbol{\alpha}^t) \\
 &= \frac{\lambda}{2}\|\mathbf{w}^t\|^2 + \frac{1}{m} \sum_{i \in \mathcal{I}} \max_{\mathbf{y} \in \mathcal{Y}} (\ell_i(\mathbf{y}) + \langle \mathbf{w}^t, \boldsymbol{\Psi}_i(\mathbf{y}) \rangle) - \sum_{i \in \mathcal{I}} \sum_{\mathbf{y} \in \mathcal{Y}} \alpha_i(\mathbf{y}) \ell_i(\mathbf{y}) + \frac{1}{2}\|\mathbf{A}\boldsymbol{\alpha}^t\|^2 \\
 &= \frac{1}{m} \sum_{i \in \mathcal{I}} \max_{\mathbf{y} \in \mathcal{Y}} (\ell_i(\mathbf{y}) + \langle \mathbf{w}^t, \boldsymbol{\Psi}_i(\mathbf{y}) \rangle) - \sum_{i \in \mathcal{I}} \sum_{\mathbf{y} \in \mathcal{Y}} \alpha_i^t(\mathbf{y}) \ell_i(\mathbf{y}) + \lambda\|\mathbf{w}^t\|^2 \\
 &= \frac{1}{m} \sum_{i \in \mathcal{I}} \max_{\mathbf{y} \in \mathcal{Y}} (\ell_i(\mathbf{y}) + \langle \mathbf{w}^t, \boldsymbol{\Psi}_i(\mathbf{y}) \rangle) - \sum_{i \in \mathcal{I}} \sum_{\mathbf{y} \in \mathcal{Y}} \alpha_i^t(\mathbf{y}) (\ell_i(\mathbf{y}) + \langle \boldsymbol{\Psi}_i(\mathbf{y}), \mathbf{w}^t \rangle) \\
 &= \frac{1}{m} \left[ \sum_{i \in \mathcal{I}} \max_{\mathbf{y} \in \mathcal{Y}} s_i(\mathbf{y}, \mathbf{w}^t) - m \sum_{\mathbf{y} \in \mathcal{Y}} \alpha_i^t(\mathbf{y}) s_i(\mathbf{y}, \mathbf{w}^t) \right] = \frac{1}{m} \sum_{i \in \mathcal{I}} G_i(\mathbf{w}^t, \boldsymbol{\alpha}^t)
 \end{aligned}$$

where  $s_i(\mathbf{y}, \mathbf{w}) = \ell_i(\mathbf{y}) + \langle \mathbf{w}, \boldsymbol{\Psi}_i(\mathbf{y}) \rangle$  is the loss augmented classification score.

- ◆ Strategy to select the block of variables  $\boldsymbol{\alpha}_i^t = (\alpha_i^t(\mathbf{y}) \mid \mathbf{y} \in \mathcal{Y})$  to be updated: go sequentially over all examples  $i \in \mathcal{I}$  and if

$$G_i(\mathbf{w}^t, \boldsymbol{\alpha}^t) > \varepsilon$$

holds then update  $\boldsymbol{\alpha}_i^t$ . Otherwise, if  $G_i(\mathbf{w}^t, \boldsymbol{\alpha}^t) \leq \varepsilon, \forall i \in \mathcal{I}$ , then  $P(\mathbf{w}^t) \leq P(\mathbf{w}^*) + \varepsilon$ .

## 10.A: SDA - steepest feasible ascent

- ◆ Assume we are going to update the  $i$ -th block of variables so that  $\beta^t$  must be from  $\mathcal{A} \cap \mathcal{B}_i$  where  $\mathcal{B}_i^t = \{\beta \in \mathbb{R}^d \mid \beta_j^t(\mathbf{y}) = \alpha_j^t(\mathbf{y}), \forall j \in \mathcal{I} \setminus \{i\}, \forall \mathbf{y} \in \mathcal{Y}\}$
- ◆ The point  $\beta^t$  can be constructed such that

$$\begin{aligned} \beta^t &= \operatorname{argmax}_{\beta \in \mathcal{A} \cap \mathcal{B}_i^t} D'_{\mathcal{A}^t}(0) = \operatorname{argmax}_{\beta \in \mathcal{A} \cap \mathcal{B}_i^t} \langle \beta - \alpha^t, \mathbf{b} - \mathbf{A}^T \mathbf{A} \alpha^t \rangle \\ &= \operatorname{argmax}_{\substack{\beta_i \geq 0 \\ \sum_{\mathbf{y} \in \mathcal{Y}} \beta_i(\mathbf{y}) = \frac{1}{m}}} \sum_{\mathbf{y} \in \mathcal{Y}} \beta_i(\mathbf{y}) \left( \ell_j(\mathbf{y}) + \langle \mathbf{w}^t, \Psi_j(\mathbf{y}) \rangle \right) \end{aligned}$$

where  $\mathbf{w}^t = -\frac{1}{\sqrt{\lambda}} \mathbf{A} \alpha^t$ .

- ◆ The solution is the vector  $\beta^t = (\beta_j(\mathbf{y}) \mid j \in \mathcal{I}, \mathbf{y} \in \mathcal{Y})$  with components

$$\beta_j^t(\mathbf{y}) = \begin{cases} \frac{1}{m} & \text{if } \mathbf{y} = \mathbf{u}^t \wedge i = j \\ 0 & \text{if } \mathbf{y} \neq \mathbf{u}^t \wedge i = j \\ \alpha_j^t(\mathbf{y}) & \text{if } j \neq i \end{cases}$$

where

$$\mathbf{u}^t \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} \left( \ell_i(\mathbf{y}) + \langle \mathbf{w}^t, \Psi_i(\mathbf{y}) \rangle \right)$$

which is the loss augmented classification problem.



## 10.A: Algorithm: SDA with steepest feasible ascent (the dual form)

**Input:** precision parameter  $\varepsilon > 0$

**Output:**  $\varepsilon$ -optimal primal-dual pair  $(\mathbf{w}^t, \boldsymbol{\alpha}^t)$

**Initialization:**

$$\alpha_i^0(\mathbf{y}) = \begin{cases} \frac{1}{m} & \text{if } \mathbf{y} = \mathbf{y}^i \\ 0 & \text{otherwise} \end{cases}, i \in \mathcal{I}$$

**repeat**

num\_updates := 0

**forall**  $i \in \mathcal{I}$  **do**

$$\mathbf{w}^t = -\frac{1}{\sqrt{\lambda}} \mathbf{A} \boldsymbol{\alpha}^t$$

$$\mathbf{u}^t \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} s_i(\mathbf{y}, \mathbf{w}^t)$$

**if**  $s_i(\mathbf{u}^t, \mathbf{w}^t) - m \sum_{\mathbf{y} \in \mathcal{Y}} \alpha_i(\mathbf{y}) s_i(\mathbf{y}, \mathbf{w}^t) > \varepsilon$  **then**

num\_updates := num\_updates + 1

$$\beta_j^t(\mathbf{y}) = \begin{cases} \frac{1}{m} & \text{if } \mathbf{y} = \mathbf{u}^t \wedge i = j \\ 0 & \text{if } \mathbf{y} \neq \mathbf{u}^t \wedge i = j \\ \alpha_j^t(\mathbf{y}) & \text{if } j \neq i \end{cases}$$

$$\tau^t = \min \left\{ 1, \frac{\langle \boldsymbol{\beta}^t - \boldsymbol{\alpha}^t, \mathbf{b} - \mathbf{A}^T \mathbf{A} \boldsymbol{\alpha}^t \rangle}{\|\mathbf{A}(\boldsymbol{\beta}^t - \boldsymbol{\alpha}^t)\|^2} \right\}$$

$$\boldsymbol{\alpha}^{t+1} = \boldsymbol{\alpha}^t (1 - \tau^t) + \boldsymbol{\beta}^t \tau^t$$

$t \leftarrow t + 1$

**until** num\_updates = 0 ;

# 10.A: SDA - steepest feasible ascent; removing the dual variables

- ◆ The key quantities can be computed without explicitly maintaining the dual variables:

$$\mathbf{w}^t = -\frac{1}{\sqrt{\lambda}}\mathbf{A}\boldsymbol{\alpha}^t = \sum_{i \in \mathcal{I}} \left( -\frac{1}{\sqrt{\lambda}}\mathbf{A}_i\boldsymbol{\alpha}_i^t \right) = \sum_{i \in \mathcal{I}} \mathbf{w}_i^t$$

$$\langle \boldsymbol{\beta}^t - \boldsymbol{\alpha}^t, \mathbf{b} - \mathbf{A}^T \mathbf{A} \boldsymbol{\alpha}^t \rangle = \frac{1}{m} s_i(\mathbf{u}^t, \mathbf{w}^t) - L_i^t + \lambda \langle \mathbf{w}^t, \mathbf{w}_i^t \rangle \quad \text{where} \quad L_i^t = \langle \boldsymbol{\alpha}_i^t, \mathbf{b}_i \rangle$$

$$\|\mathbf{A}(\boldsymbol{\beta}^t - \boldsymbol{\alpha}^t)\|^2 = \lambda \left\| \frac{1}{m\lambda} \boldsymbol{\Psi}_i(\mathbf{u}^t) + \mathbf{w}_i^t \right\|^2$$

- ◆ Starting from  $\mathbf{w}^0 = \mathbf{w}_1^0 = \dots = \mathbf{w}_m^0 = \mathbf{0}$ ,  $L_1^0 = \dots = L_m^0 = 0$ , the key quantities can be directly updated:

$$L_i^{t+1} = \langle \mathbf{b}_i, \boldsymbol{\alpha}_i^t(1 - \tau^t) + \tau^t \boldsymbol{\beta}_i^t \rangle = L_i^t(1 - \tau^t) + \frac{\tau^t}{m} \ell_i(\mathbf{u}^t)$$

$$\mathbf{w}_i^{t+1} = -\frac{1}{\sqrt{\lambda}} \mathbf{A} (\boldsymbol{\alpha}_i^t(1 - \tau^t) + \boldsymbol{\beta}_i^t \tau^t) = \mathbf{w}_i^t(1 - \tau^t) - \frac{\tau^t}{\lambda m} \boldsymbol{\Psi}_i(\mathbf{u}^t)$$

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \mathbf{w}_i^{t+1} - \mathbf{w}_i^t$$

## 10.A: Algorithm: SDA with steepest feasible ascent (the primal form)

**Input:** precision parameter  $\varepsilon > 0$

**Output:**  $\varepsilon$ -optimal primal solution  $\mathbf{w}^t$

**Initialization:**

$t := 0, \mathbf{w}^0 := \mathbf{0}, (\mathbf{w}_i^0 := \mathbf{0}, i \in \mathcal{I}), (L_i^0 := 0, i \in \mathcal{I})$

**repeat**

num\_updates := 0

**forall**  $i \in \mathcal{I}$  **do**

$\mathbf{u}^t \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} s_i(\mathbf{y}, \mathbf{w}^t)$

**if**  $s_i(\mathbf{u}^t, \mathbf{w}^t) - mL_i^t + m\lambda \langle \mathbf{w}^t, \mathbf{w}_i^t \rangle > \varepsilon$  **then**

num\_updates  $\leftarrow$  num\_updates + 1

$$\tau^t = \min \left\{ 1, \frac{\frac{1}{m}s_i(\mathbf{u}^t, \mathbf{w}^t) - L_i^t + \lambda \langle \mathbf{w}^t, \mathbf{w}_i^t \rangle}{\lambda \left\| \frac{1}{m\lambda} \Psi_i(\mathbf{u}^t) + \mathbf{w}_i^t \right\|^2} \right\}$$

$$\mathbf{w}_i^t = \mathbf{w}_i^t(1 - \tau^t) - \frac{\tau^t}{\lambda m} \Psi_i(\mathbf{u}^t)$$

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \mathbf{w}_i^{t+1} - \mathbf{w}_i^t$$

$$L_i^{t+1} = L_i^t(1 - \tau^t) + \frac{\tau^t}{m} \ell_i(\mathbf{u}^t)$$

$t := t + 1$

**until** num\_updates = 0 ;

## 10.A: SDA - improvement of the dual objective

- ◆ Provided  $\tau^{t+1} = 1$ , the improvement of the objective reads

$$\begin{aligned}
 \delta(\boldsymbol{\alpha}^t, \boldsymbol{\beta}^t) &= D(\boldsymbol{\alpha}^{t+1}) - D(\boldsymbol{\alpha}^t) \\
 &= D(\boldsymbol{\beta}^t) - D(\boldsymbol{\alpha}^t) \\
 &= \langle \boldsymbol{\alpha}^t - \boldsymbol{\beta}^t, \mathbf{b} - \mathbf{A}^T \mathbf{A} \boldsymbol{\alpha}^t \rangle - \frac{1}{2} \|\mathbf{A}(\boldsymbol{\alpha}^t - \boldsymbol{\beta}^t)\|^2
 \end{aligned}$$

- ◆ Provided  $\tau^{t+1} < 1$ , the improvement of the objective reads

$$\begin{aligned}
 \delta(\boldsymbol{\alpha}^t, \boldsymbol{\beta}^t) &= D(\boldsymbol{\alpha}^{t+1}) - D(\boldsymbol{\alpha}^t) \\
 &= D(\boldsymbol{\alpha}^t(1 - \tau^{t+1}) + \tau^{t+1}\boldsymbol{\beta}^t) - D(\boldsymbol{\alpha}^t) \\
 &= \dots \\
 &= \frac{\langle \boldsymbol{\alpha}^t - \boldsymbol{\beta}^t, \mathbf{b} - \mathbf{A}^T \mathbf{A} \boldsymbol{\alpha}^t \rangle^2}{2\|\mathbf{A}(\boldsymbol{\alpha}^t - \boldsymbol{\beta}^t)\|^2}
 \end{aligned}$$

## 10.A: SDA - variable selection based on near maximal gain

- ◆ Construct the  $\beta^t$  such that only two variables in the  $i$ -th block are updated:

$$\beta_j^t(\mathbf{y}) = \begin{cases} \alpha_i^t(\mathbf{u}^t) + \alpha_i^t(\mathbf{v}^t) & \text{if } j = i \wedge \mathbf{y} = \mathbf{u}^t \\ 0 & \text{if } j = i \wedge \mathbf{y} = \mathbf{v}^t \\ \alpha_j^t(\mathbf{y}^t) & \text{if } j \neq i \end{cases}$$

- ◆ If the labels  $\mathbf{u}^t$  and  $\mathbf{v}^t$  are selected such that  $D'_{\mathcal{A}^t}(0) > 0$ , which is equivalent to

$$\alpha_i^t(\mathbf{v}^t)(s_i(\mathbf{u}^t, \mathbf{w}^t) - s_i(\mathbf{v}^t, \mathbf{w}^t)) > 0,$$

then the improvement  $\delta^t(i, \mathbf{u}, \mathbf{v}) = D(\alpha^{t+1}) - D(\alpha^t)$  reads

$$\delta^t(i, \mathbf{u}, \mathbf{v}) = \begin{cases} \frac{\lambda(s_i(\mathbf{u}, \mathbf{w}^t) - s_i(\mathbf{v}, \mathbf{w}^t))}{\|\Psi_i(\mathbf{u}) - \Psi_i(\mathbf{v})\|^2} & \text{if } \tau^{t+1} < 1 \\ \alpha_i^t(\mathbf{v})(s_i(\mathbf{u}, \mathbf{w}^t) - s_i(\mathbf{v}, \mathbf{w}^t)) - \frac{\alpha_i^t(\mathbf{v})^2}{2\lambda} \|\Psi_i(\mathbf{u}) - \Psi_i(\mathbf{v})\|^2 & \text{if } \tau^{t+1} = 1 \end{cases}$$

- ◆ Instead of trying all possible  $|\mathcal{Y}|^2$  options we can approximate:

$$\mathbf{u}^t \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} s_i(\mathbf{y}, \mathbf{w}^t) \quad \text{and} \quad \mathbf{v}^t \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}_i} \delta^t(i, \mathbf{u}^t, \mathbf{y})$$

where  $\mathcal{Y}_i = \{\mathbf{y} \in \mathcal{Y} \mid \alpha_i^t(\mathbf{y}) > 0\}$  is a set of non-zero dual variables.

## 10.A: Algorithm: SDA with near maximal gain

**Initialization:**

$$\mathbf{w} := \mathbf{0}, \mathcal{Y}_i := \{\mathbf{y}^i\}, \alpha_i(\mathbf{y}) := \begin{cases} \frac{1}{m} & \text{if } \mathbf{y} = \mathbf{y}^i \\ 0 & \text{otherwise} \end{cases}, i \in \mathcal{I}$$

**repeat**

num\_updates := 0

**forall**  $i \in \mathcal{I}$  **do**

$\mathbf{u}_1 := \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} s_i(\mathbf{y}, \mathbf{w})$

$\mathbf{u}_2 := \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}_i} s_i(\mathbf{y}, \mathbf{w})$

**if**  $s_i(\mathbf{u}_1, \mathbf{w}) > s_i(\mathbf{u}_2, \mathbf{w})$  **then**

└  $\hat{\mathbf{u}} \leftarrow \mathbf{u}_1$

**else**

└  $\hat{\mathbf{u}} \leftarrow \mathbf{u}_2$

**if**  $s_i(\hat{\mathbf{u}}, \mathbf{w}) - m \sum_{\mathbf{y} \in \mathcal{Y}_i} \alpha_i(\mathbf{y}) s_i(\mathbf{y}, \mathbf{w}) > \varepsilon$  **then**

num\_updates := num\_updates + 1

$\hat{\mathbf{v}} := \operatorname{argmax}_{\mathbf{y} \in \{\mathbf{y}' \in \mathcal{Y}_i \mid \alpha_i(\mathbf{y}') > 0\}} \delta(i, \hat{\mathbf{u}}, \mathbf{y})$

$\tau := \min \left\{ 1, \frac{\lambda(s_i(\hat{\mathbf{u}}, \mathbf{w}) - s_i(\hat{\mathbf{v}}, \mathbf{w}))}{\alpha_i(\hat{\mathbf{v}}) \|\boldsymbol{\psi}_i(\hat{\mathbf{u}}) - \boldsymbol{\psi}_i(\hat{\mathbf{v}})\|^2} \right\}$

$\mathbf{w} := \mathbf{w} + (\boldsymbol{\psi}_i(\hat{\mathbf{v}}) - \boldsymbol{\psi}_i(\hat{\mathbf{u}})) \frac{\tau \alpha_i(\hat{\mathbf{v}})}{\lambda}$

$\alpha_i(\hat{\mathbf{u}}) := \alpha_i(\hat{\mathbf{u}}) + \tau \alpha_i(\hat{\mathbf{v}})$

$\alpha_i(\hat{\mathbf{v}}) := \alpha_i(\hat{\mathbf{v}}) - \tau \alpha_i(\hat{\mathbf{v}})$

**if**  $\hat{\mathbf{u}} = \mathbf{u}_1$  **then**

└  $\mathcal{Y}_i := \mathcal{Y}_i \cup \{\mathbf{u}_1\}$

**until** num\_updates = 0 ;

## 10.A: SDA - memory requirements

- ◆ SDA with steepest feasible ascent sel. strategy:  $\mathcal{O}(m \cdot n)$  where  $m$  is the number of examples and  $n$  the number of primal parameters.
- ◆ SDA with near maximal gain sel. strategy:  $\mathcal{O}(n' \cdot T)$  where  $n'$  is the number of non-zero elements of  $\Psi_i(\mathbf{y})$  and  $T$  is the total number of updates.

## 10.A: SDA - convergence guarantees

- ◆ The convergence theorem applies for both the SDA with the steepest feasible ascent strategy as well as the near maximal gain strategy.
- ◆ **Convergence theorem:** For any  $\varepsilon > 0$  and  $\lambda > 0$ , the SDA algorithm terminates after

$$T = \frac{8\ell_{\max}D^2}{\varepsilon^2 \lambda}$$

updates at most where  $\ell_{\max} = \max_{i \in \mathcal{I}, \mathbf{y} \in \mathcal{Y}} \ell_i(\mathbf{y})$  and  $D = \max_{i \in \mathcal{I}, \mathbf{y} \in \mathcal{Y}} \|\Psi_i(\mathbf{y})\|$ .

- ◆ Although both variants of the SDA algorithm have the same upper bound on the number of iterations their practical performance is different.



## 10.A: Experimental comparison

### Compared methods

- ◆ **CPA** Cutting Plane Algorithm: implements the Bundle Method for Risk Minimization (last lecture, 9.A)
- ◆ **BCFW** Block-Coordinate Frank-Wolfe: implements SDA with steepest feasible ascent selection strategy (this lecture, 10.A)
- ◆ **SDM** Sequential Dual Method : implements SDA with steepest feasible ascent and changing only two dual variables at time (not covered)
- ◆ **FASOLE** Fast Algorithm for Structured Output LEarning: implements SDA with near maximal gain selection strategy (this lecture, 10.A)

### Benchmark data

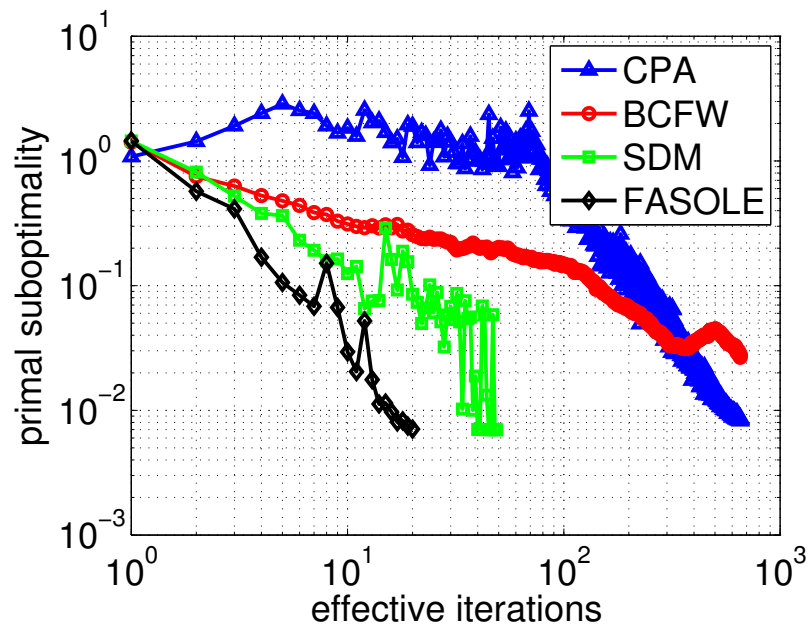
dataset		#train	#test	#params	max-sum classifier
OCR	[Taskar 2003]	5,512	1,365	4,004	HMM-chain
LANDMARK	[Uricar 2012]	5,062	3,512	232,476	HMM-tree

### Evaluation protocol

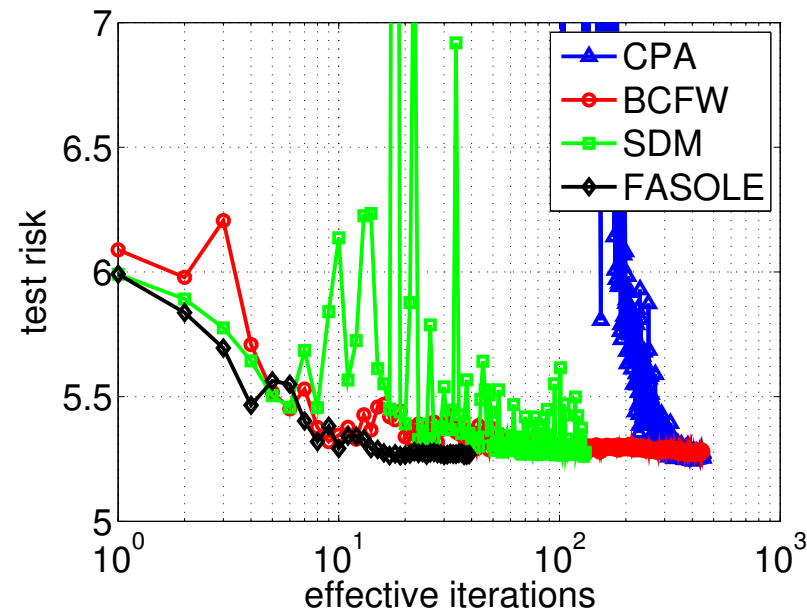
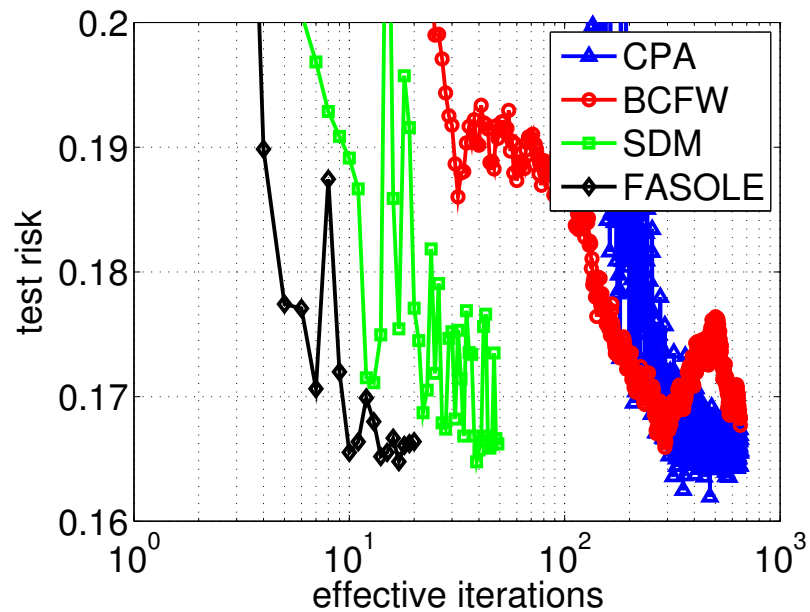
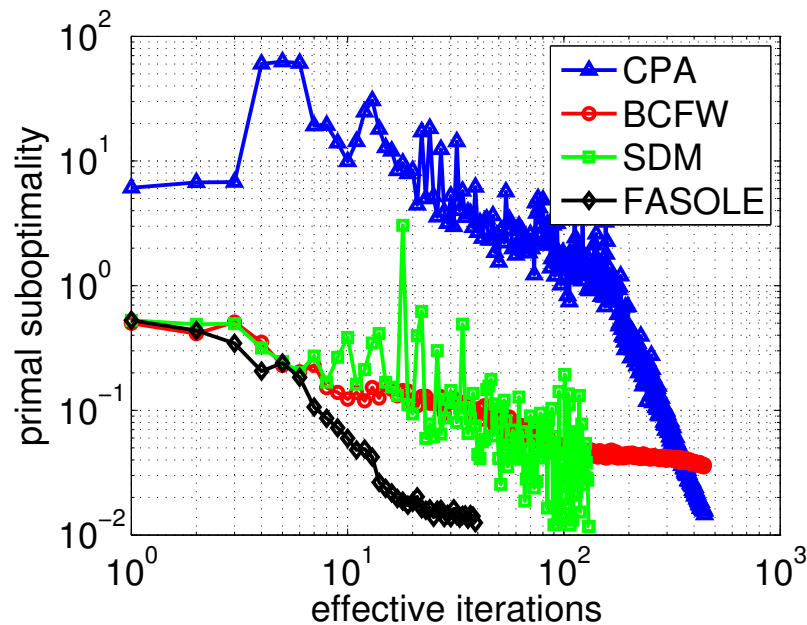
The solvers run to find the same  $\varepsilon$ -precise solution for a range of regularization parameters  $\lambda$ . We measure the convergence speed in terms of the objective value and the test risk.

# 10.A: Experimental comparison

OCR



LANDMARK



## 10.B: Stochastic Gradient Descent - formulation

- ◆ Let us consider a convex constrained problem

$$\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} F(\mathbf{w})$$

where  $\mathcal{W} \subset \mathbb{R}^n$  is a closed convex set and  $F: \mathcal{W} \rightarrow \mathbb{R}$  is a convex function.

- ◆ The stochastic gradient descent approaches the objective  $F$  via an oracle which for given  $\mathbf{w}^t$  provides a stochastic estimate  $\hat{\mathbf{g}}^t$  of the sub-gradient  $\mathbf{g}^t \in \partial F(\mathbf{w}^t)$  such that

$$\mathbb{E} \hat{\mathbf{g}}^t = \mathbf{g}^t$$

- ◆ For example, in our setting

$$F(\mathbf{w}) = \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i \in \mathcal{I}} \ell_i(\mathbf{w}) = \frac{1}{m} \sum_{i \in \mathcal{I}} \left( \frac{\lambda}{2} \|\mathbf{w}\|^2 + \ell_i(\mathbf{w}) \right) = \frac{1}{m} \sum_{i \in \mathcal{I}} F_i(\mathbf{w})$$

the oracle can pick  $i \in \mathcal{I}$  randomly with the uniform distribution over  $\mathcal{I}$  and it provides a sub-gradient

$$\hat{\mathbf{g}}^t \in \partial F_i(\mathbf{w}^t)$$

## 10.B: Stochastic Gradient Descent

- ◆ The SGD algorithm: starting from  $\mathbf{w}^1 = \mathbf{0}$ , SGD computes new iterates recursively as follows

$$\mathbf{w}^{t+1} = \Pi_{\mathcal{W}}(\mathbf{w}^t - \eta^t \hat{\mathbf{g}}^t)$$

where  $\Pi_{\mathcal{W}}: \mathbb{R}^n \rightarrow \mathcal{W}$  denotes projection on  $\mathcal{W}$ , i.e.

$$\Pi_{\mathcal{W}}(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w}' \in \mathcal{W}} \|\mathbf{w}' - \mathbf{w}\|$$

and  $\eta^t$  is a sequence of step-sizes.

- ◆ The theoretical results require a fixed step size, typically,  $\sum_{t=0}^{\infty} \eta^t = \infty$  and  $\lim_{t \rightarrow \infty} \eta^t = 0$ .
- ◆ There is no stopping condition which would provide a certificate of optimality. In practice, SGD is stopped based on monitoring a progress of the objective or a validation error.

## 10.B: Stochastic Gradient Descent - convergence guarantees

- ◆ Recall that a function  $F: \mathcal{W} \rightarrow \mathbb{R}$  is  $\lambda$ -strictly convex iff the function  $F(\mathbf{w}) - \frac{\lambda}{2}\|\mathbf{w}\|^2$  is convex.

For example,  $F(\mathbf{w}) = \frac{\lambda}{2}\|\mathbf{w}\|^2 + R(\mathbf{w})$  is  $\lambda$ -strictly convex iff  $R(\mathbf{w})$  is convex.

- ◆  **$\lambda$ -strictly convex functions:** Suppose  $F$  is  $\lambda$ -strictly convex, and that  $\mathbb{E}[\|\hat{\mathbf{g}}^t\|^2] \leq G^2$ ,  $\forall t$ . Consider SGD with step sizes  $\eta^t = \frac{1}{\lambda t}$ . Then for any  $t > 1$ , it holds that

$$\mathbb{E}[F(\mathbf{w}^t) - F(\mathbf{w}^*)] \leq \frac{17G^2(1 + \log(t))}{\lambda t}$$

- ◆ **Convex functions:** Assume that  $F$  is convex and that for some constants  $D, G$  it holds that  $\mathbb{E}[\|\hat{\mathbf{g}}^t\|] \leq G$ ,  $\forall t$ , and  $\sup_{\mathbf{w}, \mathbf{w}' \in \mathcal{W}} \|\mathbf{w} - \mathbf{w}'\| \leq D$ . Consider SGD with step size  $\eta^t = \frac{c}{\sqrt{t}}$  where  $c > 0$  is a constant. Then for any  $t > 1$  it holds that

$$\mathbb{E}[F(\mathbf{w}^t) - F(\mathbf{w}^*)] \leq \left( \frac{D^2}{c} + cG^2 \right) \frac{2 + \log(t)}{\sqrt{t}}$$

## 10.B: Stochastic Gradient Descent with averaging

- ◆ Instead of using the last iterate of the SGD, take average from the first till the last iteration

$$\bar{w}^t = \frac{1}{t} \sum_{i=1}^t w^i$$

which can be computed iteratively by setting  $\bar{w}^1 = w^1$  and update recursively

$$\bar{w}^t = \left(1 - \frac{1}{t}\right) \bar{w}^{t-1} + \frac{1}{t} w^t$$

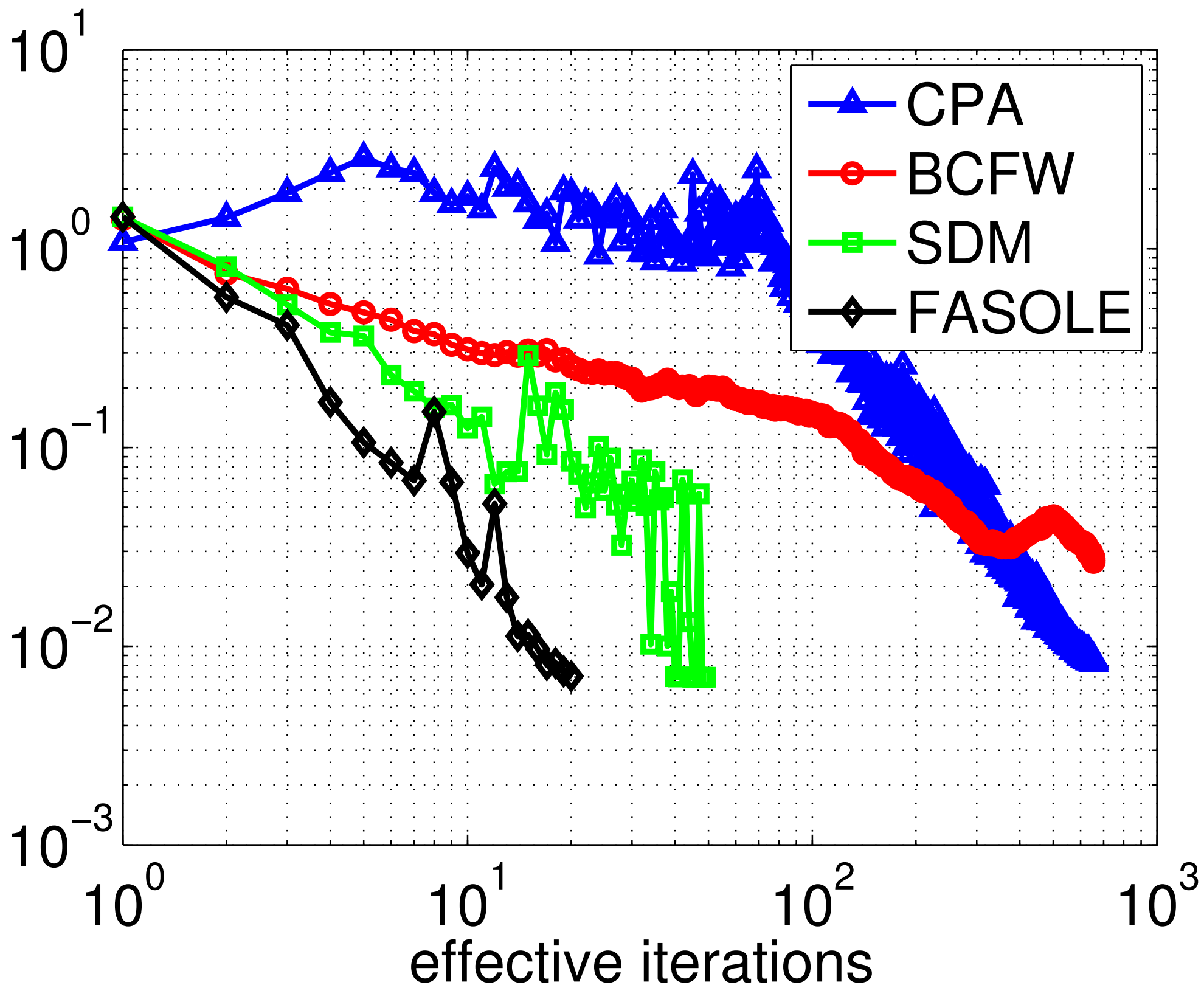
- ◆ Better results can be obtained with the polynomial decay averaging parametrized by  $\nu$ :

$$\bar{w}_\nu^t = \left(1 - \frac{\nu + 1}{t + \nu}\right) \bar{w}_\nu^{t-1} + \frac{\nu + 1}{t + \nu} w^t$$

which for  $\nu = 0$  reduces to the standard averaging while  $\nu > 0$  increases the weight of the last iterate.

- ◆ The averaging schemes provide convergence guarantees even for a fixed step-size.

primal suboptimality



primal suboptimality

