

Stochastic Dual Coordinate Ascent for training Conditional Random Fields

Rémi Le Priol

Montreal Institute of Learning Algorithms

remi.lp.17@gmail.com

September 15, 2017

- 1 Conditional Random Fields
- 2 SDCA for Max-likelihood
- 3 Non-Uniform Sampling
- 4 Leverage the Structure

- 1 Conditional Random Fields
- 2 SDCA for Max-likelihood
- 3 Non-Uniform Sampling
- 4 Leverage the Structure

data point $x \in \mathcal{X} \mapsto$ structured label $y \in \mathcal{Y}$

letter drawings \mapsto word

sentence in English \mapsto sentence in French

sentence \mapsto parsing tree

natural image \mapsto semantic segmentation

data point $x \in \mathcal{X} \mapsto$ structured label $y \in \mathcal{Y}$

letter drawings \mapsto word

sentence in English \mapsto sentence in French

sentence \mapsto parsing tree

natural image \mapsto semantic segmentation

Hypothesis

The conditional distribution $p(y|x)$ is Markov with respect to an undirected graphical model $G = (V, E)$.

Feature extractor:

$$F : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$$

Hypothesis

$$F(x, y) = \sum_{c \in \mathcal{C}} F_c(x, y_c)$$

where \mathcal{C} is the set of maximal cliques of G .

Conditional probability of y given x :

$$p(y|x; w) := \frac{\exp(w^T F(x, y))}{\sum_{y' \in \mathcal{Y}} \exp(w^T F(x, y'))}$$

Standard approach to train CRF: Maximum Likelihood

$$\min_w \mathcal{P}(w) = \frac{\lambda}{2} \|w\|^2 - \frac{1}{n} \sum_{i=1}^n \log(p(y_i|x_i; w))$$

Reformulation

New notations:

Corrected features: $\psi_i(y) := F(x_i, y_i) - F(x_i, y) \in \mathbb{R}^d$

Corrected feature matrix: $A_i := (\psi_i(1), \psi_i(2), \dots, \psi_i(|\mathcal{Y}_i|)) \in \mathbb{R}^{d \times |\mathcal{Y}_i|}$

Log-sum-exp function: $\phi_i(z) = \log \left(\sum_{y \in \mathcal{Y}_i} \exp(z_y) \right)$

Reformulation

New notations:

Corrected features: $\psi_i(y) := F(x_i, y_i) - F(x_i, y) \in \mathbb{R}^d$

Corrected feature matrix: $A_i := (\psi_i(1), \psi_i(2), \dots, \psi_i(|\mathcal{Y}_i|)) \in \mathbb{R}^{d \times |\mathcal{Y}_i|}$

Log-sum-exp function: $\phi_i(z) = \log \left(\sum_{y \in \mathcal{Y}_i} \exp(z_y) \right)$

The log-likelihood becomes :

$$\begin{aligned} -\log(p(y_i|x_i; w)) &= \log \left(\sum_y e^{w^T F(x_i, y)} \right) - w^T F(x_i, y_i) \\ &= \log \left(\sum_y e^{-w^T \psi_i(y)} \right) \\ &= \phi_i(-A_i^T w) \end{aligned}$$

Primal Objective

$$\min_{w \in \mathbb{R}^d} \mathcal{P}(w) = \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \phi_i(-A_i^T w) \quad (1)$$

$$\min_{w \in \mathbb{R}^d} \mathcal{P}(w) = \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \phi_i(-A_i^T w) \quad (1)$$

HARD ! A_i is huge.

- Exponentiated Gradient by Collins in 2008¹
- Non-Uniform Sampling – Stochastic Average Gradient (NUS-SAG) by Schmidt in 2014²

¹Collins et al., “Exponentiated Gradient Algorithms for Conditional Random Fields and Max-Margin Markov Networks”.

²Schmidt et al., “Non-Uniform Stochastic Average Gradient Method for Training Conditional Random Fields”.

³Shalev-Shwartz and Zhang, “Accelerated Proximal Stochastic Dual Coordinate Ascent for Regularized Loss Minimization”.

Anterior Work

- Exponentiated Gradient by Collins in 2008¹
- Non-Uniform Sampling – Stochastic Average Gradient (NUS-SAG) by Schmidt in 2014²
- SDCA³

¹Collins et al., “Exponentiated Gradient Algorithms for Conditional Random Fields and Max-Margin Markov Networks”.

²Schmidt et al., “Non-Uniform Stochastic Average Gradient Method for Training Conditional Random Fields”.

³Shalev-Shwartz and Zhang, “Accelerated Proximal Stochastic Dual Coordinate Ascent for Regularized Loss Minimization”.

- Exponentiated Gradient by Collins in 2008¹
- Non-Uniform Sampling – Stochastic Average Gradient (NUS-SAG) by Schmidt in 2014²
- SDCA³

Variance reduced methods.

¹Collins et al., “Exponentiated Gradient Algorithms for Conditional Random Fields and Max-Margin Markov Networks”.

²Schmidt et al., “Non-Uniform Stochastic Average Gradient Method for Training Conditional Random Fields”.

³Shalev-Shwartz and Zhang, “Accelerated Proximal Stochastic Dual Coordinate Ascent for Regularized Loss Minimization”.

- Exponentiated Gradient by Collins in 2008¹
- Non-Uniform Sampling – Stochastic Average Gradient (NUS-SAG) by Schmidt in 2014²
- SDCA³

Variance reduced methods.

Why SDCA?

¹Collins et al., “Exponentiated Gradient Algorithms for Conditional Random Fields and Max-Margin Markov Networks”.

²Schmidt et al., “Non-Uniform Stochastic Average Gradient Method for Training Conditional Random Fields”.

³Shalev-Shwartz and Zhang, “Accelerated Proximal Stochastic Dual Coordinate Ascent for Regularized Loss Minimization”.

- Exponentiated Gradient by Collins in 2008¹
- Non-Uniform Sampling – Stochastic Average Gradient (NUS-SAG) by Schmidt in 2014²
- SDCA³

Variance reduced methods.

Why SDCA?

Exact line search for cheap!

¹Collins et al., “Exponentiated Gradient Algorithms for Conditional Random Fields and Max-Margin Markov Networks”.

²Schmidt et al., “Non-Uniform Stochastic Average Gradient Method for Training Conditional Random Fields”.

³Shalev-Shwartz and Zhang, “Accelerated Proximal Stochastic Dual Coordinate Ascent for Regularized Loss Minimization”.

Overview

- 1 Conditional Random Fields
- 2 SDCA for Max-likelihood
- 3 Non-Uniform Sampling
- 4 Leverage the Structure

Dual Formulation

Primal:

$$\min_{w \in \mathbb{R}^d} \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \phi_i(-A_i^T w)$$

Dual:

$$\max_{\alpha | \forall i, \alpha_i \in \Delta_i} \mathcal{D}(\alpha) = -\frac{1}{2\lambda} \left\| \frac{1}{n} \sum_i A_i \alpha_i \right\|^2 + \frac{1}{n} \sum_{i=1}^n H_i(\alpha_i)$$

Δ_i is the simplex of dimension $|\mathcal{Y}_i|$. H_i is the entropy over Δ_i .

Conjugate variables

Primal probabilities

$$\forall i, \alpha_i(w) := \nabla \phi_i(-A_i^T w) = p(\cdot | x_i; w) \propto \exp(-w^T \psi_i(\cdot))$$

Dual weights

$$\begin{aligned} \mathbf{w}(\alpha) &= \frac{1}{\lambda n} \sum_i A_i \alpha_i = \frac{1}{\lambda n} \sum_i \mathbf{E}_{\alpha_i}[\psi_i] \\ &= \frac{1}{\lambda n} \sum_i F(x_i, y_i) - \frac{1}{\lambda n} \sum_i \mathbf{E}_{y \sim \alpha_i}[F(x_i, y)] \end{aligned}$$

Conjugate variables

Primal probabilities

$$\forall i, \alpha_i(w) := \nabla \phi_i(-A_i^T w) = p(\cdot | x_i; w) \propto \exp(-w^T \psi_i(\cdot))$$

Dual weights

$$\begin{aligned} \mathbf{w}(\alpha) &= \frac{1}{\lambda n} \sum_i A_i \alpha_i = \frac{1}{\lambda n} \sum_i \mathbf{E}_{\alpha_i}[\psi_i] \\ &= \frac{1}{\lambda n} \sum_i F(x_i, y_i) - \frac{1}{\lambda n} \sum_i \mathbf{E}_{y \sim \alpha_i}[F(x_i, y)] \end{aligned}$$

Optimality condition

$$\alpha^* = \alpha(w^*) \quad \text{and} \quad w^* = \mathbf{w}(\alpha^*)$$

$$\max_{\alpha | \forall i, \alpha_i \in \Delta_i} \mathcal{D}(\alpha) = \underbrace{-\frac{\lambda}{2} \|\mathbf{w}(\alpha)\|^2}_{\text{data fitting}} + \underbrace{\frac{1}{n} \sum_{i=1}^n H_i(\alpha_i)}_{\text{regularization}}$$

$$\max_{\alpha | \forall i, \alpha_i \in \Delta_i} \mathcal{D}(\alpha) = \underbrace{-\frac{\lambda}{2} \|\mathbf{w}(\alpha)\|^2}_{\text{data fitting}} + \underbrace{\frac{1}{n} \sum_{i=1}^n H_i(\alpha_i)}_{\text{regularization}}$$

HARD ! α is huge.

Principle of SDCA⁴

- Store dual probabilities α and $\mathbf{w}(\alpha)$.
- Sample $i \in \{1, \dots, n\}$
- Update $\alpha_i^+ \leftarrow (1 - \gamma)\alpha_i + \gamma\alpha_i(\mathbf{w}(\alpha))$ with $\gamma \in [0, 1]$

⁴Shalev-Shwartz and Zhang, “Accelerated Proximal Stochastic Dual Coordinate Ascent for Regularized Loss Minimization”.

Algorithm 1 SDCA for Logistic Regression

$\forall i$, initialize $\alpha_i^{(0)}$ at random in Δ_K

Let $w^{(0)} = \frac{1}{\lambda n} A \alpha$

Let $\forall i, g_i = 1$ (optional)

for $k = 0 \dots K$ **do**

 Pick i at random in $\{1, \dots, n\}$ (optionally, proportional to g_i)

 Let $\beta_i := p(\cdot | x; w^{(k)}) = \alpha_i(\mathbf{w}(\alpha_i^{(k)}))$

 Let $g_i = D_{KL}(\alpha_i || \beta_i)$ (optional)

 Let $d_i = \beta_i - \alpha_i^{(k)}$ (dual ascent direction)

 Let $v_i = \frac{1}{\lambda n} A_i d_i$ (primal descent direction)

 Solve $\gamma^* = \arg \max_{\gamma \in [0,1]} H_i(\alpha_i^{(k)} + \gamma d_i) - \frac{\lambda n}{2} \|w^{(k)} + \gamma v_i\|^2$ (Line Search)

 Update $\alpha_i^{(k+1)} := \alpha_i^{(k)} + \gamma^* d_i$

 Update $w^{(k+1)} := w^{(k)} + \gamma^* v_i$

Overview

- 1 Conditional Random Fields
- 2 SDCA for Max-likelihood
- 3 Non-Uniform Sampling**
- 4 Leverage the Structure

On duality gaps

Duality gap:

$$g(w, \alpha) = \mathcal{P}(w) - \mathcal{D}(\alpha)$$

Primal model gap:

$$g(w, \alpha(w)) = \frac{\lambda}{2} \|w - \mathbf{w}(\alpha(w))\|^2 = \frac{1}{2\lambda} \|\nabla \mathcal{P}(w)\|^2$$

Dual model gap:

$$g(\mathbf{w}(\alpha), \alpha) = \frac{1}{n} \sum_i D_{KL}(\alpha_i || \alpha_i(\mathbf{w}(\alpha)))$$

Our Scheme

- At each step, update individual duality gap $g_i = D_{KL}(\alpha_i || \alpha_i(\mathbf{w}(\alpha)))$
- Sample i proportionally to g_i .

^aOsokin et al., “Minding the Gaps for Block Frank-Wolfe Optimization of Structured SVMs”.

Our Scheme

- At each step, update individual duality gap $g_i = D_{KL}(\alpha_i || \alpha_i(\mathbf{w}(\alpha)))$
- Sample i proportionally to g_i .

Scheme adapted from Block-Coordinate Frank-Wolfe^a.

Transposable to the exponentiated gradient.

^aOsokin et al., “Minding the Gaps for Block Frank-Wolfe Optimization of Structured SVMs”.

Our Scheme

- At each step, update individual duality gap $g_i = D_{KL}(\alpha_i || \alpha_i(\mathbf{w}(\alpha)))$
- Sample i proportionally to g_i .

Scheme adapted from Block-Coordinate Frank-Wolfe^a.

Transposable to the exponentiated gradient.

^aOsokin et al., “Minding the Gaps for Block Frank-Wolfe Optimization of Structured SVMs”.

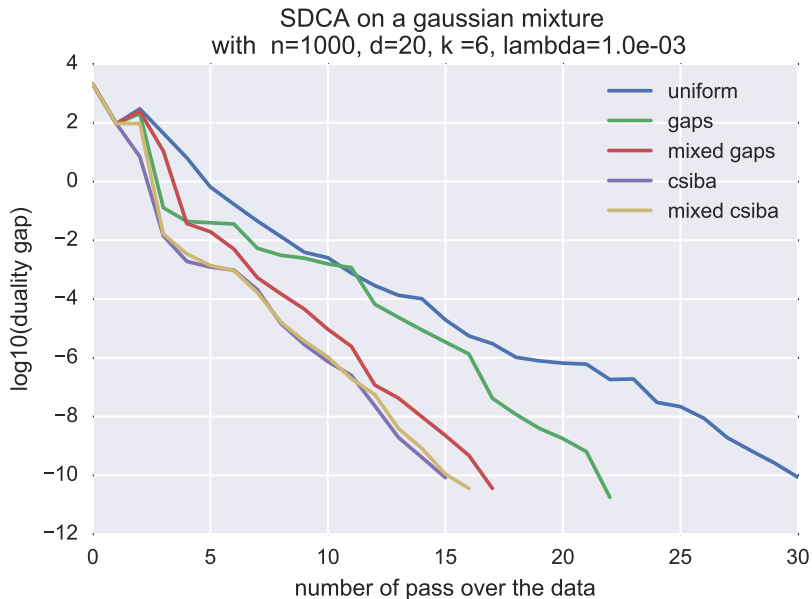
Competing scheme

$$g_i = \|\alpha_i(\mathbf{w}(\alpha)) - \alpha_i\| \sqrt{R_i^2 + 2\lambda n}$$

where R_i is the operator norm of A_i .^a

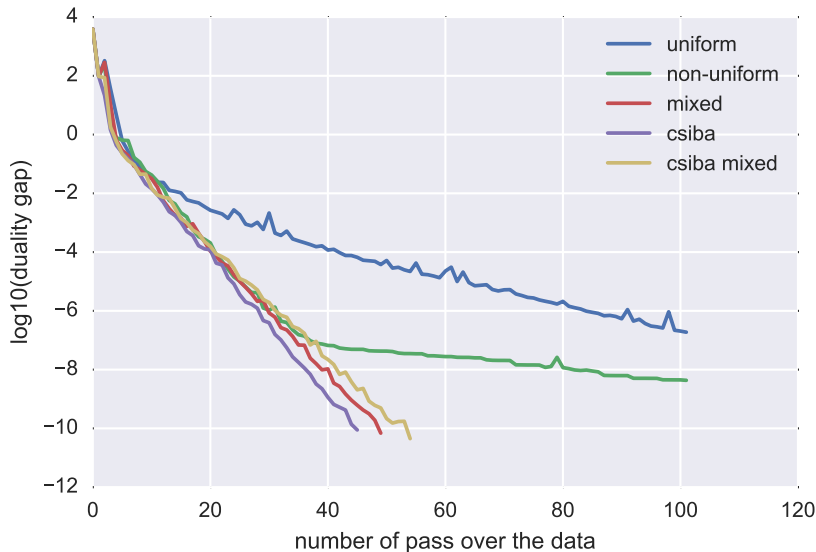
^aCsiba, Qu, and Richtarik, “Stochastic Dual Coordinate Ascent with Adaptive Probabilities”.

Results 1



Results 2

Training of SDCA on Covertypes with $n=10000$ and $\lambda=1.0e-04$



Overview

- 1 Conditional Random Fields
- 2 SDCA for Max-likelihood
- 3 Non-Uniform Sampling
- 4 Leverage the Structure

Joints and Marginals⁵

Marginal probability on the clique $c \in \mathcal{C}$.

$$\mu_{i,c}(y_c) := \sum_{y' \in \mathcal{Y}_i \mid y'_c = y_{i,c}} \alpha_i(y')$$

The marginals of the sample i live in the local consistency polytope L_i .

If $\mathcal{T} = (\mathcal{C}, \mathcal{S})$ is a junction tree of G :

$$\alpha(y) = \frac{\prod_{c \in \mathcal{C}_{max}} \mu_c(y_c)}{\prod_{s \in \mathcal{S}} \mu_s(y_s)} \quad (2)$$

Junction Tree algorithm = marginalization oracle.

We can compute $p_c(y_c | x_i; w) = \mu'_{i,c}(y_c)$.

⁵Taskar, Guestrin, and Koller, “Max-Margin Markov Networks”.

Marginals to the Weights

$$\begin{aligned}\mathbf{w}(\alpha) &= \frac{1}{\lambda n} \sum_i \mathbf{E}_{y \sim \alpha_i} [\psi_i(y)] \\ &= \frac{1}{\lambda n} \sum_i \sum_{c \in \mathcal{C}_i} \mathbf{E}_{y \sim \alpha_i} [\psi_{i,c}(y_c)] \\ &= \frac{1}{\lambda n} \sum_i \sum_{c \in \mathcal{C}_i} \mathbf{E}_{y_c \sim \mu_{i,c}} [\psi_{i,c}(y_c)]\end{aligned}$$

Marginal weights

$$\mathcal{W}(\mu) := \frac{1}{\lambda n} \sum_i \sum_c B_{i,c} \mu_{i,c}$$

where $B_{i,c}$ has size $d \times |\mathcal{Y}_c|$, with $\psi_{i,c}(y_c)$ in the column y_c .

Marginals to the Entropy

We express the entropy of the joint probability as a function of the marginals with equation 2.

$$\mathcal{H}(\mu) := H_{|\mathcal{Y}|}(\alpha) = \sum_c H_{|c|}(\mu_c) - \sum_s H_{|s|}(\mu_s)$$
$$\mathcal{D}(\mu||\mu') := D_{KL}(\alpha||\alpha') = \sum_c D(\mu_c||\mu'_c) - \sum_s D(\mu_s||\mu'_s)$$

Remark: We can directly transpose our non-uniform sampling scheme. This is not true for Cisba's scheme.

A New Dual Objective

$$\max_{\forall i, \mu_i \in \mathcal{L}_i} -\frac{\lambda}{2} \|\mathcal{W}(\mu)\|^2 + \frac{1}{n} \sum_i \mathcal{H}_i(\mu_i) \quad (3)$$

We apply the coordinate ascent directly on this objective

Algorithm 2 SDCA for CRF

Let $\forall i, c, \mu_{i,c}^{(0)} := \frac{1}{|c|}$ and $w^{(0)} := \frac{1}{\lambda n} B \mu^{(0)}$

Let $\forall i g_i = 1$ (optional)

for $k = 0 \dots K$ **do**

 Pick i at random in $\{1, \dots, n\}$ (optionally, proportional to g_i)

 Compute $\forall c, \mu'_{i,c}(y_c) := p(y_c | x; w^{(k)})$ (marginalization oracle)

 Let $g_i = \mathcal{D}(\mu_i || \mu'_i)$ (optional)

 Let $d_i = \mu'_i - \mu_i^{(k)}$ (ascent direction)

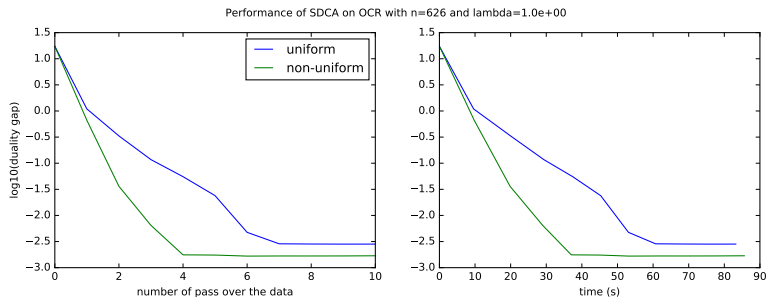
 Let $v_i = \frac{1}{\lambda} B_i d_i$ (primal direction)

 Solve $\gamma^* = \arg \max_{\gamma \in [0,1]} \mathcal{H}_i(\mu_i^{(k)} + \gamma d_i) - \frac{\lambda n}{2} \|w^{(k)} + \frac{\gamma}{n} v_i\|^2$ (Line Search)

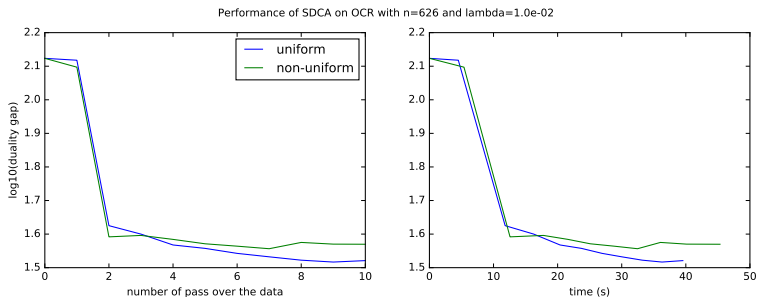
 Update $\mu_i^{(k+1)} := \mu_i^{(k)} + \gamma^* d_i$

 Update $w^{(k+1)} := w^{(k)} + \frac{\gamma^*}{n} v_i$

No convergence yet.

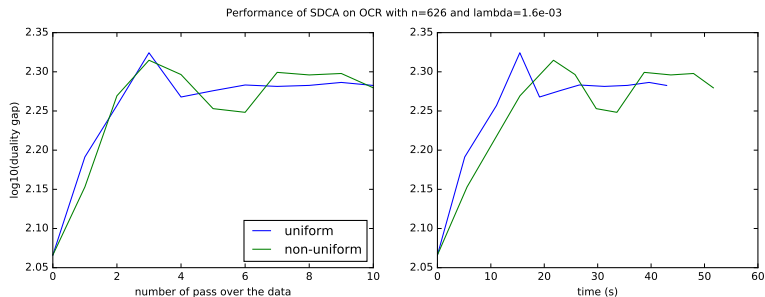


No convergence yet.



Results 1

No convergence yet.



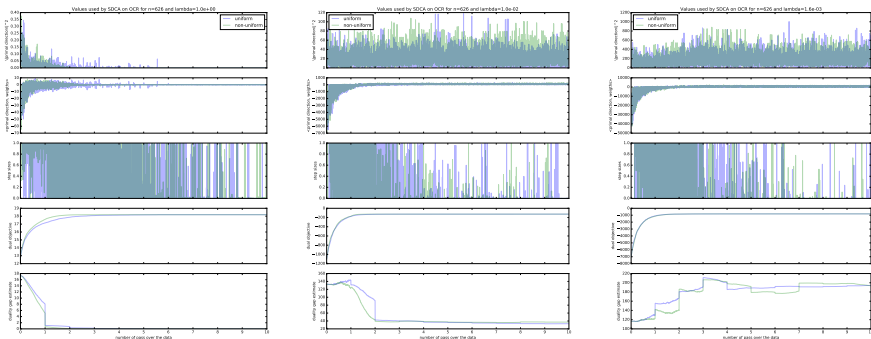


Figure: Some values of interest tracked along the run of SDCA.

Questions?

Appendix 1

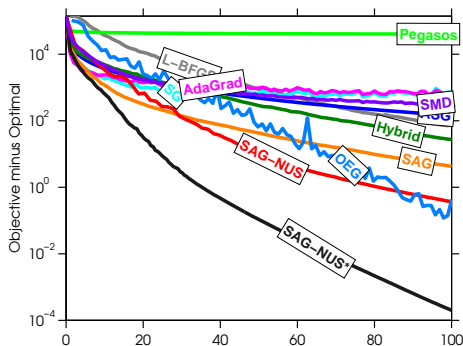


Figure: Training curves for various optimization algorithms on the OCR dataset. The x axis is the number of epochs, while the y axis is the primal suboptimality.

Appendix 1

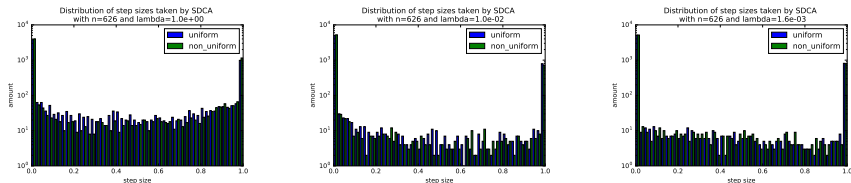


Figure: Distributions of step sizes taken by SDCA. The y axis is a log-scale. A large majority of steps are either taken with full size 1, either not taken at all. When the algorithm works better, with λ large, there are more intermediate step sizes.

Appendix 1

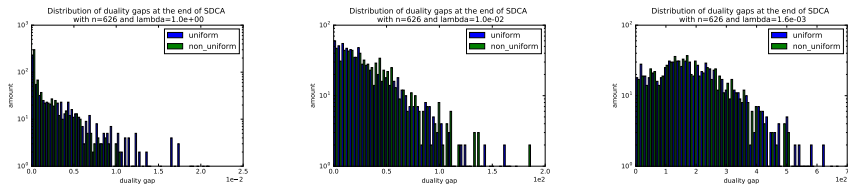


Figure: Distribution of individual duality gaps after the run of SDCA.