

Block-Coordinate Frank-Wolfe for Structural SVMs

Martin
Jaggi



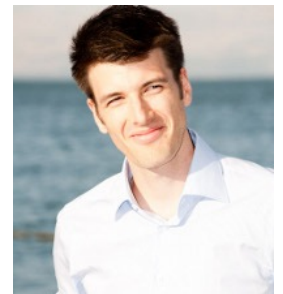
Simon
Lacoste-Julien



Mark
Schmidt



Patrick
Pletscher

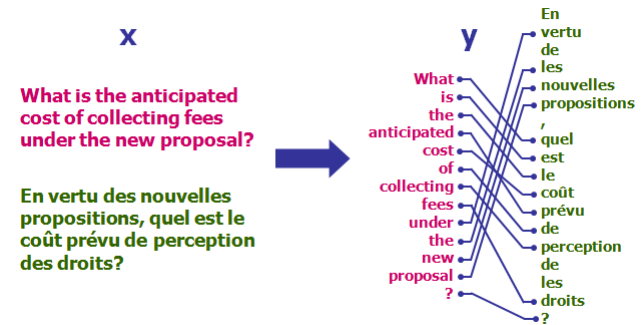


NIPS OPT2012 Workshop – Dec. 8th 2012

Structural SVM optimization

Structural SVM optimization

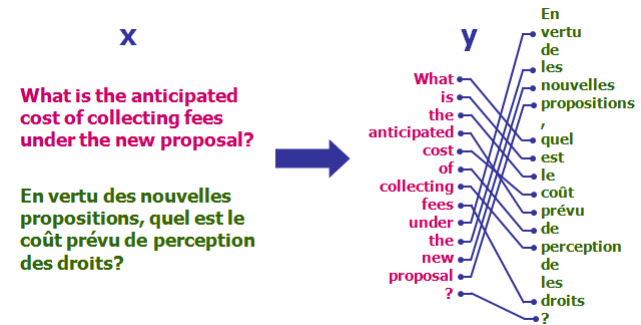
- structured prediction:



Structural SVM optimization

- structured prediction:
- learn linear classifier:

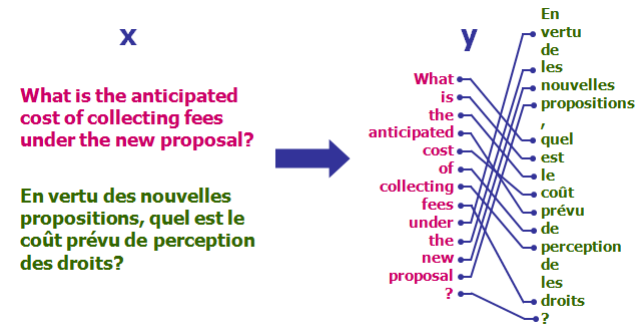
$$h_w(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle w, \phi(x, y) \rangle$$



Structural SVM optimization

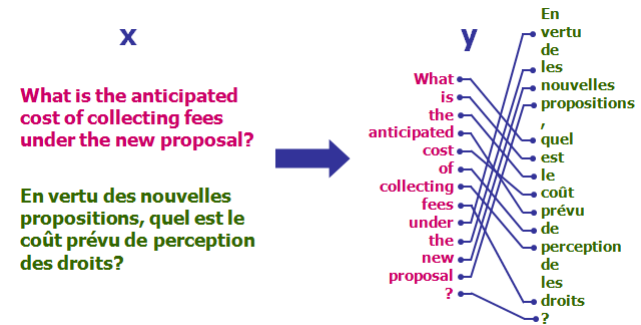
- structured prediction:
- learn linear classifier:

$$h_w(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle w, \phi(x, y) \rangle \quad \leftarrow \text{decoding}$$



Structural SVM optimization

- structured prediction:
- learn linear classifier:



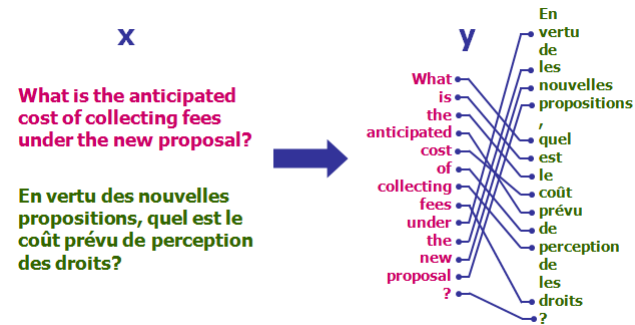
$$h_w(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle w, \phi(x, y) \rangle \quad \leftarrow \text{decoding}$$

- structural SVM primal:

$$\min_w \quad \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}} \left\{ L(y_i, y) + \langle w, \phi(x_i, y) \rangle \right\} - \langle w, \phi(x_i, y_i) \rangle$$

Structural SVM optimization

- structured prediction:
- learn linear classifier:



$$h_w(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle w, \phi(x, y) \rangle \quad \leftarrow \text{decoding}$$

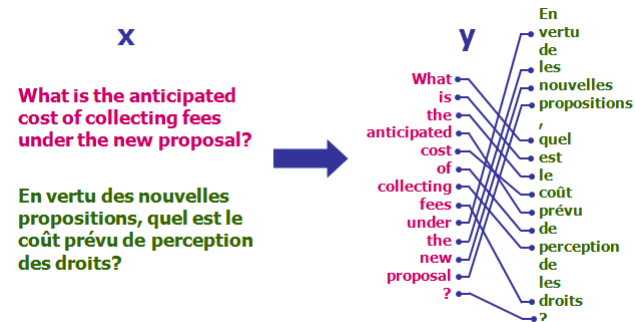
- structural SVM primal:

$$\min_w \quad \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}} \left\{ L(y_i, y) + \langle w, \phi(x_i, y) \rangle \right\} - \langle w, \phi(x_i, y_i) \rangle$$

$$\text{vs. binary hinge loss:} \quad \max \left\{ 0, 1 - \langle w, \phi(x_i) y_i \rangle \right\}$$

Structural SVM optimization

- structured prediction:
- learn linear classifier:



$$h_w(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle w, \phi(x, y) \rangle \quad \leftarrow \text{decoding}$$

structured hinge loss:

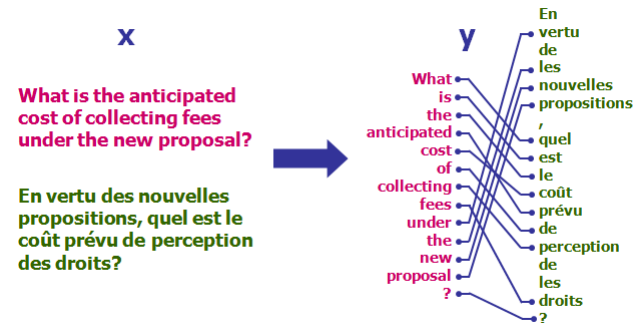
- structural SVM primal:

$$\min_w \quad \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}} \left\{ L(y_i, y) + \langle w, \phi(x_i, y) \rangle \right\} - \langle w, \phi(x_i, y_i) \rangle$$

vs. binary hinge loss: $\max \left\{ 0, 1 - \langle w, \phi(x_i) y_i \rangle \right\}$

Structural SVM optimization

- structured prediction:
- learn linear classifier:



$$h_w(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle w, \phi(x, y) \rangle \quad \leftarrow \text{decoding}$$

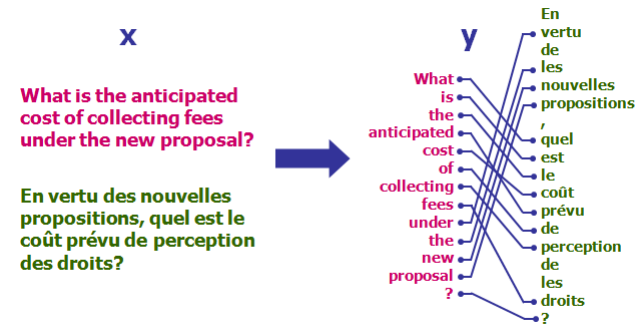
- structural SVM primal: \rightarrow **loss-augmented decoding**

$$\min_w \quad \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}} \left\{ L(y_i, y) + \langle w, \phi(x_i, y) \rangle \right\} - \langle w, \phi(x_i, y_i) \rangle$$

vs. binary hinge loss: $\max \left\{ 0, 1 - \langle w, \phi(x_i) y_i \rangle \right\}$

Structural SVM optimization

- structured prediction:
- learn linear classifier:



$$h_w(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle w, \phi(x, y) \rangle \quad \leftarrow \text{decoding}$$

- structural SVM primal: \rightarrow **loss-augmented decoding**

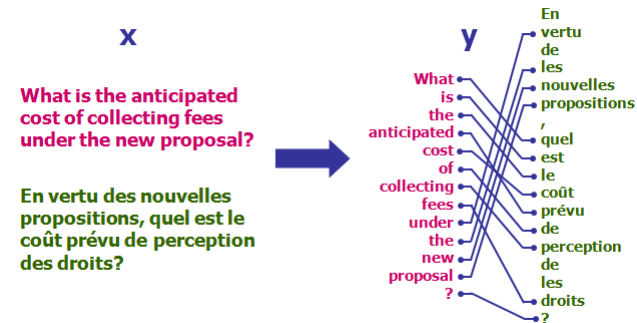
$$\min_w \quad \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}} \left\{ L(y_i, y) + \langle w, \phi(x_i, y) \rangle \right\} - \langle w, \phi(x_i, y_i) \rangle$$

vs. binary hinge loss: $\max \left\{ 0, 1 - \langle w, \phi(x_i) y_i \rangle \right\}$

- structural SVM dual: $\max_{\alpha \in \mathcal{M}} \quad b^T \alpha - \frac{\lambda}{2} \|A\alpha\|^2$ **primal-dual pair: $w = A\alpha$**
- $$\mathcal{M} := \Delta_{|y_1|} \times \dots \times \Delta_{|y_n|}$$

Structural SVM optimization

- structured prediction:
- learn linear classifier:



$$h_w(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \langle w, \phi(x, y) \rangle \quad \leftarrow \text{decoding}$$

- structural SVM primal: \rightarrow loss-augmented decoding

$$\min_w \quad \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}} \left\{ L(y_i, y) + \langle w, \phi(x_i, y) \rangle \right\} - \langle w, \phi(x_i, y_i) \rangle$$

vs. binary hinge loss: $\max \left\{ 0, 1 - \langle w, \phi(x_i) y_i \rangle \right\}$

- structural SVM dual:

\rightarrow exp. number of variables!

$$\max_{\alpha \in \mathcal{M}} \quad b^T \alpha - \frac{\lambda}{2} \|A\alpha\|^2$$

$$\mathcal{M} := \Delta_{|y_1|} \times \dots \times \Delta_{|y_n|}$$

primal-dual
pair: $w = A\alpha$

Structural SVM optimization (2)

$$\min_w \quad \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}} \left\{ L(y_i, y) + \langle w, \phi(x_i, y) \rangle \right\} - \langle w, \phi(x_i, y_i) \rangle$$

rate: after T passes
through data:

$$\tilde{O} \left(\frac{1}{nT} \right)$$

$$O \left(\frac{1}{T} \right)$$

$$O \left(\frac{1}{nT} \right)$$

Structural SVM optimization (2)

$$\min_w \quad \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}} \left\{ L(y_i, y) + \langle w, \phi(x_i, y) \rangle \right\} - \langle w, \phi(x_i, y_i) \rangle$$

- popular approaches:

- stochastic subgradient descent

rate: after T passes
through data:

$$\tilde{O}\left(\frac{1}{nT}\right)$$

- cutting plane method (SVMstruct)

$$O\left(\frac{1}{T}\right)$$

$$O\left(\frac{1}{nT}\right)$$

Structural SVM optimization (2)

$$\min_w \quad \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}} \left\{ L(y_i, y) + \langle w, \phi(x_i, y) \rangle \right\} - \langle w, \phi(x_i, y_i) \rangle$$

- popular approaches:

- stochastic subgradient descent

- pros: online!

- cons: sensitive to step-size; don't know when to stop

- cutting plane method (SVMstruct)

rate: after T passes
through data:

$$\tilde{O}\left(\frac{1}{nT}\right)$$

$$O\left(\frac{1}{T}\right)$$

$$O\left(\frac{1}{nT}\right)$$

Structural SVM optimization (2)

$$\min_w \quad \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}} \left\{ L(y_i, y) + \langle w, \phi(x_i, y) \rangle \right\} - \langle w, \phi(x_i, y_i) \rangle$$

- popular approaches:

rate: after T passes
through data:

$$\tilde{O}\left(\frac{1}{nT}\right)$$

- stochastic subgradient descent

- pros: online!
 - cons: sensitive to step-size; don't know when to stop

- cutting plane method (SVMstruct)

- pros: automatic step-size; duality gap
 - cons: batch! -> slow for large n

$$O\left(\frac{1}{T}\right)$$

$$O\left(\frac{1}{nT}\right)$$

Structural SVM optimization (2)

$$\min_w \quad \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}} \left\{ L(y_i, y) + \langle w, \phi(x_i, y) \rangle \right\} - \langle w, \phi(x_i, y_i) \rangle$$

- popular approaches:

rate: after T passes
through data:

$$\tilde{O}\left(\frac{1}{nT}\right)$$

- stochastic subgradient descent

- pros: online!

- cons: sensitive to step-size; don't know when to stop

- cutting plane method (SVMstruct)

- pros: automatic step-size; duality gap

- cons: batch! -> slow for large n

$$O\left(\frac{1}{T}\right)$$

- **our approach:** block-coordinate Frank-Wolfe on dual

-> combines best of both worlds:

$$O\left(\frac{1}{nT}\right)$$

Structural SVM optimization (2)

$$\min_w \quad \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \max_{y \in \mathcal{Y}} \left\{ L(y_i, y) + \langle w, \phi(x_i, y) \rangle \right\} - \langle w, \phi(x_i, y_i) \rangle$$

■ popular approaches:

rate: after T passes
through data:

$$\tilde{O}\left(\frac{1}{nT}\right)$$

■ stochastic subgradient descent

- pros: online!

- cons: sensitive to step-size; don't know when to stop

■ cutting plane method (SVMstruct)

- pros: automatic step-size; duality gap

- cons: batch! -> slow for large n

$$O\left(\frac{1}{T}\right)$$

■ **our approach:** block-coordinate Frank-Wolfe on dual

-> combines best of both worlds:

- **online!**

- automatic step-size via analytic line search

- duality gap

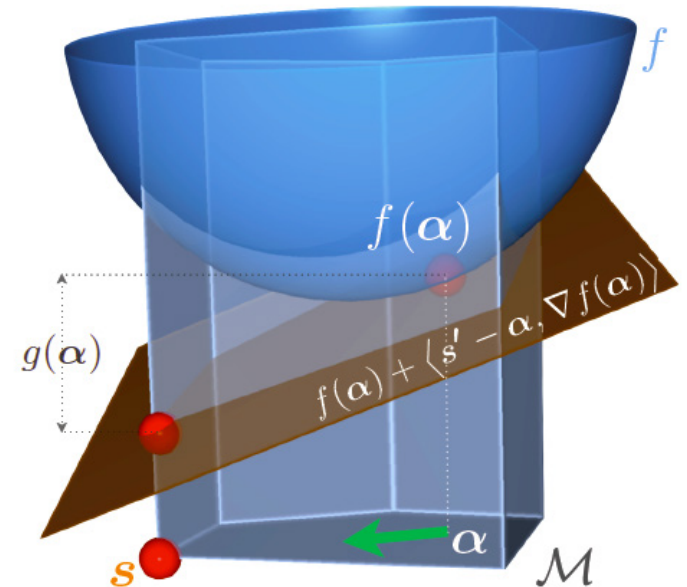
- rates also hold for approximate oracles

$$O\left(\frac{1}{nT}\right)$$

Frank-Wolfe algorithm [Frank, Wolfe 1956]

(aka conditional gradient)

- alg. for constrained opt.: $\min_{\alpha \in \mathcal{M}} f(\alpha)$
where:
 - f convex & cts. differentiable
 - \mathcal{M} convex & compact
- FW algorithm – repeat:



Frank-Wolfe algorithm [Frank, Wolfe 1956]

(aka conditional gradient)

- alg. for constrained opt.: $\min_{\alpha \in \mathcal{M}} f(\alpha)$

where:

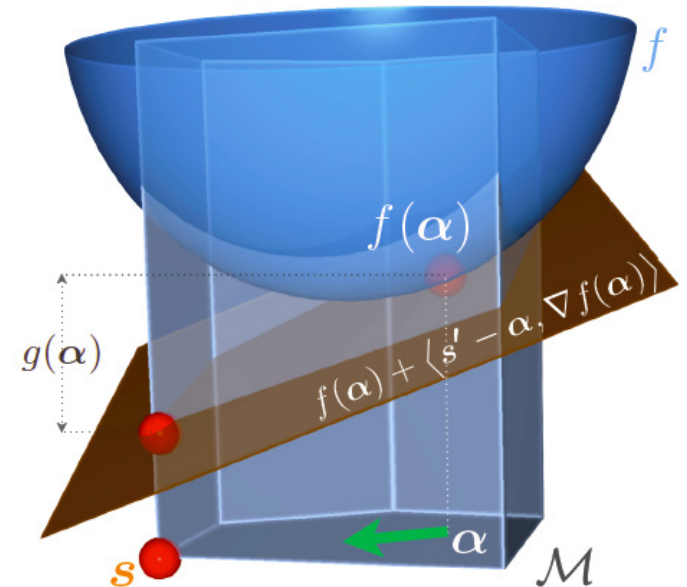
f convex & cts. differentiable

\mathcal{M} convex & compact

- FW algorithm – repeat:

1) Find good feasible direction by minimizing linearization of f :

$$s_{t+1} \in \arg \min_{s' \in \mathcal{M}} \langle s', \nabla f(\alpha_t) \rangle$$



Frank-Wolfe algorithm [Frank, Wolfe 1956]

(aka conditional gradient)

- alg. for constrained opt.: $\min_{\alpha \in \mathcal{M}} f(\alpha)$

where:

f convex & cts. differentiable

\mathcal{M} convex & compact

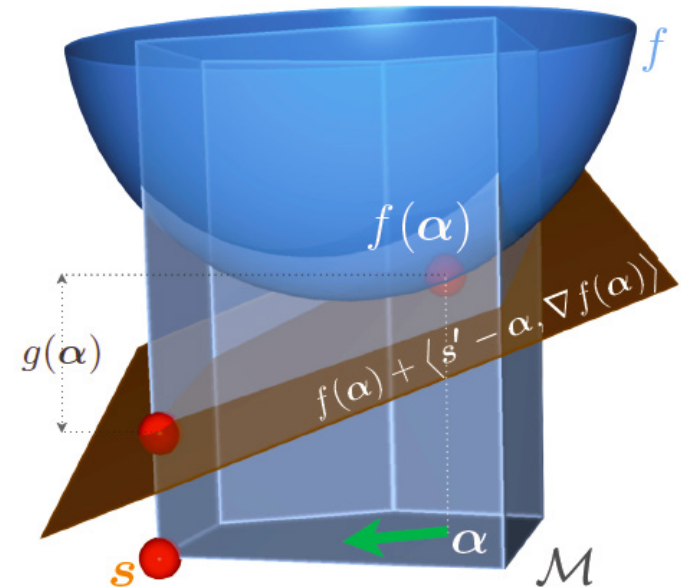
- FW algorithm – repeat:

1) Find good feasible direction by minimizing linearization of f :

$$s_{t+1} \in \arg \min_{s' \in \mathcal{M}} \langle s', \nabla f(\alpha_t) \rangle$$

2) Take convex step in direction:

$$\alpha_{t+1} = (1 - \gamma_t) \alpha_t + \gamma_t s_{t+1}$$



Frank-Wolfe algorithm [Frank, Wolfe 1956]

(aka conditional gradient)

- alg. for constrained opt.: $\min_{\alpha \in \mathcal{M}} f(\alpha)$

where:

f convex & cts. differentiable

\mathcal{M} convex & compact

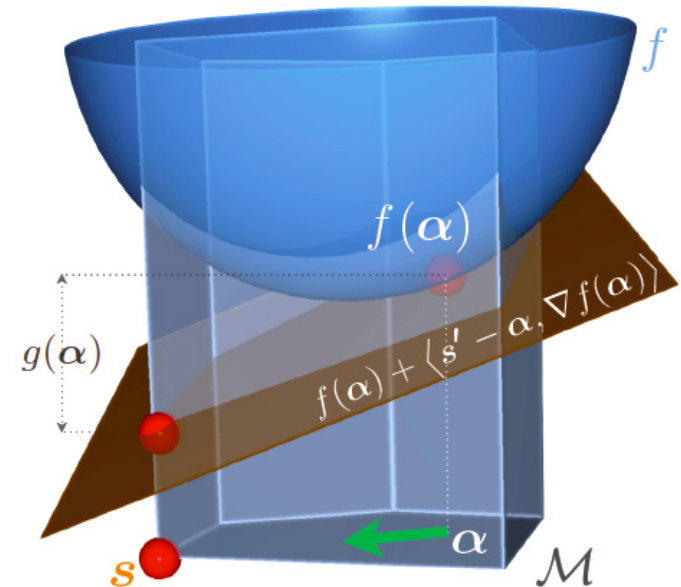
- FW algorithm – repeat:

1) Find good feasible direction by minimizing linearization of f :

$$s_{t+1} \in \arg \min_{s' \in \mathcal{M}} \langle s', \nabla f(\alpha_t) \rangle$$

2) Take convex step in direction:

$$\alpha_{t+1} = (1 - \gamma_t) \alpha_t + \gamma_t s_{t+1}$$



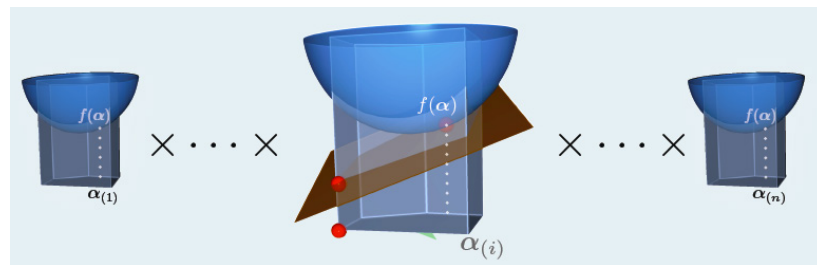
- Properties: $O(1/T)$ rate
 - sparse iterates
 - get duality gap $g(\alpha)$ for free
 - rate holds even if linear subproblem solved **approximately**

Block-Coordinate Frank-Wolfe (new!)

- for constrained optimization over compact **product domain**:

$$\min_{\alpha \in \mathcal{M}^{(1)} \times \dots \times \mathcal{M}^{(n)}} f(\alpha)$$

$$\alpha = (\alpha_{(1)}, \dots, \alpha_{(n)})$$

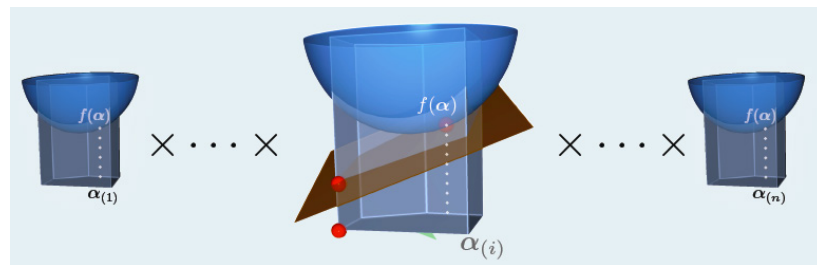


- Properties: $O(1/T)$ rate
 - sparse iterates
 - duality gap guarantees
 - rate holds even if linear subproblem solved **approximately**

Block-Coordinate Frank-Wolfe (new!)

- for constrained optimization over compact **product domain**:

$$\min_{\alpha \in \mathcal{M}^{(1)} \times \dots \times \mathcal{M}^{(n)}} f(\alpha)$$
$$\alpha = (\alpha_{(1)}, \dots, \alpha_{(n)})$$



- pick i at random; update only block i with a FW step:

$$s_{(i)} = \operatorname{argmin}_{s'_{(i)} \in \mathcal{M}^{(i)}} \langle s'_{(i)}, \nabla_{(i)} f(\alpha^{(k)}) \rangle$$

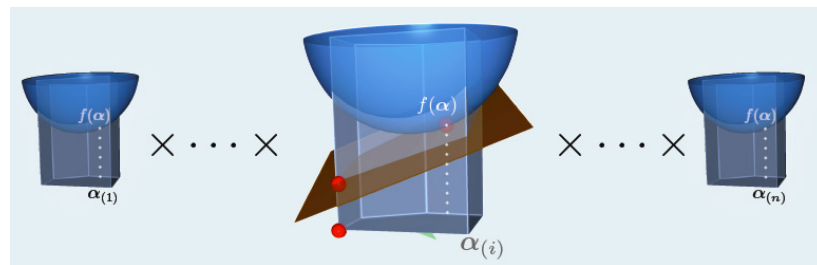
$$\alpha_{(i)}^{(k+1)} = (1 - \gamma) \alpha_{(i)}^{(k)} + \gamma s_{(i)}$$

- Properties: $O(1/T)$ rate
 - sparse iterates
 - duality gap guarantees
 - rate holds even if linear subproblem solved **approximately**

Block-Coordinate Frank-Wolfe (new!)

- for constrained optimization over compact **product domain**:

$$\min_{\alpha \in \mathcal{M}^{(1)} \times \dots \times \mathcal{M}^{(n)}} f(\alpha)$$
$$\alpha = (\alpha_{(1)}, \dots, \alpha_{(n)})$$



- pick i at random; update only block i with a FW step:

$$s_{(i)} = \operatorname{argmin}_{s'_{(i)} \in \mathcal{M}^{(i)}} \langle s'_{(i)}, \nabla_{(i)} f(\alpha^{(k)}) \rangle$$

$$\alpha_{(i)}^{(k+1)} = (1 - \gamma) \alpha_{(i)}^{(k)} + \gamma s_{(i)}$$

- we proved **same** $O(1/T)$ rate

as batch FW

-> each step **n times cheaper** though

-> constant can be the same (SVM e.g.)

- Properties: $O(1/T)$ rate

- sparse iterates
- duality gap guarantees
- rate holds even if linear subproblem solved **approximately**

Block-Coordinate Frank-Wolfe (new!)

- for constrained optimization over compact **product domain**:

$$\min_{\alpha \in \mathcal{M}^{(1)} \times \dots \times \mathcal{M}^{(n)}} f(\alpha)$$

$$\alpha = (\alpha_{(1)}, \dots, \alpha_{(n)})$$

- pick i at random; update only block i with a FW step:

$$s_{(i)} = \operatorname{argmin}_{s'_{(i)} \in \mathcal{M}^{(i)}} \langle s'_{(i)}, \nabla_{(i)} f(\alpha^{(k)}) \rangle$$

$$\alpha_{(i)}^{(k+1)} = (1 - \gamma) \alpha_{(i)}^{(k)} + \gamma s_{(i)}$$

- we proved **same** $O(1/T)$ rate

as batch FW

-> each step **n times cheaper** though

-> constant can be the same (SVM e.g.)

- Properties: $O(1/T)$ rate

- sparse iterates
- duality gap guarantees
- rate holds even if linear subproblem solved **approximately**

Block-Coordinate Frank-Wolfe (new!)

- for constrained optimization over compact **product domain**:

$$\min_{\alpha \in \mathcal{M}^{(1)} \times \dots \times \mathcal{M}^{(n)}} f(\alpha)$$
$$\alpha = (\alpha_{(1)}, \dots, \alpha_{(n)})$$

structural SVM:

$$\max_{\alpha \in \mathcal{M}} b^T \alpha - \frac{\lambda}{2} \|A\alpha\|^2$$

$$\mathcal{M} := \Delta_{|y_1|} \times \dots \times \Delta_{|y_n|}$$

- pick i at random; update only block i with a FW step:

$$s_{(i)} = \operatorname{argmin}_{s'_{(i)} \in \mathcal{M}^{(i)}} \langle s'_{(i)}, \nabla_{(i)} f(\alpha^{(k)}) \rangle$$

$$\alpha_{(i)}^{(k+1)} = (1 - \gamma) \alpha_{(i)}^{(k)} + \gamma s_{(i)}$$

- we proved **same** $O(1/T)$ rate

as batch FW

-> each step **n times cheaper** though

-> constant can be the same (SVM e.g.)

- Properties: $O(1/T)$ rate

- sparse iterates
- duality gap guarantees
- rate holds even if linear subproblem solved **approximately**

Block-Coordinate Frank-Wolfe (new!)

- for constrained optimization over compact **product domain**:

$$\min_{\alpha \in \mathcal{M}^{(1)} \times \dots \times \mathcal{M}^{(n)}} f(\alpha)$$
$$\alpha = (\alpha_{(1)}, \dots, \alpha_{(n)})$$

structural SVM:

$$\max_{\alpha \in \mathcal{M}} b^T \alpha - \frac{\lambda}{2} \|A\alpha\|^2$$

$$\mathcal{M} := \Delta_{|y_1|} \times \dots \times \Delta_{|y_n|}$$

- pick i at random; update only block i with a FW step:

$$s_{(i)} = \operatorname{argmin}_{s'_{(i)} \in \mathcal{M}^{(i)}} \langle s'_{(i)}, \nabla_{(i)} f(\alpha^{(k)}) \rangle$$

key insight:

$$\alpha_{(i)}^{(k+1)} = (1 - \gamma) \alpha_{(i)}^{(k)} + \gamma s_{(i)}$$

- we proved **same** $O(1/T)$ rate

as batch FW

- > each step **n times cheaper** though
- > constant can be the same (SVM e.g.)

- Properties: $O(1/T)$ rate

- sparse iterates
- duality gap guarantees
- rate holds even if linear subproblem solved **approximately**

Block-Coordinate Frank-Wolfe (new!)

- for constrained optimization over compact **product domain**:

$$\min_{\alpha \in \mathcal{M}^{(1)} \times \dots \times \mathcal{M}^{(n)}} f(\alpha)$$
$$\alpha = (\alpha_{(1)}, \dots, \alpha_{(n)})$$

structural SVM:

$$\max_{\alpha \in \mathcal{M}} b^T \alpha - \frac{\lambda}{2} \|A\alpha\|^2$$
$$\mathcal{M} := \Delta_{|Y_1|} \times \dots \times \Delta_{|Y_n|}$$

- pick i at random; update only block i with a FW step:

$$s_{(i)} = \operatorname{argmin}_{s'_{(i)} \in \mathcal{M}^{(i)}} \langle s'_{(i)}, \nabla_{(i)} f(\alpha^{(k)}) \rangle \Leftrightarrow$$

key insight:

$$\max_{y \in \mathcal{Y}} \left\{ L(y_i, y) + \langle w, \phi(x_i, y) \rangle \right\}$$

loss-augmented decoding

$$\alpha_{(i)}^{(k+1)} = (1 - \gamma) \alpha_{(i)}^{(k)} + \gamma s_{(i)}$$

- we proved **same** $O(1/T)$ rate

as batch FW

- > each step **n times cheaper** though
- > constant can be the same (SVM e.g.)

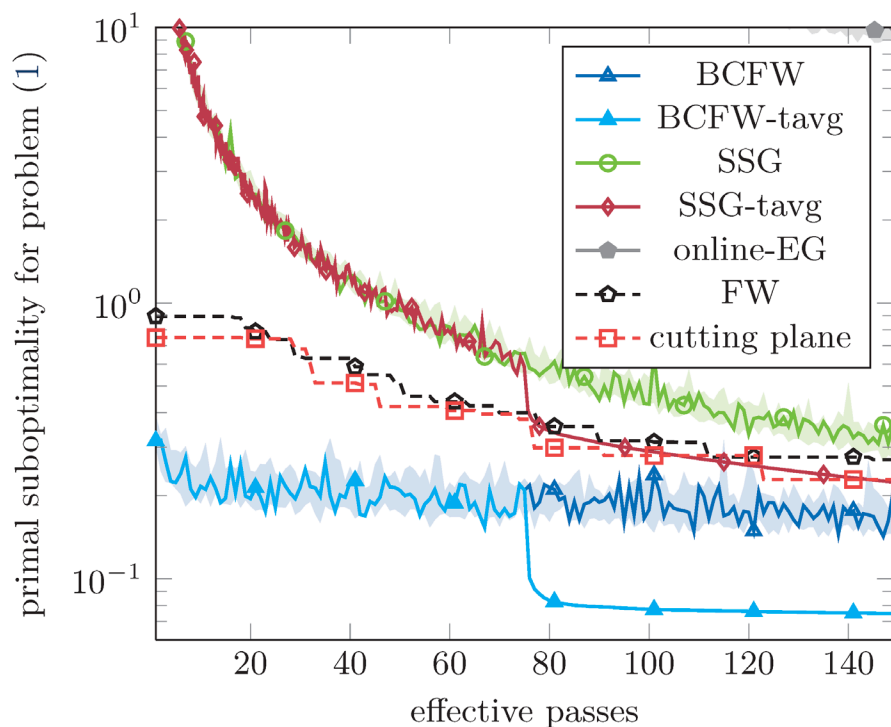
- Properties: $O(1/T)$ rate

- sparse iterates
- duality gap guarantees
- rate holds even if linear subproblem solved **approximately**

Experiments

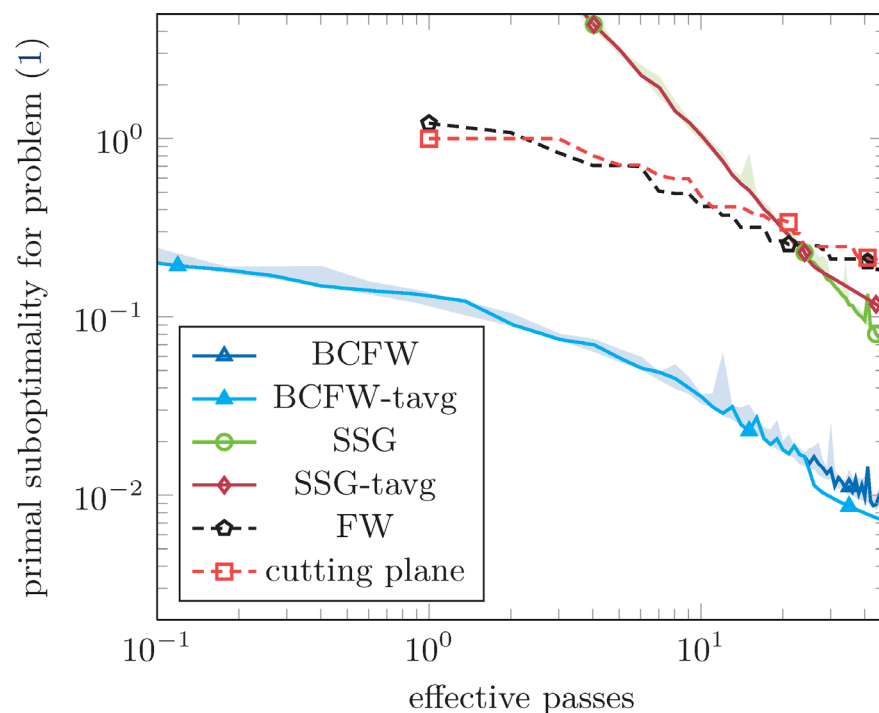
$$\lambda = 1/n$$

OCR dataset



$$n = 6k, d = 4k$$

CoNLL dataset



$$n = 9k, d = 1.6M$$

Conclusion

- new block-coordinate variant of Frank-Wolfe algorithm
 - same convergence rate but with cheaper iteration cost
- applied to structural SVM, yields:
 - online algorithm
 - optimal step-size computed in close form
 - duality gap
 - rates hold with approximate oracles