

Block-Coordinate Frank-Wolfe for Structural SVMs

Simon Lacoste-Julien^b

Martin Jaggi^a

Mark Schmidt^b

Patrick Pletscher^c



informatics

mathematics

ETH

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich



^a CMAP - ERC SIPA project,
École Polytechnique,
CNRS UMR 7641, Paris, France

^b INRIA - SIERRA project-team,
École Normale Supérieure,
CNRS UMR 8548, Paris, France

^c Machine Learning Laboratory,
ETH Zurich, Switzerland

Short Summary

Motivation

Despite their wider applicability, optimization of structural SVMs remains challenging.

Contributions

New **block-coordinate** variant of the classic **Frank-Wolfe algorithm**
(for convex optim. with block-separable constraints)

Giving a new simple **online** algorithm for structural SVMs, with primal-dual convergence rate, outperforming existing solvers in practice

Advantages

- The **optimal step-size** can be computed in closed-form (no parameter tuning)
- Duality gap** guarantee, (e.g. as a stopping criterion)
- Allows use of **approximate maximization oracles** (weakest / most general oracle)

Frank-Wolfe (or conditional gradient)

Constrained Convex Optimization

over a compact domain

$$\min_{\alpha \in \mathcal{M}} f(\alpha)$$

Algorithm 1 Frank-Wolfe

```
Let  $\alpha^{(0)} \in \mathcal{M}$ 
for  $k = 0 \dots K$  do
    Compute  $s := \operatorname{argmin}_{s' \in \mathcal{M}} \langle s', \nabla f(\alpha^{(k)}) \rangle$ 
    Let  $\gamma := \frac{2}{k+2}$  or find the optimal  $\gamma$ 
    Update  $\alpha^{(k+1)} := (1 - \gamma)\alpha^{(k)} + \gamma s$ 
end for
```

Idea: Minimize a linear approximation

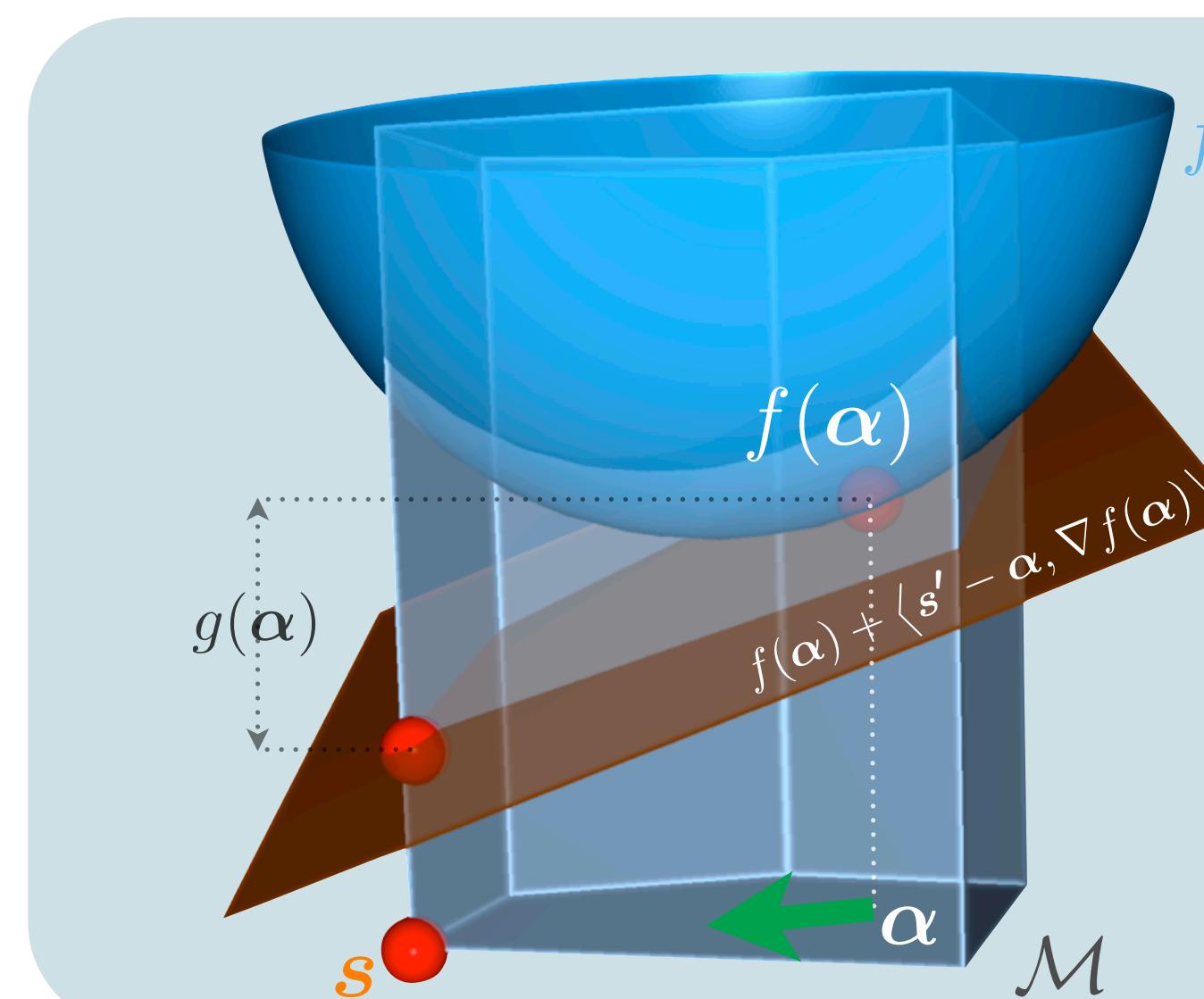
Convergence:
Error $\leq \frac{2C_f}{k+2}$
after k steps.

(also in **duality gap**,
and with **inexact subproblems**)

Duality Gap

$g(\alpha)$ = efficient certificate
for approximation quality

Sparse Iterates!



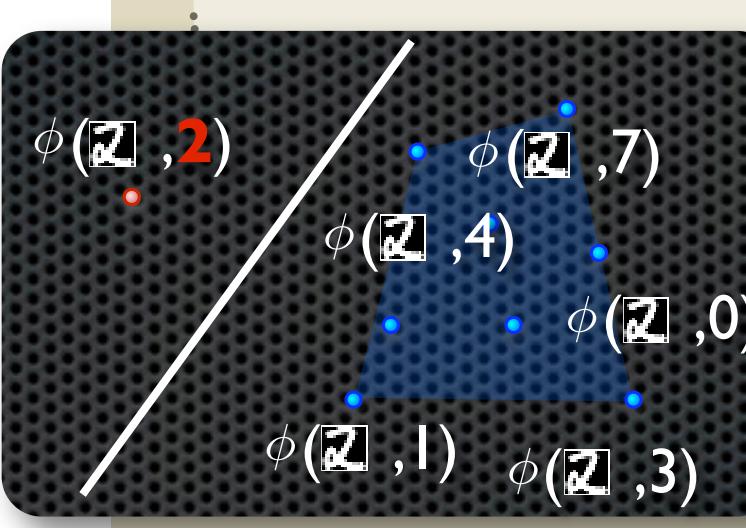
Constant bounded by the Lipschitz constant L_f of the gradient, $C_f \leq L_f \operatorname{diam}(\mathcal{M})^2$

Structural SVM

Structured Prediction

Goal: Given a joint “structured” feature map $\phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$, construct a good linear classifier of the form

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



Large margin separation

Primal

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

= structured hinge loss $=: \psi_i(y)$

primal-dual correspondence
 $w = A\alpha$

Dual

$$\min_{\alpha \in \mathbb{R}^{n \times |\mathcal{Y}|}} f(\alpha) := \frac{\lambda}{2} \|A\alpha\|^2 - b^T \alpha$$

s.t. $\sum_{y \in \mathcal{Y}} \alpha_i(y) = 1 \quad \forall i \in [n]$
and $\alpha_i(y) \geq 0 \quad \forall i \in [n], \forall y \in \mathcal{Y}$

Challenge: exponential # of variables

$$A := \left\{ \frac{1}{\lambda n} \psi_i(y) \in \mathbb{R}^d \mid i \in [n], y \in \mathcal{Y} \right\}$$

$$b := \left(\frac{1}{n} L_i(y) \right)_{i \in [n], y \in \mathcal{Y}}$$

Optimization of the Structural SVM Dual

Key Insight: Frank-Wolfe step = Maximization oracle

Batch Frank-Wolfe:

Duality gap $\leq \varepsilon$ after $O\left(\frac{R^2}{\lambda\varepsilon}\right)$ iterations
(iteration cost: **n** oracle calls)

Relation with Batch Subgradient

Can interpret batch subgradient (in the primal) as classic Frank-Wolfe (in the dual)

Relation with Cutting Plane

Can interpret cutting plane (SVM^{struct}, bundle methods) as a Frank-Wolfe variant, giving a simpler convergence proof

Related Work

Table 1. Convergence rates given in the number of calls to the oracles for different optimization algorithms for the structural SVM objective (1) in the case of a Markov random field structure, to reach a specific accuracy ε measured for different types of gaps, in term of the number of training examples n , regularization parameter λ , size of the label space $|\mathcal{Y}|$, maximum feature norm $R := \max_{i,y} \|\psi_i(y)\|_2$ (some minor terms were ignored for succinctness). Table inspired from (Zhang et al., 2011). Notice that only stochastic subgradient and our proposed algorithm have rates independent of n .

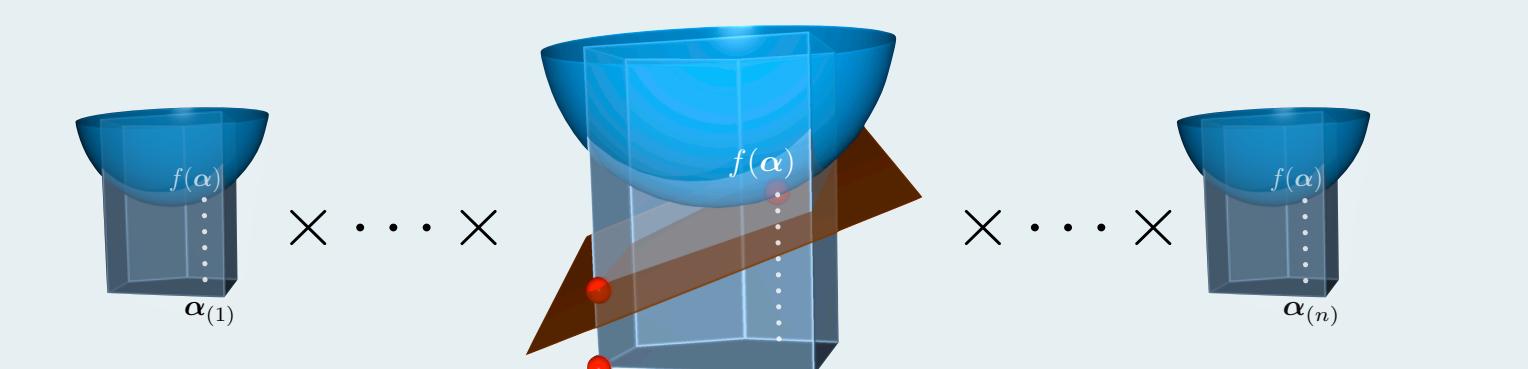
Optimization algorithm	Online	Primal/Dual	Type of guarantee	Oracle type	# Oracle calls	
dual extragradient (Taskar et al., 2006)		no	primal “dual”	saddle point gap	Bregman projection	$O\left(\frac{nR \log \mathcal{Y} }{\lambda\varepsilon}\right)$
online exponentiated gradient (Collins et al., 2008)	yes	dual	expected dual error	expectation	$O\left(\frac{(n+\log \mathcal{Y})R^2}{\lambda\varepsilon}\right)$	
excessive gap reduction (Zhang et al., 2011)	no	primal-dual	duality gap	expectation	$O\left(nR\sqrt{\frac{\log \mathcal{Y} }{\lambda\varepsilon}}\right)$	
BMRM (Teo et al., 2010)	no	primal	\geq primal error	maximization	$O\left(\frac{nR^2}{\lambda\varepsilon}\right)$	
1-slack SVM-Struct (Joachims et al., 2009)	no	primal-dual	duality gap	maximization	$O\left(\frac{nR^2}{\lambda\varepsilon}\right)$	
stochastic subgradient (Shalev-Shwartz et al., 2010)	yes	primal	primal error w.h.p.	maximization	$\tilde{O}\left(\frac{R^2}{\lambda\varepsilon}\right)$	
this paper: stochastic block-coordinate Frank-Wolfe	yes	primal-dual	expected duality gap	maximization	$O\left(\frac{R^2}{\lambda\varepsilon}\right)$ Thm. 3	

Block-Coordinate Frank-Wolfe

Problem: Minimize a convex function over block-separable compact constraints

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

Idea: Combine Coordinate Descent with cheaper Frank-Wolfe steps



(pick one single block at random, and perform a Frank-Wolfe step affecting only this block)

Algorithm 3 Block-Coordinate Frank-Wolfe

```
Let  $\alpha^{(0)} \in \mathcal{M} = \mathcal{M}^{(1)} \times \dots \times \mathcal{M}^{(n)}$ 
for  $k = 0 \dots K$  do
    Pick  $i \in_{u.a.r.} [n]$ 
    Find  $s_{(i)} := \operatorname{argmin}_{s'_{(i)} \in \mathcal{M}^{(i)}} \langle s'_{(i)}, \nabla_i f(\alpha^{(k)}) \rangle$ 
    Let  $\gamma := \frac{2n}{k+2n}$ , or find the optimal  $\gamma$ 
    Update  $\alpha_{(i)}^{(k+1)} := \alpha_{(i)}^{(k)} + \gamma(s_{(i)} - \alpha_{(i)}^{(k)})$ 
end for
```

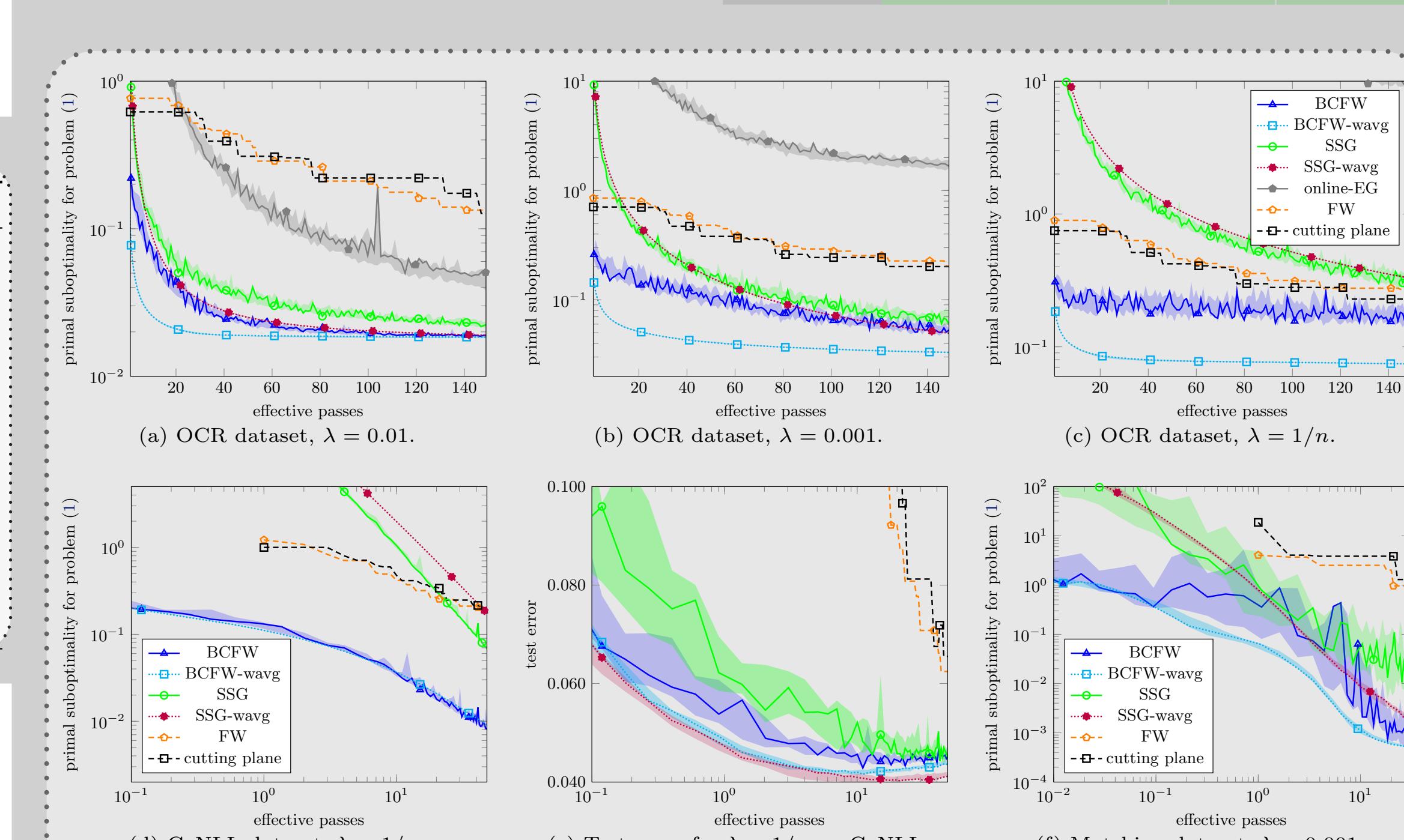
The constant C_f^{prod} can be much smaller than C_f . (For structural SVM, $nC_f^{prod} \approx C_f$)

Convergence:
Error $\leq \frac{2nC_f^{prod}}{k+2n}$
after k steps.

(also in **duality gap**,
and with **inexact subproblems**)

Experimental Results

dataset	sequence labeling	n	d
OCR		6251	4028
CoNLL	POS sequence labeling	8936	1643026
Matching	word alignment	5000	82



comparing block-coordinate Frank-Wolfe (BCFW) to stochastic subgradient (SSG), online exponentiated gradient (EG), batch Frank-Wolfe (FW) and cutting plane (wavg = weighted averaging of the iterates)