

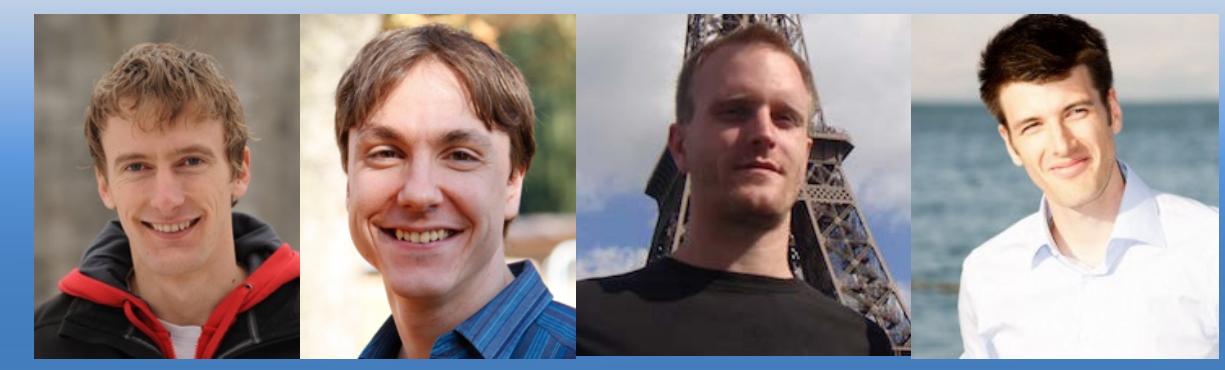
# Block-Coordinate Frank-Wolfe for Structural SVMs

Martin Jaggi<sup>a</sup>

Simon Lacoste-Julien<sup>b</sup>

Mark Schmidt<sup>b</sup>

Patrick Pletscher<sup>c</sup>



## 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:

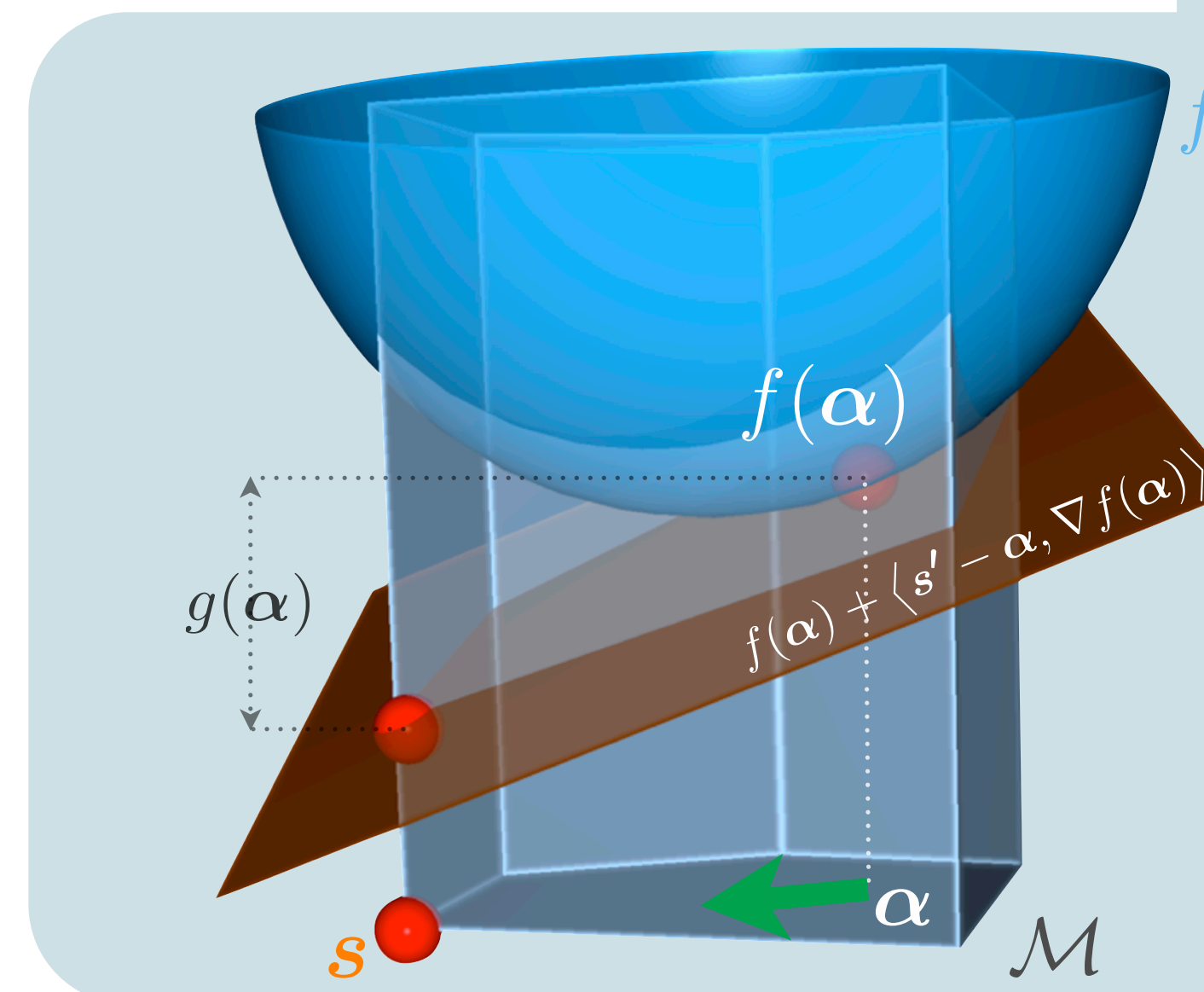
$$\text{Error} \leq \frac{2C_f}{k+2} \text{ after } k \text{ steps.}$$

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

### Duality Gap

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

**Sparse Iterates!**



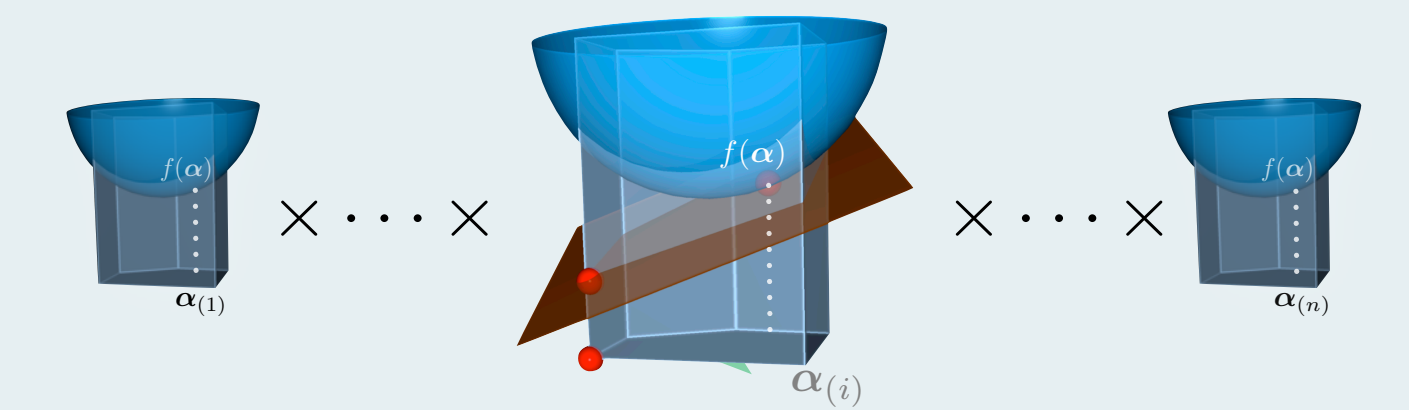
Constant bounded by the Lipschitz constant  $L_f$  of the gradient,  $C_f \leq L_f \operatorname{diam}(\mathcal{M})^2$

## 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
    
```

**new**

### Convergence:

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

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

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

## 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**

maximization oracle

(loss augmented decoding)

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

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

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

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

$$\text{and} \quad \alpha_i(y) \geq 0 \quad \forall i \in [n], \forall y \in \mathcal{Y}$$

Challenge: exponential # of variables

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

primal-dual correspondence

block-structure!

## Optimization of the Structural SVM Dual

### Batch Frank-Wolfe:

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

### Block-Coordinate Frank-Wolfe:

Duality gap  $\leq \varepsilon$  after  $O\left(\frac{R^2}{\lambda \varepsilon}\right)$  iter. (iteration cost: **one** oracle call)

### 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<sup>struct</sup>, bundle methods) as a Frank-Wolfe variant, giving a simpler convergence proof

### Relation with Stochastic Subgradient (SGD)

Same cheap iteration cost, but we have stronger primal-dual guarantees, more robustness, no step-size tuning, and faster in experiments

## 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  $\varepsilon$  measured for different types of gaps, in term of the number of training examples  $n$ , regularization parameter  $\lambda$ , size of the label space  $|\mathcal{Y}|$ , maximum feature norm  $R := \max_{y \in \mathcal{Y}} \|\psi_i(y)\|_2$  (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

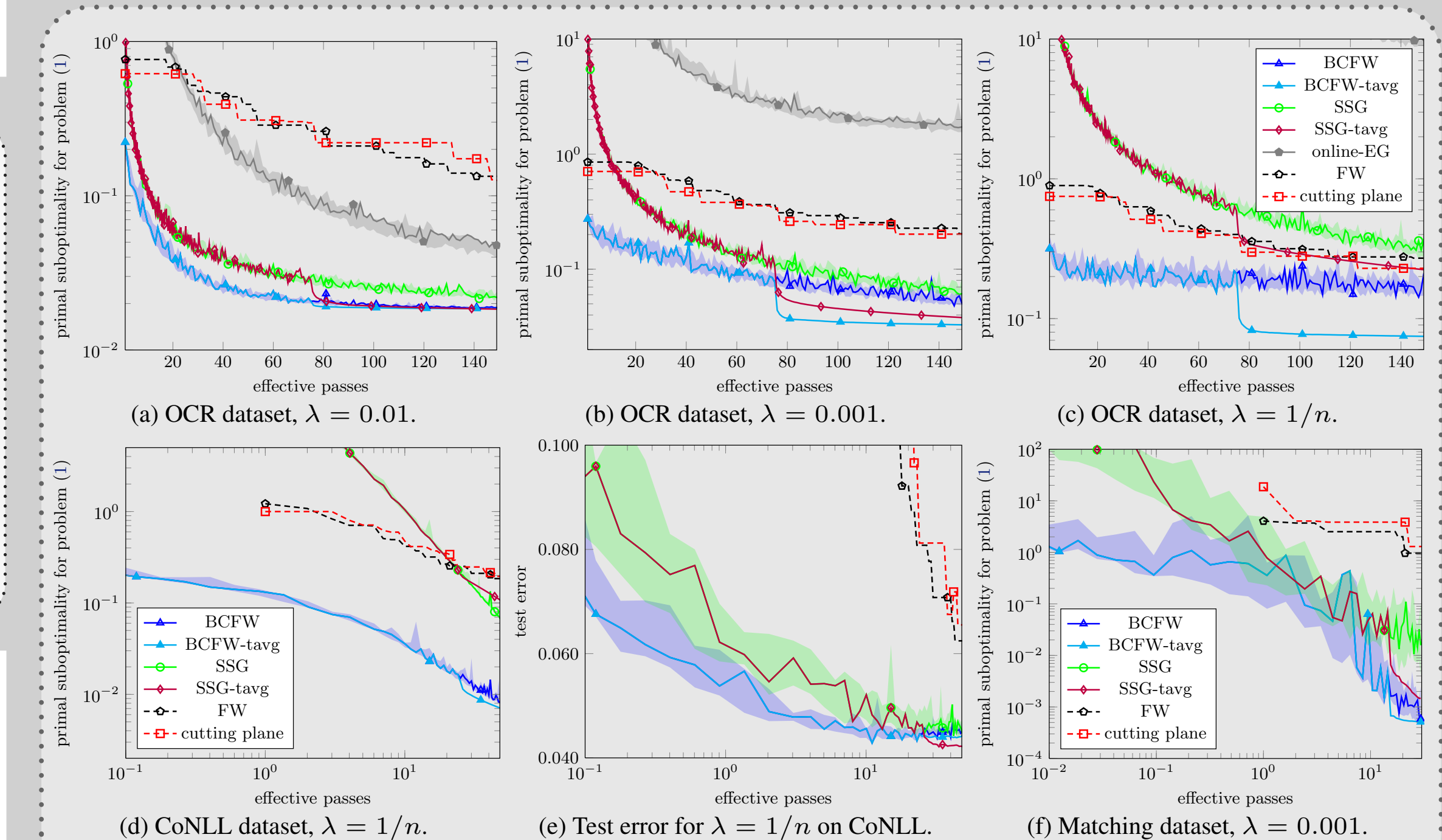
#### Algorithm 4 BCFW for Structural SVM

```

Let  $w^{(0)} := w_i^{(0)} := 0$ 
for  $k = 0 \dots K$  do
  Pick  $i \in_{u.a.r.} [n]$ 
  Solve  $y_i^* := \operatorname{argmax}_{y \in \mathcal{Y}_i} H_i(y; w^{(k)})$ 
  Let  $w_s := \frac{1}{\lambda n} \psi_i(y_i^*)$ 
  Let  $\gamma := \frac{2n}{k+2n}$ , or find the optimal  $\gamma$ 
  Update  $w_i^{(k+1)} := (1 - \gamma)w_i^{(k)} + \gamma w_s$ 
  Update  $w^{(k+1)} := w^{(k)} + w_i^{(k+1)} - w_i^{(k)}$ 
end for
    
```

## Experimental Results

dataset		n	d
OCR	sequence labeling	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 (tavg = tail averaging for the second half of the iterations)