

Blended Conditional Gradients: The unconditioning of conditional gradients

Joint work with Gabor Braun, Sebastian Pokutta, Steve Wright

Dan Tu

CONDITIONAL GRADIENTS: PROJECTION-FREE

Given a polytope P , solve the optimization problem:

$$\min f(x) \quad \text{s.t. } x \in P$$

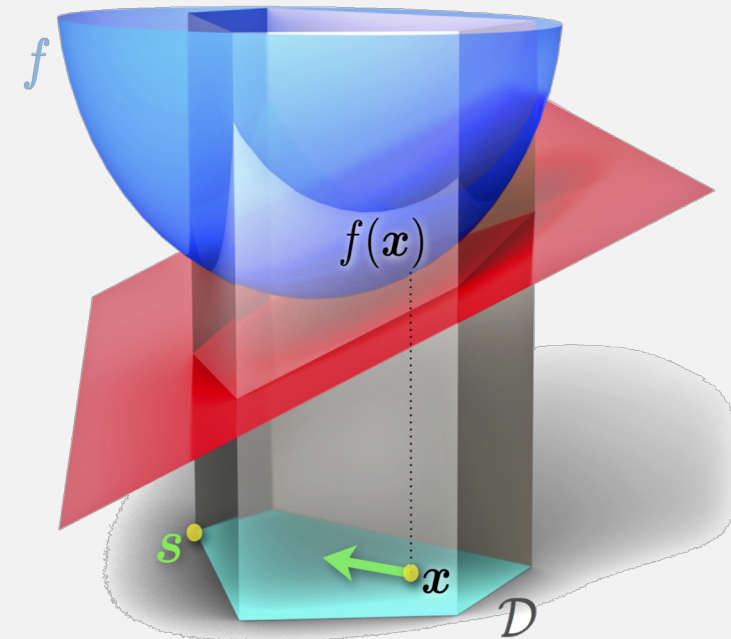
where the objective function f is smooth and strongly convex

Frank-Wolfe Algorithm

- 1: **Input:** initial point $x_0 \in P$, convex function f
- 2: **for** $t = 1$ to T **do**
- 3: $s_{t+1} = \operatorname{argmin}_{s \in P} \nabla f(x_t)^\top s$
- 4: $x_{t+1} = (1 - \gamma)x_t + \gamma s_{t+1}$
- 5: **end for**

Find a vertex through LP Oracle.

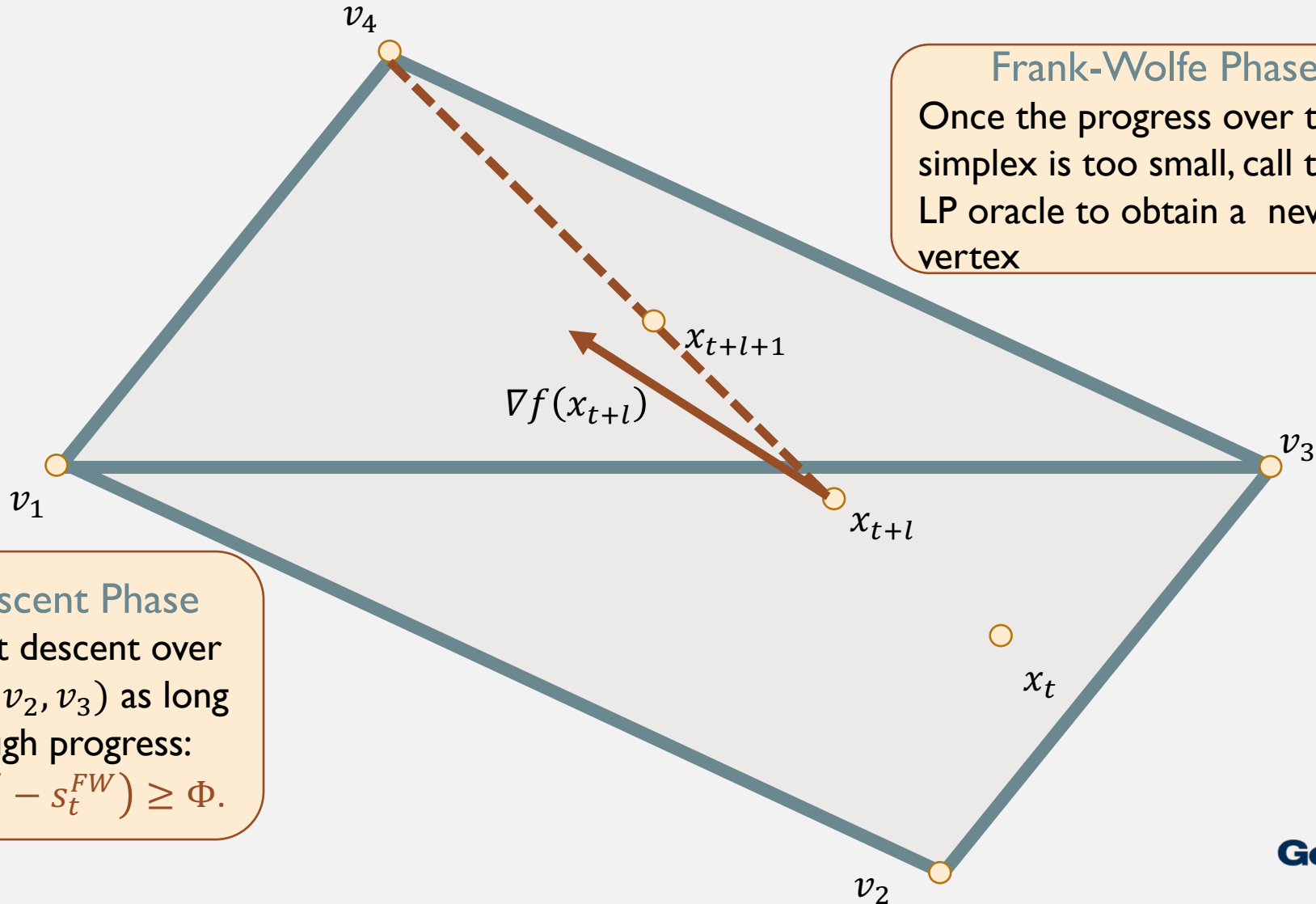
Walk along the conditional gradient direction.



Problems:

- LP Oracle can be computationally expensive.
- The conditional gradient direction, as an approximation of the negative gradient, can be inefficient.

BLENDED CONDITIONAL GRADIENT



Frank-Wolfe Phase

Once the progress over the simplex is too small, call the LP oracle to obtain a new vertex

Gradient Descent Phase

Perform gradient descent over the simplex (v_1, v_2, v_3) as long as it makes enough progress:

$$\nabla f(x_t)^T (v_t^{Away} - s_t^{FW}) \geq \Phi.$$

GRADIENT DESCENT PHASE

SIMPLEX GRADIENT DESCENT ORACLE

For a general simplex, decompose x as a convex combination

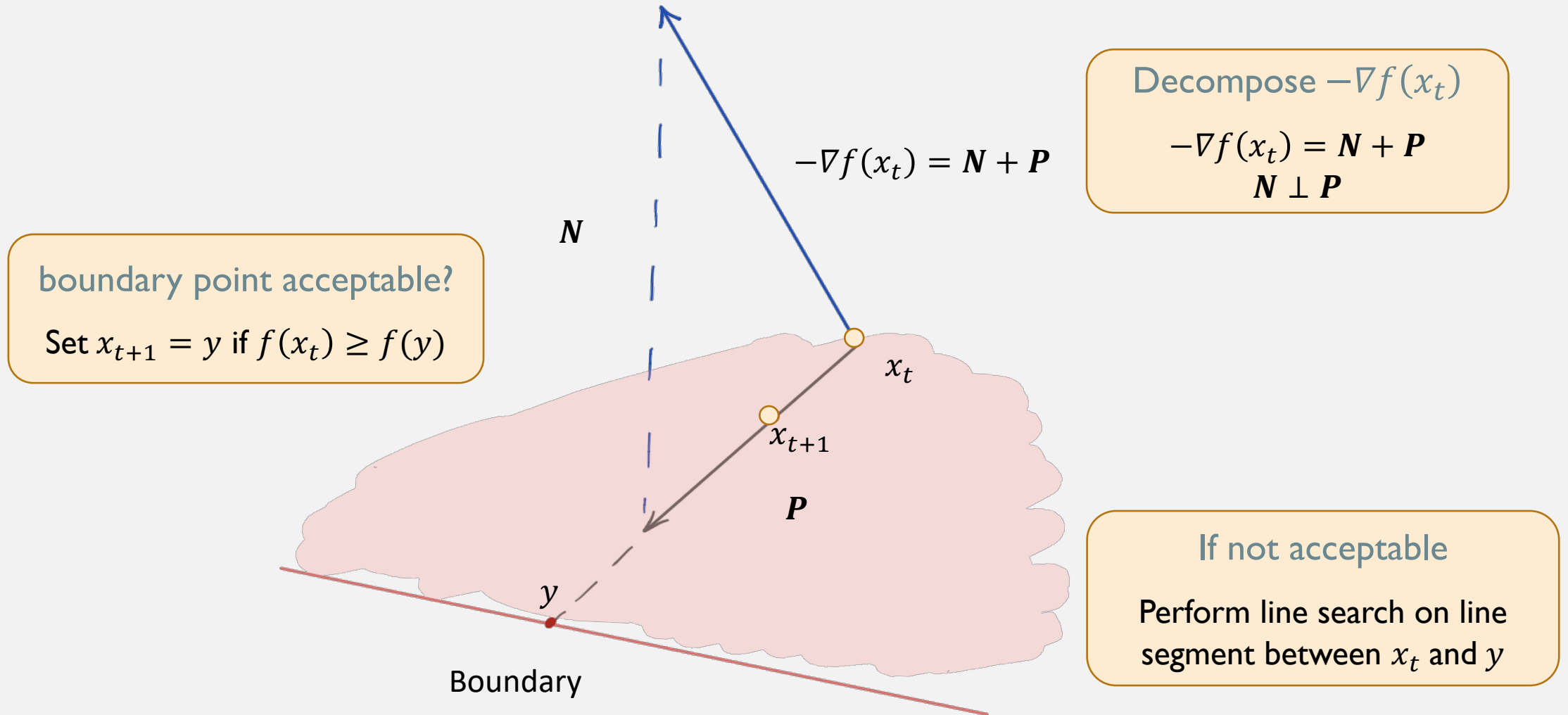
$$x = \sum_{i=1}^t \lambda_i v_i, \text{ with } \sum_{i=1}^t \lambda_i = 1 \text{ and } \lambda_i \geq 0, i = 1, 2, \dots, t$$

Treat λ_i as variables $\rightarrow x$ in a standard simplex with normal vector:

$$N = (1, 1, \dots, 1) / \sqrt{t}$$

GRADIENT DESCENT PHASE

SIMPLEX GRADIENT DESCENT ORACLE



BCG ALGORITHM

Algorithm Blended Conditional Gradients (BCG)

Require: smooth convex function f , start vertex $x_0 \in P$, weak separation oracle LPsep_P , accuracy $K \geq 1$

Ensure: points x_t in P for $t = 1, \dots, T$

```
1:  $\Phi_0 \leftarrow \max_{v \in P} \nabla f(x_0)(x_0 - v)/2$  {Initial dual gap estimate}
2:  $S_0 \leftarrow \{x_0\}$ 
3: for  $t = 0$  to  $T - 1$  do
4:    $v_t^A \leftarrow \operatorname{argmax}_{v \in S_t} \nabla f(x_t)v$ 
5:    $v_t^{FW-S} \leftarrow \operatorname{argmin}_{v \in S_t} \nabla f(x_t)v$ 
6:   if  $\nabla f(x_t)(v_t^A - v_t^{FW-S}) \geq \Phi_t$  then
7:      $x_{t+1}, S_{t+1} \leftarrow \text{Simplex Gradient Descent}(x_t, S_t)$ 
8:      $\Phi_{t+1} \leftarrow \Phi_t$ 
9:   else
10:     $v_t \leftarrow \text{LPsep}_P(\nabla f(x_t), x_t, \Phi_t, K)$ 
11:    if  $v_t = \text{false}$  then
12:       $x_{t+1} \leftarrow x_t$ 
13:       $\Phi_{t+1} \leftarrow \Phi_t/2$  {update dual gap estimate}
14:       $S_{t+1} \leftarrow S_t$ 
15:    else
16:       $x_{t+1} \leftarrow \operatorname{argmin}_{x \in [x_t, v_t]} f(x)$ 
17:      Choose  $S_{t+1} \subseteq S_t \cup \{v_t\}$  minimal such that  $x_{t+1} \in S_{t+1}$ .
18:       $\Phi_{t+1} \leftarrow \Phi_t$ 
19:    end if
20:  end if
21: end for
```

Gradient Descent Phase

Frank-Wolfe Phase

COMPUTATIONAL RESULTS

BCG outperforms several recent variants of Frank-Wolfe algorithm

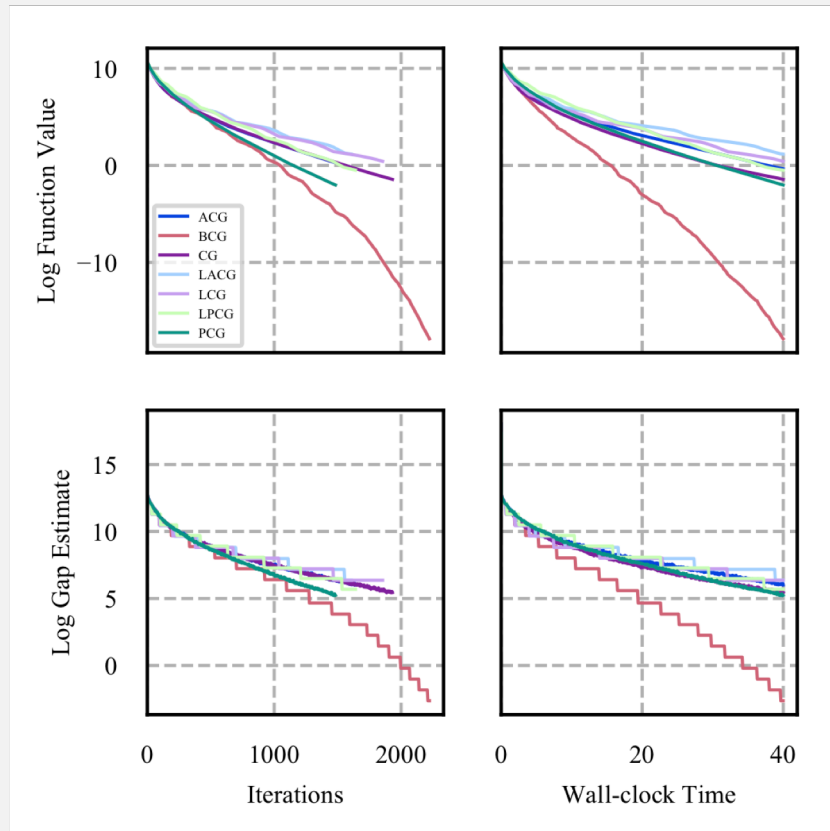


Fig 1: Lasso Regression

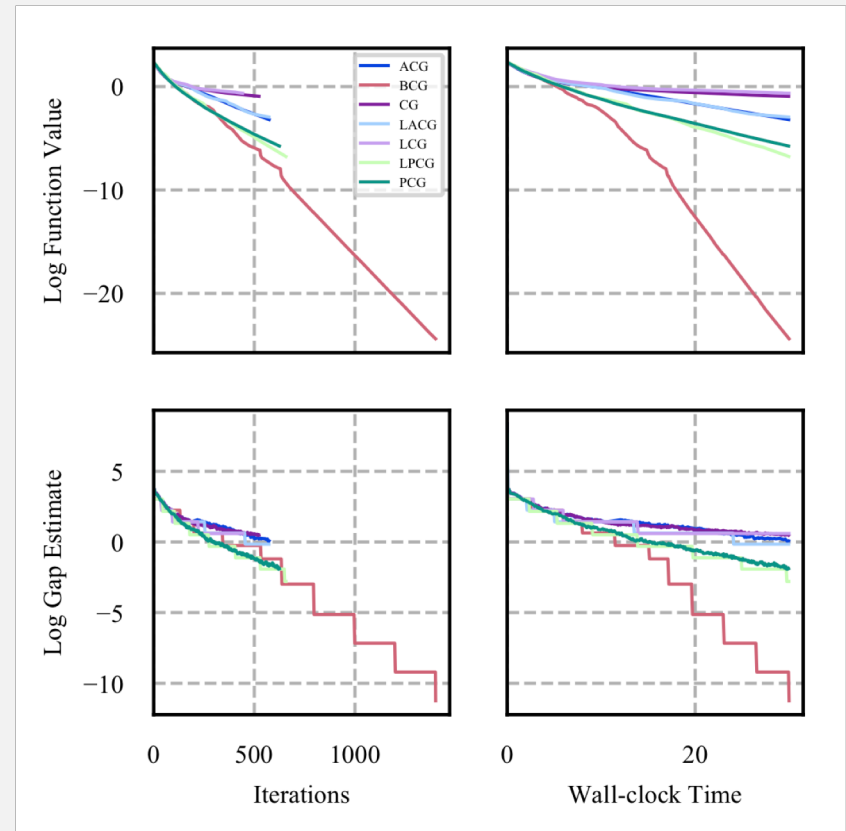


Fig 2: Sparse Signal Recovery

CONVERGENCE

Theorem

If f is a strongly convex and smooth function over the polytope P with geometric strong convexity μ and simplicial curvature, then BCG algorithm ensures $f(x_t) - f(x^*) \leq \epsilon$ for some T that satisfies:

$$T \leq \left\lceil \log \frac{2\Phi_0}{\epsilon} \right\rceil + 8 \left\lceil \log \frac{\Phi_0}{2C^\Delta} \right\rceil + \frac{64C^\Delta}{\mu} \left\lceil \log \frac{4C^\Delta}{\epsilon} \right\rceil = O\left(\frac{C^\Delta}{\mu} \log \frac{\Phi_0}{\epsilon}\right)$$

THANKS!