

Block-Coordinate Frank-Wolfe Optimization

with applications to structured prediction

Martin Jaggi

CMAP, Ecole Polytechnique

XRCE Seminar 2012 / 11 / 13

Co-Authors: Simon Lacoste-Julien, Mark Schmidt and Patrick Pletscher

Outline

- Two Old First-Order Optimizers
 - Coordinate Descent
 - The Frank-Wolfe Algorithm
- Duality for Constrained Convex Optimization
- Block-Separable Problems
 - A new block-coordinate variant of Frank-Wolfe
- Applications: Large Margin Prediction
 - binary SVMs
 - structural SVMs

Coordinate Descent

(for snow-
avalanche
rescue)



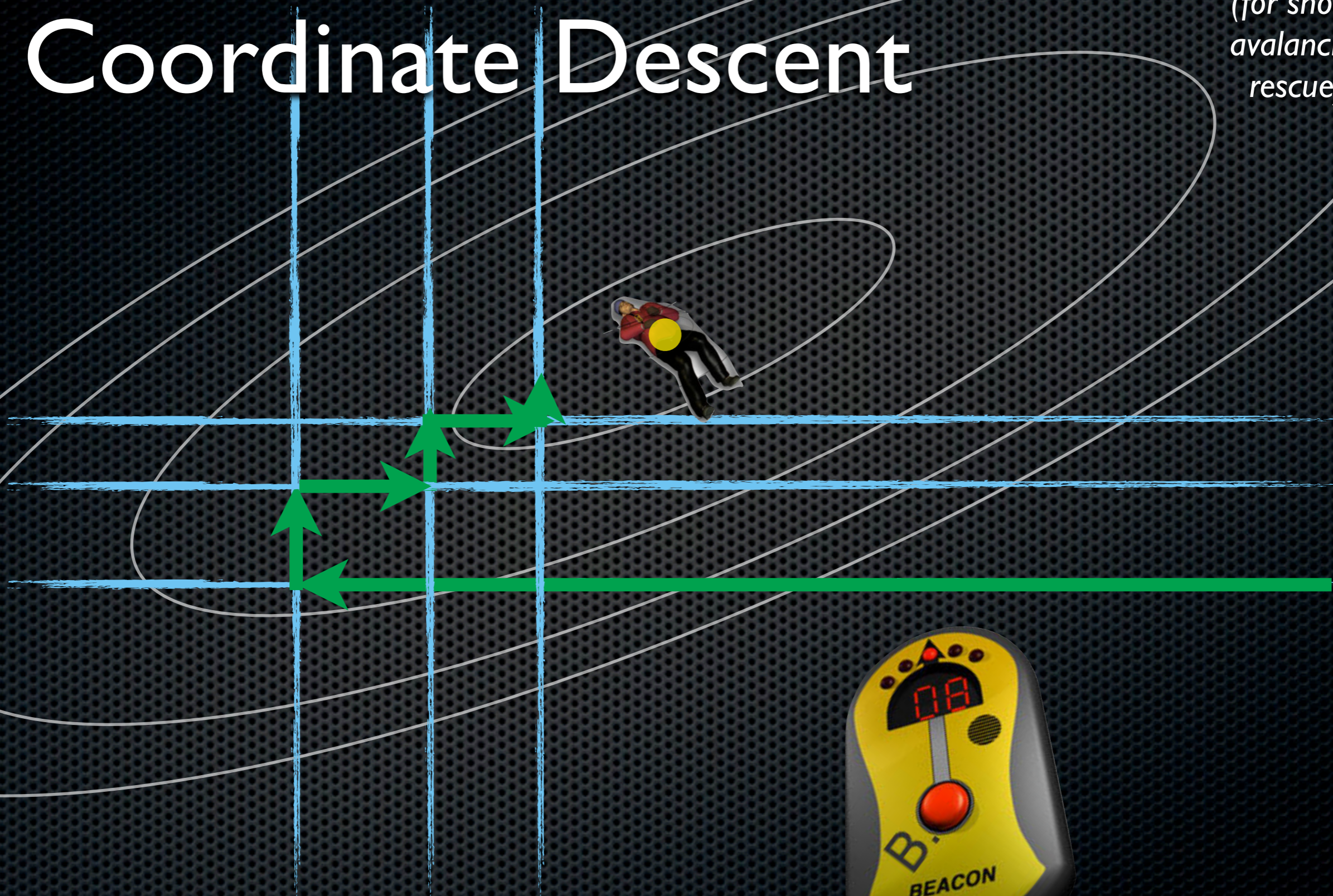
Coordinate Descent

(for snow-
avalanche
rescue)

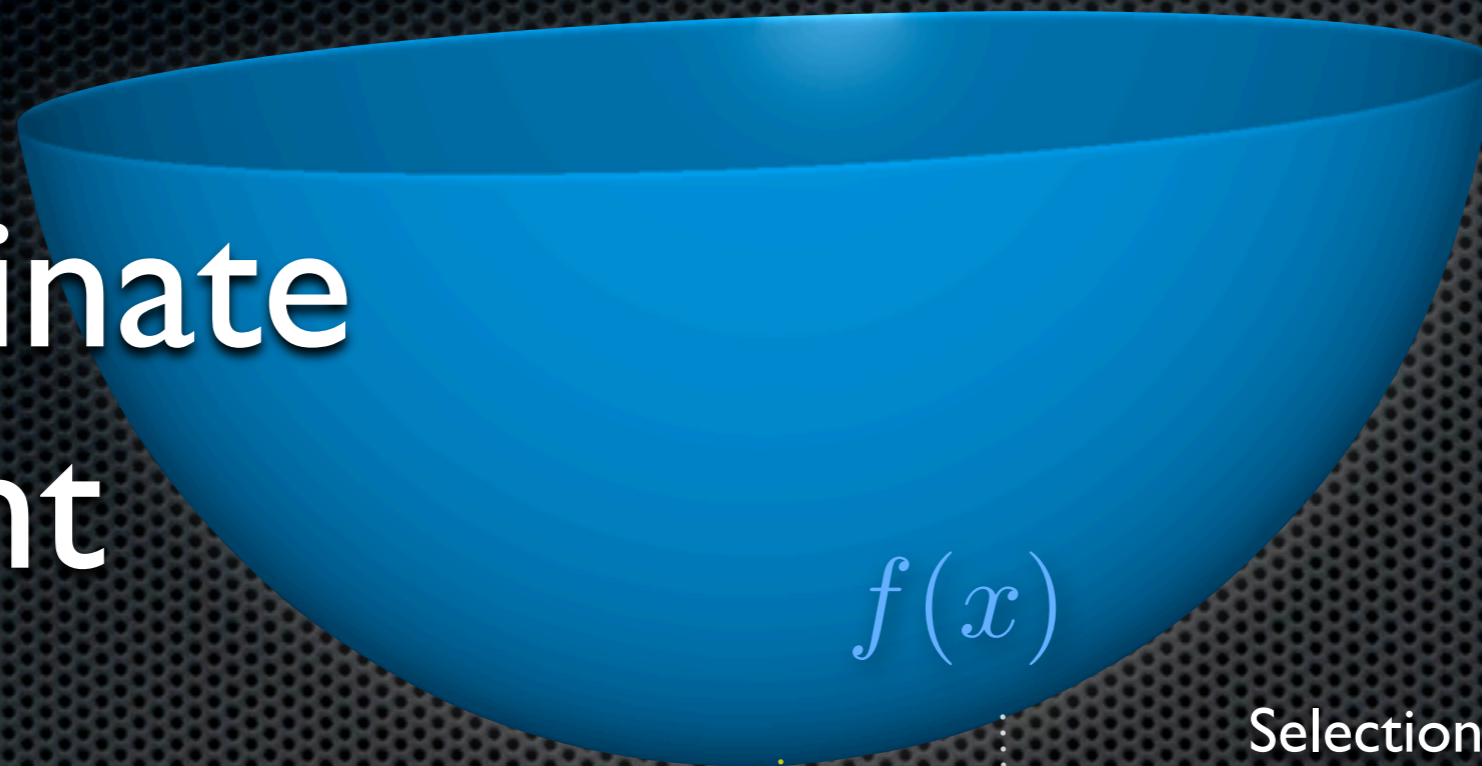


Coordinate Descent

(for snow-
avalanche
rescue)

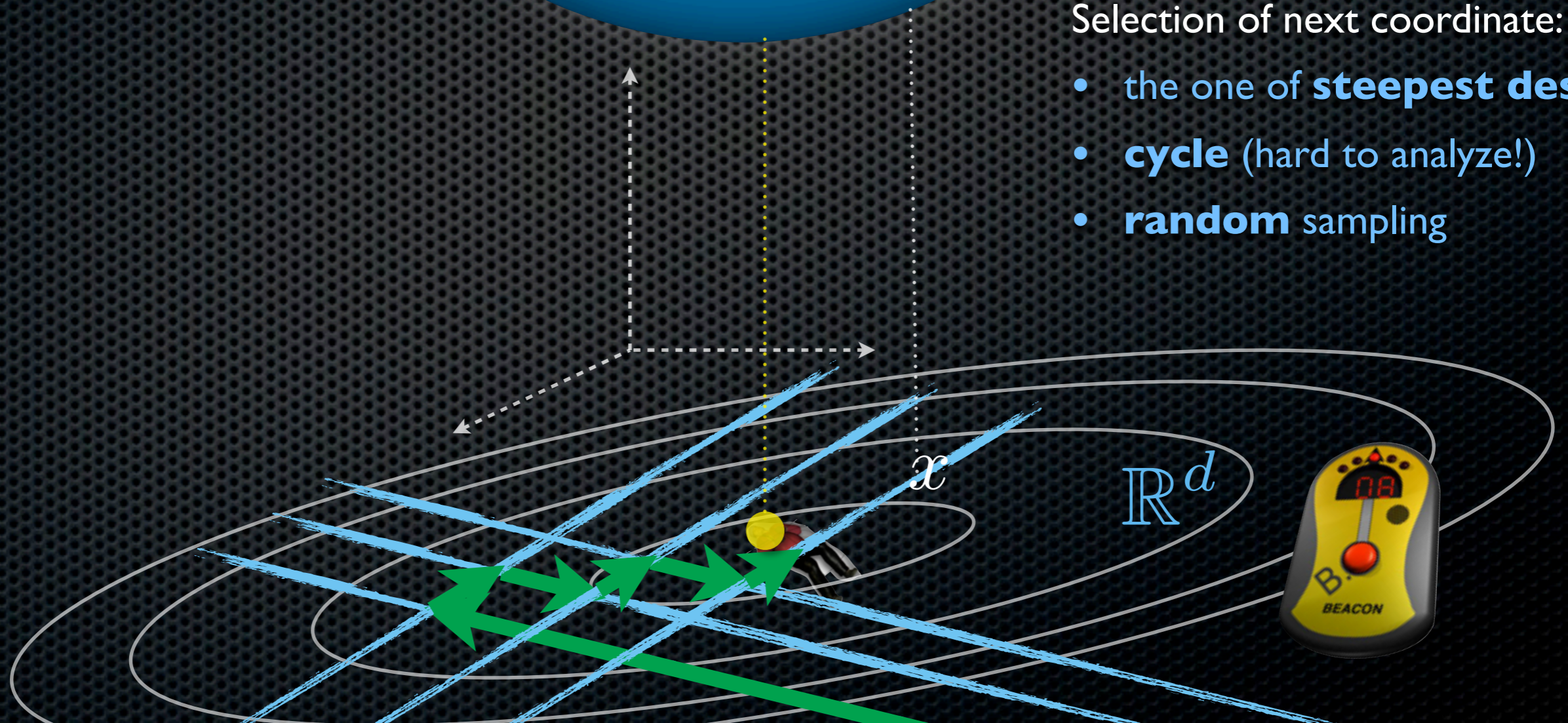


Coordinate Descent



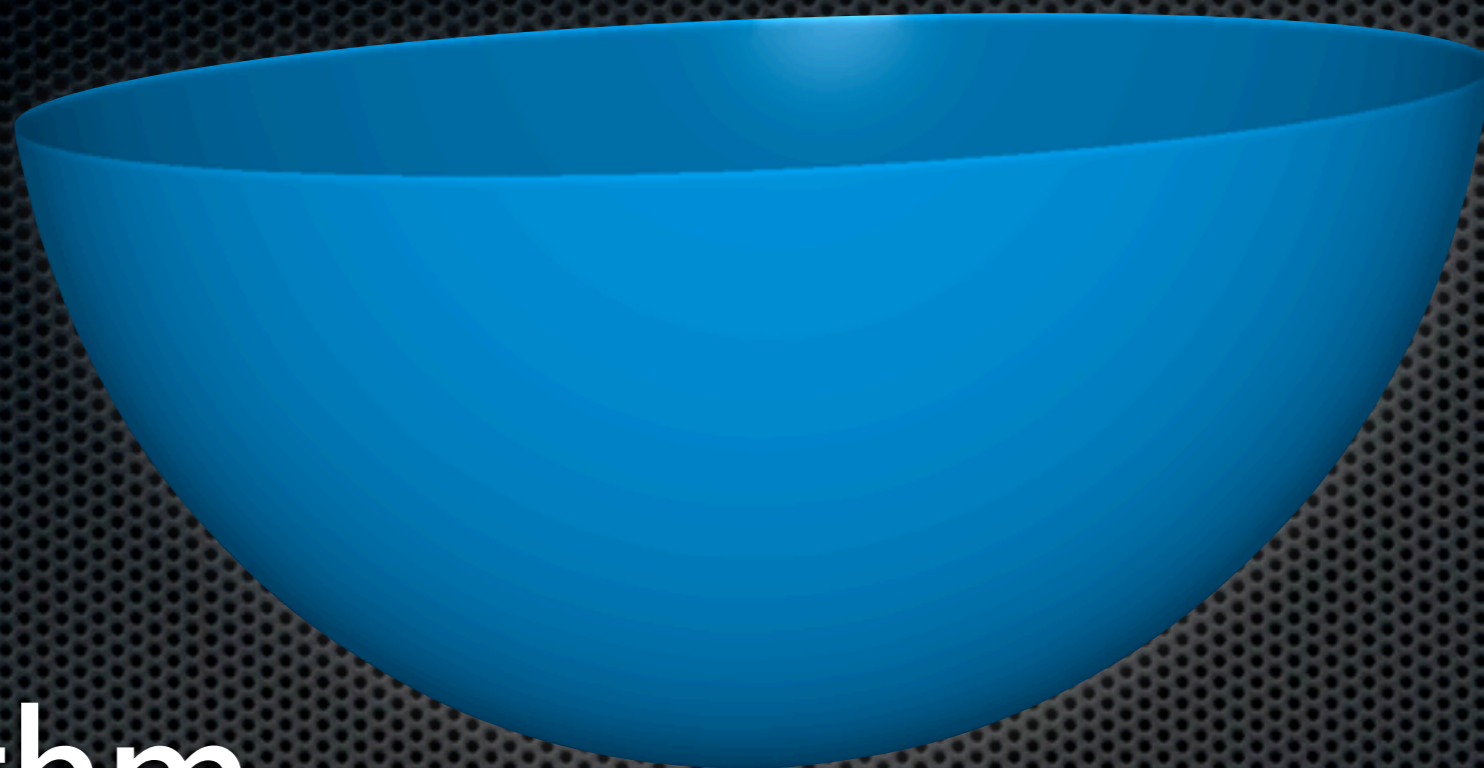
Selection of next coordinate:

- the one of **steepest desc.**
- **cycle** (hard to analyze!)
- **random** sampling

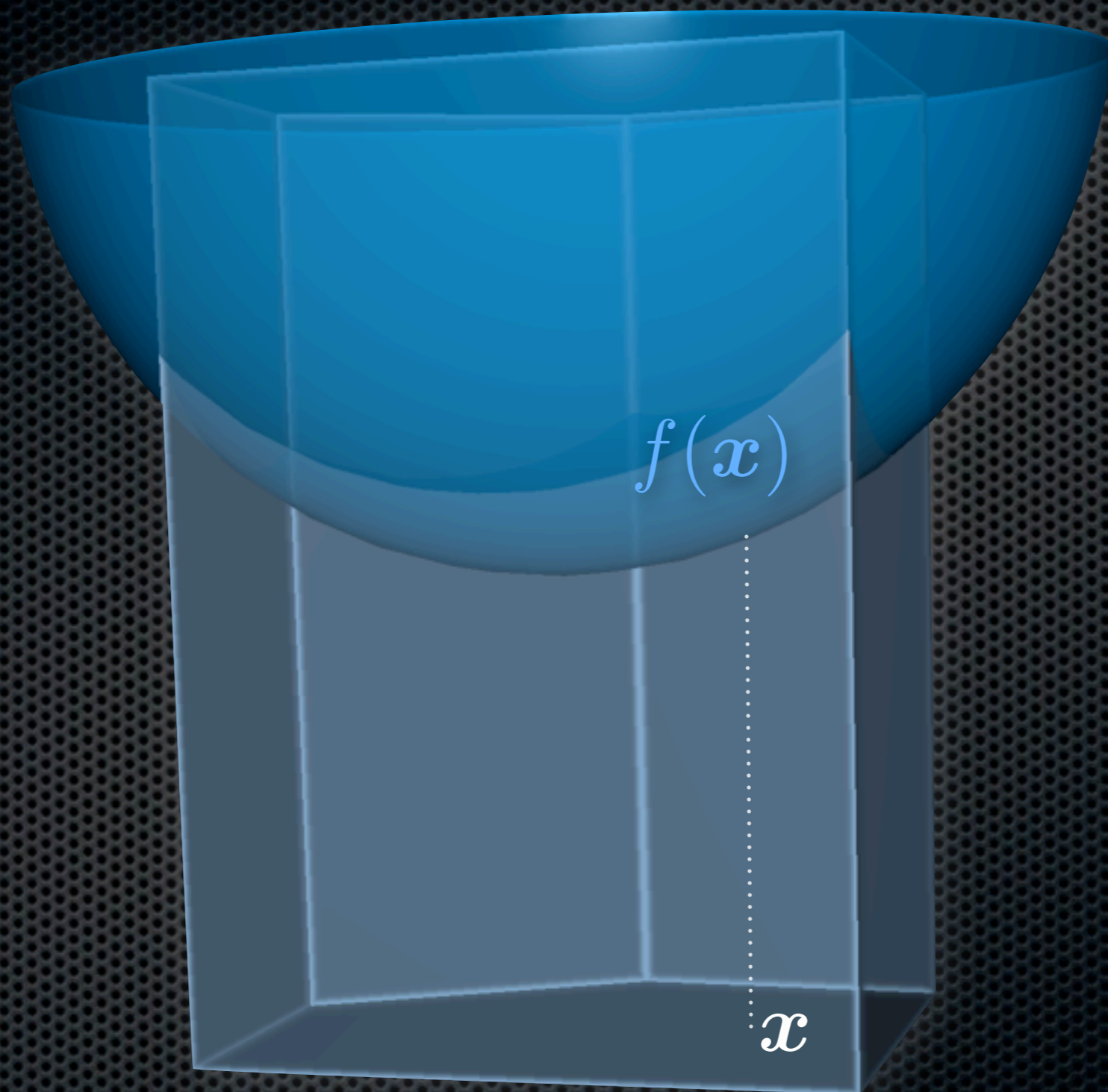


The Frank- Wolfe Algorithm

Frank and Wolfe (1956)

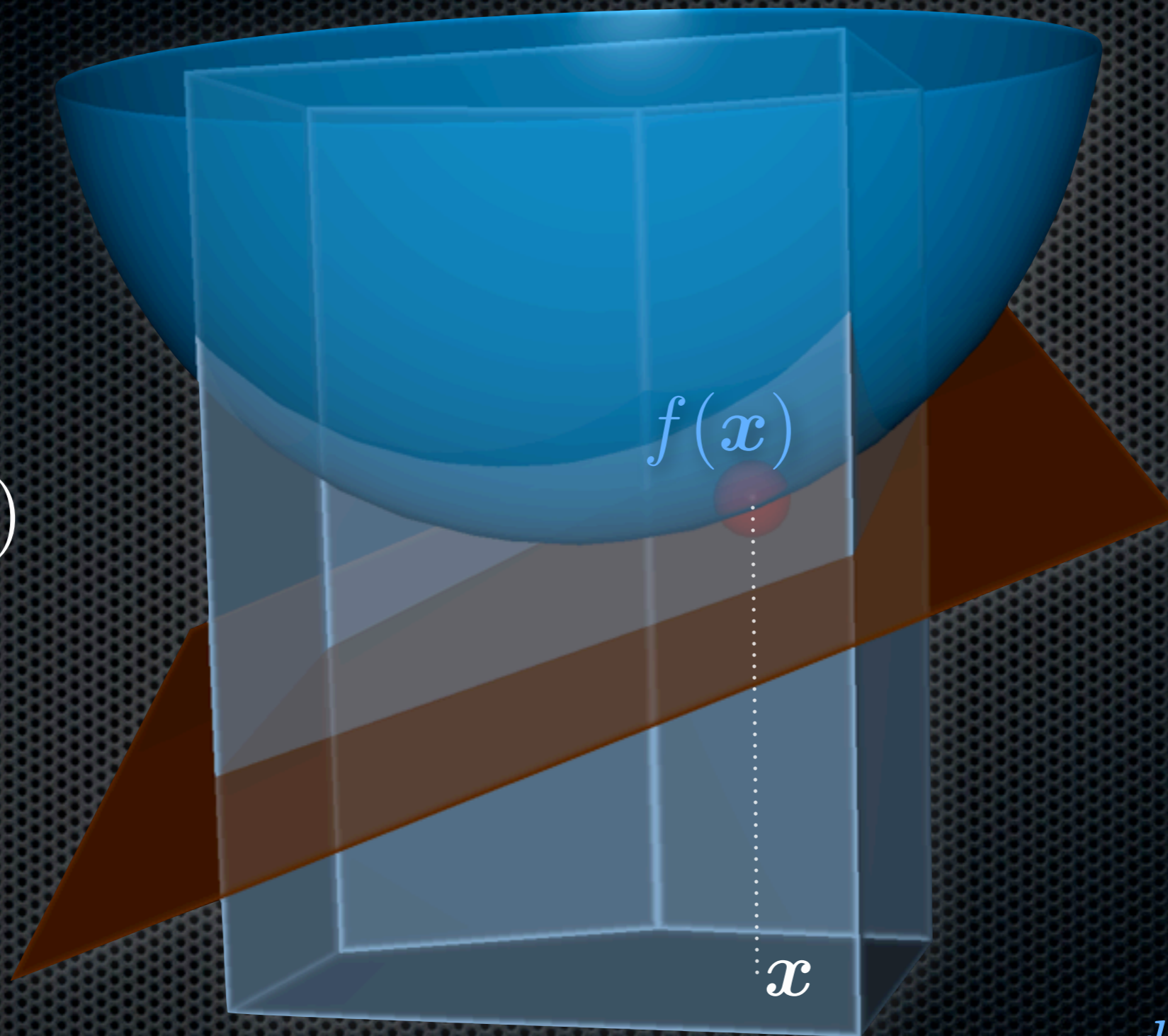


$$\min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x})$$



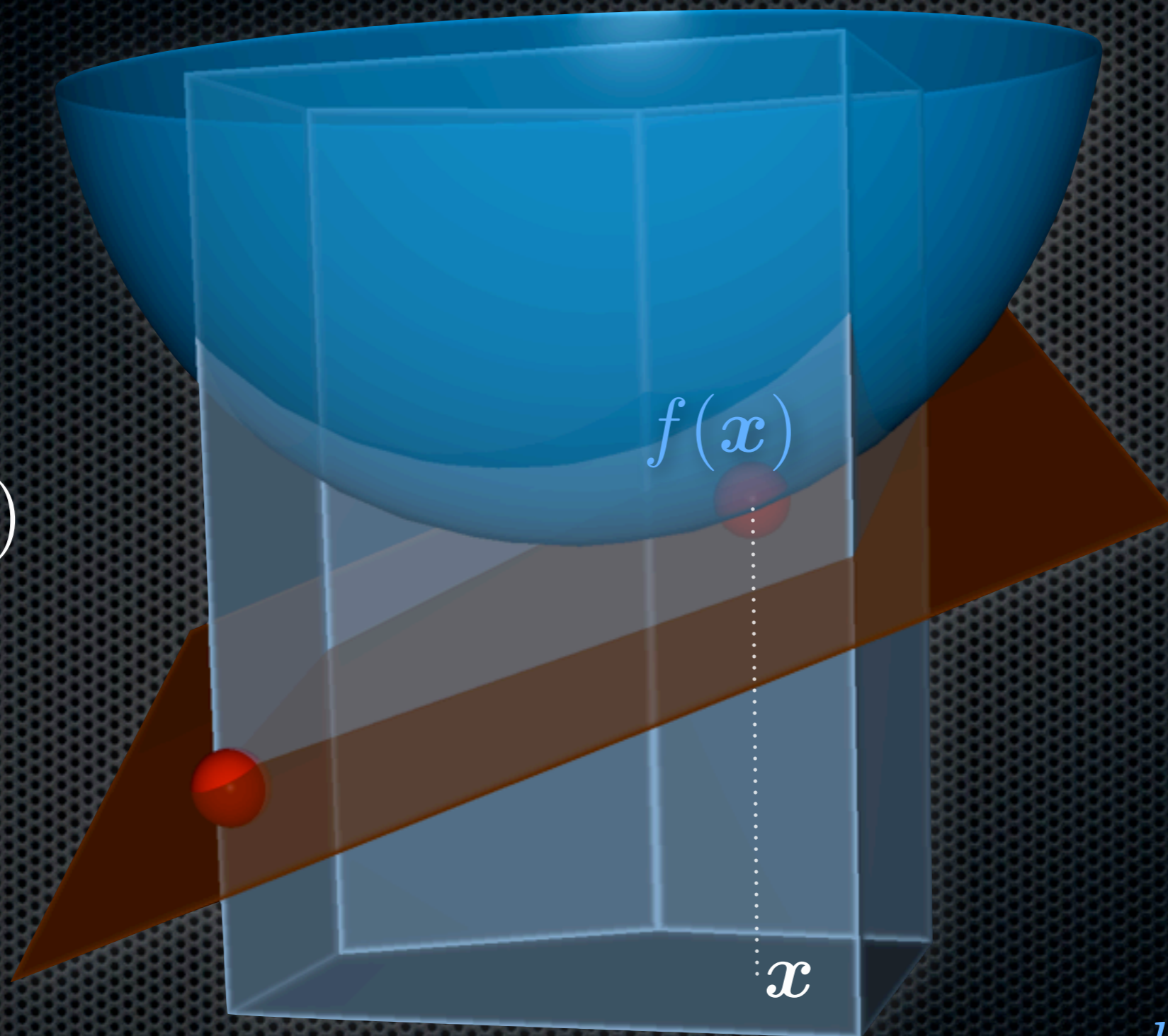
$$\mathcal{D} \subset \mathbb{R}^d$$

$$\min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x})$$



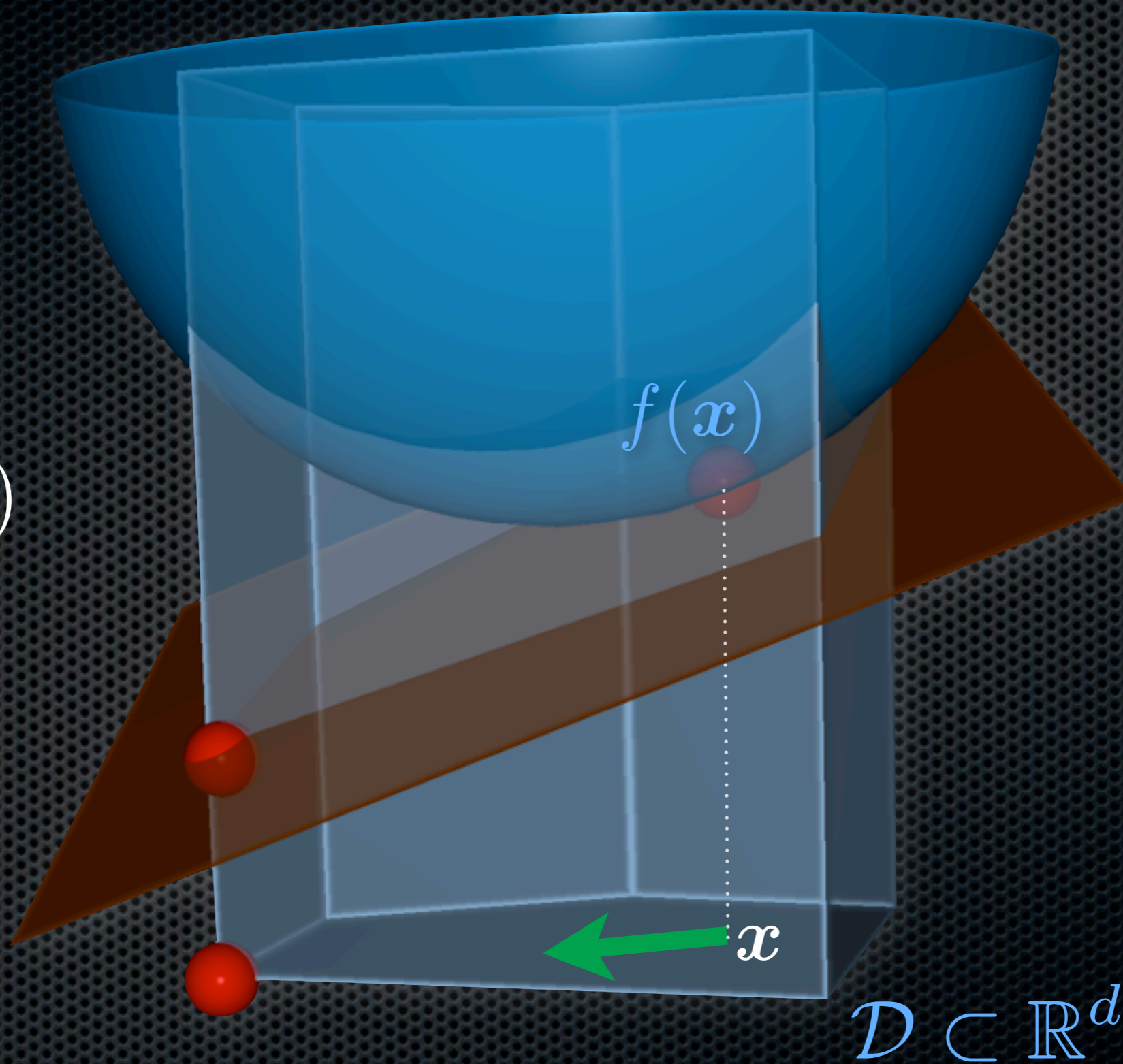
$$\mathcal{D} \subset \mathbb{R}^d$$

$$\min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x})$$



$$\mathcal{D} \subset \mathbb{R}^d$$

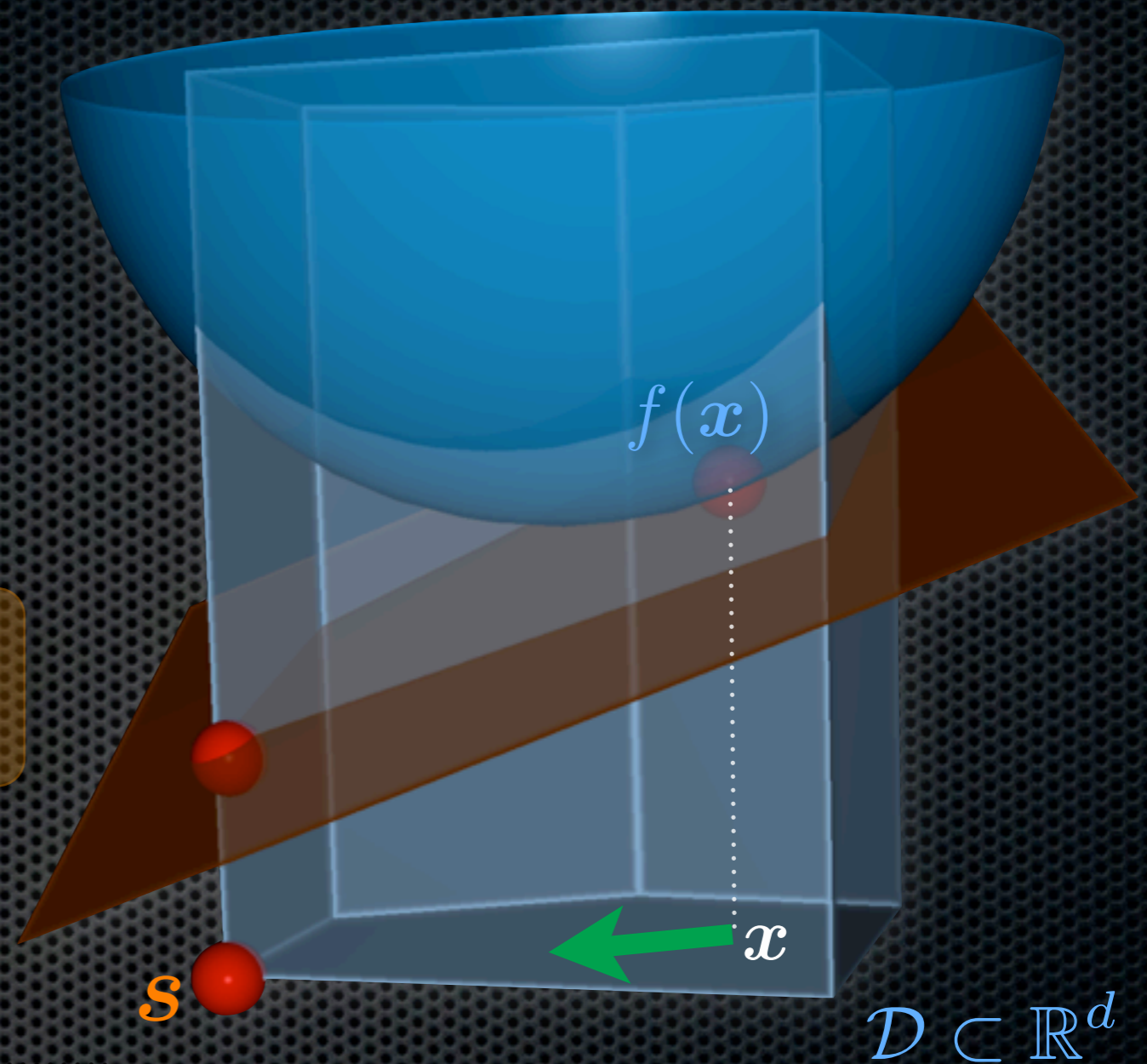
$$\min_{x \in D} f(x)$$



$$D \subset \mathbb{R}^d$$

The Linearized Problem

$$\min_{s' \in \mathcal{D}} f(x) + \langle s' - x, \nabla f(x) \rangle$$



Algorithm 1: Frank-Wolfe

Let $x^{(0)} \in \mathcal{D}$

for $k = 0 \dots K$ do

Compute $s := \arg \min_{s' \in \mathcal{D}} \langle s', \nabla f(x^{(k)}) \rangle$

Let $\gamma := \frac{2}{k+2}$, or optimize γ by line-search

Update $x^{(k+1)} := (1 - \gamma)x^{(k)} + \gamma s$

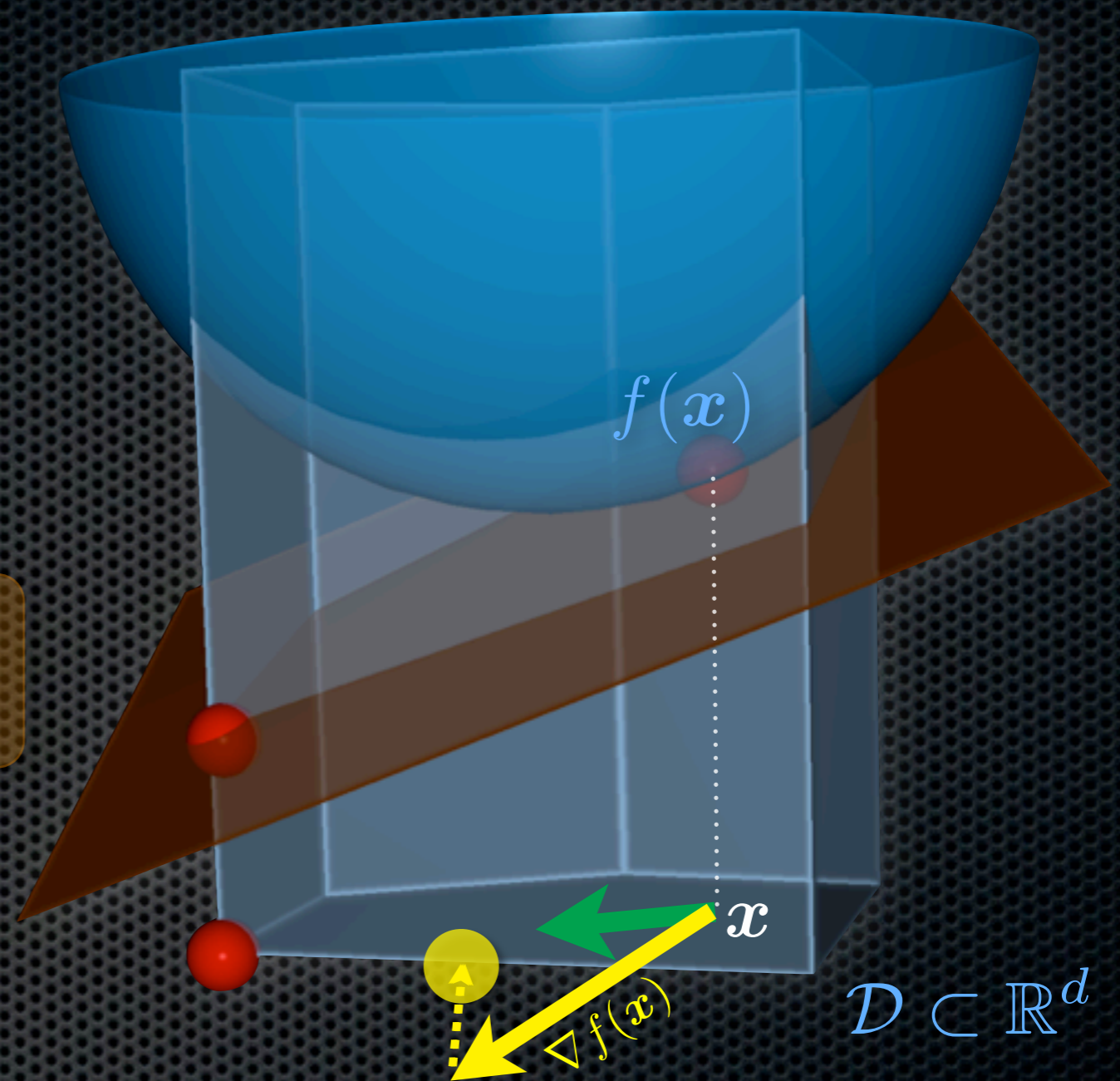
end

Theorem:

Algorithm obtains
accuracy $O\left(\frac{1}{k}\right)$
after k **steps**.

The Linearized Problem

$$\min_{s' \in \mathcal{D}} f(x) + \langle s' - x, \nabla f(x) \rangle$$



	Frank-Wolfe	Gradient Descent
Cost per step	Approx. solve linearized problem on \mathcal{D}	Projection back to \mathcal{D}
Sparse / Low Rank Solutions	✓ (depending on the domain)	✗
Convergence	$1/k$	$1/k$

Some Examples of Atomic Domains Suitable for Frank-Wolfe

\mathcal{X}	Optimization Domain		Complexity of one Frank-Wolfe Iteration	
	Atoms \mathcal{A}	$\mathcal{D} = \text{conv}(\mathcal{A})$	$\sup_{s \in \mathcal{D}} \langle s, \mathbf{y} \rangle$	Complexity
\mathbb{R}^n	Sparse Vectors	$\ \cdot\ _1$ -ball	$\ \mathbf{y}\ _\infty$	$O(n)$
\mathbb{R}^n	Sign-Vectors	$\ \cdot\ _\infty$ -ball	$\ \mathbf{y}\ _1$	$O(n)$
\mathbb{R}^n	ℓ_p -Sphere	$\ \cdot\ _p$ -ball	$\ \mathbf{y}\ _q$	$O(n)$
\mathbb{R}^n	Sparse Non-neg. Vectors	Simplex Δ_n	$\max_i \{y_i\}$	$O(n)$
\mathbb{R}^n	Latent Group Sparse Vec.	$\ \cdot\ _{\mathcal{G}}$ -ball	$\max_{g \in \mathcal{G}} \ \mathbf{y}_{(g)}\ _g^*$	$\sum_{g \in \mathcal{G}} g $
$\mathbb{R}^{m \times n}$	Matrix Trace Norm	$\ \cdot\ _{tr}$ -ball	$\ \mathbf{y}\ _{op} = \sigma_1(\mathbf{y})$	$\tilde{O}(N_f / \sqrt{\varepsilon'})$ (Lanczos)
$\mathbb{R}^{m \times n}$	Matrix Operator Norm	$\ \cdot\ _{op}$ -ball	$\ \mathbf{y}\ _{tr} = \ \sigma_i(\mathbf{y})\ _1$	SVD
$\mathbb{R}^{m \times n}$	Schatten Matrix Norms	$\ \sigma_i(\cdot)\ _p$ -ball	$\ \sigma_i(\mathbf{y})\ _q$	SVD
$\mathbb{R}^{m \times n}$	Matrix Max-Norm	$\ \cdot\ _{\max}$ -ball		$\tilde{O}(N_f (n+m)^{1.5} / \varepsilon'^{2.5})$
$\mathbb{R}^{n \times n}$	Permutation Matrices	Birkhoff polytope		$O(n^3)$
$\mathbb{R}^{n \times n}$	Rotation Matrices			SVD (Procrustes prob.)
$\mathbb{S}^{n \times n}$	Rank-1 PSD matrices of unit trace	$\{\mathbf{x} \succeq 0, \text{Tr}(\mathbf{x})=1\}$	$\lambda_{\max}(\mathbf{y})$	$\tilde{O}(N_f / \sqrt{\varepsilon'})$ (Lanczos)
$\mathbb{S}^{n \times n}$	PSD matrices of bounded diagonal	$\{\mathbf{x} \succeq 0, \mathbf{x}_{ii} \leq 1\}$		$\tilde{O}(N_f n^{1.5} / \varepsilon'^{2.5})$

Table 1: Some examples of atomic domains suitable for optimization using the Frank-Wolfe algorithm. Here SVD refers to the complexity of computing a singular value decomposition, which is $O(\min\{mn^2, m^2n\})$. N_f is the number of non-zero entries in the gradient of the objective function f , and $\varepsilon' = \frac{2\delta C_f}{k+2}$ is the required accuracy for the linear subproblems. For any $p \in [1, \infty]$, the conjugate value q is meant to satisfy $\frac{1}{p} + \frac{1}{q} = 1$, allowing $q = \infty$ for $p = 1$ and $q = 1$ for $p = \infty$.

A Simple Alternative Optimization Duality

The Problem

$$\min_{x \in D} f(x)$$

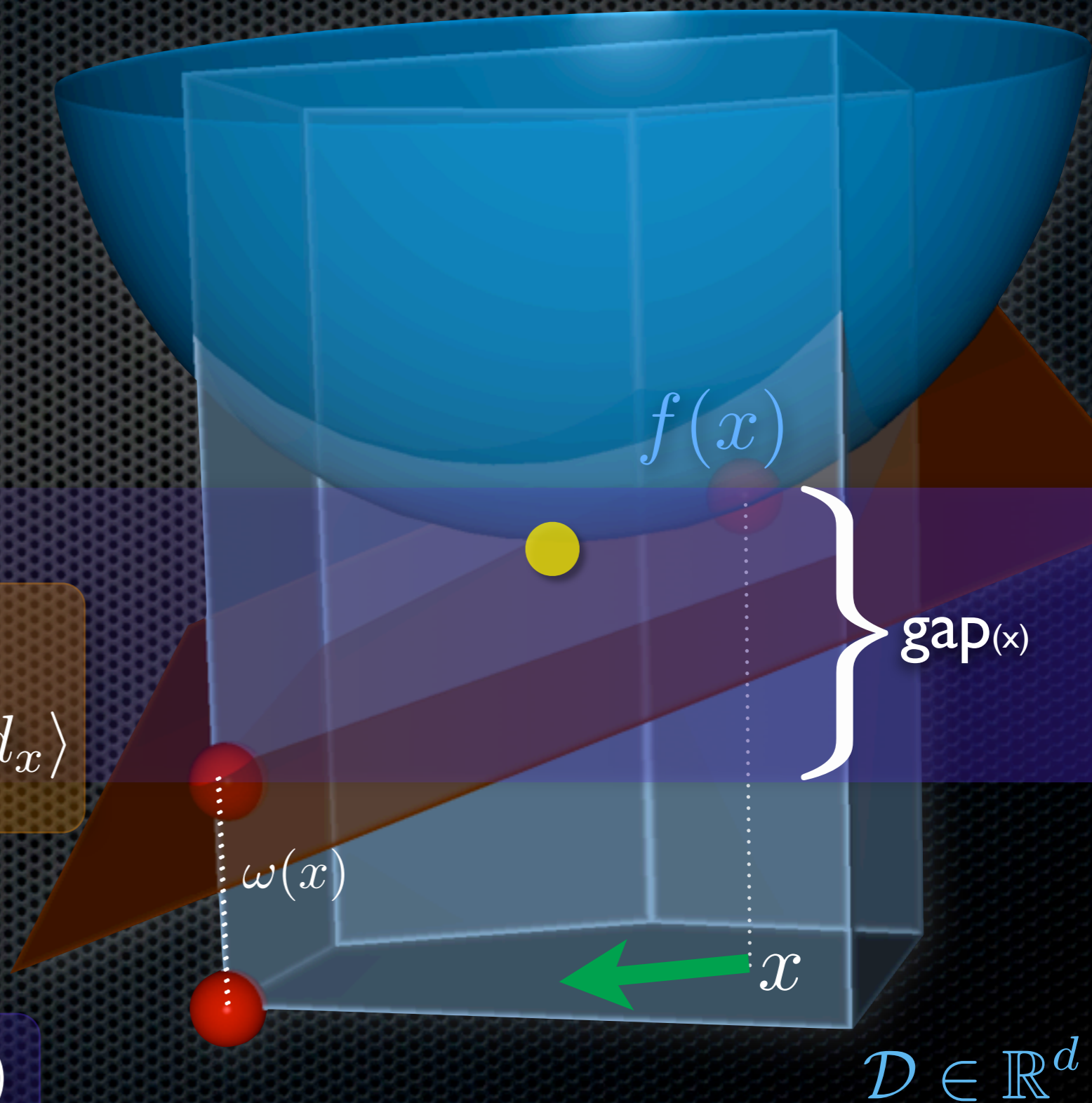
The Dual

$$\omega(x) :=$$

$$\min_{y \in D} f(x) + \langle y - x, d_x \rangle$$

Weak Duality

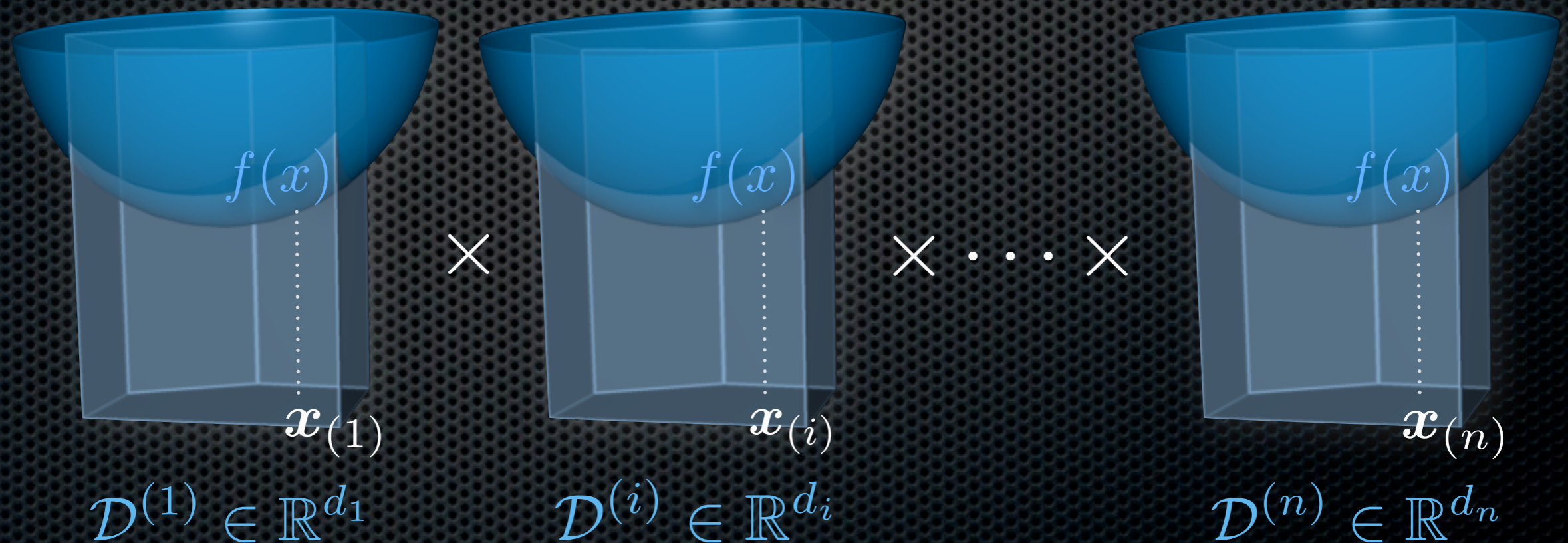
$$\omega(x) \leq f(x^*) \leq f(x')$$

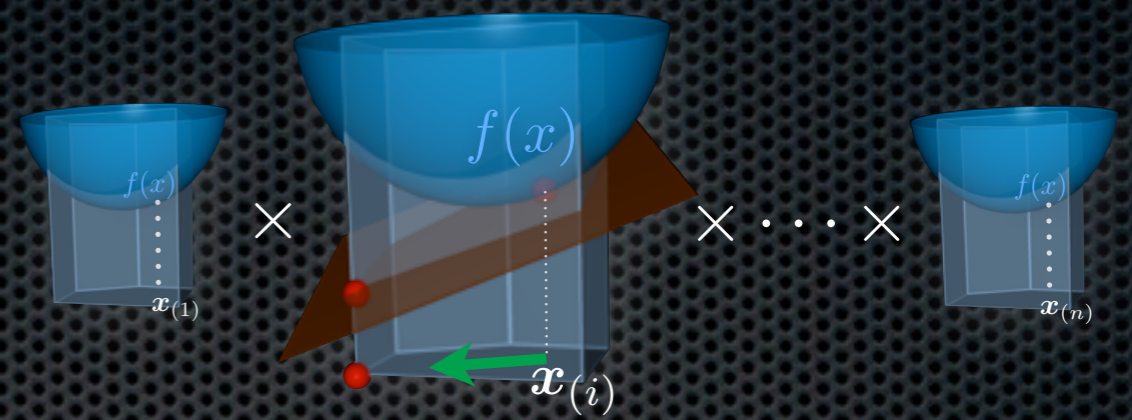
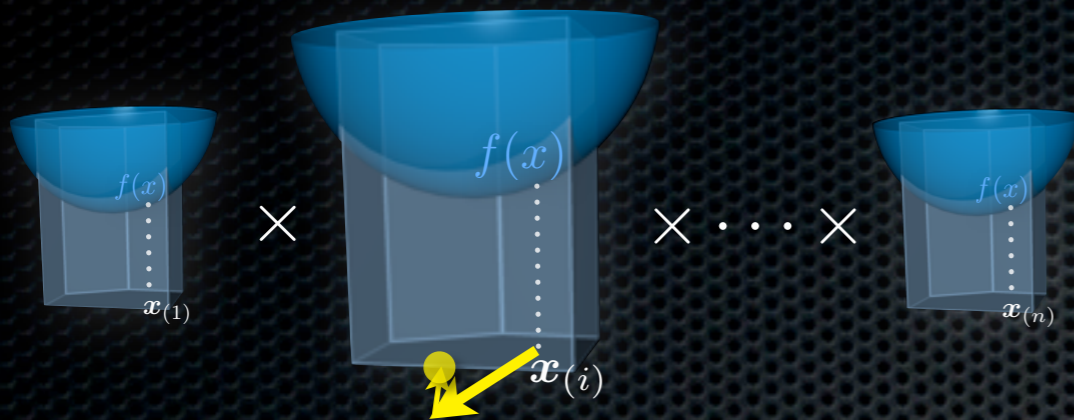


Block-Separable Optimization Problems

domain is a
product of
 n blocks

$$\min_{\mathbf{x} \in \mathcal{D}^{(1)} \times \dots \times \mathcal{D}^{(n)}} f(\mathbf{x})$$
$$\mathbf{x} = (\mathbf{x}_{(1)}, \dots, \mathbf{x}_{(n)})$$





Algorithm 2: Uniform Coordinate Descent

Let $\mathbf{x}^{(0)} \in \mathcal{D}$

for $k = 0 \dots K$ do

Pick $i \in_{u.a.r.} [n]$

Compute $\mathbf{s}_{(i)} := \arg \min_{\mathbf{s}_{(i)} \in \mathcal{D}^{(i)}} \left\langle \mathbf{s}_{(i)}, \nabla_{(i)} f(\mathbf{x}^{(k)}) \right\rangle + \frac{L_i}{2} \|\mathbf{s}_{(i)} - \mathbf{x}_{(i)}\|^2$

Update $\mathbf{x}_{(i)}^{(k+1)} := \mathbf{x}_{(i)}^{(k)} + (\mathbf{s}_{(i)} - \mathbf{x}_{(i)}^{(k)})$

end

Algorithm 3: Block-Coordinate "Frank-Wolfe"

Let $\mathbf{x}^{(0)} \in \mathcal{D}$

for $k = 0 \dots K$ do

Pick $i \in_{u.a.r.} [n]$

Compute $\mathbf{s}_{(i)} := \arg \min_{\mathbf{s}_{(i)} \in \mathcal{D}^{(i)}} \left\langle \mathbf{s}_{(i)}, \nabla_{(i)} f(\mathbf{x}^{(k)}) \right\rangle$

Let $\gamma := \frac{2n}{k+2n}$, or optimize γ by line-search

Update $\mathbf{x}_{(i)}^{(k+1)} := \mathbf{x}_{(i)}^{(k)} + \gamma(\mathbf{s}_{(i)} - \mathbf{x}_{(i)}^{(k)})$

end

Theorem:

Algorithm obtains

accuracy

$$O\left(\frac{2n}{k+2n}\right)$$

after k steps.

Hidden constant:

Curvature

$$\leq \sum_i L_f \text{diam}^2(\mathcal{D}^{(i)})$$

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

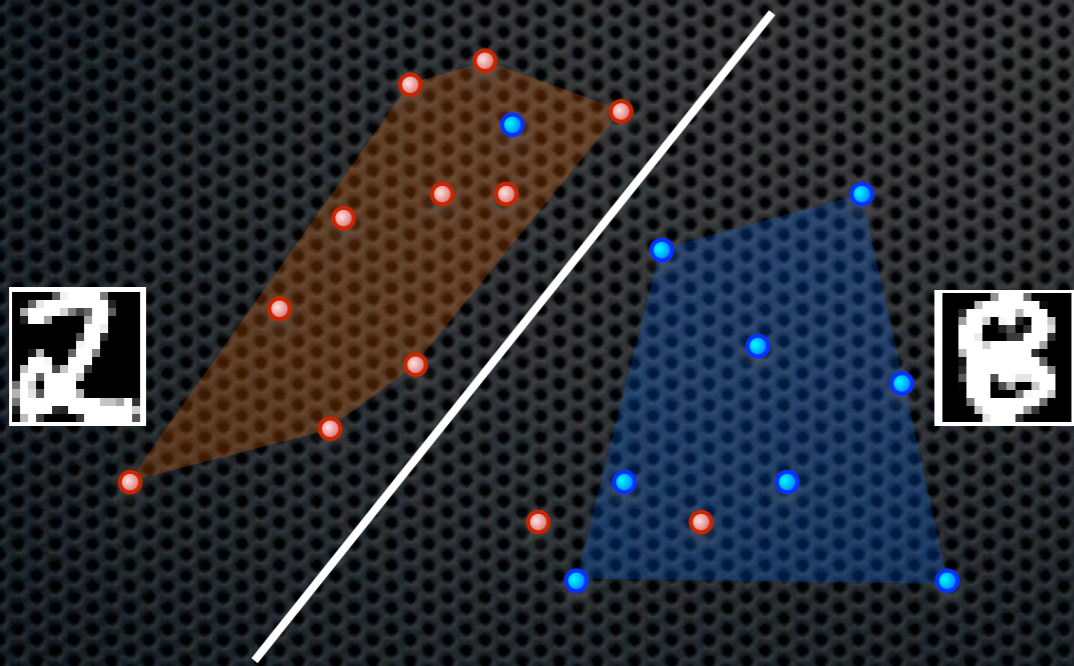
our
arXiv
paper

Nesterov (2012)

“Huge-Scale” Coordinate
Descent. *J. Opt*

Applications: Large Margin Prediction

- Binary Support Vector Machine
(no bias)
- also: Ranking SVM

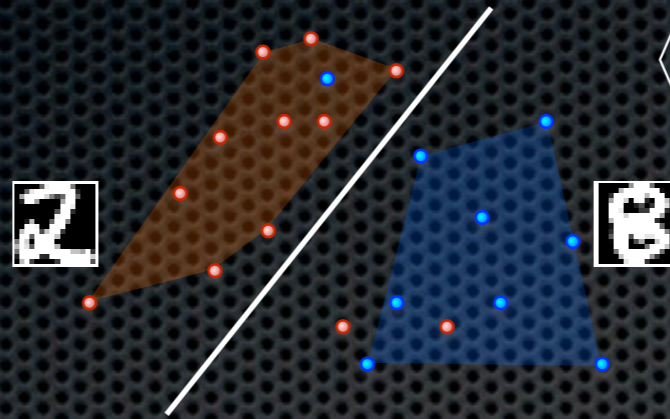


$$\langle \mathbf{w}, \phi(\mathbf{x}_i) \mathbf{y}_i \rangle \geq 1 - \xi_i$$

primal problem:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{\lambda}{2} \|\mathbf{w}\|^2 \\ & + \frac{1}{n} \sum_{i=1}^n \max \{0, 1 - \langle \mathbf{w}, \phi(\mathbf{x}_i) \mathbf{y}_i \rangle\} \end{aligned}$$

Binary SVM



$$\langle \mathbf{w}, \phi(\mathbf{x}_i) \mathbf{y}_i \rangle \geq 1 - \xi_i$$

primal

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{\lambda}{2} \|\mathbf{w}\|^2 \\ & + \frac{1}{n} \sum_{i=1}^n \max \left\{ 0, 1 - \langle \mathbf{w}, \phi(\mathbf{x}_i) \mathbf{y}_i \rangle \right\} \end{aligned}$$

- d -dim
- non-smooth, strongly convex
- unconstrained

dual

$$\begin{aligned} \min_{\alpha \in \mathbb{R}^n} \quad & f(\alpha) := \\ & \frac{\lambda}{2} \left\| \underbrace{\sum_{i \in [n]} \alpha_i \frac{\phi(\mathbf{x}_i) \mathbf{y}_i}{\lambda n}}_{=: \mathbf{w} = A\alpha} \right\|^2 - \underbrace{\sum_{i \in [n]} \frac{\alpha_i}{n}}_{=: \mathbf{b}^T \alpha} \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq 1 \quad \forall i \in [n]. \end{aligned}$$

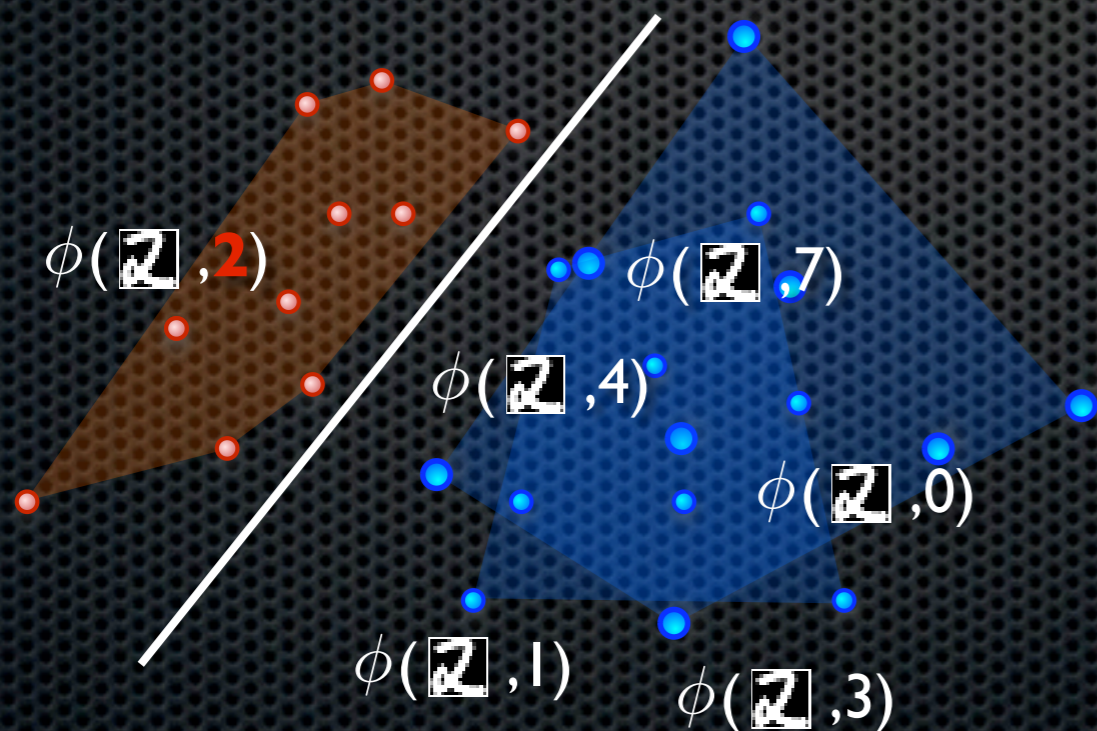
- n -dim
- smooth
- box-constrained

Structural SVM

“joint” feature map $\phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$

large margin “separation”

$$\langle \mathbf{w}, \phi(\mathbf{x}_i, \mathbf{y}_i) - \phi(\mathbf{x}_i, \mathbf{y}) \rangle \geq L(\mathbf{y}, \mathbf{y}_i) - \xi_i \quad \forall \mathbf{y}$$



primal problem:

$$\min_{\mathbf{w}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \max_{\mathbf{y} \in \mathcal{Y}} \left\{ L(\mathbf{y}, \mathbf{y}_i) - \langle \mathbf{w}, \underbrace{\phi(\mathbf{x}_i, \mathbf{y}_i) - \phi(\mathbf{x}_i, \mathbf{y})}_{=: \psi_i(\mathbf{y})} \rangle \right\}$$

loss-augmented decoding

Binary SVM

primal

$$\min_{\mathbf{w}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \max \left\{ 0, 1 - \langle \mathbf{w}, \phi(\mathbf{x}_i) \mathbf{y}_i \rangle \right\}$$

dual

$$\min_{\alpha \in \mathbb{R}^n} f(\alpha) := \frac{\lambda}{2} \left\| \underbrace{\sum_{i \in [n]} \alpha_i \frac{\phi(\mathbf{x}_i) \mathbf{y}_i}{\lambda n}}_{=: \mathbf{w} = A\alpha} \right\|^2 - \underbrace{\sum_{i \in [n]} \frac{\alpha_i}{n}}_{=: \mathbf{b}^T \alpha}$$

s.t. $0 \leq \alpha_i \leq 1 \quad \forall i \in [n].$

Structural SVM

primal

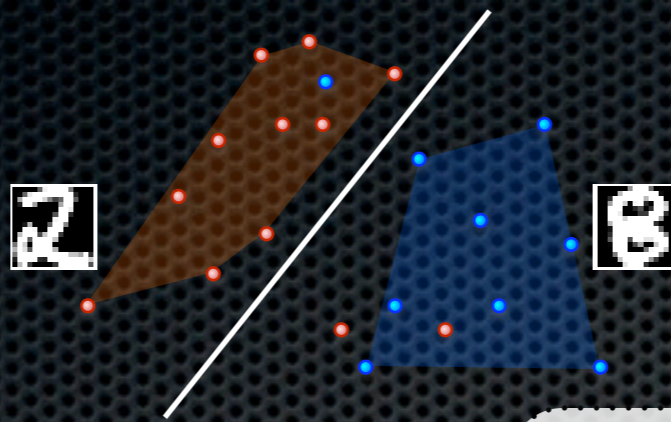
$$\min_{\mathbf{w}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \max_{\mathbf{y} \in \mathcal{Y}} \left\{ L(\mathbf{y}, \mathbf{y}_i) - \langle \mathbf{w}, \underbrace{\phi(\mathbf{x}_i, \mathbf{y}_i) - \phi(\mathbf{x}_i, \mathbf{y})}_{=: \psi_i(\mathbf{y})} \rangle \right\}$$

dual

$$\min_{\alpha \in \mathbb{R}^{n \cdot |\mathcal{Y}|}} f(\alpha) := \frac{\lambda}{2} \left\| \underbrace{\sum_{\substack{i \in [n], \\ \mathbf{y} \in \mathcal{Y}}} \alpha_i(\mathbf{y}) \frac{\psi_i(\mathbf{y})}{\lambda n}}_{=: \mathbf{w} = A\alpha} \right\|^2 - \underbrace{\sum_{\substack{i \in [n], \\ \mathbf{y} \in \mathcal{Y}}} \alpha_i(\mathbf{y}) \frac{L(\mathbf{y}, \mathbf{y}_i)}{n}}_{=: \mathbf{b}^T \alpha}$$

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

Binary SVM



$$\langle \mathbf{w}, \phi(\mathbf{x}_i) \mathbf{y}_i \rangle \geq 1 - \xi_i$$

primal

$$\min_{\mathbf{w}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \max \left\{ 0, 1 - \langle \mathbf{w}, \phi(\mathbf{x}_i) \mathbf{y}_i \rangle \right\}$$

- d -dim
- unconstrained
- non-smooth, strongly convex

dual

$$\min_{\alpha \in \mathbb{R}^n} f(\alpha) := \frac{\lambda}{2} \left\| \underbrace{\sum_{i \in [n]} \alpha_i \frac{\phi(\mathbf{x}_i) \mathbf{y}_i}{\lambda n}}_{=: \mathbf{w} = A\alpha} \right\|^2 - \underbrace{\sum_{i \in [n]} \frac{\alpha_i}{n}}_{=: \mathbf{b}^T \alpha}$$

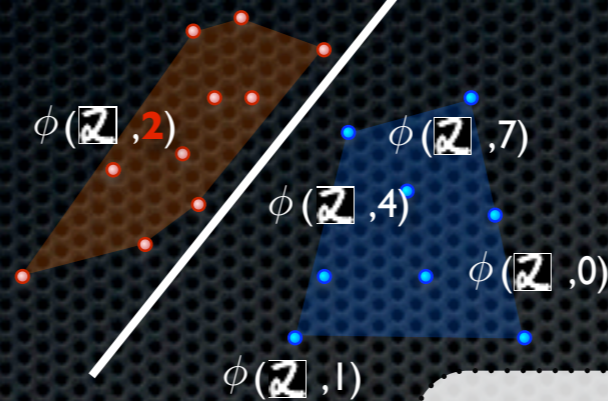
s.t. $0 \leq \alpha_i \leq 1 \quad \forall i \in [n].$

- n -dim
- box-constrained
- smooth

Optimization Algorithms

	primal	dual
batch	<ul style="list-style-type: none"> • subgradient descent • bundle methods 	<ul style="list-style-type: none"> • Frank-Wolfe • cutting planes (<i>SVM-light</i>)
online	<ul style="list-style-type: none"> • stochastic subgradient (<u><i>SGD, Pegasos</i></u>) 	<ul style="list-style-type: none"> • coordinate descent (<u><i>Hsieh, LibLinear</i></u>) • block-coordinate Frank-Wolfe

Structural SVM



primal

$$\min_{\mathbf{w}} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \max_{\mathbf{y} \in \mathcal{Y}} \left\{ L(\mathbf{y}, \mathbf{y}_i) - \langle \mathbf{w}, \underbrace{\phi(\mathbf{x}_i, \mathbf{y}_i) - \phi(\mathbf{x}_i, \mathbf{y})}_{=: \psi_i(\mathbf{y})} \rangle \right\}$$

- d -dim
- unconstrained
- non-smooth, strongly convex

dual

$$\min_{\alpha \in \mathbb{R}^{n \cdot |\mathcal{Y}|}} f(\alpha) := \frac{\lambda}{2} \left\| \underbrace{\sum_{\substack{i \in [n], \\ \mathbf{y} \in \mathcal{Y}}} \alpha_i(\mathbf{y}) \frac{\psi_i(\mathbf{y})}{\lambda n}}_{=: \mathbf{w} = A\alpha} \right\|^2 - \underbrace{\sum_{\substack{i \in [n], \\ \mathbf{y} \in \mathcal{Y}}} \alpha_i(\mathbf{y}) \frac{L(\mathbf{y}, \mathbf{y}_i)}{n}}_{=: \mathbf{b}^T \alpha}$$

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

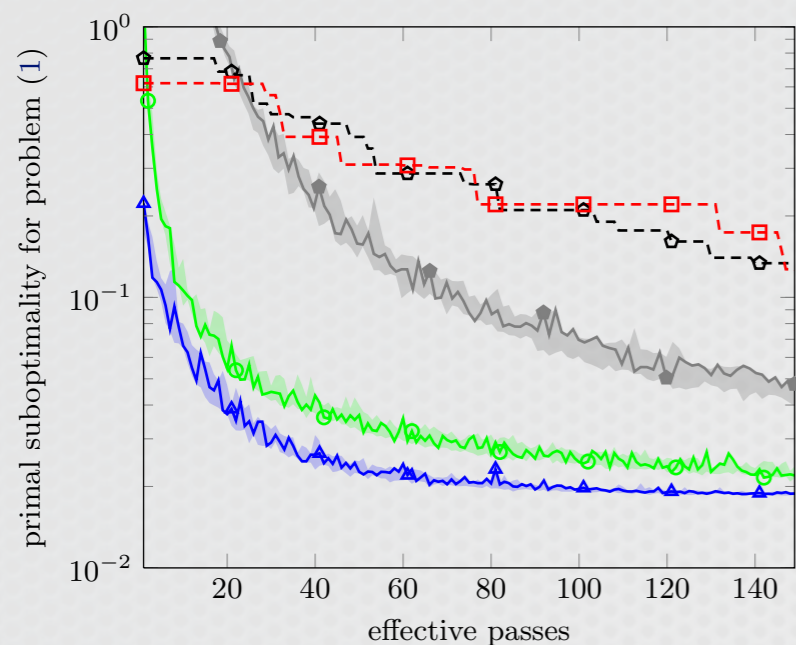
- $n |\mathcal{Y}|$ - dim
- block-constrained
- smooth

Optimization Algorithms

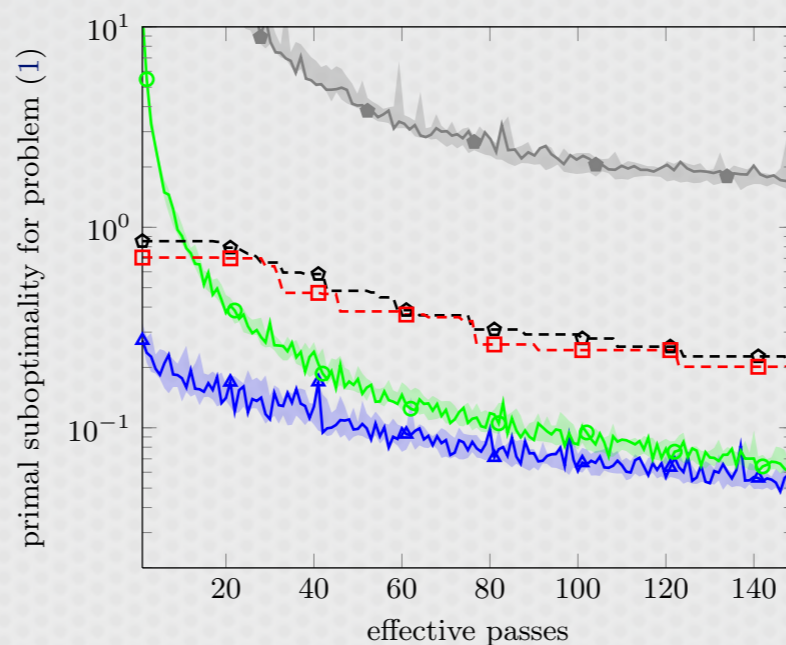
	primal	dual
batch	<ul style="list-style-type: none"> • subgradient descent • bundle methods 	<ul style="list-style-type: none"> • Frank-Wolfe • cutting planes (<i>SVM-struct</i>)
online	<ul style="list-style-type: none"> • stochastic subgradient (<i>SGD, Pegasos</i>) 	<ul style="list-style-type: none"> • block-coordinate descent (Nesterov) • block-coordinate Frank-Wolfe

Experimental Results

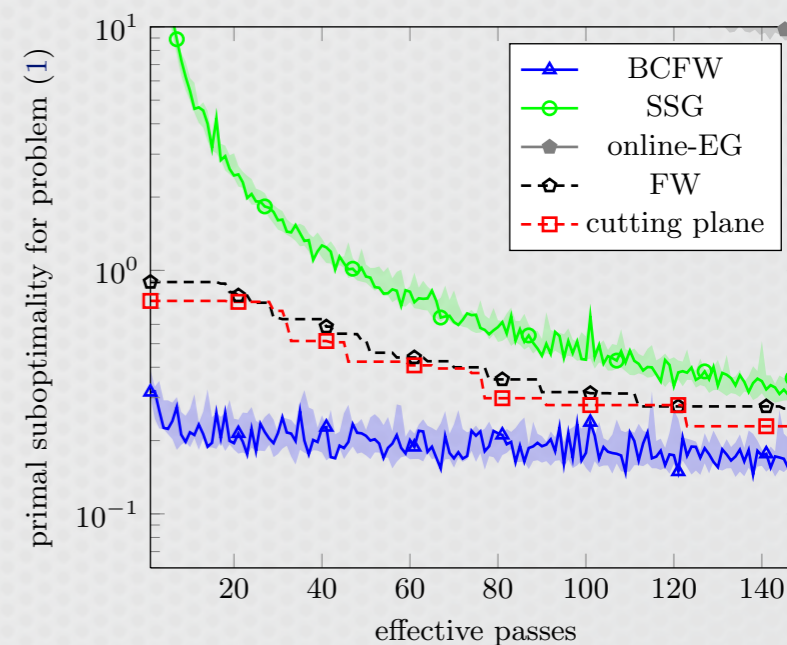
dataset		n	d
OCR	sequence labeling	6251	4028
CoNLL	POS sequence labeling	8936	1643026
Matching	word alignment	5000	82



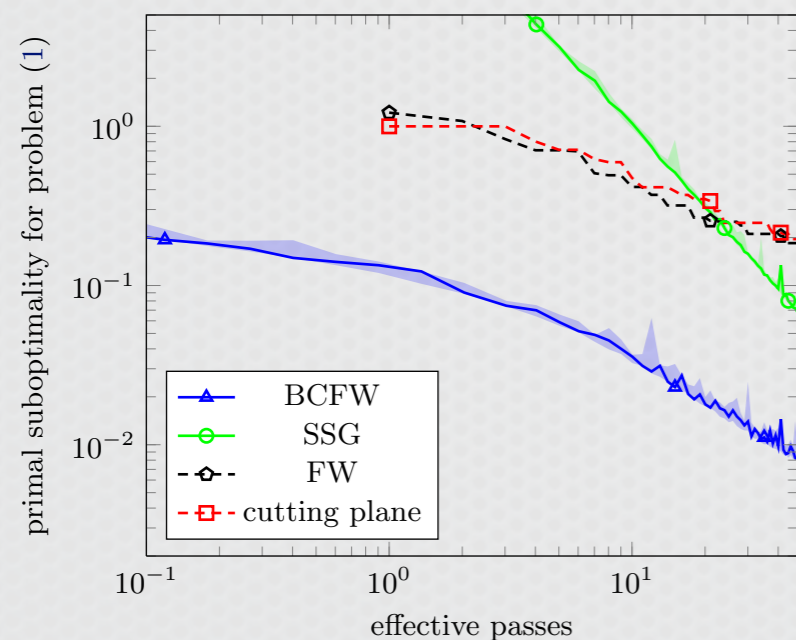
(a) OCR dataset, $\lambda = 0.01$.



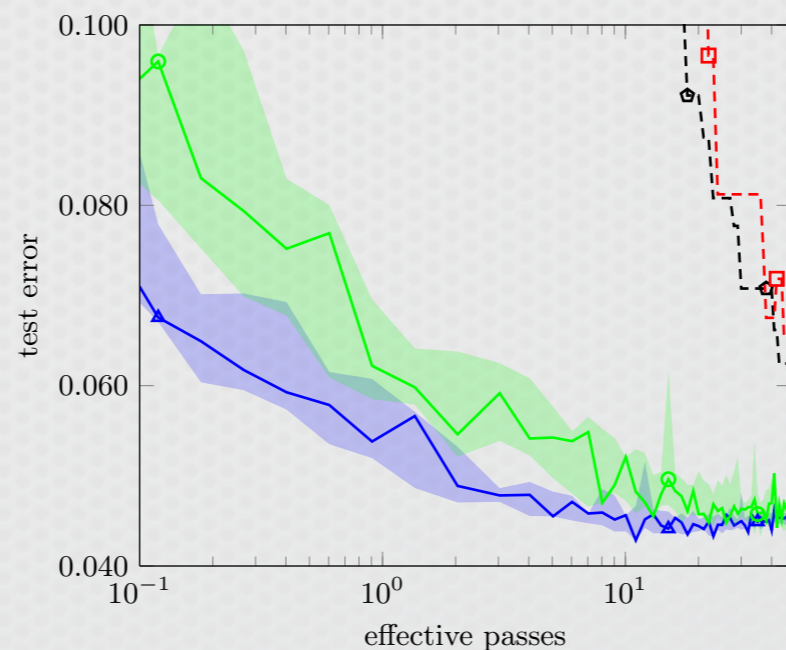
(b) OCR dataset, $\lambda = 0.001$.



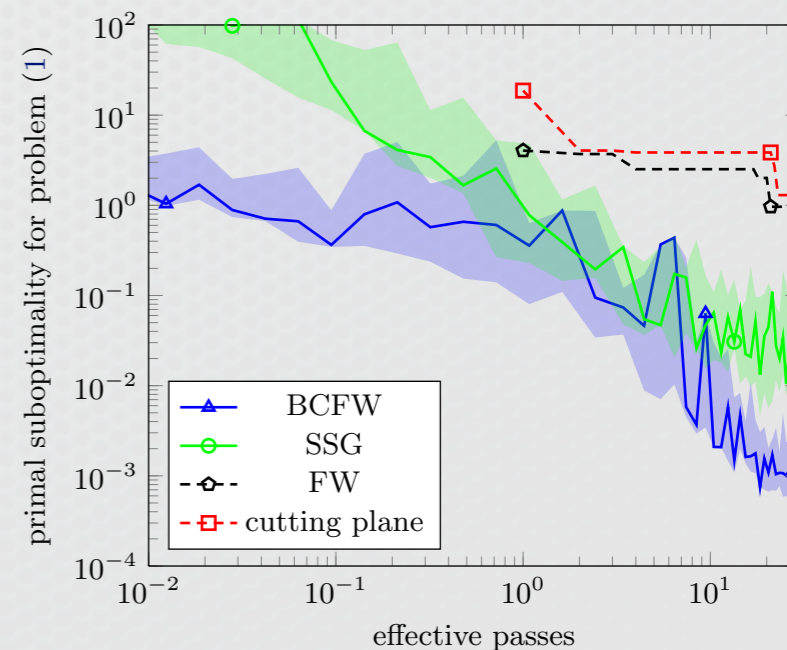
(c) OCR dataset, $\lambda = 1/n$.



(d) CoNLL dataset, $\lambda = 1/n$.



(e) Test error for $\lambda = 1/n$ on CoNLL.



(f) Matching dataset, $\lambda = 0.001$.

Thanks!

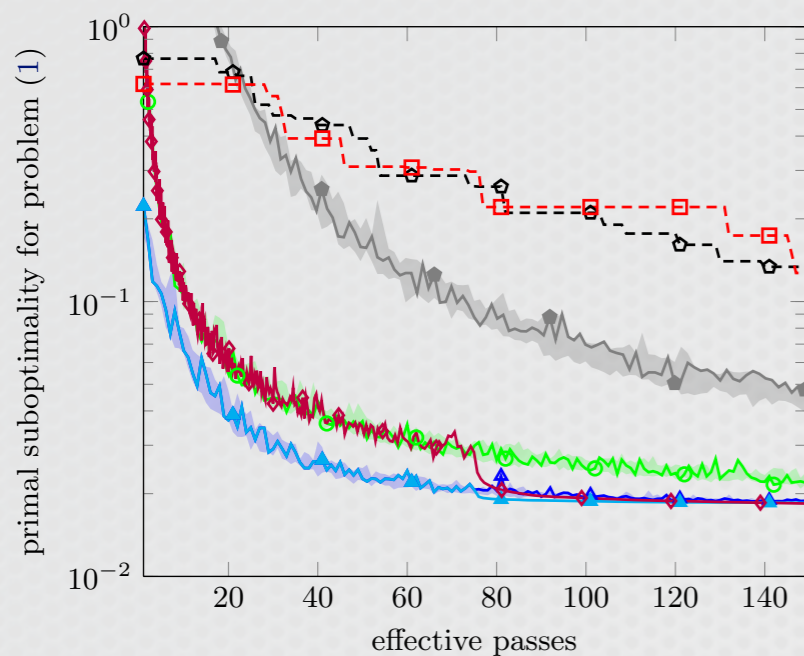
Related Work

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 $R := \max_{i,\mathbf{y}} \|\psi_i(\mathbf{y})\|_2$ (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 .

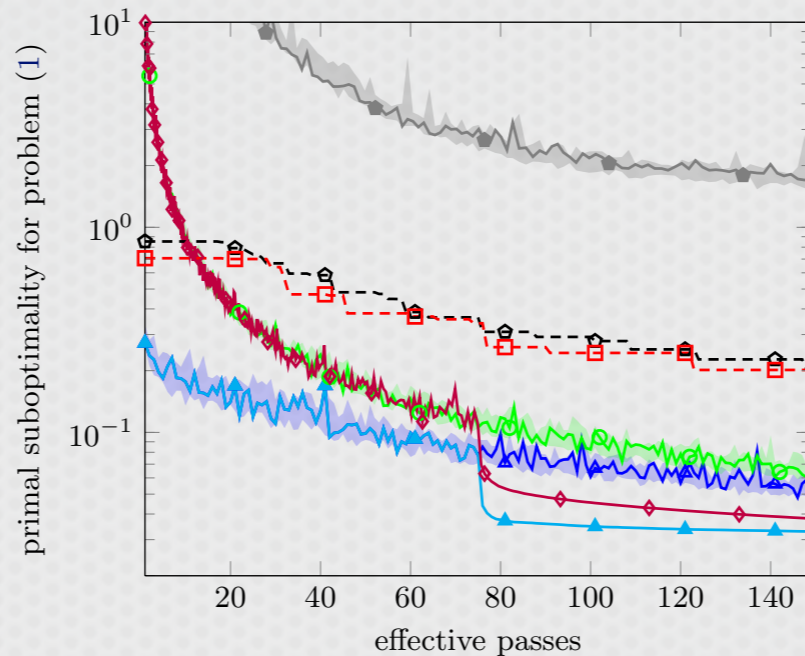
Optimization algorithm	Online	Primal/Dual	Type of guarantee	Oracle type	# Oracle calls
dual extragradient (Taskar et al., 2006)	no	primal-“dual”	saddle point gap	Bregman projection	$O\left(\frac{nR \log \mathcal{Y} }{\lambda\varepsilon}\right)$
online exponentiated gradient (Collins et al., 2008)	yes	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

Experimental Results (w/ averaging)

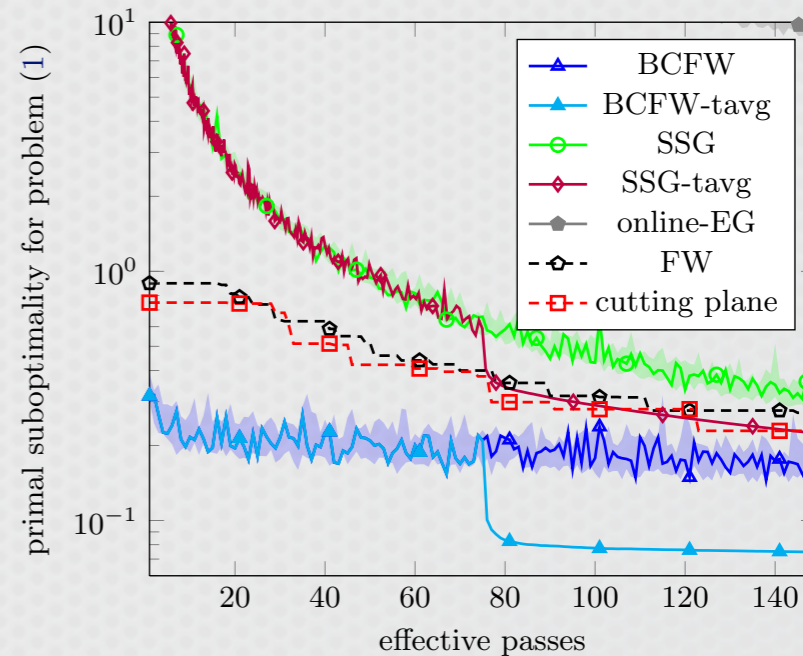
dataset		n	d
OCR	sequence labeling	6251	4028
CoNLL	POS sequence labeling	8936	1643026
Matching	word alignment	5000	82



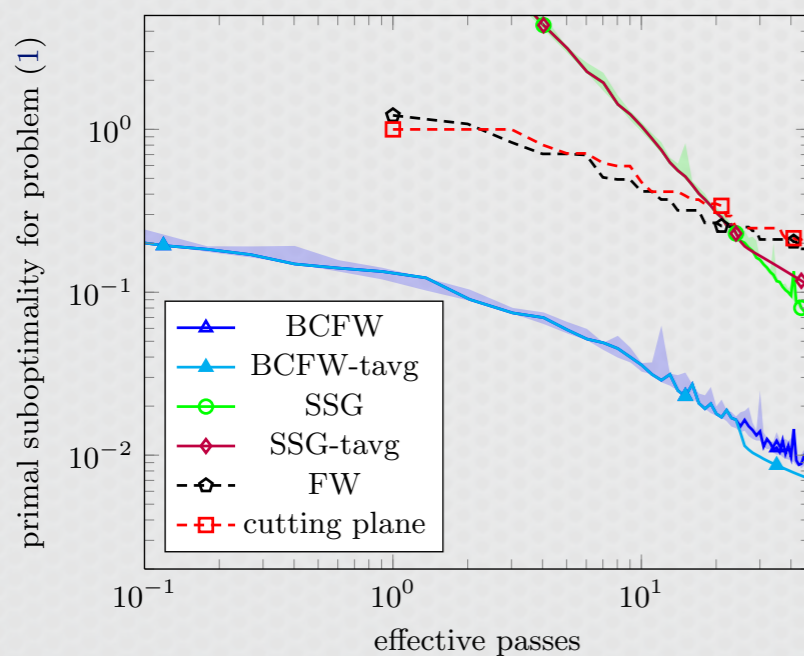
(a) OCR dataset, $\lambda = 0.01$.



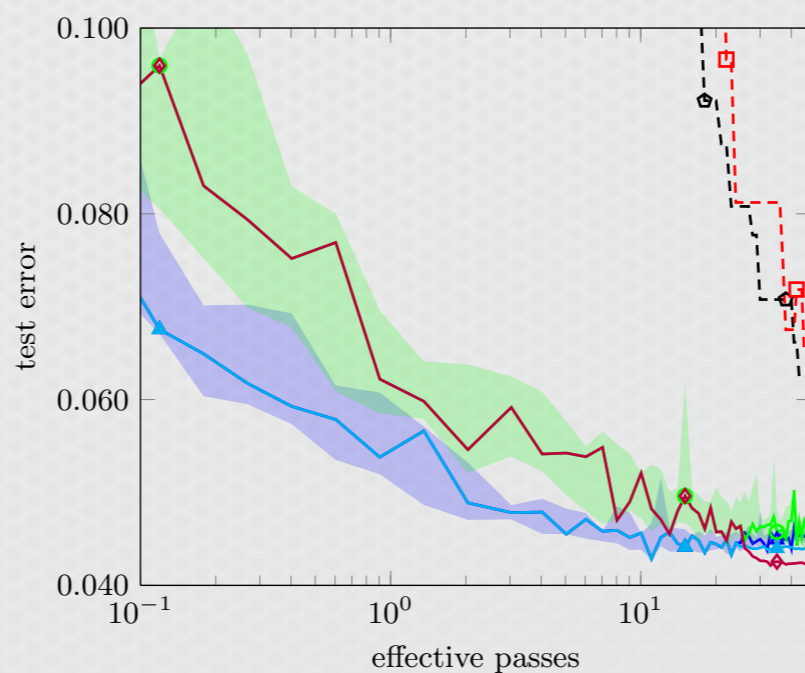
(b) OCR dataset, $\lambda = 0.001$.



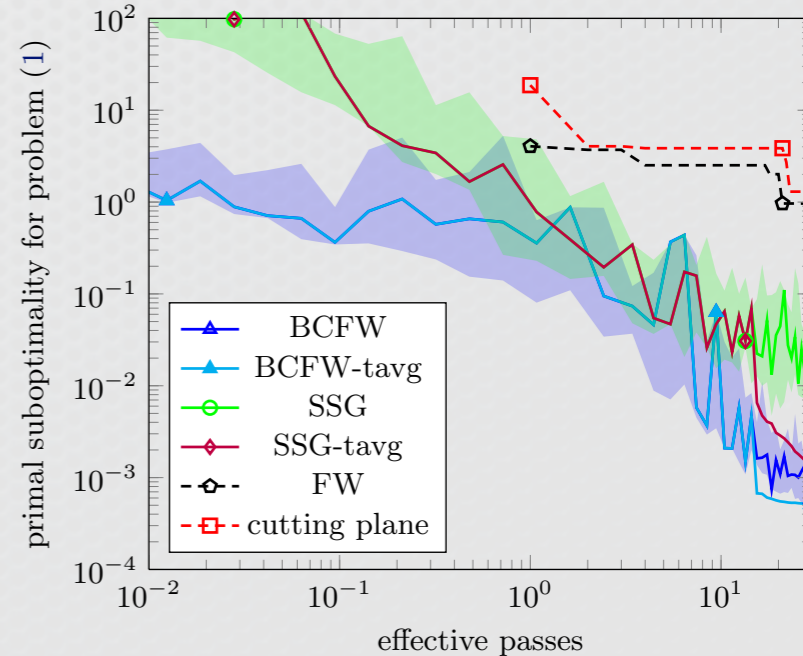
(c) OCR dataset, $\lambda = 1/n$.



(d) CoNLL dataset, $\lambda = 1/n$.



(e) Test error for $\lambda = 1/n$ on CoNLL.



(f) Matching dataset, $\lambda = 0.001$.

Frank-Wolfe: History & Related Work

	Domain	Known Stepsize	Approx. Subproblem	Primal-Dual Guarantee
Frank & Wolfe 1956	linear inequality constraints	✗	✗	✗
Dunn 1978, 1980	general bounded convex domain	✗	✓	✗
Zhang 2003	convex hulls	✗	✓	✗
Clarkson 2008, 2010	unit simplex	✓	✗	✓
Hazan 2008	semidefinite matrices of bounded trace	✓	✓	✓
J. PhD Thesis	general bounded convex domain	✓	✓	✓