

# Improved Regret Guarantees for OCO in Bandit Setting

Ankan Saha <sup>1</sup>   Ambuj Tewari <sup>2</sup>

AISTATS 2011

---

<sup>1</sup>University of Chicago

<sup>2</sup>University of Texas at Austin

# Introduction and Motivation

- Sequential decision making: important question in machine learning, Economics, Operations Research and related fields.
- Modeled as a sequential game between learner and adversary.

- At every time step  $t$ :

- At every time step  $t$ :
- Player plays a point  $\mathbf{x}_t \in \mathcal{K} \subseteq \mathbb{R}^d$ ,  $\mathcal{K}$  is convex and compact.

- At every time step  $t$ :
- Player plays a point  $\mathbf{x}_t \in \mathcal{K} \subseteq \mathbb{R}^d$ ,  $\mathcal{K}$  is convex and compact.
- Adversary responds with a function  $f_t \in \mathcal{F}$ .
- Player suffers a loss  $f_t(\mathbf{x}_t)$ .

- At every time step  $t$ :
- Player plays a point  $\mathbf{x}_t \in \mathcal{K} \subseteq \mathbb{R}^d$ ,  $\mathcal{K}$  is convex and compact.
- Adversary responds with a function  $f_t \in \mathcal{F}$ .
- Player suffers a loss  $f_t(\mathbf{x}_t)$ .
- *Online Convex Optimization*:  $f_t$  are convex functions.

- At every time step  $t$ :
- Player plays a point  $\mathbf{x}_t \in \mathcal{K} \subseteq \mathbb{R}^d$ ,  $\mathcal{K}$  is convex and compact.
- Adversary responds with a function  $f_t \in \mathcal{F}$ .
- Player suffers a loss  $f_t(\mathbf{x}_t)$ .
- *Online* Convex Optimization:  $f_t$  are convex functions.
- Full information vs. Bandit Setting.

- Goal: To minimize the **Regret** *i.e.* the player's performance with respect to the best performance in hindsight.

$$R_T = \sum_{t=1}^T f_t(\mathbf{x}_t) - \min_{\mathbf{x}^* \in \mathcal{K}} \sum_{t=1}^T f_t(\mathbf{x}^*)$$



- Goal: To minimize the **Regret** *i.e.* the player's performance with respect to the best performance in hindsight.

$$R_T = \sum_{i=1}^T f_t(\mathbf{x}_t) - \min_{\mathbf{x}^* \in \mathcal{K}} \sum_{t=1}^T f_t(\mathbf{x}^*)$$

- Full information Setting: Zinkevich's algorithm.

$$\mathbf{x}_{t+1} = \mathbf{Proj}_{\mathcal{K}} \{ \mathbf{x}_t - \eta_t \nabla f_t(\mathbf{x}_t) \}$$

- Goal: To minimize the **Regret** *i.e.* the player's performance with respect to the best performance in hindsight.

$$R_T = \sum_{i=1}^T f_t(\mathbf{x}_t) - \min_{\mathbf{x}^* \in \mathcal{K}} \sum_{t=1}^T f_t(\mathbf{x}^*)$$

- Full information Setting: Zinkevich's algorithm.

$$\mathbf{x}_{t+1} = \mathbf{Proj}_{\mathcal{K}} \{ \mathbf{x}_t - \eta_t \nabla f_t(\mathbf{x}_t) \}$$

- Projected gradient descent strategy incurs  $O(\sqrt{T})$  regret.
- Zinkevich's algorithm just requires information about the gradients of  $f_t$  at every time step.

# Zinkevich's Algorithm

- Key lies in the fact that we minimize the regret and not the actual loss.

## Theorem

If the updates are given by  $\mathbf{x}_{t+1} = \mathbf{Proj}_{\mathcal{K}}(\mathbf{x}_t - \eta_t \nabla f_t(\mathbf{x}_t))$ , choosing  $\eta_t = t^{-1/2}$  gives the following bound on regret after  $T$  steps

$$R_T \leq \frac{\text{diam}(\mathcal{K})^2 \sqrt{T}}{2} + \left(\sqrt{T} - \frac{1}{2}\right) \|\nabla f\|^2$$

where  $\|\nabla f\| = \max_{\mathbf{x} \in \mathcal{K}, t \in [T]} \|\nabla f_t(\mathbf{x})\|$

# Zinkevich's Algorithm

- Key lies in the fact that we minimize the regret and not the actual loss.

## Theorem

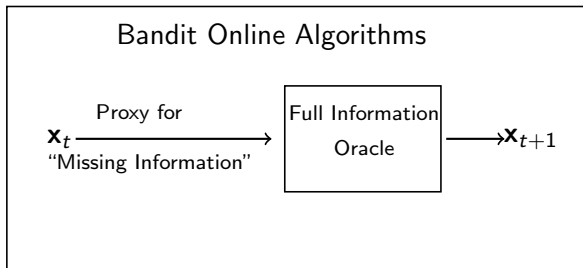
If the updates are given by  $\mathbf{x}_{t+1} = \mathbf{Proj}_{\mathcal{K}}(\mathbf{x}_t - \eta_t \nabla f_t(\mathbf{x}_t))$ , choosing  $\eta_t = t^{-1/2}$  gives the following bound on regret after  $T$  steps

$$R_T \leq \frac{\text{diam}(\mathcal{K})^2 \sqrt{T}}{2} + \left(\sqrt{T} - \frac{1}{2}\right) \|\nabla f\|^2$$

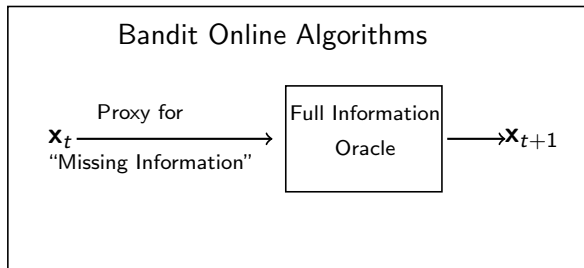
where  $\|\nabla f\| = \max_{\mathbf{x} \in \mathcal{K}, t \in [T]} \|\nabla f_t(\mathbf{x})\|$

- Convergence rate depends on  $\|\nabla f\|$ .

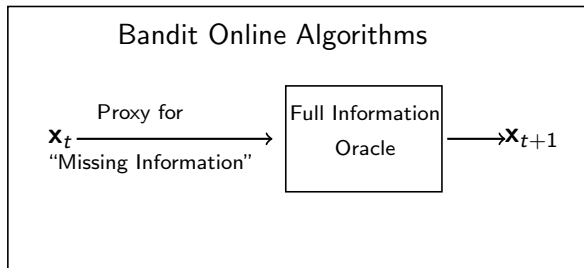
# Bandit Online Optimization



# Bandit Online Optimization



- Missing information  $\implies$  Gradient of the function.



- Missing information  $\implies$  Gradient of the function.
- How to evaluate the gradient from a single point evaluation?

# Unbiased gradient estimation

- [FKM05] provides such a scheme.
- Introduce randomness.

$$\hat{f}(\mathbf{x}) = \mathbb{E}_{\mathbf{v} \in \mathbb{B}^d} [f(\mathbf{x} + \delta \mathbf{v})]$$



# Unbiased gradient estimation

- [FKM05] provides such a scheme.
- Introduce randomness.

$$\hat{f}(\mathbf{x}) = \mathbb{E}_{\mathbf{v} \in \mathbb{B}^d} [f(\mathbf{x} + \delta \mathbf{v})]$$

- Then

$$\nabla \hat{f}(\mathbf{x}) = \mathbb{E}_{\mathbf{u} \in \mathbb{S}^d} \left[ \frac{d}{\delta} f(\mathbf{x} + \delta \mathbf{u}) \cdot \mathbf{u} \right]$$

# Problems with existing approach

- Key trick in [FKM05] is to choose the size of the ball for evaluating the unbiased gradient estimate.

# Problems with existing approach

- Key trick in [FKM05] is to choose the size of the ball for evaluating the unbiased gradient estimate.
- Too small  $\delta$  blows up the  $\|\nabla f\|$ , depending on the location of  $\mathbf{x}_t$  and weakens the bounds.

# Problems with existing approach

- Key trick in [FKM05] is to choose the size of the ball for evaluating the unbiased gradient estimate.
- Too small  $\delta$  blows up the  $\|\nabla f\|$ , depending on the location of  $\mathbf{x}_t$  and weakens the bounds.
- Too large  $\delta \implies$  Gradient estimates are not very accurate.

# Problems with existing approach

- Key trick in [FKM05] is to choose the size of the ball for evaluating the unbiased gradient estimate.
- Too small  $\delta$  blows up the  $\|\nabla f\|$ , depending on the location of  $\mathbf{x}_t$  and weakens the bounds.
- Too large  $\delta \implies$  Gradient estimates are not very accurate.
- Trading off the ball size along with Zinkevich's scheme gives  $O(T^{3/4})$  regret in the OCO bandit setting.

# Convex Analysis Basics

- Convex functions: Lower Bounded by linear approximation.

$$f(\mathbf{x}) \geq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$$

- Convex functions: Lower Bounded by linear approximation.

$$f(\mathbf{x}) \geq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$$

- Strongly Convex functions: Lower Bounded by a quadratic.  
Eg. 2-norm squared.

$$f(\mathbf{x}) \geq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{\sigma}{2} \|\mathbf{x} - \mathbf{y}\|^2$$



- Convex functions: Lower Bounded by linear approximation.

$$f(\mathbf{x}) \geq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$$

- Strongly Convex functions: Lower Bounded by a quadratic.  
Eg. 2-norm squared.

$$f(\mathbf{x}) \geq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{\sigma}{2} \|\mathbf{x} - \mathbf{y}\|^2$$

- *l.c.g* functions: Upper Bounded by a quadratic.

$$f(\mathbf{x}) \leq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2$$

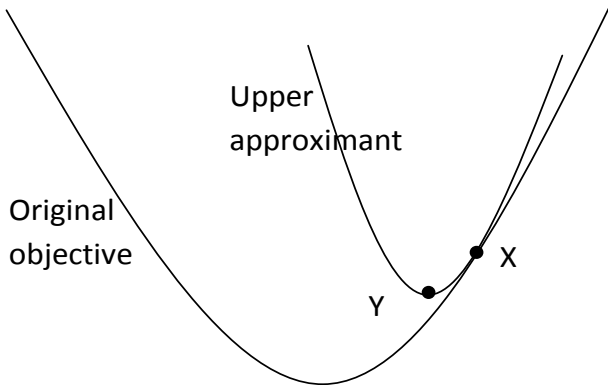
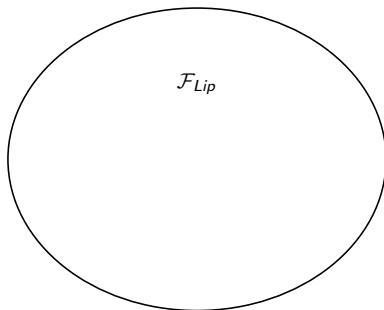


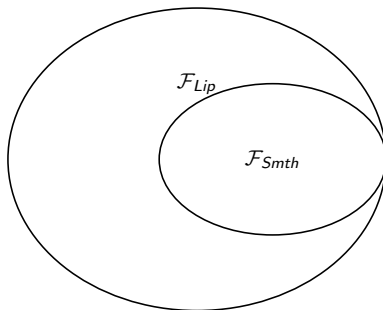
Figure: *l.c.g* functions

# Results for Subclasses of Convex functions(Full Information Setting)



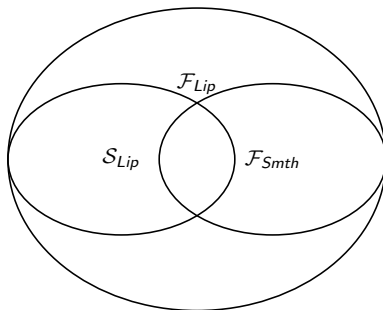
Function Class	Best Rates on Regret
Lipschitz ( $\mathcal{F}_{Lip}$ )	$\Theta(T^{1/2})$ [Zin03]

# Results for Subclasses of Convex functions(Full Information Setting)



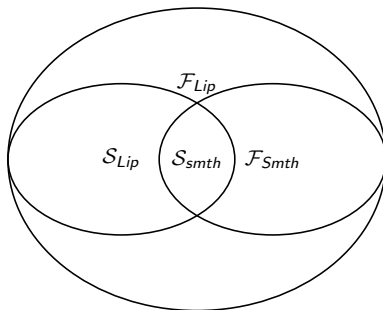
Function Class	Best Rates on Regret
Lipschitz ( $\mathcal{F}_{Lip}$ )	$\Theta(T^{1/2})$ [Zin03]
Smooth ( $\mathcal{F}_{Smth}$ )	$\Theta(T^{1/2})$ [Zin03]

# Results for Subclasses of Convex functions(Full Information Setting)



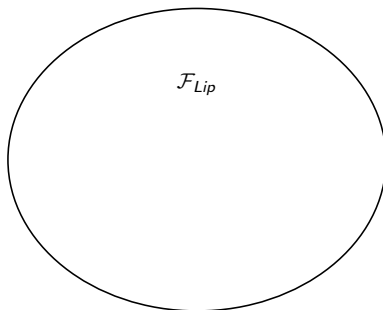
Function Class	Best Rates on Regret
Lipschitz ( $\mathcal{F}_{Lip}$ )	$\Theta(T^{1/2})$ [Zin03]
Smooth ( $\mathcal{F}_{Smth}$ )	$\Theta(T^{1/2})$ [Zin03]
Lipschitz and Strongly Convex ( $\mathcal{S}_{Lip}$ )	$\Theta(\log T)$ [HAK07]

# Results for Subclasses of Convex functions(Full Information Setting)



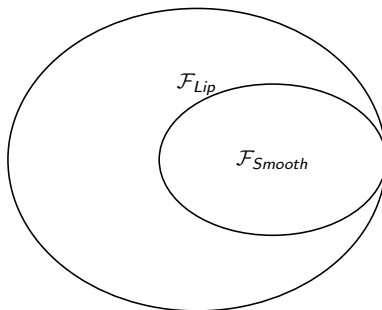
Function Class	Best Rates on Regret
Lipschitz ( $\mathcal{F}_{Lip}$ )	$\Theta(T^{1/2})$ [Zin03]
Smooth ( $\mathcal{F}_{Smth}$ )	$\Theta(T^{1/2})$ [Zin03]
Lipschitz and Strongly Convex ( $\mathcal{S}_{Lip}$ )	$\Theta(\log T)$ [HAK07]
Smooth and Strongly Convex ( $\mathcal{S}_{Smth}$ )	$\Theta(\log T)$ [HAK07]

# Results for Bandit Setting



Function Class	Best Rates on Regret	Lower Bounds
Lipschitz ( $\mathcal{F}_{Lip}$ )	$O(T^{3/4})$ [FKM05]	$\Omega(T^{1/2})$

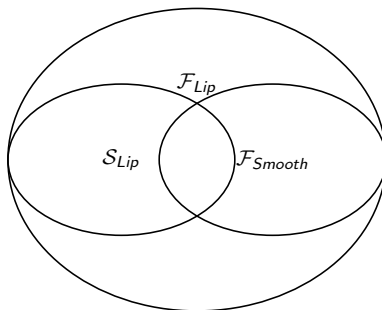
# Results for Bandit Setting



Function Class	Best Rates on Regret	Lower Bounds
Lipschitz ( $\mathcal{F}_{Lip}$ )	$O(T^{3/4})$ [FKM05]	$\Omega(T^{1/2})$
<b>Smooth (<math>\mathcal{F}_{Smooth}</math>)</b>	<b><math>O^*(T^{2/3})</math> [This paper]</b>	$\Omega(T^{1/2})$

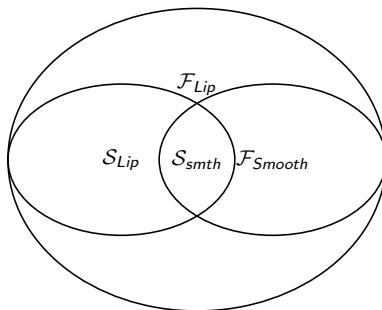


# Results for Bandit Setting



Function Class	Best Rates on Regret	Lower Bounds
Lipschitz ( $\mathcal{F}_{Lip}$ )	$O(T^{3/4})$ [FKM05]	$\Omega(T^{1/2})$
<b>Smooth</b> ( $\mathcal{F}_{Smooth}$ )	$O^*(T^{2/3})$ <b>[This paper]</b>	$\Omega(T^{1/2})$
Lipschitz and S.C. ( $\mathcal{S}_{Lip}$ )	$O^*(T^{2/3})$ [ADX10]	$\Omega(\log T)$

# Results for Bandit Setting



Function Class	Best Rates on Regret	Lower Bounds
Lipschitz ( $\mathcal{F}_{Lip}$ )	$O(T^{3/4})$ [FKM05]	$\Omega(T^{1/2})$
<b>Smooth</b> ( $\mathcal{F}_{Smth}$ )	$O^*(T^{2/3})$ <b>[This paper]</b>	$\Omega(T^{1/2})$
Lipschitz and S.C. ( $S_{Lip}$ )	$O^*(T^{2/3})$ [ADX10]	$\Omega(\log T)$
Smooth and S.C. ( $S_{Smth}$ )	$O^*(T^{2/3})$ [ADX10]	$\Omega(\log T)$

# New Approach

- Tighter bounds obtained by calculating  $\|\nabla f\|$  with respect to a changing *local* norm, demonstrated for Bandit OLO [AHR08].

## Definition

A function  $R : \mathcal{K} \rightarrow \mathbb{R}$  is a self-concordant barrier if

- a)  $R \rightarrow \infty$  near the boundary of  $\mathcal{K}$  and
- b)  $R$  and  $\nabla^2 R$  are Lipschitz continuous with respect to the local norm  $\|\cdot\|_{R,\mathbf{w}}$ , given by  $\|f\|_{R,\mathbf{w}} = \sqrt{\langle f, \nabla^2 R(\mathbf{w})f \rangle}$ .

# New Approach

- Tighter bounds obtained by calculating  $\|\nabla f\|$  with respect to a changing *local* norm, demonstrated for Bandit OLO [AHR08].

## Definition

A function  $R : \mathcal{K} \rightarrow \mathbb{R}$  is a self-concordant barrier if

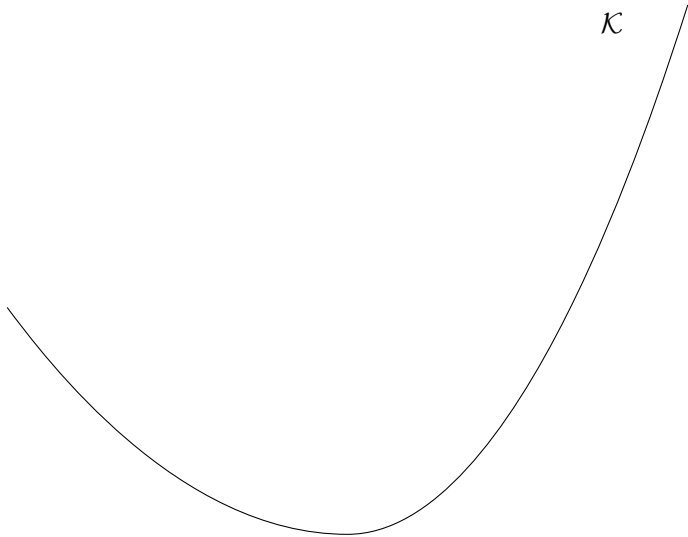
- a)  $R \rightarrow \infty$  near the boundary of  $\mathcal{K}$  and
- b)  $R$  and  $\nabla^2 R$  are Lipschitz continuous with respect to the local norm  $\|\cdot\|_{R,\mathbf{w}}$ , given by  $\|f\|_{R,\mathbf{w}} = \sqrt{\langle f, \nabla^2 R(\mathbf{w})f \rangle}$ .

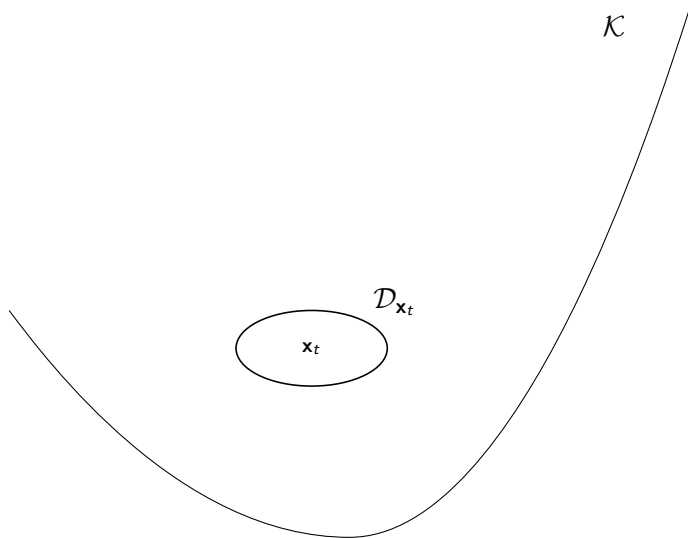
- Given a self concordant barrier  $R$ , for every  $\mathbf{w} \in \mathcal{K}$ , the Dikin ellipsoid centered at  $\mathbf{w}$

$$\mathcal{D}_{\mathbf{w}} = \left\{ \mathbf{w}' : \|\mathbf{w} - \mathbf{w}'\|_{R,\mathbf{w}} \leq 1 \right\}$$

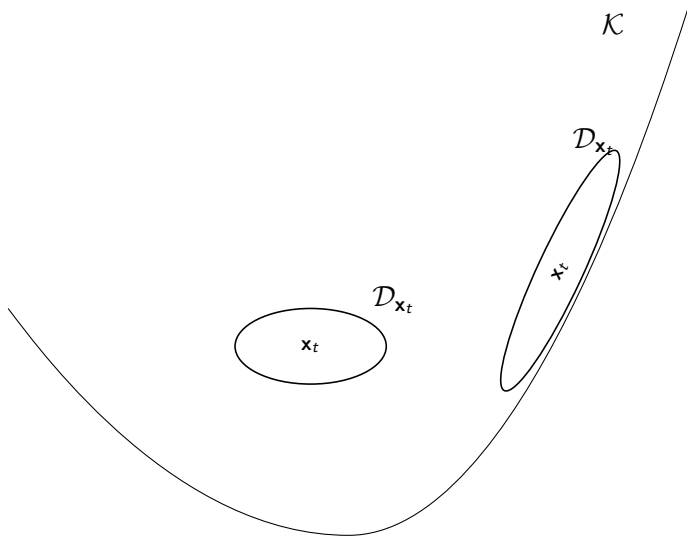
is always contained in  $\mathcal{K}$ .

$\kappa$

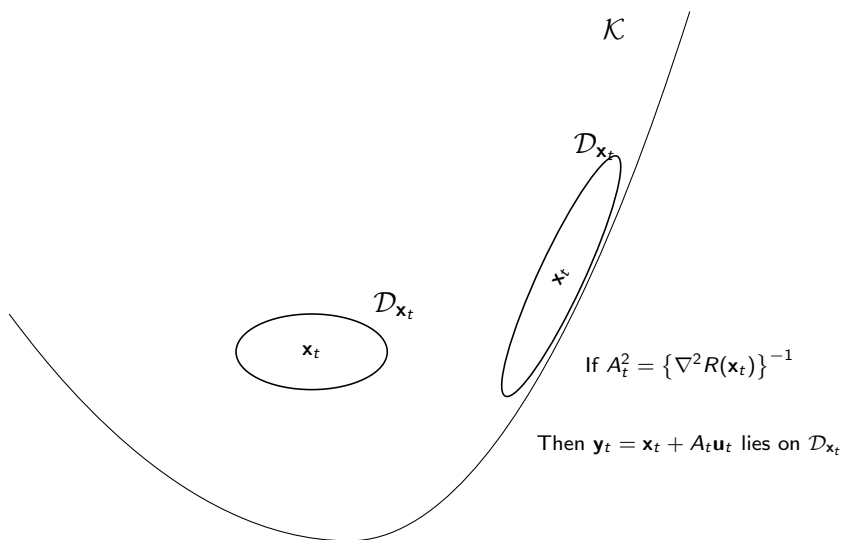




- Dikin ellipsoid always lies inside  $\mathcal{K}$ .



- Dikin ellipsoid always lies inside  $\mathcal{K}$ .
- $R$  needs to curve strongly near the boundary.



- Dikin ellipsoid always lies inside  $\mathcal{K}$ .
- $R$  needs to curve strongly near the boundary.



# Combining Ideas

- Combine single point gradient estimate with the idea of self concordant barriers.

# Combining Ideas

- Combine single point gradient estimate with the idea of self concordant barriers.
- Removes the need for projections.

# Combining Ideas

- Combine single point gradient estimate with the idea of self concordant barriers.
- Removes the need for projections.
- Generate gradient estimate  $\mathbf{g}_t$  at step  $t$ .
- Feed it to the full information algorithm blackbox from [AHR08] to obtain

$$\mathbf{x}_{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}} \sum_{i=1}^t \eta \langle \mathbf{g}_i, \mathbf{x} \rangle + R(\mathbf{x})$$

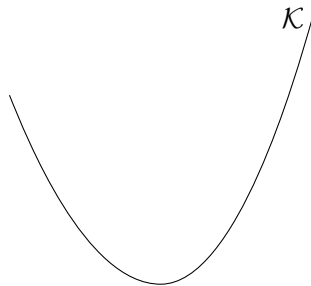
# Algorithm

- Given  $R$ : a self concordant barrier for  $\mathcal{K}$ .

# Algorithm

- Given  $R$ : a self concordant barrier for  $\mathcal{K}$ .
- Pick  $\mathbf{x}_1 \in \mathcal{K}$ .

At every step  $t = 1, 2, 3, \dots$ ,

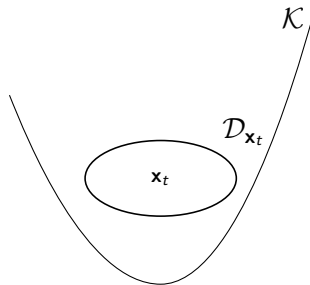


# Algorithm

- Given  $R$ : a self concordant barrier for  $\mathcal{K}$ .
- Pick  $\mathbf{x}_1 \in \mathcal{K}$ .

At every step  $t = 1, 2, 3, \dots$ ,

- Evaluate  $A_t = (\nabla^2 R(\mathbf{x}_t))^{-1/2}$ .

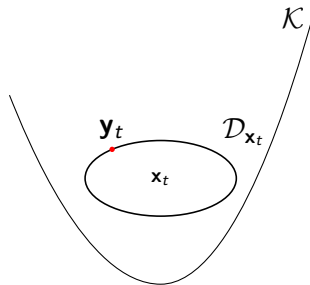


# Algorithm

- Given  $R$ : a self concordant barrier for  $\mathcal{K}$ .
- Pick  $\mathbf{x}_1 \in \mathcal{K}$ .

At every step  $t = 1, 2, 3, \dots$ ,

- Evaluate  $A_t = (\nabla^2 R(\mathbf{x}_t))^{-1/2}$ .
  - Sample  $\mathbf{u}_t \in \mathbb{S}^d$ .
- $$\mathbf{y}_t = \mathbf{x}_t + \delta A_t \mathbf{u}_t.$$

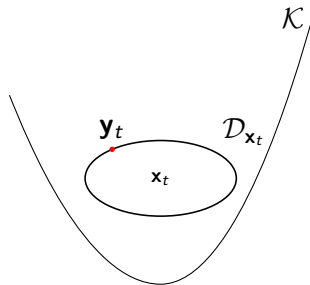


# Algorithm

- Given  $R$ : a self concordant barrier for  $\mathcal{K}$ .
- Pick  $\mathbf{x}_1 \in \mathcal{K}$ .

At every step  $t = 1, 2, 3, \dots$ ,

- Evaluate  $A_t = (\nabla^2 R(\mathbf{x}_t))^{-1/2}$ .
- Sample  $\mathbf{u}_t \in \mathbb{S}^d$ .  
 $\mathbf{y}_t = \mathbf{x}_t + \delta A_t \mathbf{u}_t$ .
- Player plays  $\mathbf{y}_t$  and receives loss  $f_t(\mathbf{y}_t)$ .



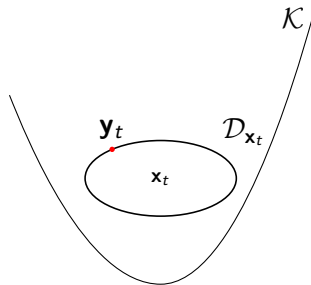


# Algorithm

- Given  $R$ : a self concordant barrier for  $\mathcal{K}$ .
- Pick  $\mathbf{x}_1 \in \mathcal{K}$ .

At every step  $t = 1, 2, 3, \dots$ ,

- Evaluate  $A_t = (\nabla^2 R(\mathbf{x}_t))^{-1/2}$ .
- Sample  $\mathbf{u}_t \in \mathbb{S}^d$ .  
 $\mathbf{y}_t = \mathbf{x}_t + \delta A_t \mathbf{u}_t$ .
- Player plays  $\mathbf{y}_t$  and receives loss  $f_t(\mathbf{y}_t)$ .
- $\mathbf{g}_t = \frac{d}{\delta} f_t(\mathbf{y}_t) A_t^{-1} \mathbf{u}_t$ .

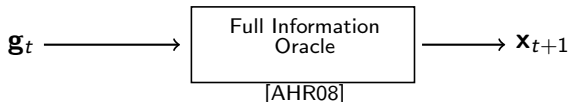
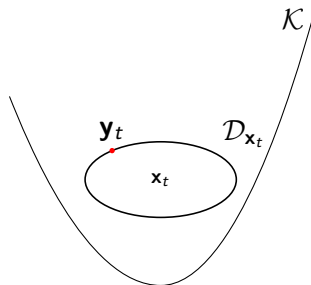


# Algorithm

- Given  $R$ : a self concordant barrier for  $\mathcal{K}$ .
- Pick  $\mathbf{x}_1 \in \mathcal{K}$ .

At every step  $t = 1, 2, 3, \dots$ ,

- Evaluate  $A_t = (\nabla^2 R(\mathbf{x}_t))^{-1/2}$ .
- Sample  $\mathbf{u}_t \in \mathbb{S}^d$ .  
 $\mathbf{y}_t = \mathbf{x}_t + \delta A_t \mathbf{u}_t$ .
- Player plays  $\mathbf{y}_t$  and receives loss  $f_t(\mathbf{y}_t)$ .
- $\mathbf{g}_t = \frac{d}{\delta} f_t(\mathbf{y}_t) A_t^{-1} \mathbf{u}_t$ .



- Define the approximation

$$\hat{f}_t(\mathbf{x}) = \mathbb{E}_{\mathbf{v} \in \mathbb{B}^d} [f_t(\mathbf{x} + \delta A_t \mathbf{v}_t)]$$

- Regret written as a telescopic sum.
- Using approximation + Randomization incurs  $O(T\delta^2)$  cost.
- [AHR08] blackbox incurs  $O^*(\sqrt{T}/\delta)$ .
- Trading off  $\delta$  gives  $O^*(T^{2/3})$  regret.

# Conclusions and Future Work

- Improve bounds on regret for bandit OCO from  $O(T^{3/4})$  to  $O^*(T^{2/3})$ , when  $f_t$  are smooth.

# Conclusions and Future Work

- Improve bounds on regret for bandit OCO from  $O(T^{3/4})$  to  $O^*(T^{2/3})$ , when  $f_t$  are smooth.
- Lower Bounds still correspond to full information setting.

# Conclusions and Future Work

- Improve bounds on regret for bandit OCO from  $O(T^{3/4})$  to  $O^*(T^{2/3})$ , when  $f_t$  are smooth.
- Lower Bounds still correspond to full information setting.
- How to bridge the gap?

# Conclusions and Future Work

- Improve bounds on regret for bandit OCO from  $O(T^{3/4})$  to  $O^*(T^{2/3})$ , when  $f_t$  are smooth.
- Lower Bounds still correspond to full information setting.
- How to bridge the gap?
- Better bounds on “price of bandit information” for more specific function classes?

# Conclusions and Future Work

- Improve bounds on regret for bandit OCO from  $O(T^{3/4})$  to  $O^*(T^{2/3})$ , when  $f_t$  are smooth.
- Lower Bounds still correspond to full information setting.
- How to bridge the gap?
- Better bounds on “price of bandit information” for more specific function classes?

Thank You!





Alekh Agarwal, Ofer Dekel, and Lin Xiao.

Optimal algorithms for online convex optimization with multi-point bandit feedback, 2010.

longer version available at

<http://www.cs.berkeley.edu/~alekh/bandit-colt.pdf>.



Jacob Abernethy, Elad Hazan, and Alexander Rakhlin.

Competing in the dark: An efficient algorithm for bandit linear optimization.

In *COLT*, pages 263–274, 2008.



Abraham Flaxman, Adam Tauman Kalai, and H. Brendan McMahan.

Online convex optimization in the bandit setting: gradient descent without a gradient.

In *SODA*, pages 385–394, 2005.



Elad Hazan, Amit Agarwal, and Satyen Kale.

Logarithmic regret algorithms for online convex optimization.

*Machine Learning*, 69(2-3):169–192, 2007.



M. Zinkevich.

Online convex programming and generalised infinitesimal gradient ascent.

In *ICML*, pages 928–936, 2003.