Algorithms for sustainable data centers

Adam Wierman (Caltech) Minghong Lin (Caltech) Zhenhua Liu (Caltech) Lachlan Andrew (Swinburne) and many others...











This talk: An overview of algorithmic questions & challenges

Two examples
 A general model & some results
 A case study

Example 1: Dynamic resizing in a data center





Goal: Adapt provisioning to minimize cost. **Challenges:** Switching is costly & prediction is hard.

Example 1: Dynamic resizing in a data center





Example 1: Dynamic resizing in a data center





Example 2: Geographical load balancing





<u>Data centers:</u>

Diverse, time varying electricity prices, renewable availabilities, cooling efficiencies, etc.



Example 2: Geographical load balancing



Goal: Adapt routing & provisioning to minimize cost. **Challenges:** Switching is costly & prediction is hard.

Dynamic resizing





[Chase et al. 01] [Pinheiro et al. 01] ...[Chen et al. 05] ...[Ghandi et al. 09] [Khargharia et al. 10] [Kansal et al. 10] ...[LWAT10] ...



[Pakbaznia et al. 09] [Qureshi et al. 09] [Rao et al. 10] [Stanojevic et al. 10] [Wendell et al 10] [Le et al 2010] [LLWLA 11] [LLWAL11] [LAW12], ...



A (familiar) model





Data center goal: $\min_{x^t \in F} \sum_{t} (\text{operating cost}) + (\text{switching cost})$ (stability constraints)

```
Data center goal:

\min_{x^t \in F} \sum_{t} (\text{operating cost}) + (\text{switching cost})

stability constraints

\frac{x_t \in \mathbb{N}}{N_{max} \ge x_t \ge 0}

x_t \ge SLA(\lambda_t)

...
```

Data center goal: $\min_{x^t \in F} \sum_t (\text{operating cost}) + (\text{switching cost})$



Data center goal: $\min_{x^{t} \in F} \sum_{t} c^{t}(x^{t}) + ||x^{t} - x^{t-1}||$ $\sum_{t} \sum_{t} c^{t}(x^{t}) + ||x^{t} - x^{t-1}||$



Smoothed online convex optimization (SOCO)





SOCO comes up in many applications ...

geographical load balancing dynamic capacity management video streaming electricity generation planning product/service selection portfolio management labor markets penalized estimation







Can an algorithm maintain sub-linear regret? [KV02] [BBCM03] [Z03] ... [HAK07]... [<u>LRAMW12</u>]

Competitive ratio(Alg) = $\frac{\text{Cost(Alg)}}{\text{Cost(Offline_Opt)}}$

Can an algorithm maintain sub-linear regret? [KV02] [BBCM03] [Z03] ... [HAK07]... [LRAMW12]

Can an algorithm maintain a constant competitive ratio? [BKRS92] [BLS92]...[BB00]...[LWAT11][LRAMW12]

Can an algorithm maintain sub-linear regret and a constant competitive ratio?





$$x^{1}, c^{1}, x^{2}, c^{2}, x^{3}, c^{3}, ...$$

$$\lim_{x^{t} \in F} \sum_{t} c^{t}(x^{t})$$

$$\|\nabla c^{t}(\cdot)\|_{2} < C$$
Goal: Algorithms with sub-linear regret
$$name{formula}$$
Online gradient descent (0GD) [Z03]
$$x^{t+1} = \operatorname{Proj}(x^{t} - \eta_{t} \nabla c^{t}(x^{t}))$$

<u>Theorem [Z03]</u> When $\eta_t = \Theta(1/\sqrt{t})$, online gradient descent has $O(\sqrt{T})$ -regret for OCO.

Theorem [Z03,H06] Any algorithm for OCO must incur $\Omega(\sqrt{T})$ -regret on linear cost functions.

 $\begin{aligned} x^{1}, c^{1}, x^{2}, c^{2}, x^{3}, c^{3}, \dots \\ \min_{x^{t} \in F} \sum_{t} c^{t}(x^{t}) \\ \|\nabla c^{t}(\cdot)\|_{2} < C \end{aligned}$ **Online convex optimization** Goal: Algorithms with sub-linear regret - Online gradient descent (OGD) [Z03] - Multiplicative weights Many algorithms achieve sub-linear regret [AHK05] [FS99] - Newton's method based [HKKA06][HAK07]

 $\begin{aligned} x^{1}, c^{1}, x^{2}, c^{2}, x^{3}, c^{3}, \dots \\ \min_{x^{t} \in F} \sum_{t} c^{t}(x^{t}) \\ \| \nabla c^{t}(\cdot) \|_{2} < C \end{aligned}$ Goal: Algorithms with sub-linear regret

Online convex optimization

Do OCO algorithms maintain sub-linear regret for SOCO?

 $x^{1}, c^{1}, x^{2}, c^{2}, x^{3}, c^{3}, \dots$ $\min_{x^{t} \in F} \sum_{t} c^{t}(x^{t}) + \|x^{t} - x^{t-1}\|$ $\|\nabla c^{t}(\cdot)\|_{2} < C$ smoothed
Smoo Goal: Algorithms with sub-linear regret

Do OCO algorithms maintain sub-linear regret for SOCO?
$\frac{\text{Corollary:}}{\text{When } \eta_t = \Theta(1/\sqrt{t}), \text{ online gradient descent still has } O(\sqrt{T}) \text{-regret for SOCO.}$

Proof:

 $\sum_{t} \|x^{t} - x^{t-1}\| \le M \sum \|x^{t} - x^{t-1}\|_{2}$ (Equivalence of norms) $\leq M \sum \eta_t \| \nabla c^{t-1}(x^{t-1}) \|_2$ (Projection is non-expansive) $\leq MC \sum \eta_t = O(\sqrt{T})$ (Bounded gradient)

Can an algorithm maintain sub-linear regret? <u>Yes</u>: Online gradient descent has $O(\sqrt{T})$ -regret.

Can an algorithm maintain a constant competitive ratio? [BKRS92] [BLS92]...[BB00]...[LWAT11][LRAMW12]

Can an algorithm maintain sub-linear regret and a constant competitive ratio?







 $\frac{n \text{ is the size of the action set}}{\sum_{n \in \mathbb{N}}}$ $\frac{1}{\sum_{n \in \mathbb{N}}} = \frac{1}{\sum_{n \in \mathbb{N}}} = \frac{1}{\sum_{n \in \mathbb{N}}}$ Any deterministic algorithm is $\Omega(n)$ -competitive.
Any randomized algorithm is $\Omega(\sqrt{\log n}/\log \log n)$ -competitive







"Lazy Capacity Provisioning" (LCP) is 3-competitive. Further $Cost(LCP) \le Cost(OPT) + 2 SwitchingCost(OPT)$





Theorem:

"Lazy Capacity Provisioning" (LCP) is 3-competitive. Further $Cost(LCP) \le Cost(OPT) + 2 SwitchingCost(OPT)$

<u>Lemma:</u>

The offline optimal solution is "lazy in reverse time".







Classic Approach: Receding Horizon Control

Classic Approach: Receding Horizon Control



Receding Horizon Control (RHC):

Choose x_t to minimize cost over [t, t + w], given prediction window and x_{t-1} .



Receding Horizon Control (RHC):

Choose x_t to minimize cost over [t, t + w], given prediction window and x_{t-1} .





<u>Averaging Fixed Horizon Control (AFHC):</u>

Choose x_t as the average of w + 1 fixed horizon control algorithms.





1 Can an algorithm maintain sub-linear regret? <u>Yes</u>: Online gradient descent has $\Theta(\sqrt{T})$ -regret.

Can an algorithm maintain a constant competitive ratio? Yes. (in the scalar case): LCP is 3-competitive. Yes. (in the vector case): AFHC is O(||1||/w)-comp

3 Can an algorithm maintain sub-linear regret and a constant competitive ratio? <u>NO.</u>







- Is ϵT -regret "good enough"?
- ls log log *T* competitive "good enough"?
- What about under stochastic assumptions?

<u>Theorem:</u>

Given a 1-dimensional SOCO, for arbitrary, increasing f(T), "Randomly Biased Greedy" (RBG) is: (1 + f(T)/||1||)-competitive with $O(\max\{T/f(T), f(T)\})$ -regret Randomly Biased Greedy (f(T)) $x^{t} = \operatorname{argmin} w^{t}(x^{t}) + x^{t}r$ Where: $\|1\|_{RBG} = f(T)$. $w^{t}(x) = \min_{y} \{w^{t-1}(y) + c^{t}(y) + \|x - y\|_{RBG}\}$ $r = \operatorname{UnifRand}(-\|1\|_{RBG}, \|1\|_{RBG})$

<u>Theorem:</u>

Given a 1-dimensional SOCO, for arbitrary, increasing f(T), "Randomly Biased Greedy" (RBG) is: (1 + f(T)/||1||)-competitive with $O(\max\{T/f(T), f(T)\})$ -regret **1** Can an algorithm maintain sub-linear regret? <u>Yes</u>: Online gradient descent has $\Theta(\sqrt{T})$ -regret.

Can an algorithm maintain a constant competitive ratio? <u>Yes. (in the scalar case)</u>: LCP is 3-competitive. <u>Yes. (in the vector case)</u>: AFHC is O(||1||/w)-comp

Can an algorithm maintain sub-linear regret and a constant competitive ratio? No! But (in the scalar case): RBG is O(f(T))-competitive with $O(\max\{T/f(T), f(T)\})$ -regret.

Back to the applications...

An implementation: Dynamic resizing







THE WALL STREET JOURNAL.	Subscribe now and get 4 weeks free SUBSCRIBE >
Europe Edition Home • Today's Paper • Video • Blogs • Emails • Journal Community • Mobile • Tablet	Subscribe Log In
World • Europe • U.K. • U.S. • Business • Markets • Market Data • Tech • Life & Style •	Opinion • Real Estate • Jobs •
May 30, 2012, 12:58 p.m. ET	
HP Unveils Architecture for First Net Zero Energy Data Center	
Article	
Email Errinter Share: Share:	- Text +
The things that matter to us.	Make it matter.
HP Unveils Architecture for First Net Zero Energy Data Center	
PALO ALTO, CA (MARKETWIRE) 05/30/12	
HP (NYSE: HPQ) today unveiled research from HP Labs, the company's central research arm, that illustrates the architecture for a data center that requires no net energy from traditional power grids.	
The research shows how the architecture, combined with holistic energy-management techniques, enables organizations to cut total power usage by 30 percent, as well as dependence on grid power and costs by more than 80 percent.(1)	
With the HP Net-Zero Energy Data Center research, HP aims to provide businesses and societies around the world the potential to operate data centers using local renewable resources, removing dependencies such as location, energy supply and costs. This opens up the possibility of introducing IT services to organizations of all sizes.	
"Information technology has the power to be an equalizer across societies globally, but the cost of IT services, and by extension the cost of energy, is prohibitive and inhibits widespread adoption," said Cullen Bash, distinguished technologist, HP, and interim director, Sustainable Ecosystems Research Group, HP Labs. "The HP Net-Zero Energy Data Center not only aims to minimize the environmental impact of computing, but also has a goal of reducing energy costs associated with data- center operations to extend the reach of IT accessibility globally."	

HP collaborators: Yuan Chen, Cullen Bash, Martin Arlitt, Daniel Gmach, Zhikui Wang, Manish Marwah and Chris Hyser

An implementation: Dynamic resizing



HP collaborators: Yuan Chen, Cullen Bash, Martin Arlitt, Daniel Gmach, Zhikui Wang, Manish Marwah and Chris Hyser

A case study: Geographical load balancing





<u>Data centers:</u>

Google data center locations Real traces for wind, solar, electricity price, etc.



A case study: Geographical load balancing



<u>GLB:</u> geographical load balancing & dynamic resizing

VS.

LOCAL: Route to the closest data center & dynamic resizing



Follow the renewables routing emerges.



The good

Follow the renewables routing emerges. Huge reductions in grid usage become possible.








Final thoughts





Algorithms for sustainable data centers Algorithms for sustainable data centers

- Minghong Lin, Adam Wierman, Lachlan Andrew, and Eno Thereska. Dynamic right-sizing for power proportional data centers. Infocom, 2011. Best Paper award winner.
- Zhenhua Liu, Minghong Lin, Adam Wierman, Steven Low, and Lachlan Andrew. Greening geographical load balancing. Sigmetrics 2011.
- Zhenhua Liu, Minghong Lin, Adam Wierman, Steven Low, and Lachlan Andrew. Geographical load balancing with renewables. Greenmetrics, 2011. Best Student Paper award winner.
- Zhenhua Liu, Yuan Chen, Cullen Bash, Adam Wierman, et al. Renewable and cooling aware workload management for sustainable data centers. Sigmetrics 2012.
- Minghong Lin, Lachlan Andrew, and Adam Wierman. Online Algorithms for Geographical Load Balancing. Green Computing Conference, 2012. Best Paper award winner.