

MetaGrad: Faster Convergence **Without Curvature** in Online Convex Optimization



Wouter M. Koolen Tim van Erven



Centrum Wiskunde & Informatica



Universiteit Leiden

Inria Lille
Friday 15th April, 2016

This talk



- ▶ Online Convex Optimization
- ▶ Learning the Learning rate
- ▶ Second-order (variance) bounds (individual sequence)
- ▶ Fast rates without curvature

Outline

Online Convex Optimization

A New Type of Guarantee

Fast Rates

MetaGrad Algorithm

Fundamental Learning Model: Online Convex Optimization

- ▶ In round $t = 1, 2, \dots$
 - ▶ Learner predicts \mathbf{w}_t (from unit ball)
 - ▶ Encounter convex loss function $f_t(\mathbf{u}) : \mathbb{R}^d \rightarrow \mathbb{R}$



- ▶ Learner
 - ▶ observes gradient $\mathbf{g}_t := \nabla f_t(\mathbf{w}_t)$ (from unit ball)
 - ▶ incurs loss $f_t(\mathbf{w}_t)$

Fundamental Learning Model: Online Convex Optimization

- ▶ In round $t = 1, 2, \dots$
 - ▶ Learner predicts \mathbf{w}_t (from unit ball)
 - ▶ Encounter convex loss function $f_t(\mathbf{u}) : \mathbb{R}^d \rightarrow \mathbb{R}$



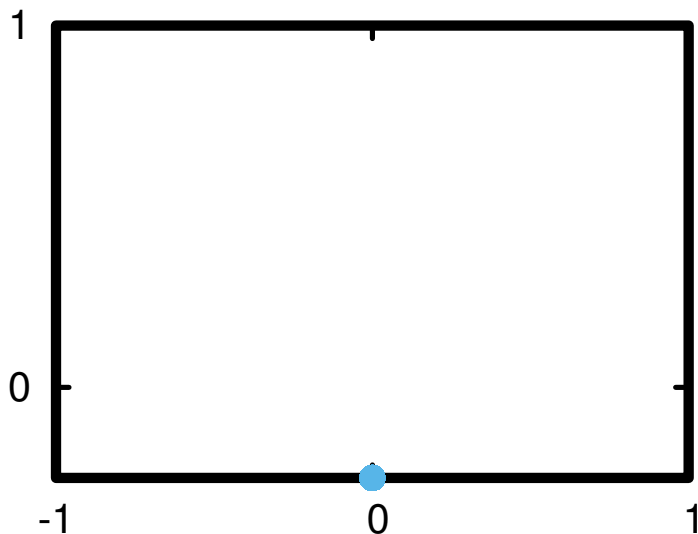
- ▶ Learner
 - ▶ observes gradient $\mathbf{g}_t := \nabla f_t(\mathbf{w}_t)$ (from unit ball)
 - ▶ incurs loss $f_t(\mathbf{w}_t)$
- ▶ The goal is to have small **regret**

$$R_T^u := \underbrace{\sum_{t=1}^T f_t(\mathbf{w}_t)}_{\text{Learner}} - \underbrace{\sum_{t=1}^T f_t(\mathbf{u})}_{\text{Point } \mathbf{u}}$$

with respect to **every** point \mathbf{u} .

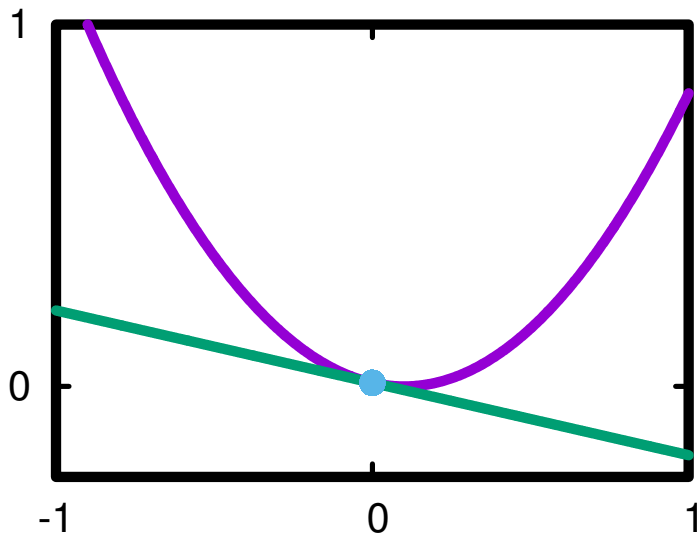
The Learner's Perspective

Round 1: Learner plays $w_1 = 0$



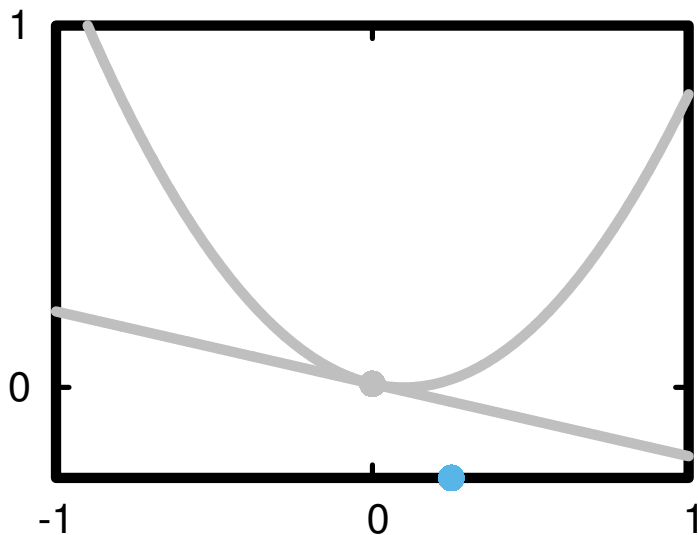
The Learner's Perspective

Round 1: Learner incurs $f_1(w_1)$ and sees $g_1 = \nabla f_1(w_1)$



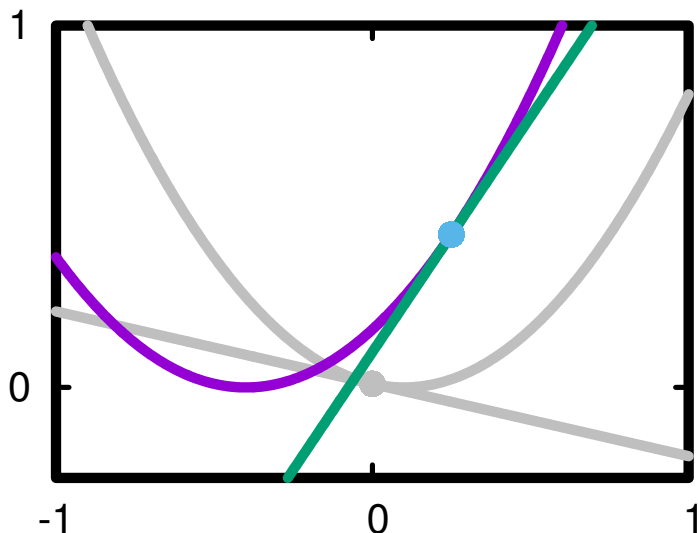
The Learner's Perspective

Round 2: Learner plays $w_1 = 1/4$



The Learner's Perspective

Round 2: Learner incurs $f_2(w_2)$ and sees $g_2 = \nabla f_2(w_2)$

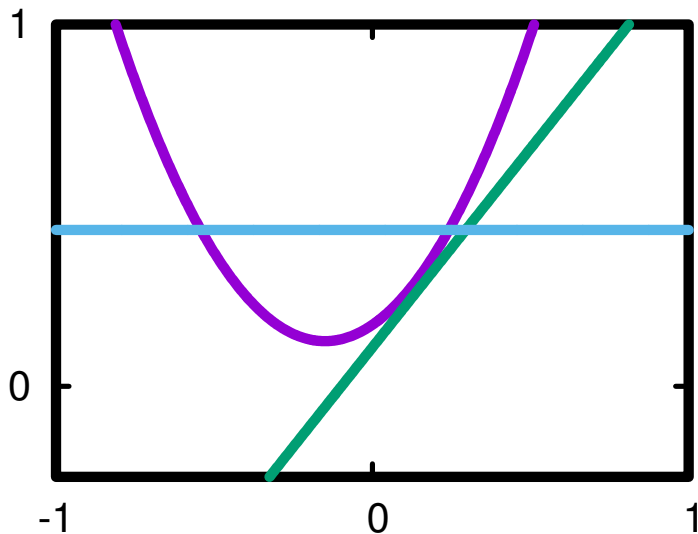


The Learner's Perspective

...

The Learner's Perspective

Evaluate Learner using **regret**: $R_T^u = \sum_{t=1}^T (f_t(w_t) - f_t(u))$



State of the Art



Online gradient descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{g}_t$$

recall $\mathbf{g}_t = \nabla f_t(\mathbf{w}_t)$

State of the Art



Online gradient descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{g}_t$$

recall $\mathbf{g}_t = \nabla f_t(\mathbf{w}_t)$

OGD bound: After T rounds,

$$R_T^{\mathbf{u}} \leq O\left(\sqrt{\sum_{t=1}^T \|\mathbf{g}_t\|^2}\right) \quad \text{for all } \mathbf{u} \text{ with } \|\mathbf{u}\| \leq 1.$$

Bounds Reveal Our Dearest Hopes



Always have **worst-case guarantee**

$$R_T^u \leq O\left(\sqrt{\sum_{t=1}^T \|g_t\|^2}\right) \leq O(\sqrt{T}).$$

Yet bound says we **might get lucky**

For smooth functions f_t with common optimum u^* , as $w_t \rightarrow u^*$, we have $g_t \rightarrow \mathbf{0}$, and

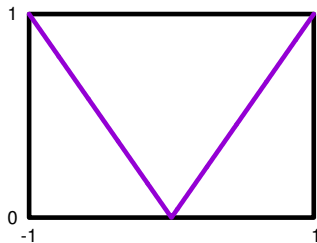
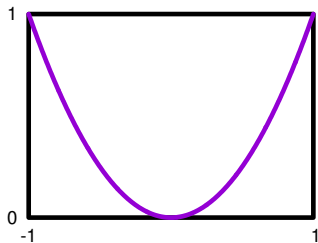
$$\sqrt{\sum_{t=1}^T \|g_t\|^2} \ll \sqrt{T}$$

grows **much slower** than \sqrt{T} .

What We Hope Happens



$$R_T^u \leq O \left(\sqrt{\sum_{t=1}^T \|g_t\|^2} \right)$$



Can We Do Better?



No **in general**: matching lower bound.

$$R_T^u \geq \Omega(\sqrt{T})$$

Yes, **with curvature**:

$$R_T^u \leq O(\ln T)$$

- ▶ Strongly convex: $\mathbf{I} \preceq \nabla^2 f(\mathbf{u})$, e.g.

$$f_t(\mathbf{u}) = \|\mathbf{u} - \mathbf{y}_t\|^2$$

⇒ gradient descent with small η

- ▶ Exp-concave: $\nabla f(\mathbf{u}) \nabla f(\mathbf{u})^\top \preceq \nabla^2 f(\mathbf{u})$, e.g.

$$f_t(\mathbf{u}) = -\ln(1 + \mathbf{y}_t^\top \mathbf{u})$$

⇒ Online Newton Step

Can We Do Better?



But do we really need **curvature**?

This talk: no, **stability** is enough.

New algorithm MetaGrad:

Separate learning rate η for each point u

Outline

Online Convex Optimization

A New Type of Guarantee

Fast Rates

MetaGrad Algorithm

Refined Bound



Recall bound for gradient descent:

$$R_T^u \leq O\left(\sqrt{\sum_{t=1}^T \|g_t\|^2}\right)$$

New bound for MetaGrad:

$$R_T^u \leq O\left(\sqrt{V_T^u d \ln T}\right) \quad \text{where} \quad V_T^u := \sum_{t=1}^T ((w_t - u)^\top g_t)^2$$

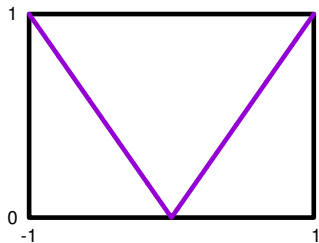
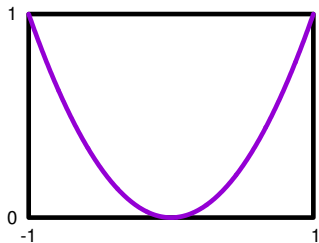
Data-dependent. Whoa! Ouroboric.

Always improvement:

$$((w_t - u)^\top g_t)^2 \leq \|w_t - u\|^2 \|g_t\|^2$$

Now What We Hope Happens

$$R_T^u \leq O \left(\sqrt{\sum_{t=1}^T ((w_t - u)^\top g_t)^2} \right)$$



Outline

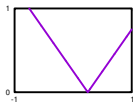
Online Convex Optimization

A New Type of Guarantee

Fast Rates

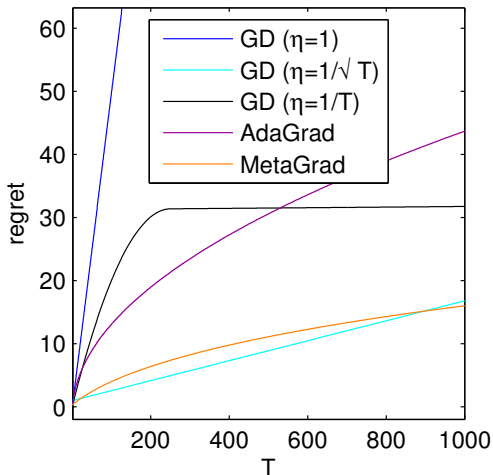
MetaGrad Algorithm

Does it Really Work?



Offline optimization (fixed function):

$$f_t(u) = |u - 1/4|$$



The “Fast Rates” Pipeline



Combine

- ▶ refined individual-sequence regret bound

$$R_T^u \leq \sqrt{V_T^u d \ln T} \quad \forall u$$

- ▶ Special-purpose argument that for best u^*

$$V_T^{u^*} \leq R_T^{u^*}$$

- ▶ Profit!

$$R_T^{u^*} \leq \sqrt{R_T^{u^*} d \ln T} \quad \text{so} \quad R_T^{u^*} \leq d \ln T$$

Significant Improvement: Fixed Function



Any fixed $f_t(\mathbf{u}) = f(\mathbf{u})$.

Let $\mathbf{u}^* = \arg \min_{\mathbf{u}} f(\mathbf{u})$ be the offline minimiser.

Crux: $(\mathbf{w}_t - \mathbf{u}^*)^\top \mathbf{g}_t \in [0, 2]$.

Now from the regret bound

$$R_T^{\mathbf{u}^*} \leq \sum_{t=1}^T (\mathbf{w}_t - \mathbf{u}^*)^\top \mathbf{g}_t \leq \sqrt{V_T^{\mathbf{u}^*} d \ln T}$$

and **special-purpose** observation

$$V_T^{\mathbf{u}^*} = \sum_{t=1}^T ((\mathbf{u}^* - \mathbf{w}_t)^\top \mathbf{g}_t)^2 \leq 2 \sum_{t=1}^T (\mathbf{w}_t - \mathbf{u}^*)^\top \mathbf{g}_t$$

we can solve for $V_T^{\mathbf{u}^*}$ to find $V_T^{\mathbf{u}^*} \leq d \ln T$ and hence

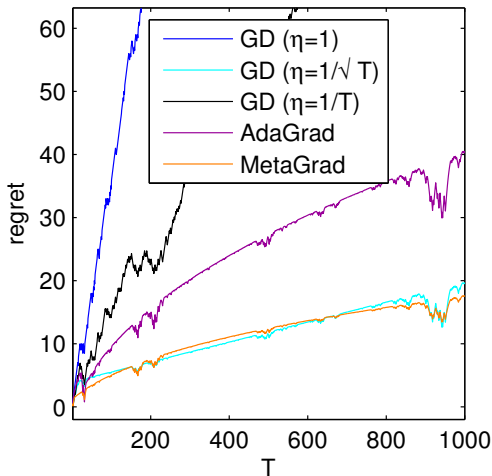
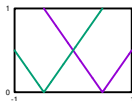
$$R_T^{\mathbf{u}^*} \leq \sqrt{2d \ln T}$$

Does It Really Actually Work?

Stochastic optimization:

$$f_t(u) = |u - x_t|$$

where $x_t = \pm \frac{1}{2}$ i.i.d. with probability 0.4 and 0.6.

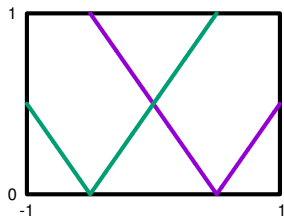


What's Going On, Really?

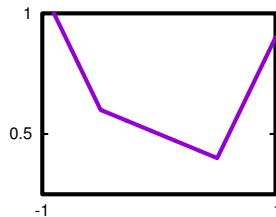
Stochastic optimization:

$$f_t(u) = |u - x_t|$$

where $x_t = \pm \frac{1}{2}$ i.i.d. with probability 0.4 and 0.6.



Individual functions



Long-term average

Stable minimum easy to converge to

Significant Improvement: Stochastic Case



Consider i.i.d.

$$f \sim \mathbb{P} \quad \text{with} \quad \mathbf{u}^* = \arg \min_{\mathbf{u}} \mathbb{E}[f(\mathbf{u})]$$

Condition: there is a $c > 0$ such that

$$\forall \mathbf{w} : (\mathbf{w} - \mathbf{u}^*)^\top \mathbb{E} [\nabla f(\mathbf{w}) \nabla f(\mathbf{w})^\top] (\mathbf{w} - \mathbf{u}^*) \leq c (\mathbf{w} - \mathbf{u}^*)^\top \mathbb{E} [\nabla f(\mathbf{w})]$$

Now from the **special-case** condition

$$\mathbb{E}[V_T^{\mathbf{u}^*}] \leq c \mathbb{E} \left[\sum_{t=1}^T (\mathbf{w}_t - \mathbf{u}^*)^\top \nabla f(\mathbf{w}_t) \right]$$

and by the **generic** regret bound, in expectation,

$$\mathbb{E}[R_T^{\mathbf{u}^*}] \leq \mathbb{E} \left[\sum_{t=1}^T (\mathbf{w}_t - \mathbf{u}^*)^\top \nabla f(\mathbf{w}_t) \right] \leq \mathbb{E} \left[\sqrt{V_T^{\mathbf{u}^*} d \ln T} \right]$$

and by Jensen's inequality $\mathbb{E} \left[\sqrt{V_T^{\mathbf{u}^*}} \right] \leq \sqrt{\mathbb{E} [V_T^{\mathbf{u}^*}]}$, so that

$$\mathbb{E}[R_T^{\mathbf{u}^*}] \leq \sqrt{cd \ln T}$$

Outline

Online Convex Optimization

A New Type of Guarantee

Fast Rates

MetaGrad Algorithm

Our Approach in a Nutshell



1. Replace actual loss $f_t(\mathbf{u})$ by **surrogate loss** $\ell_t^\eta(\mathbf{u})$
 - ▶ parametrised by learning rate η
 - ▶ exp-concave in \mathbf{u}
 - ▶ So can get good bound for surrogate regret
2. Exponentially spaced **grid** $\eta_1, \eta_2, \dots, \eta_{\log(T)}$ ($\eta_i = 2^{-i}$).
3. Off-the-shelf exp-concave **Slave** for grid point η_i predicts

$$\mathbf{w}_1^{\eta_i}, \mathbf{w}_2^{\eta_i}, \dots$$

4. At each round t , **Master** aggregates $\mathbf{w}_t^{\eta_1}, \mathbf{w}_t^{\eta_2}, \dots$ into \mathbf{w}_t .

Surrogate Loss



Real loss

$$f_t(\mathbf{u}) \leq f_t(\mathbf{w}_t) + (\mathbf{u} - \mathbf{w}_t)^\top \mathbf{g}_t$$

Surrogate loss

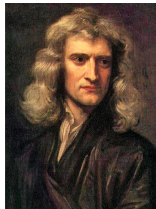
$$\ell_t^\eta(\mathbf{u}) := \eta(\mathbf{u} - \mathbf{w}_t)^\top \mathbf{g}_t + (\eta(\mathbf{u} - \mathbf{w}_t)^\top \mathbf{g}_t)^2$$

Exp-concave! In particular:

$$e^{-\ell_t^\eta(\mathbf{u})} \leq 1 + \eta(\mathbf{w}_t - \mathbf{u})^\top \mathbf{g}_t.$$

Excellent bound $O(\ln T)$ for wrong loss.

MetaGrad Slave



η -Slave (variant of Online Newton Step) predicts

$$\mathbf{w}_{t+1}^\eta = \mathbf{w}_t^\eta - \eta \Sigma_{t+1}^\eta \mathbf{g}_t$$

where the **covariance matrix** is given by

$$\Sigma_{t+1}^\eta = \left(\frac{1}{4} \mathbf{I} + 2\eta^2 \sum_{s=1}^t \mathbf{g}_s \mathbf{g}_s^\top \right)^{-1}$$

η -Slave guarantees

$$\begin{aligned} \sum_{t=1}^T (\ell_t^\eta(\mathbf{w}_t^\eta) - \ell_t^\eta(\mathbf{u})) &\leq \frac{1}{8} \|\mathbf{u}\|^2 + \frac{1}{2} \ln \det \left(\mathbf{I} + 8\eta^2 \sum_{t=1}^T \mathbf{g}_t \mathbf{g}_t^\top \right) \\ &\leq O(d \ln T) \quad \forall \mathbf{u} \end{aligned}$$

MetaGrad Master



Input: Grid points $\eta_i = 2^{-i}$ with weights $\pi_i = \frac{1}{i(i+1)}$.

Goal: aggregate $w_t^{\eta_1}, w_t^{\eta_2}, \dots$

Idea: Potential

$$\Phi_t := \sum_i \pi_i e^{-\sum_{s=1}^t \ell_s^{\eta_i}(w_s^{\eta_i})}.$$

Two steps:

- ▶ Find predictions w_t that ensure $1 \geq \Phi_1 \geq \Phi_2 \geq \dots$
- ▶ Derive regret bound from $1 \geq \Phi_T$.



Tilted exponentially weighted average

$$\mathbf{w}_{t+1} = \frac{\sum_i \pi_i e^{-\sum_{s=1}^t \ell_s^{\eta_i}(\mathbf{w}_s^{\eta_i})} \eta_i \mathbf{w}_{t+1}^{\eta_i}}{\sum_i \pi_i e^{-\sum_{s=1}^t \ell_s^{\eta_i}(\mathbf{w}_s^{\eta_i})} \eta_i}$$

ensures potential shrinks:

$$\begin{aligned} \Phi_{t+1} - \Phi_t &= \sum_i \pi_i e^{-\sum_{s=1}^t \ell_s^{\eta_i}(\mathbf{w}_s^{\eta_i})} \left(e^{-\ell_{t+1}^{\eta_i}(\mathbf{w}_{t+1}^{\eta_i})} - 1 \right) \\ &\stackrel{\text{exp-con}}{\leq} \sum_i \pi_i e^{-\sum_{s=1}^t \ell_s^{\eta_i}(\mathbf{w}_s^{\eta_i})} \eta_i (\mathbf{w}_{t+1} - \mathbf{w}_{t+1}^{\eta_i})^\top \mathbf{g}_{t+1} \stackrel{\text{weights}}{=} 0 \end{aligned}$$

and hence $\Phi_t \leq 1$.

MetaGrad Master, Small Potential is Good



The Master achieves for all t :

$$1 \geq \Phi_t = \sum_i \pi_i e^{-\sum_{s=1}^t \ell_s^{\eta_i}(\mathbf{w}_s^{\eta_i})}.$$

It follows that

$$\sum_{t=1}^T (0 - \ell_t^{\eta_i}(\mathbf{w}_t^{\eta_i})) \leq -\ln \pi_i \quad \forall i \text{ in grid}$$

(Master has zero surrogate loss)

MetaGrad Analysis



Now combine the Master and Slave guarantee.

For **each** grid point η and comparator \mathbf{u}

$$\sum_{t=1}^T (0 - \ell_t^\eta(\mathbf{w}_t^\eta)) \leq -\ln \pi_i \leq \ln \ln T$$

$$\sum_{t=1}^T (\ell_t^\eta(\mathbf{w}_t^\eta) - \ell_t^\eta(\mathbf{u})) \leq O(d \ln T)$$

so

$$\sum_{t=1}^T (0 - \ell_t^\eta(\mathbf{u})) \leq O(d \ln T).$$

Unpacking $\ell_t^\eta(\mathbf{u}) = \eta(\mathbf{u} - \mathbf{w}_t)^\top \mathbf{g}_t + (\eta(\mathbf{u} - \mathbf{w}_t)^\top \mathbf{g}_t)^2$ yields

$$\eta R_T^{\mathbf{u}} \leq \eta^2 V_T^{\mathbf{u}} + O(d \ln T).$$

MetaGrad Analysis (ctd.)



Reorganise the bound to:

$$R_T^u \leq \eta V_T^u + \frac{O(d \ln T)}{\eta}$$

Now pick the best grid point

$$\hat{\eta} = \sqrt{\frac{O(d \ln T)}{V_T^u}}$$

to find

$$R_T^u \leq O\left(\sqrt{V_T^u d \ln T}\right)$$

of course we need a grid point close to $\hat{\eta}$ and we need to deal with off-grid $\hat{\eta} \gg 1$ and $\hat{\eta} \ll \frac{1}{\sqrt{T}}$.

MetaGrad Outlook



- ▶ Run-time $O(d^2)$ per round
- ▶ Projections (avoid $O(d^3)$ per round!)
- ▶ We design and analyze two versions of Slave
 - ▶ Full covariance (quadratic)
 - ▶ Diagonal approximation (linear)
- ▶ Very welcome to discuss further

MetaGrad Outlook



- ▶ Run-time $O(d^2)$ per round
- ▶ Projections (avoid $O(d^3)$ per round!)
- ▶ We design and analyze two versions of Slave
 - ▶ Full covariance (quadratic)
 - ▶ Diagonal approximation (linear)
- ▶ Very welcome to discuss further

Learn more:

- ▶ Paper submitted to COLT 2016, preprint available
- ▶ Code is available



<http://bitbucket.org/wmkoolen/metagrad>

- ▶ Experiments coming soon.



<http://blog.wouterkoolen.info>

Summary



Low regret through stability, even without curvature.

- ▶ New MetaGrad algorithm.
- ▶ Hierarchical Master-Slave construction.
- ▶ Learns the learning rate.
- ▶ Refined (adaptive) regret bound.
- ▶ Stochastic condition for logarithmic regret (fast rates)

Thank you!