

# MetaGrad: Multiple Learning Rates in Online Learning

Tim van Erven Wouter M. Koolen

## Abstract

To get good performance in online convex optimization you need to **select and tune your algorithm** based on **lots of technical stuff**.

Grand goal: single algorithm that works well in all cases.

Multiple Eta Gradient (MetaGrad) algorithm learns **optimal learning rate** from data.

Provable Guarantees:

- Robust to worst-case convex losses
- Adapts to curvature (strong-convex, exp-concave)
- Exploits stochastic data (Bernstein)

## Online Convex Optimization Setting

- 1: **for**  $t = 1, 2, \dots, T$  **do**
- 2: Learner plays  $w_t$  in convex domain  $\mathcal{U}$
- 3: Environment reveals convex loss function  $f_t : \mathcal{U} \rightarrow \mathbb{R}$
- 4: Learner incurs loss  $f_t(w_t)$ , observes gradient  $g_t = \nabla f_t(w_t)$
- 5: **end for**

Measure **regret** w.r.t.  $u \in \mathcal{U}$ :  $\text{Regret}_T^u = \sum_{t=1}^T f_t(w_t) - \sum_{t=1}^T f_t(u)$ .

## Standard Theory

Rates based on curvature:

Convex $f_t$	$\sqrt{T}$	GD with $\eta_t \propto \frac{1}{\sqrt{t}}$
Strongly convex $f_t$	$\ln T$	GD with $\eta_t \propto \frac{1}{t}$
Exp-concave $f_t$	$d \ln T$	ONS with $\eta_t = \text{constant}$

[Bartlett, Hazan, and Rakhlin, 2007], [Do et al., 2009] handle two cases: **strongly convex + worst-case convex**

## MetaGrad Covers Many Cases

Convex $f_t$	$\sqrt{T \ln \ln T}$
Exp-concave, strongly convex $f_t$	$d \ln T$
$\beta$ -Bernstein i.i.d. $f_t$	$(d \ln T)^{\frac{1}{2-\beta}} T^{\frac{1-\beta}{2-\beta}}$

Bernstein distributions with  $\beta = 1$  very common:

Absolute loss*	$f_t(u) =  u - X_t $	$\ln T$
Hinge loss*	$\max\{0, 1 - Y_t \langle u, X_t \rangle\}$	$d \ln T$

## Main Theorem

**Theorem 1.** MetaGrad's regret is bounded by

$$\text{Regret}_T^u \leq \sum_{t=1}^T (w_t - u)^\top g_t \leq \min \left\{ \begin{array}{l} O(\sqrt{V_T^u d \ln T} + d \ln T) \\ O(\sqrt{T \ln \ln T}), \end{array} \right.$$

where  $V_T^u = \sum_{t=1}^T ((u - w_t)^\top g_t)^2$ .

## Fast Rates: Directional Derivative Condition

**Theorem 2.** If there exist  $a, b > 0$  such that all  $f_t$  satisfy

$$f_t(u) \geq f_t(w) + a(u - w)^\top \nabla f_t(w) + b((u - w)^\top \nabla f_t(w))^2 \quad \forall w \in \mathcal{U},$$

then  $\text{Regret}_T^u \leq O(d \ln T)$ .

- Satisfied by **exp-concave** and **strongly convex** functions [Hazan, Agarwal, and Kale, 2007] with  $a = 1$ .
- Satisfied for any **fixed convex** function  $f_t = f$  with minimizer  $u$ , even without any curvature, with  $a = 2$ .

## Fast Rates: Stochastic Bernstein Condition

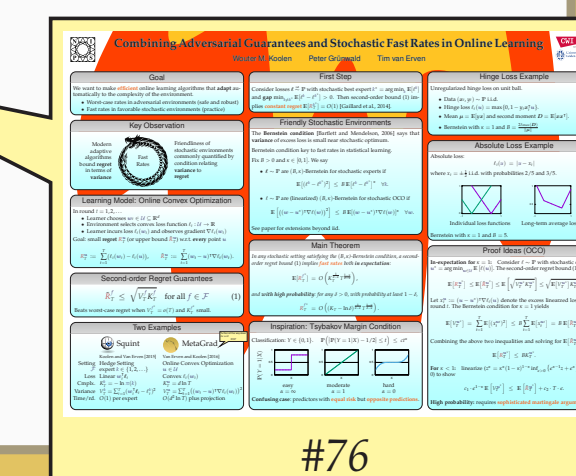
Consider  $f_t \stackrel{\text{iid}}{\sim} \mathbb{P}$  with stochastic optimum  $u^* = \arg \min_{u \in \mathcal{U}} \mathbb{E}_f[f(u)]$  satisfying the (linearized) **(B,  $\beta$ )-Bernstein condition**

$$\mathbb{E} \left[ ((w - u^*)^\top \nabla f(w))^2 \right] \leq B \mathbb{E} [(w - u^*)^\top \nabla f(w)]^\beta \quad \forall w \in \mathcal{U}.$$

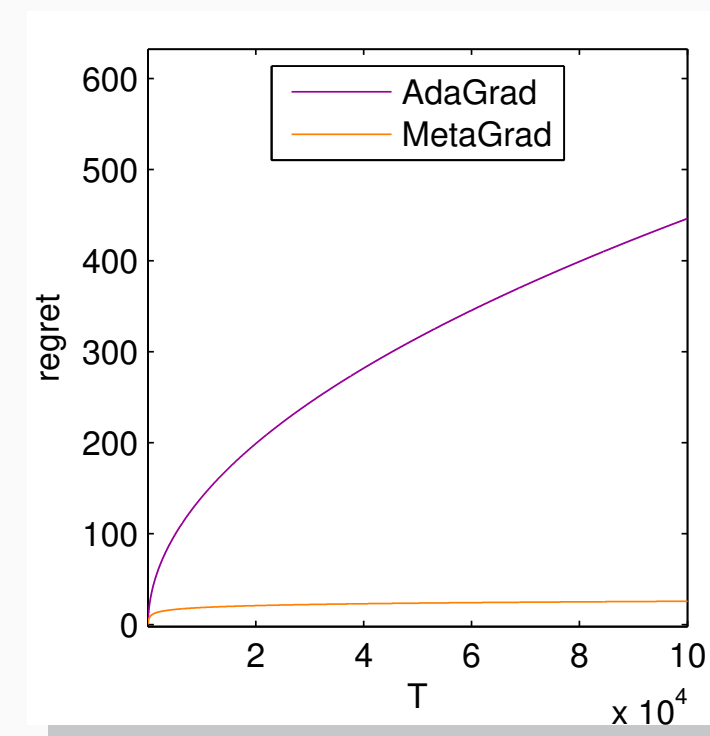
**Example:** Hinge loss (unit ball):  $\beta = 1, B = \frac{2\lambda_{\max}(\mathbb{E}[XX^\top])}{\|\mathbb{E}[YX]\|}$

**Theorem 3** (Koolen, Grünwald, Van Erven, 2016).

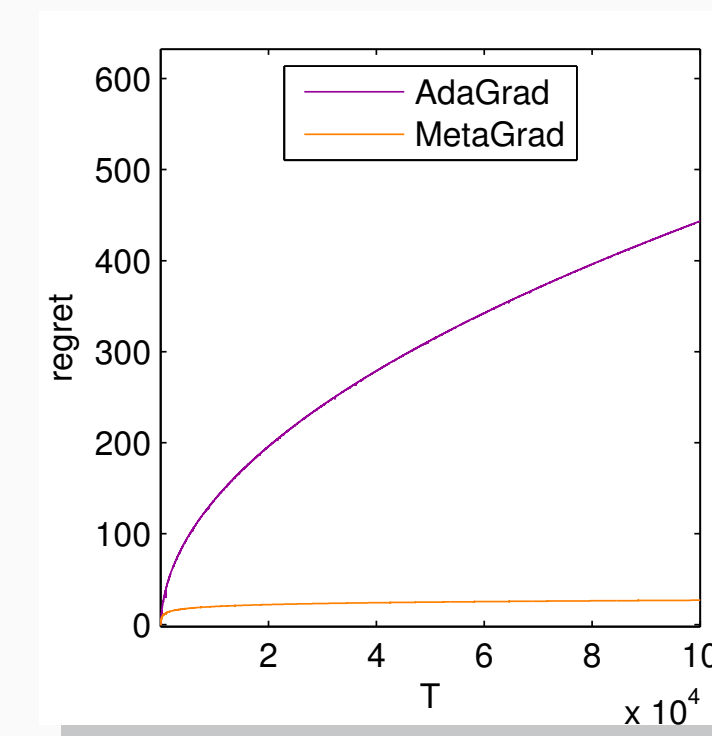
$$\mathbb{E}[\text{Regret}_T^{u^*}] = O\left( (Bd \ln T)^{\frac{1}{2-\beta}} T^{\frac{1-\beta}{2-\beta}} \right)$$



## Experiments (Proof-of-Concept)



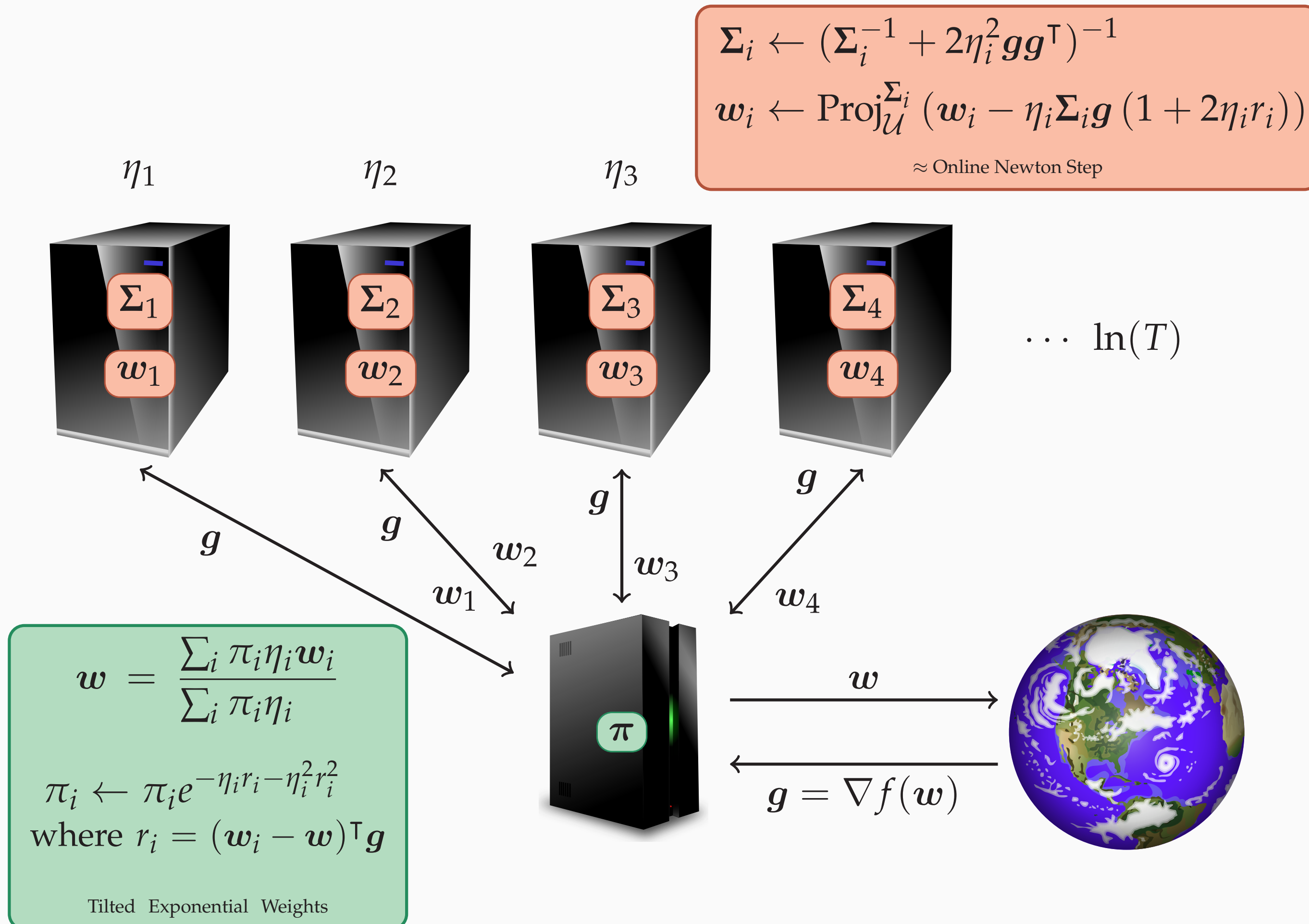
Offline: fixed  $f_t(u) = |u - 1/4|$



Stochastic Online:  $f_t(u) = |u - X_t|$  where  $X_t = \pm 1/2$  i.i.d. w.p. 2/5 and 3/5.

- MetaGrad:  $O(\ln T)$  regret, AdaGrad:  $O(\sqrt{T})$ , match bounds
- Functions neither strongly convex nor smooth

## MetaGrad Algorithm



## Proof Ideas

Analysis based on second-order **surrogate loss**. For each  $\eta$ :

$$\ell_t^\eta(u) := \eta(u - w_t)^\top g_t + \eta^2((u - w_t)^\top g_t)^2$$

Since surrogate is **exp-concave** for each fixed  $\eta$ , we can use **online quasi-Newton method** like Online Newton Step [Hazan et al., 2007] to get predictions  $w_t^\eta$  that achieve logarithmic regret:

$$\sum_{t=1}^T \ell_t^\eta(w_t^\eta) - \sum_{t=1}^T \ell_t^\eta(u) \leq O(d \ln T) \quad \forall u \in \mathcal{U}$$

To learn the **best  $\eta$**  we combine the predictions  $w_t^\eta$  for multiple  $\eta$  into a single master prediction  $w_t$  using an **experts algorithm for combining multiple learning rates** similar to Squint [Koolen and Van Erven, 2015], to get:

$$\sum_{t=1}^T \ell_t^\eta(w_t) - \sum_{t=1}^T \ell_t^\eta(w_t^\eta) \leq O(\ln \ln T) \quad \forall \eta$$

**Difficulty:** Master has to perform well under multiple loss functions simultaneously. No standard experts algorithm works!

Together:  $-\sum_{t=1}^T \ell_t^\eta(u) \leq O(d \ln T)$  for each  $\eta$  and  $u$ , resulting in

$$\sum_{t=1}^T (w_t - u)^\top g_t \leq \frac{O(d \ln T)}{\eta} + \eta V_T^u \Rightarrow O\left(\sqrt{V_T^u d \ln T}\right).$$