

Leiden Lorentz Workshop, November 9, 2016

# MetaGrad: Multiple Learning Rates in Online Learning

**Tim van Erven**

Wouter Koolen



Universiteit  
Leiden



Centrum Wiskunde & Informatica

# Online Convex Optimization

Parameters  $w$  take values in a convex domain  $\mathcal{U}$

- 1: **for**  $t = 1, 2, \dots, T$  **do**
- 2:   Learner plays  $w_t \in \mathcal{U}$
- 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.  $u \in \mathcal{U}$ :

$$\text{Regret}_T^u = \sum_{t=1}^T f_t(w_t) - \sum_{t=1}^T f_t(u)$$

Assumptions:  $\text{diameter}(\mathcal{U}) \leq D$ ,  $\|g_t\|_2 \leq G$ .

# The Standard Picture

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]:  
Adaptive GD: **strongly convex + general convex**

# The Standard Picture

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]:  
Adaptive GD: **strongly convex + general convex**

Our goals:

- ▶ Adaptivity to more types of functions  $f_t$
- ▶ Fast rates **without curvature** for 'easy' stochastic data

## Other Types of Adaptivity

- ▶ [Orabona, 2014, Orabona and Pál, 2016]: adapt to size  $\|\mathbf{u}\|_2$  of comparator
- ▶ AdaGrad [Duchi et al., 2011]: box-like domain ( $\ell_\infty$ -ball) instead of  $\ell_2$ -ball
- ▶ [Hazan and Kale, 2010], [Chiang et al., 2012]: linear functions  $f_t$  that vary little over time
- ▶ [Orabona, Crammer, and Cesa-Bianchi, 2015]: data-dependent time-varying regularizers
- ▶ ...

### Key techniques:

- ▶ Adaptive tuning of learning rate  $\eta_t$
- ▶ Use second-order information about covariance of features in time-varying regularizer

# MetaGrad: Multiple Eta Gradient Algorithm

## Theorem

MetaGrad's Regret $_T^u$  is bounded by

$$\sum_{t=1}^T (\mathbf{w}_t - \mathbf{u})^\top \mathbf{g}_t = \begin{cases} O(\sqrt{V_T^u d \ln T} + d \ln T) \\ O(\sqrt{T \ln \ln T}), \end{cases}$$

where

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

- ▶ By convexity,  $f_t(\mathbf{w}_t) - f_t(\mathbf{u}) \leq (\mathbf{w}_t - \mathbf{u})^\top \mathbf{g}_t$ .

# MetaGrad: Multiple Eta Gradient Algorithm

## Theorem

MetaGrad's Regret $_T^u$  is bounded by

$$\sum_{t=1}^T (\mathbf{w}_t - \mathbf{u})^\top \mathbf{g}_t = \begin{cases} O(\sqrt{V_T^u d \ln T} + d \ln T) \\ O(\sqrt{T \ln \ln T}), \end{cases}$$

where

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

- ▶ By convexity,  $f_t(\mathbf{w}_t) - f_t(\mathbf{u}) \leq (\mathbf{w}_t - \mathbf{u})^\top \mathbf{g}_t$ .
- ▶ Covariance:  $\mathbf{g}_t \mathbf{g}_t^\top \propto \mathbf{X}_t \mathbf{X}_t^\top$  when  $f_t(\mathbf{w}) = \text{loss}(Y_t \langle \mathbf{w}, \mathbf{X}_t \rangle)$

# MetaGrad: Multiple Eta Gradient Algorithm

## Theorem

MetaGrad's Regret $_T^u$  is bounded by

$$\sum_{t=1}^T (\mathbf{w}_t - \mathbf{u})^\top \mathbf{g}_t = \begin{cases} O(\sqrt{V_T^u d \ln T} + d \ln T) \\ O(\sqrt{T \ln \ln T}), \end{cases}$$

where

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

- ▶ By convexity,  $f_t(\mathbf{w}_t) - f_t(\mathbf{u}) \leq (\mathbf{w}_t - \mathbf{u})^\top \mathbf{g}_t$ .
- ▶ Covariance:  $\mathbf{g}_t \mathbf{g}_t^\top \propto \mathbf{X}_t \mathbf{X}_t^\top$  when  $f_t(\mathbf{w}) = \text{loss}(Y_t \langle \mathbf{w}, \mathbf{X}_t \rangle)$
- ▶ Optimal learning rate depends on  $V_T^u$ , but  $\mathbf{u}$  unknown!  
Solution: aggregate **multiple learning rates**.



## Consequences

Non-stochastic adaptation:

Convex $f_t$	$\sqrt{T \ln \ln T}$
Exp-concave $f_t$	$d \ln T$
Fixed convex $f_t = f$	$d \ln T$

## Consequences

Non-stochastic adaptation:

Convex $f_t$	$\sqrt{T \ln \ln T}$
Exp-concave $f_t$	$d \ln T$
Fixed convex $f_t = f$	$d \ln T$

**Loose end:** strongly convex  $\Rightarrow$  exp-concave gives  $d \ln T$

# Consequences

Non-stochastic adaptation:

Convex $f_t$	$\sqrt{T \ln \ln T}$
Exp-concave $f_t$	$d \ln T$
Fixed convex $f_t = f$	$d \ln T$

**Loose end:** strongly convex  $\Rightarrow$  exp-concave gives  $d \ln T$

Stochastic without curvature [joint work with Grünwald]

Suppose  $f_t$  i.i.d. with stochastic optimum

$\mathbf{u}^* = \arg \min_{\mathbf{u} \in \mathcal{U}} \mathbb{E}_f[f(\mathbf{u})]$ . Then expected regret  $\mathbb{E}[\text{Regret}_T^{\mathbf{u}^*}]$ :

Absolute loss* $f_t(\mathbf{u}) =  \mathbf{u} - \mathbf{X}_t $	$\ln T$
Hinge loss* $\max\{0, 1 - Y_t \langle \mathbf{u}, \mathbf{X}_t \rangle\}$	$d \ln T$
<b><math>(B, \beta)</math>-Bernstein</b>	$(Bd \ln T)^{1/(2-\beta)} T^{(1-\beta)/(2-\beta)}$

# Outline

## Two General Fast Rate Conditions

### MetaGrad Algorithm

### Exponential Weights Interpretation of Online Learning Algorithms

# 1. Directional Derivative Condition

## Theorem

If there exist  $a, b > 0$  such that all  $f_t$  satisfy

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

then  $O(d \ln T)$  regret w.r.t.  $\mathbf{u}$ .

$a = 1$

- ▶ Satisfied by **exp-concave** functions [Hazan, Agarwal, and Kale, 2007]
- ▶ Requires quadratic curvature in direction of minimizer  $\mathbf{u}$ .

General  $a$

- ▶ Satisfied for any **fixed convex** function  $f_t = f$  with minimizer  $\mathbf{u}$ , even without any curvature, with  $a = 2$  and  $b = 1/(DG)$ .

## 2. Bernstein Condition for Online Learning

Suppose  $f_t$  i.i.d. with stochastic optimum  $\mathbf{u}^* = \arg \min_{\mathbf{u} \in \mathcal{U}} \mathbb{E}[f(\mathbf{u})]$ .

**Standard Bernstein condition:**

$$\mathbb{E}(f(\mathbf{w}) - f(\mathbf{u}^*))^2 \leq B(\mathbb{E}[f(\mathbf{w}) - f(\mathbf{u}^*)])^\beta \quad \text{for all } \mathbf{w} \in \mathcal{U}.$$

## 2. Bernstein Condition for Online Learning

Suppose  $f_t$  i.i.d. with stochastic optimum  $\mathbf{u}^* = \arg \min_{\mathbf{u} \in \mathcal{U}} \mathbb{E}[f(\mathbf{u})]$ .

**Standard Bernstein condition:**

$$\mathbb{E}(f(\mathbf{w}) - f(\mathbf{u}^*))^2 \leq B(\mathbb{E}[f(\mathbf{w}) - f(\mathbf{u}^*)])^\beta \quad \text{for all } \mathbf{w} \in \mathcal{U}.$$

Replace by **weaker linearized version:**

- ▶ Apply with  $\tilde{f}(\mathbf{u}) = \langle \mathbf{u}, \nabla f(\mathbf{w}) \rangle$  instead of  $f$ !
- ▶ By convexity,  $f(\mathbf{w}) - f(\mathbf{u}^*) \leq \tilde{f}(\mathbf{w}) - \tilde{f}(\mathbf{u}^*)$ .

$$\mathbb{E}((\mathbf{w} - \mathbf{u}^*) \nabla f(\mathbf{w}))^2 \leq B(\mathbb{E}[(\mathbf{w} - \mathbf{u}^*) \nabla f(\mathbf{w})])^\beta \quad \text{for all } \mathbf{w} \in \mathcal{U}.$$

## 2. Bernstein Condition for Online Learning

Suppose  $f_t$  i.i.d. with stochastic optimum  $\mathbf{u}^* = \arg \min_{\mathbf{u} \in \mathcal{U}} \mathbb{E}[f(\mathbf{u})]$ .

**Standard Bernstein condition:**

$$\mathbb{E}(f(\mathbf{w}) - f(\mathbf{u}^*))^2 \leq B(\mathbb{E}[f(\mathbf{w}) - f(\mathbf{u}^*)])^\beta \quad \text{for all } \mathbf{w} \in \mathcal{U}.$$

Replace by **weaker linearized version:**

- ▶ Apply with  $\tilde{f}(\mathbf{u}) = \langle \mathbf{u}, \nabla f(\mathbf{w}) \rangle$  instead of  $f$ !
- ▶ By convexity,  $f(\mathbf{w}) - f(\mathbf{u}^*) \leq \tilde{f}(\mathbf{w}) - \tilde{f}(\mathbf{u}^*)$ .

$$\mathbb{E}((\mathbf{w} - \mathbf{u}^*) \nabla f(\mathbf{w}))^2 \leq B(\mathbb{E}[(\mathbf{w} - \mathbf{u}^*) \nabla f(\mathbf{w})])^\beta \quad \text{for all } \mathbf{w} \in \mathcal{U}.$$

Hinge loss (with  $G = D = 1$ ):  $\beta = 1$ ,  $B = \frac{2\lambda_{\max}(\mathbb{E}[\mathbf{X}\mathbf{X}^\top])}{\|\mathbb{E}[\mathbf{Y}\mathbf{X}]\|}$



## 2. Bernstein Condition for Online Learning

Suppose  $f_t$  i.i.d. with stochastic optimum  $\mathbf{u}^* = \arg \min_{\mathbf{u} \in \mathcal{U}} \mathbb{E}[f(\mathbf{u})]$ .

**Standard Bernstein condition:**

$$\mathbb{E}(f(\mathbf{w}) - f(\mathbf{u}^*))^2 \leq B(\mathbb{E}[f(\mathbf{w}) - f(\mathbf{u}^*)])^\beta \quad \text{for all } \mathbf{w} \in \mathcal{U}.$$

Replace by **weaker linearized version:**

- ▶ Apply with  $\tilde{f}(\mathbf{u}) = \langle \mathbf{u}, \nabla f(\mathbf{w}) \rangle$  instead of  $f$ !
- ▶ By convexity,  $f(\mathbf{w}) - f(\mathbf{u}^*) \leq \tilde{f}(\mathbf{w}) - \tilde{f}(\mathbf{u}^*)$ .

$$\mathbb{E}((\mathbf{w} - \mathbf{u}^*) \nabla f(\mathbf{w}))^2 \leq B(\mathbb{E}[(\mathbf{w} - \mathbf{u}^*) \nabla f(\mathbf{w})])^\beta \quad \text{for all } \mathbf{w} \in \mathcal{U}.$$

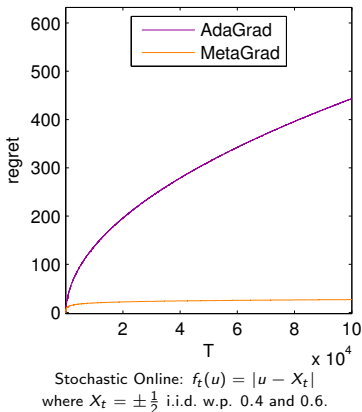
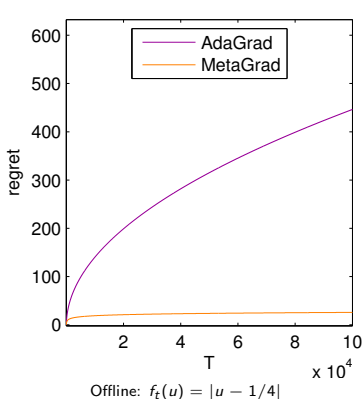
Hinge loss (with  $G = D = 1$ ):  $\beta = 1$ ,  $B = \frac{2\lambda_{\max}(\mathbb{E}[\mathbf{X}\mathbf{X}^\top])}{\|\mathbb{E}[\mathbf{Y}\mathbf{X}]\|}$

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

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

$$\text{Regret}_T^{\mathbf{u}^*} = O\left((Bd \ln T - \ln \delta)^{1/(2-\beta)} T^{(1-\beta)/(2-\beta)}\right) \quad \text{w.p.} \geq 1 - \delta$$

## Difference in Rates Not Just Theoretical



- ▶ MetaGrad:  $O(\ln T)$  regret, AdaGrad:  $O(\sqrt{T})$ , match bounds
- ▶ Functions neither strongly convex nor smooth
- ▶ **Caveat:** comparison more complicated for higher dimensions, unless we run a separate copy of MetaGrad per dimension, like the diagonal version of AdaGrad runs GD per dimension

# Outline

Two General Fast Rate Conditions

**MetaGrad Algorithm**

Exponential Weights Interpretation of Online Learning Algorithms

# MetaGrad Algorithm

Second-order **surrogate loss** for each  $\eta$  of interest (from a grid):

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

# MetaGrad Algorithm

Second-order **surrogate loss** for each  $\eta$  of interest (from a grid):

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

One **Slave** algorithm per  $\eta$  produces  $\mathbf{w}_t^\eta$  such that

$$\sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t^\eta) - \sum_{t=1}^T \ell_t^\eta(\mathbf{u}) \leq R_{\text{slave}}^{\mathbf{u}}(\eta)$$

# MetaGrad Algorithm

Second-order **surrogate loss** for each  $\eta$  of interest (from a grid):

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

One **Slave** algorithm per  $\eta$  produces  $\mathbf{w}_t^\eta$  such that

$$\sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t^\eta) - \sum_{t=1}^T \ell_t^\eta(\mathbf{u}) \leq R_{\text{slave}}^{\mathbf{u}}(\eta)$$

Single **Master** algorithm produces  $\mathbf{w}_t$  such that

$$\underbrace{\sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t)}_{=0} - \sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t^\eta) \leq R_{\text{master}}(\eta) \quad \forall \eta$$

# MetaGrad Algorithm

Second-order **surrogate loss** for each  $\eta$  of interest (from a grid):

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

One **Slave** algorithm per  $\eta$  produces  $\mathbf{w}_t^\eta$  such that

$$\sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t^\eta) - \sum_{t=1}^T \ell_t^\eta(\mathbf{u}) \leq R_{\text{slave}}^{\mathbf{u}}(\eta)$$

Single **Master** algorithm produces  $\mathbf{w}_t$  such that

$$\underbrace{\sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t)}_{=0} - \sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t^\eta) \leq R_{\text{master}}(\eta) \quad \forall \eta$$

Together:  $-\sum_{t=1}^T \ell_t^\eta(\mathbf{u}) \leq R_{\text{slave}}^{\mathbf{u}}(\eta) + R_{\text{master}}(\eta) \quad \forall \eta$

# MetaGrad Algorithm

Second-order **surrogate loss** for each  $\eta$  of interest (from a grid):

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

One **Slave** algorithm per  $\eta$  produces  $\mathbf{w}_t^\eta$  such that

$$\sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t^\eta) - \sum_{t=1}^T \ell_t^\eta(\mathbf{u}) \leq R_{\text{slave}}^u(\eta)$$

Single **Master** algorithm produces  $\mathbf{w}_t$  such that

$$\underbrace{\sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t) - \sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t^\eta)}_{=0} \leq R_{\text{master}}(\eta) \quad \forall \eta$$

Together:  $-\sum_{t=1}^T \ell_t^\eta(\mathbf{u}) \leq R_{\text{slave}}^u(\eta) + R_{\text{master}}(\eta) \quad \forall \eta$

$$\sum_{t=1}^T (\mathbf{w}_t - \mathbf{u})^\top \mathbf{g}_t \leq \frac{R_{\text{slave}}^u(\eta) + R_{\text{master}}(\eta)}{\eta} + \eta V_T^u$$



# MetaGrad Algorithm

Second-order **surrogate loss** for each  $\eta$  of interest (from a grid):

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

One **Slave** algorithm per  $\eta$  produces  $\mathbf{w}_t^\eta$  such that

$$\sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t^\eta) - \sum_{t=1}^T \ell_t^\eta(\mathbf{u}) \leq R_{\text{slave}}^u(\eta)$$

Single **Master** algorithm produces  $\mathbf{w}_t$  such that

$$\underbrace{\sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t) - \sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t^\eta)}_{=0} \leq R_{\text{master}}(\eta) \quad \forall \eta$$

Together:  $-\sum_{t=1}^T \ell_t^\eta(\mathbf{u}) \leq R_{\text{slave}}^u(\eta) + R_{\text{master}}(\eta) \quad \forall \eta$

$$\sum_{t=1}^T (\mathbf{w}_t - \mathbf{u})^\top \mathbf{g}_t \leq \frac{O(d \ln T) + O(\ln \ln T)}{\eta} + \eta V_T^u$$

## MetaGrad Algorithm

Second-order **surrogate loss** for each  $\eta$  of interest (from a grid):

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

One **Slave** algorithm per  $\eta$  produces  $\mathbf{w}_t^\eta$  such that

$$\sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t^\eta) - \sum_{t=1}^T \ell_t^\eta(\mathbf{u}) \leq R_{\text{slave}}^u(\eta)$$

Single **Master** algorithm produces  $\mathbf{w}_t$  such that

$$\underbrace{\sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t) - \sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t^\eta)}_{=0} \leq R_{\text{master}}(\eta) \quad \forall \eta$$

Together:  $-\sum_{t=1}^T \ell_t^\eta(\mathbf{u}) \leq R_{\text{slave}}^u(\eta) + R_{\text{master}}(\eta) \quad \forall \eta$

$$\sum_{t=1}^T (\mathbf{w}_t - \mathbf{u})^\top \mathbf{g}_t \leq \frac{O(d \ln T) + O(\ln \ln T)}{\eta} + \eta V_T^u \Rightarrow O\left(\sqrt{V_T^u d \ln T}\right)$$

# MetaGrad Master

**Goal:** aggregate slave predictions  $\mathbf{w}_t^\eta$  for all  $\eta$  in exponentially spaced grid  $\frac{2^{-0}}{5DG}, \frac{2^{-1}}{5DG}, \dots, \frac{2^{-\lceil \frac{1}{2} \log_2 T \rceil}}{5DG}$

**Difficulty:** master's predictions must be good w.r.t. different loss functions  $\ell_t^\eta$  for all  $\eta$  simultaneously

Compute **exponential weights** with performance of each  $\eta$  measured by its own surrogate loss:

$$\pi_t(\eta) = \frac{\pi_1(\eta) e^{-\sum_{s < t} \ell_s^\eta(\mathbf{w}_s^\eta)}}{Z}$$

Then predict with **tilted** exponentially weighted average:

$$\mathbf{w}_t = \frac{\sum_{\eta} \pi_t(\eta) \eta \mathbf{w}_t^\eta}{\sum_{\eta} \pi_t(\eta) \eta}$$

# MetaGrad Master Analysis

**Potential**  $\Phi_T = \sum_{\eta} \pi_1(\eta) e^{-\sum_{t=1}^T \ell_t^{\eta}(\mathbf{w}_t^{\eta})}$

Proof outline:

$$\Phi_T \leq \Phi_{T-1} \leq \dots \leq \Phi_0 = 1$$

$$\pi_1(\eta) e^{-\sum_{t=1}^T \ell_t^{\eta}(\mathbf{w}_t^{\eta})} \leq 1 \quad \forall \eta$$

$$\underbrace{\sum_{t=1}^T \ell_t^{\eta}(\mathbf{w}_t)}_{=0} - \sum_{t=1}^T \ell_t^{\eta}(\mathbf{w}_t^{\eta}) \leq -\ln \pi_1(\eta)$$

# MetaGrad Master Analysis

**Potential**  $\Phi_T = \sum_{\eta} \pi_1(\eta) e^{-\sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t^\eta)}$

Proof outline:

$$\Phi_T \leq \Phi_{T-1} \leq \dots \leq \Phi_0 = 1$$

$$\pi_1(\eta) e^{-\sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t^\eta)} \leq 1 \quad \forall \eta$$

$$\underbrace{\sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t)}_{=0} - \sum_{t=1}^T \ell_t^\eta(\mathbf{w}_t^\eta) \leq -\ln \pi_1(\eta)$$

Grid has  $\lceil \frac{1}{2} \log_2 T \rceil + 1$  learning rates, so for heavy-tailed prior:

$$-\ln \pi_1(\eta) = O(\ln \ln T)$$

## MetaGrad Master Analysis: Decreasing Potential

Surrogate loss  $\ell_t^\eta(\mathbf{u}) = \eta(\mathbf{u} - \mathbf{w}_t)^\top \mathbf{g}_t + \eta^2(\mathbf{u} - \mathbf{w}_t)^\top \mathbf{g}_t \mathbf{g}_t^\top (\mathbf{u} - \mathbf{w}_t)$  is **exp-concave**, even if  $f_t$  is not.

Upper bound by tangent at  $\mathbf{u} = \mathbf{w}_t$ :

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

## MetaGrad Master Analysis: Decreasing Potential

Surrogate loss  $\ell_t^\eta(\mathbf{u}) = \eta(\mathbf{u} - \mathbf{w}_t)^\top \mathbf{g}_t + \eta^2(\mathbf{u} - \mathbf{w}_t)^\top \mathbf{g}_t \mathbf{g}_t^\top (\mathbf{u} - \mathbf{w}_t)$  is **exp-concave**, even if  $f_t$  is not.

Upper bound by tangent at  $\mathbf{u} = \mathbf{w}_t$ :

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

Choose master's weights to ensure decreasing potential:

$$\begin{aligned} \Phi_T - \Phi_{T-1} &= \sum_{\eta} \pi_1(\eta) e^{-\sum_{t < T} \ell_t^\eta(\mathbf{w}_t^\eta)} \left( e^{-\ell_T^\eta(\mathbf{w}_T^\eta)} - 1 \right) \\ &\leq \sum_{\eta} \pi_1(\eta) e^{-\sum_{t < T} \ell_t^\eta(\mathbf{w}_t^\eta)} \eta (\mathbf{w}_T - \mathbf{w}_T^\eta)^\top \mathbf{g}_T \\ &= 0 \quad \text{for any } \mathbf{g}_T \end{aligned}$$

## MetaGrad Slave

**Goal:** Given  $\eta$ , minimize regret w.r.t. exp-concave surrogate  $\ell_t^\eta$ .

**Update:**

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

where

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

**Project onto domain:**

$$\mathbf{w}_{t+1}^\eta = \arg \min_{\mathbf{u} \in \mathcal{U}} (\mathbf{u} - \tilde{\mathbf{w}}_{t+1}^\eta)^\top (\Sigma_{t+1}^\eta)^{-1} (\mathbf{u} - \tilde{\mathbf{w}}_{t+1}^\eta)$$



# MetaGrad Slave

**Goal:** Given  $\eta$ , minimize regret w.r.t. exp-concave surrogate  $\ell_t^\eta$ .

**Update:**

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

where

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

**Project onto domain:**

$$\mathbf{w}_{t+1}^\eta = \arg \min_{\mathbf{u} \in \mathcal{U}} (\mathbf{u} - \tilde{\mathbf{w}}_{t+1}^\eta)^\top (\Sigma_{t+1}^\eta)^{-1} (\mathbf{u} - \tilde{\mathbf{w}}_{t+1}^\eta)$$

- ▶ If master = slave, i.e.  $\mathbf{w}_t^\eta = \mathbf{w}_t$ , then is **Online Newton Step**

# MetaGrad Slave as Continuous Exponential Weights

## Exponential Weights

- ▶ Continuous set of experts  $\mathbf{u} \in \mathbb{R}^d$
- ▶ Gaussian prior  $P_1 = \mathcal{N}(\mathbf{0}, D^2 \mathbf{I})$

$$d\tilde{P}_{t+1}(\mathbf{u}) = \frac{e^{-\ell_t^\eta(\mathbf{u})} dP_t(\mathbf{u})}{Z} \quad (\text{update})$$

$$P_{t+1} = \min_{P: \mu_P \in \mathcal{U}} \text{KL}(P \| \tilde{P}_{t+1}) \quad (\text{project})$$

Play the mean of exponential weights:

$$\tilde{P}_t = \mathcal{N}(\tilde{\mathbf{w}}_t^\eta, \Sigma_t^\eta)$$

$$P_t = \mathcal{N}(\mathbf{w}_t^\eta, \Sigma_t^\eta)$$

- ▶ Can understand **Online Newton Step** and many other algorithms this way

## MetaGrad Slave Analysis

Standard exponential weights analysis gives **regret bound in space of distributions** for all  $Q$ :

$$\begin{aligned} \text{KL}(Q \| P_1) &\geq \sum_{t=1}^T -\ln \mathbb{E}_{P_t} [e^{-\ell_t^\eta(\mathbf{u})}] - \sum_{t=1}^T \mathbb{E}_Q [\ell_t^\eta(\mathbf{u})] \\ &\stackrel{\text{exp-conc}}{\geq} \sum_{t=1}^T \ell_t^\eta(\boldsymbol{\mu}_{P_t}) - \sum_{t=1}^T \mathbb{E}_Q [\ell_t^\eta(\mathbf{u})] \end{aligned}$$

## MetaGrad Slave Analysis

Standard exponential weights analysis gives **regret bound in space of distributions** for all  $Q$ :

$$\begin{aligned} \text{KL}(Q \| P_1) &\geq \sum_{t=1}^T -\ln \mathbb{E}_{P_t} [e^{-\ell_t^\eta(\mathbf{u})}] - \sum_{t=1}^T \mathbb{E}_Q [\ell_t^\eta(\mathbf{u})] \\ &\stackrel{\text{exp-conc}}{\geq} \sum_{t=1}^T \ell_t^\eta(\boldsymbol{\mu}_{P_t}) - \sum_{t=1}^T \mathbb{E}_Q [\ell_t^\eta(\mathbf{u})] \end{aligned}$$

Specialize to  $Q = \mathcal{N}(\mathbf{u}^*, D^2\Sigma) + \text{algebra}$ :

$$\begin{aligned} &\frac{1}{2D^2} \|\mathbf{u}^*\|^2 + \frac{1}{2} (-\ln \det(\Sigma) + \text{tr}(\Sigma) - d) \\ &\geq \sum_{t=1}^T \ell_t^\eta(\boldsymbol{\mu}_{P_t}) - \sum_{t=1}^T \ell_t^\eta(\mathbf{u}^*) - \sum_{t=1}^T \eta^2 D^2 \text{tr}(\Sigma \mathbf{g}_t \mathbf{g}_t^\top) \end{aligned}$$

## MetaGrad Slave Analysis

Standard exponential weights analysis gives **regret bound in space of distributions** for all  $Q$ :

$$\begin{aligned} \text{KL}(Q \| P_1) &\geq \sum_{t=1}^T -\ln \mathbb{E}_{P_t} [e^{-\ell_t^n(\mathbf{u})}] - \sum_{t=1}^T \mathbb{E}_Q [\ell_t^n(\mathbf{u})] \\ &\stackrel{\text{exp-conc}}{\geq} \sum_{t=1}^T \ell_t^n(\boldsymbol{\mu}_{P_t}) - \sum_{t=1}^T \mathbb{E}_Q [\ell_t^n(\mathbf{u})] \end{aligned}$$

Specialize to  $Q = \mathcal{N}(\mathbf{u}^*, D^2\Sigma) + \text{algebra}$ :

$$\begin{aligned} &\frac{1}{2D^2} \|\mathbf{u}^*\|^2 + \frac{1}{2} (-\ln \det(\Sigma) + \text{tr}(\Sigma) - d) \\ &\geq \sum_{t=1}^T \ell_t^n(\boldsymbol{\mu}_{P_t}) - \sum_{t=1}^T \ell_t^n(\mathbf{u}^*) - \sum_{t=1}^T \eta^2 D^2 \text{tr}(\Sigma \mathbf{g}_t \mathbf{g}_t^\top) \end{aligned}$$

**Optimize  $\Sigma$ :**

$$R_{\text{slave}}^u(\eta) \leq \frac{1}{2D^2} \|\mathbf{u}^*\|^2 + \frac{1}{2} \ln \det \left( \mathbf{I} + 2\eta^2 D^2 \sum_{t=1}^T \mathbf{g}_t \mathbf{g}_t^\top \right) = O(d \ln T)$$

# Outline

Two General Fast Rate Conditions

MetaGrad Algorithm

**Exponential Weights Interpretation of Online Learning Algorithms**

# Many Algorithms as Exponential Weights

## Exponentiated Gradient [Kivinen and Warmuth, 1997]

- ▶ Continuous set of experts  $\mathbf{u} \in \mathbb{R}^d$  with loss  $\ell_t(\mathbf{u}) = \langle \mathbf{u}, \mathbf{g}_t \rangle$
- ▶ Prior  $P_1$  puts point masses on corners of probability simplex

$$dP_{t+1}(\mathbf{u}) = \frac{e^{-\eta \ell_t(\mathbf{u})} dP_t(\mathbf{u})}{Z} \quad \implies \quad \mathbb{E}_{P_{t+1}}[\mathbf{u}] = P_{t+1} = \text{EG}$$

# Many Algorithms as Exponential Weights

## Exponentiated Gradient [Kivinen and Warmuth, 1997]

- ▶ Continuous set of experts  $\mathbf{u} \in \mathbb{R}^d$  with loss  $\ell_t(\mathbf{u}) = \langle \mathbf{u}, \mathbf{g}_t \rangle$
- ▶ Prior  $P_1$  puts point masses on corners of probability simplex

$$dP_{t+1}(\mathbf{u}) = \frac{e^{-\eta \ell_t(\mathbf{u})} dP_t(\mathbf{u})}{Z} \quad \implies \quad \mathbb{E}_{P_{t+1}}[\mathbf{u}] = P_{t+1} = \text{EG}$$

## Gradient Descent

$$\text{Gaussian prior } P_1 = \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad \implies \quad P_{t+1}(\mathbf{u}) = \mathcal{N}\left(-\eta \sum_{t=1}^T \mathbf{g}_t, \mathbf{I}\right)$$



# Many Algorithms as Exponential Weights

## Exponentiated Gradient [Kivinen and Warmuth, 1997]

- ▶ Continuous set of experts  $\mathbf{u} \in \mathbb{R}^d$  with loss  $\ell_t(\mathbf{u}) = \langle \mathbf{u}, \mathbf{g}_t \rangle$
- ▶ Prior  $P_1$  puts point masses on corners of probability simplex

$$dP_{t+1}(\mathbf{u}) = \frac{e^{-\eta \ell_t(\mathbf{u})} dP_t(\mathbf{u})}{Z} \quad \implies \quad \mathbb{E}_{P_{t+1}}[\mathbf{u}] = P_{t+1} = \text{EG}$$

## Gradient Descent

$$\text{Gaussian prior } P_1 = \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad \implies \quad P_{t+1}(\mathbf{u}) = \mathcal{N}\left(-\eta \sum_{t=1}^T \mathbf{g}_t, \mathbf{I}\right)$$

## Mirror Descent [Van der Hoeven, 2016]

- ▶ Generalize to arbitrary prior  $P_1$
- ▶  $\Phi(\theta) = \ln \int e^{\langle \theta, \mathbf{u} \rangle} dP_1(\mathbf{u})$ :  
CGF for exponential family with carrier  $P_1$

$$\mu_{P_{t+1}} = \underbrace{\nabla \Phi(\nabla \Phi^*(\mu_{P_t}) - \eta \mathbf{g}_t)}_{\text{update in natural parameters}}$$



Dirk

# Summary and Last Remarks

## MetaGrad

- ▶  $\tilde{O}(\sqrt{T})$  regret
- ▶ **Fast rates** (often  $O(d \ln T)$ ) for:
  - ▶ Adversarial: exp-concave, strongly convex, fixed functions
  - ▶ Stochastic: under Bernstein condition (including for hinge loss)

# Summary and Last Remarks

## MetaGrad

- ▶  $\tilde{O}(\sqrt{T})$  regret
- ▶ **Fast rates** (often  $O(d \ln T)$ ) for:
  - ▶ Adversarial: exp-concave, strongly convex, fixed functions
  - ▶ Stochastic: under Bernstein condition (including for hinge loss)

## MetaGrad Algorithm

- ▶ Master:
  - ▶ Pays only  $O(\ln \ln T)$  for learning the  $\eta$  that is empirically optimal on the data
  - ▶ Almost exponential weights for surrogate loss, but need to **tilt towards larger learning rates**
- ▶ Slave:
  - ▶ Continuous exponential weights on surrogate loss
  - ▶ Matrix updates take  $O(d^2)$  work, projections often  $O(d^3)$
  - ▶ Open problem: add sketching like [Luo et al., 2016]?

# Summary and Last Remarks

## MetaGrad

- ▶  $\tilde{O}(\sqrt{T})$  regret
- ▶ **Fast rates** (often  $O(d \ln T)$ ) for:
  - ▶ Adversarial: exp-concave, strongly convex, fixed functions
  - ▶ Stochastic: under Bernstein condition (including for hinge loss)

## MetaGrad Algorithm

- ▶ Master:
  - ▶ Pays only  $O(\ln \ln T)$  for learning the  $\eta$  that is empirically optimal on the data
  - ▶ Almost exponential weights for surrogate loss, but need to **tilt towards larger learning rates**
- ▶ Slave:
  - ▶ Continuous exponential weights on surrogate loss
  - ▶ Matrix updates take  $O(d^2)$  work, projections often  $O(d^3)$
  - ▶ Open problem: add sketching like [Luo et al., 2016]?

**Cool Aside:** View many online algorithms as continuous exponential weights

# References

- P. L. Bartlett, E. Hazan, and A. Rakhlin. Adaptive online gradient descent. In *NIPS 20*, pages 65–72, 2007.
- C.-K. Chiang, T. Yang, C.-J. Le, M. Mahdavi, C.-J. Lu, R. Jin, and S. Zhu. Online optimization with gradual variations. In *Proc. of the 25th Annual Conf. on Learning Theory (COLT)*, pages 6.1–6.20, 2012.
- C. B. Do, Q. V. Le, and C.-S. Foo. Proximal regularization for online and batch learning. In *Proc. of the 26th Annual International Conf. on Machine Learning (ICML)*, pages 257–264, 2009.
- J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011.
- E. Hazan and S. Kale. Extracting certainty from uncertainty: Regret bounded by variation in costs. *Machine learning*, 80(2-3):165–188, 2010.
- E. Hazan, A. Agarwal, and S. Kale. Logarithmic regret algorithms for online convex optimization. *Machine Learning*, 69(2-3):169–192, 2007.
- J. Kivinen and M. K. Warmuth. Exponentiated gradient versus gradient descent for linear predictors. *Information and Computation*, 132(1):1–63, 1997.
- W. M. Koolen, P. Grünwald, and T. van Erven. Combining adversarial guarantees and stochastic fast rates in online learning. *NIPS*, 2016.
- H. Luo, A. Agarwal, N. Cesa-Bianchi, and J. Langford. Efficient second order online learning by sketching. In *NIPS 29*, 2016.
- F. Orabona. Simultaneous model selection and optimization through parameter-free stochastic learning. In *NIPS 27*, pages 1116–1124, 2014.
- F. Orabona and D. Pál. Coin betting and parameter-free online learning. In *NIPS 29*, 2016.
- F. Orabona, K. Crammer, and N. Cesa-Bianchi. A generalized online mirror descent with applications to classification and regression. *Machine Learning*, 99(3):411–435, 2015.
- D. van der Hoeven. Is mirror descent a special case of exponential weights? Master's thesis, Leiden University, 2016. Available from <https://www.math.leidenuniv.nl/en/theses/year/2016/>.
- T. van Erven and W. M. Koolen. Metagrad: Multiple learning rates in online learning. *NIPS*, 2016.