

MetaGrad: Adapting to the Distribution of the Data in Online Sequential Prediction

Tim van Erven



Universiteit
Leiden

Joint work with: Wouter Koolen, Peter Grünwald

Rennes, March 17, 2017

Offline vs Online Learning

I. (Offline) Statistical Learning

- ▶ Limits of minimax analysis pretty well understood (margin/Bernstein condition)
- ▶ Adaptive algorithms known to exploit 'easy' data distributions

II. Online Sequential Prediction

- ▶ Most results still about minimax analysis
- ▶ We develop theory of adaptive analysis:
 - ▶ New online version of the Bernstein condition
 - ▶ New adaptive algorithm: MetaGrad
 - ▶ Adaptively achieving (optimal) fast rates

Part I: Statistical Learning Background

Statistical Learning

$$\begin{pmatrix} Y_1 \\ \mathbf{X}_1 \end{pmatrix}, \dots, \begin{pmatrix} Y_N \\ \mathbf{X}_N \end{pmatrix} \quad \text{independently distributed } \sim P$$

↓

$$\hat{f} \in \mathcal{F}$$

↓

small risk $R(\hat{f}) = \mathbb{E}_{(\mathbf{X}, Y) \sim P} [\text{loss}(\mathbf{X}, Y, \hat{f})]$

compared to minimizer f^* of risk in model \mathcal{F}

Minimax Rate:

Rate for most difficult possible P

$$\min_{\hat{f}} \max_P \mathbb{E}[R(\hat{f})] - R(f^*)$$

Classification

Given $\mathbf{X} \in \mathbb{R}^d$, predict binary label $Y \in \{0, 1\}$

$$\text{loss}(\mathbf{X}, Y, f) = \begin{cases} 0 & \text{if } f(\mathbf{X}) = Y, \\ 1 & \text{if } f(\mathbf{X}) \neq Y \end{cases}$$

$$R(f) = P(f(\mathbf{X}) \neq Y)$$

Minimax Rate:

For worst-case P , learning is slow:

$$\mathbb{E}[R(\hat{f})] - R(f^*) \asymp \sqrt{\frac{\text{complexity}(\mathcal{F})}{N}}$$

But Faster Rates Are Common

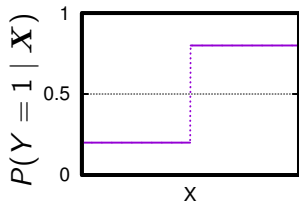
- ▶ Worst-case distribution: $P(Y = 1 | \mathbf{X})$ very close to $\frac{1}{2}$
- ▶ But then learning is (almost) useless!

The Margin Condition: [Tsybakov, 2004]

- ▶ Common case: $P(Y = 1 | \mathbf{X})$ not too close to $\frac{1}{2}$
- ▶ Assume $f^*(\mathbf{X}) = f_B(\mathbf{X}) = \arg \max_y P(Y = y | \mathbf{X})$
- ▶ Then learning is much faster, up to

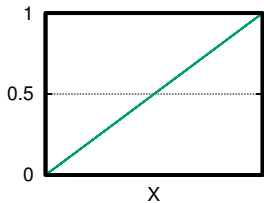
$$\mathbb{E}[R(\hat{f})] - R(f^*) \preccurlyeq \frac{\text{complexity}(\mathcal{F})}{N}$$

The Margin Condition



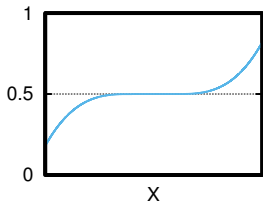
easy

$$\alpha = \infty$$



moderate

$$\alpha = 1$$



hard

$$\alpha = 0$$

$$P_{\mathbf{X}}(|P(Y | \mathbf{X}) - \frac{1}{2}| \leq t) \leq ct^{\alpha}$$

Fast Rates Because Margin Reduces Variance

Important source of excess risk $R(\hat{f}) - R(f^*)$ is **variance in excess loss**:

$$V(\hat{f}, f^*) = \mathbb{E} \left(\text{loss}(\mathbf{X}, Y, \hat{f}) - \text{loss}(\mathbf{X}, Y, f^*) \right)^2$$

Lemma (Tsybakov)

If $f^* = f_B$. Then α -margin is equivalent to the (B, β) -Bernstein condition

$$V(\hat{f}, f^*) \leq B \left(R(\hat{f}) - R(f^*) \right)^\beta$$

for $\beta = \frac{\alpha}{1+\alpha} \in [0, 1]$ and some constant $B \geq 0$.

Fast Rates Because Margin Reduces Variance

Important source of excess risk $R(\hat{f}) - R(f^*)$ is **variance in excess loss**:

$$V(\hat{f}, f^*) = \mathbb{E} \left(\text{loss}(\mathbf{X}, Y, \hat{f}) - \text{loss}(\mathbf{X}, Y, f^*) \right)^2$$

Lemma (Tsybakov)

If $f^* = f_B$. Then α -margin is equivalent to the (B, β) -Bernstein condition

$$V(\hat{f}, f^*) \leq B \left(R(\hat{f}) - R(f^*) \right)^\beta$$

for $\beta = \frac{\alpha}{1+\alpha} \in [0, 1]$ and some constant $B \geq 0$.

Smaller excess risk

Smaller variance



Fast Rates Because Margin Reduces Variance

Important source of excess risk $R(\hat{f}) - R(f^*)$ is **variance in excess loss**:

$$V(\hat{f}, f^*) = \mathbb{E} \left(\text{loss}(\mathbf{X}, Y, \hat{f}) - \text{loss}(\mathbf{X}, Y, f^*) \right)^2$$

Lemma (Tsybakov)

If $f^* = f_B$. Then α -margin is equivalent to the (B, β) -Bernstein condition

$$V(\hat{f}, f^*) \leq B \left(R(\hat{f}) - R(f^*) \right)^\beta$$

for $\beta = \frac{\alpha}{1+\alpha} \in [0, 1]$ and some constant $B \geq 0$.



Adaptive Statistical Learning

- ▶ For simplicity: **prior** π on **countable model** $\mathcal{F} = \{f_1, f_2, \dots\}$.
- ▶ **Penalized ERM** \hat{f} minimizes

$$\sum_{i=1}^N \text{loss}(\mathbf{X}_i, Y_i, f) + \lambda \log \frac{1}{\pi(f)}$$

Crucial question: **how to tune λ ?**

Adaptive Statistical Learning

- ▶ For simplicity: **prior** π on **countable model** $\mathcal{F} = \{f_1, f_2, \dots\}$.
- ▶ **Penalized ERM** \hat{f} minimizes

$$\sum_{i=1}^N \text{loss}(\mathbf{X}_i, Y_i, f) + \lambda \log \frac{1}{\pi(f)}$$

Proposition (Bernstein Condition Rate)

(Assuming bounded loss.) Knowing β , $\lambda = \left(\frac{N}{\log \frac{1}{\pi(f^*)}}\right)^{\frac{1-\beta}{2-\beta}}$ achieves

$$R(\hat{f}) - R(f^*) \preceq \left(\frac{\log \frac{1}{\delta \pi(f^*)}}{N}\right)^{\frac{1}{2-\beta}} \quad \text{w.p.} \geq 1 - \delta.$$

Adaptive Statistical Learning

- ▶ For simplicity: **prior** π on **countable model** $\mathcal{F} = \{f_1, f_2, \dots\}$.
- ▶ **Penalized ERM** \hat{f} minimizes

$$\sum_{i=1}^N \text{loss}(\mathbf{X}_i, Y_i, f) + \lambda \log \frac{1}{\pi(f)}$$

Proposition (Bernstein Condition Rate)

(Assuming bounded loss.) Knowing β , $\lambda = \left(\frac{N}{\log \frac{1}{\pi(f^*)}}\right)^{\frac{1-\beta}{2-\beta}}$ achieves

$$R(\hat{f}) - R(f^*) \preceq \left(\frac{\log \frac{1}{\delta \pi(f^*)}}{N}\right)^{\frac{1}{2-\beta}} \quad \text{w.p.} \geq 1 - \delta.$$

- ▶ **Simple adaptive** method: estimate λ on hold out set
- ▶ Or **very sophisticated adaptive approaches**:
 - ▶ Slope heuristic (Birgé, Massart)
 - ▶ Lepski's method
 - ▶ Safe Bayes (Grünwald)

Part II: Online Sequential Prediction

- ▶ Same problem of tuning $\lambda \leftrightarrow 1/\eta$
- ▶ Hold-out estimation not possible
- ▶ Very limited adaptive methods and theory

Part II: Online Sequential Prediction

- ▶ Same problem of tuning $\lambda \leftrightarrow 1/\eta$
- ▶ Hold-out estimation not possible
- ▶ Very limited adaptive methods and theory
- ▶ Our contribution: develop theory and solve tuning problem

Online Sequential Prediction

Online Data

- ▶ Data $\begin{pmatrix} Y_1 \\ \mathbf{X}_1 \end{pmatrix}, \dots, \begin{pmatrix} Y_T \\ \mathbf{X}_T \end{pmatrix}$ arrive one by one, sequentially
- ▶ No assumption that data follow any distribution

Online Estimation

- ▶ Focus on parametric models $\mathcal{F} = \{f_{\mathbf{w}} \mid \mathbf{w} \in \mathbb{R}^d\}$, but often high dimensional
- ▶ Estimate \mathbf{w}_t from $t - 1$ examples, predict Y_t by $f_{\mathbf{w}_t}(\mathbf{X}_t)$

Notation

- ▶ Abbreviate $\ell_t(\mathbf{w}) = \text{loss}(\mathbf{X}_t, Y_t, f_{\mathbf{w}})$

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 $\ell_t : \mathcal{U} \rightarrow \mathbb{R}$
- 4: Learner incurs loss $\ell_t(w_t)$, observes gradient $g_t = \nabla \ell_t(w_t)$
- 5: **end for**

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 $\ell_t : \mathcal{U} \rightarrow \mathbb{R}$
- 4: Learner incurs loss $\ell_t(w_t)$, observes gradient $g_t = \nabla \ell_t(w_t)$
- 5: **end for**

Example: Classification with Convex Surrogate Losses

Given bounded $X_t \in [-1, +1]^d$, predict label $Y_t \in \{-1, 1\}$

$$\ell_t(w) = \max\{0, 1 - Y_t \langle w, X_t \rangle\} \quad (\text{hinge loss})$$

$$\ell_t(w) = \ln \left(1 + e^{-Y_t \langle w, X_t \rangle} \right) \quad (\text{logistic loss})$$

$$\ell_t(w) = (Y_t - \langle w, X_t \rangle)^2 \quad (\text{squared loss})$$

for w in Euclidean ball with diameter D .

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 $\ell_t : \mathcal{U} \rightarrow \mathbb{R}$
- 4: Learner incurs loss $\ell_t(w_t)$, observes gradient $g_t = \nabla \ell_t(w_t)$
- 5: **end for**

Minimize **regret** w.r.t. oracle parameters $u \in \mathcal{U}$:

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

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

Standard Methods

Online Gradient Descent (OGD)

Move into direction of steepest descent:

$$\begin{aligned} \mathbf{w}_{t+1} &= \arg \min_{\mathbf{w} \in \mathcal{U}} \langle \mathbf{w}, \mathbf{g}_t \rangle + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|_2^2 \\ &= \mathbf{w}_t - \eta_t \mathbf{g}_t \quad (\text{and project onto } \mathcal{U} \text{ if go outside}) \end{aligned}$$

Regularization/step size determined by **learning rate** $\eta_t > 0$.

Standard Methods

Online Gradient Descent (OGD)

Move into direction of steepest descent:

$$\begin{aligned} \mathbf{w}_{t+1} &= \arg \min_{\mathbf{w} \in \mathcal{U}} \langle \mathbf{w}, \mathbf{g}_t \rangle + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|_2^2 \\ &= \mathbf{w}_t - \eta_t \mathbf{g}_t \quad (\text{and project onto } \mathcal{U} \text{ if go outside}) \end{aligned}$$

Regularization/step size determined by **learning rate** $\eta_t > 0$.

Online Newton Step (ONS)

Make less sensitive to parametrization by running OGD on pre-conditioned functions $\ell_t(\mathbf{H}_{t+1}^{-1/2} \tilde{\mathbf{w}})$:

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

where $\mathbf{H}_{t+1} = I + 2 \sum_{s=1}^t \mathbf{g}_s \mathbf{g}_s^\top$ and η fixed.

The Standard Picture

Minimax rates based on curvature [Hazan, 2016]:

Convex ℓ_t	\sqrt{T}	GD with $\eta_t \propto \frac{1}{\sqrt{t}}$
Strongly convex ℓ_t	$\ln T$	GD with $\eta_t \propto \frac{1}{t}$
Exp-concave ℓ_t	$d \ln T$	ONS with $\eta \propto 1$

- ▶ **Strongly convex:** second derivative at least $\alpha > 0$, implies exp-concave
- ▶ **Exp-concave:** $e^{-\alpha \ell_t}$ concave
Satisfied by logistic loss, squared loss, but not hinge loss

The Standard Picture

Minimax rates based on curvature [Hazan, 2016]:

Convex ℓ_t	\sqrt{T}	GD with $\eta_t \propto \frac{1}{\sqrt{t}}$
Strongly convex ℓ_t	$\ln T$	GD with $\eta_t \propto \frac{1}{t}$
Exp-concave ℓ_t	$d \ln T$	ONS with $\eta \propto 1$

Limitations:

- ▶ Different method in each case! (Requires sophisticated users.)
- ▶ Theoretical tuning of η_t **very conservative**
- ▶ What if curvature varies between rounds?
- ▶ In many applications data are **stochastic** (i.i.d.) Should be easier than worst case...

Need Adaptive Methods!

Existing Adaptivity Results

- ▶ [Bartlett, Hazan, and Rakhlin, 2007], [Do et al., 2009]: Adaptive GD: **strongly convex + general convex**

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 (ℓ_∞ -ball) instead of ℓ_2 -ball
- ▶ [Hazan and Kale, 2010], [Chiang et al., 2012]: linear functions ℓ_t that vary little over time
- ▶ [Orabona, Crammer, and Cesa-Bianchi, 2015]: data-dependent time-varying regularizers
- ▶ ...

Key techniques:

- ▶ Adaptive tuning of learning rate η_t

Main difficulty:

- ▶ Overhead for learning η_t dominates benefits of tuning η_t

MetaGrad: Multiple Eta Gradient Algorithm

Theorem (Van Erven, Koolen, 2016)

MetaGrad's Regret_T^u is bounded by

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

where

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

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

MetaGrad: Multiple Eta Gradient Algorithm

Theorem (Van Erven, Koolen, 2016)

MetaGrad's Regret_T^u is bounded by

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

where

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

- ▶ By convexity, $\ell_t(\mathbf{w}_t) - \ell_t(\mathbf{u}) \leq (\mathbf{w}_t - \mathbf{u})^\top \mathbf{g}_t$.
- ▶ Covariance: $\mathbf{g}_t \mathbf{g}_t^\top \propto \mathbf{X}_t \mathbf{X}_t^\top$ when $\ell_t(\mathbf{w}) = h_t(\langle \mathbf{w}, \mathbf{X}_t \rangle)$
e.g. hinge, logistic, squared loss

MetaGrad: Multiple Eta Gradient Algorithm

Theorem (Van Erven, Koolen, 2016)

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

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

where

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

- ▶ By convexity, $\ell_t(\mathbf{w}_t) - \ell_t(\mathbf{u}) \leq (\mathbf{w}_t - \mathbf{u})^\top \mathbf{g}_t$.
- ▶ Covariance: $\mathbf{g}_t \mathbf{g}_t^\top \propto \mathbf{X}_t \mathbf{X}_t^\top$ when $\ell_t(\mathbf{w}) = h_t(\langle \mathbf{w}, \mathbf{X}_t \rangle)$
e.g. hinge, logistic, squared loss
- ▶ Optimal learning rate η depends on V_T^u , but \mathbf{u} unknown!
Solution: aggregate **multiple learning rates** at almost no cost

Consequences

1. Non-stochastic adaptation to curvature:

Convex l_t	$\sqrt{T \ln \ln T}$
Exp-concave l_t	$d \ln T$
Fixed convex $l_t = l$	$d \ln T$

Consequences

1. Non-stochastic adaptation to curvature:

Convex l_t	$\sqrt{T \ln \ln T}$
Exp-concave l_t	$d \ln T$
Fixed convex $l_t = l$	$d \ln T$

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

Consequences

1. Non-stochastic adaptation to curvature:

Convex l_t	$\sqrt{T \ln \ln T}$
Exp-concave l_t	$d \ln T$
Fixed convex $l_t = l$	$d \ln T$

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

2. Stochastic without curvature

Suppose l_t i.i.d. with stochastic optimum

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

Absolute loss* $l_t(u) = u - X_t $	$\ln T$
Hinge loss* $\max\{0, 1 - Y_t \langle u, X_t \rangle\}$	$d \ln T$
(B, β)-Bernstein	$(Bd \ln T)^{1/(2-\beta)} T^{(1-\beta)/(2-\beta)}$

*Conditions apply

1. Directional Derivative Condition

Corollary (Van Erven, Koolen, 2016)

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

$$\ell_t(\mathbf{u}) \geq \ell_t(\mathbf{w}) + a(\mathbf{u} - \mathbf{w})^\top \nabla \ell_t(\mathbf{w}) + b((\mathbf{u} - \mathbf{w})^\top \nabla \ell_t(\mathbf{w}))^2$$

for all $\mathbf{w} \in \mathcal{U}$, then $O(d \ln T)$ regret w.r.t. \mathbf{u} .

$a = 1$

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

General a

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

2. Bernstein Condition for Online Learning

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

Standard Bernstein condition:

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

2. Bernstein Condition for Online Learning

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

Standard Bernstein condition:

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

Replace by **weaker linearized version:**

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

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

2. Bernstein Condition for Online Learning

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

Standard Bernstein condition:

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

Replace by **weaker linearized version:**

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

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

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

Standard Bernstein condition:

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

Replace by **weaker linearized version:**

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

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

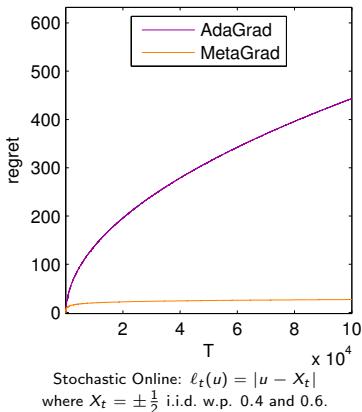
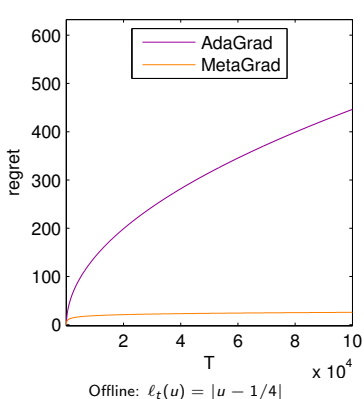
Hinge loss (with $G = D = 1$): $\beta = 1$, $B = \frac{2\lambda_{\max}(\mathbb{E}[\mathbf{X}\mathbf{X}^\top])}{\|\mathbb{E}[\mathbf{Y}\mathbf{X}]\|}$

Theorem (Koolen, Grünwald, Van Erven, 2016)

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

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

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

MetaGrad Algorithm

Second-order **surrogate loss** for each η 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 η 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 η produces \mathbf{w}_t^η 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 η 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 η produces \mathbf{w}_t^η 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

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

MetaGrad Algorithm

Second-order **surrogate loss** for each η 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 η produces \mathbf{w}_t^η 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

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

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

MetaGrad Algorithm

Second-order **surrogate loss** for each η 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 η produces \mathbf{w}_t^η 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 η 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 η produces \mathbf{w}_t^η 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 η 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 η produces \mathbf{w}_t^η 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

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

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

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

MetaGrad Master

Goal: aggregate slave predictions \mathbf{w}_t^η for all η 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 ℓ_t^η for all η simultaneously

Compute **exponential weights** with performance of each η 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 ℓ_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 ℓ_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 η , minimize regret w.r.t. exp-concave surrogate ℓ_t^η .

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 η , minimize regret w.r.t. exp-concave surrogate ℓ_t^η .

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**

Summary

1. Statistical learning

- ▶ Limits of minimax analysis pretty well understood
- ▶ Margin/Bernstein condition relates variance to excess risk
- ▶ Adaptive algorithms known to exploit 'easy data'

2. Online sequential prediction

- ▶ Most results still about minimax analysis
- ▶ MetaGrad:
 - ▶ new adaptive algorithm
 - ▶ with new variance bound
- ▶ Variance bound implies fast rates in:
 - ▶ all known cases: exp-concave, strong convex
 - ▶ new cases characterized by new online version of Bernstein condition

Future Work

Computation

- ▶ Online learning often applied in high dimensions d
- ▶ Gradient descent: $O(d)$ work per step
- ▶ MetaGrad: $O(d^2)$ work per step + projection on domain
- ▶ Need to speed up MetaGrad to work in high dimensions

Important Application: Deep Learning

- ▶ Online learning (e.g. gradient descent) is used to train large neural networks
- ▶ But loss is not convex, so theory breaks down
- ▶ Make it work experimentally

Papers

- ▶ T. van Erven and W. M. Koolen. **Metagrad: Multiple learning rates in online learning**. In Advances in Neural Information Processing Systems 29 (NIPS), pages 3666–3674, 2016.
- ▶ W. M. Koolen, P. Grünwald, and T. van Erven. **Combining adversarial guarantees and stochastic fast rates in online learning**. In Advances in Neural Information Processing Systems 29 (NIPS), pages 4457–4465, 2016.

References

- P. L. Bartlett, E. Hazan, and A. Rakhlin. Adaptive online gradient descent. In *Advances in Neural Information Processing Systems 20 (NIPS)*, 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 *Proceedings of the 25th Annual Conference 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 *Proceedings of the 26th Annual International Conference 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. Introduction to online optimization. Draft, April 10, 2016, available from ocobook.cs.princeton.edu, 2016.
- 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.
- H. Luo, A. Agarwal, N. Cesa-Bianchi, and J. Langford. Efficient second order online learning by sketching. In *Advances in Neural Information Processing Systems 29 (NIPS)*, 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.
- A. B. Tsybakov. Optimal aggregation of classifiers in statistical learning. *The Annals of Statistics*, 32(1):135–166, 2004.