

STATS 701 – Theory of Reinforcement Learning

Online Learning with Full Information

Ambuj Tewari

Associate Professor, Department of Statistics, University of Michigan
tewaria@umich.edu

<https://ambujtewari.github.io/stats701-winter2021/>

Slide Credits: Wouter Koolen @ CWI Amsterdam, The Netherlands

Winter 2021

1. The Experts Problem; Exponential Weights

2. Two Peeks Beyond the Basics

- Follow the Regularized Leader and Mirror Descent
- Online Quadratic Optimization; Online Newton Step

3. Conclusion and Extensions

From Learning Parameters to Picking Actions

We now turn to the second elementary online learning task.

- Decision Theoretic Online Learning
- Experts setting (also: Hedge setting)
- Prediction with Expert Advice

From Learning Parameters to Picking Actions

We now turn to the second elementary online learning task.

- Decision Theoretic Online Learning
- Experts setting (also: Hedge setting)
- Prediction with Expert Advice

Protocol: Prediction With Expert Advice

Given: game length T , number K of experts

For $t = 1, 2, \dots, T$,

- Learner chooses a distribution $w_t \in \Delta_K$ on K “experts”.
- Adversary reveals loss vector $\ell_t \in [0, 1]^K$.
- Learner’s loss is the **dot loss** $w_t^\top \ell_t = \sum_{k=1}^K w_t^k \ell_t^k$

From Learning Parameters to Picking Actions

We now turn to the second elementary online learning task.

- Decision Theoretic Online Learning
- Experts setting (also: Hedge setting)
- Prediction with Expert Advice

Protocol: Prediction With Expert Advice

Given: game length T , number K of experts

For $t = 1, 2, \dots, T$,

- Learner chooses a distribution $w_t \in \Delta_K$ on K “experts”.
- Adversary reveals loss vector $\ell_t \in [0, 1]^K$.
- Learner’s loss is the **dot loss** $w_t^\top \ell_t = \sum_{k=1}^K w_t^k \ell_t^k$

The goal: control the **regret** (w.r.t. the best expert after T rounds)

$$\mathcal{R}_T = \sum_{t=1}^T w_t^\top \ell_t - \min_{k \in [K]} \sum_{t=1}^T \ell_t^k$$

using a computationally **efficient** algorithm for learner.

Let's apply what we know

Observations:

- Dot loss $\mathbf{u} \mapsto \mathbf{u}^\top \ell_t$ is *linear* (hence convex).
- Gradient $\ell_t \in [0, 1]^K$ bounded by $\|\ell_t\| \leq \sqrt{K}$.
- Probability simplex Δ_K is contained in unit ball.

So: Instance of Online Convex Optimization.

OGD with $D = 1$ and $G = \sqrt{K}$ gives $\mathcal{R}_T \leq \sqrt{KT}$.

Let's apply what we know

Observations:

- Dot loss $\mathbf{u} \mapsto \mathbf{u}^\top \ell_t$ is *linear* (hence convex).
- Gradient $\ell_t \in [0, 1]^K$ bounded by $\|\ell_t\| \leq \sqrt{K}$.
- Probability simplex Δ_K is contained in unit ball.

So: Instance of Online Convex Optimization.

OGD with $D = 1$ and $G = \sqrt{K}$ gives $\mathcal{R}_T \leq \sqrt{KT}$.

Q: **Optimal?**

Let's apply what we know

Observations:

- Dot loss $\mathbf{u} \mapsto \mathbf{u}^\top \ell_t$ is *linear* (hence convex).
- Gradient $\ell_t \in [0, 1]^K$ bounded by $\|\ell_t\| \leq \sqrt{K}$.
- Probability simplex Δ_K is contained in unit ball.

So: Instance of Online Convex Optimization.

OGD with $D = 1$ and $G = \sqrt{K}$ gives $\mathcal{R}_T \leq \sqrt{KT}$.

Q: **Optimal?**

Maybe not. There are no points with loss difference \sqrt{K} in the simplex ...

Exponential Weights / Hedge Algorithm

Algorithm: Exponential Weights (EW)

EW with *learning rate* $\eta > 0$ plays weights in round t :

$$w_t^k = \frac{e^{-\eta \sum_{s=1}^{t-1} \ell_s^k}}{\sum_{j=1}^K e^{-\eta \sum_{s=1}^{t-1} \ell_s^j}}. \quad (\text{EW})$$

Exponential Weights / Hedge Algorithm

Algorithm: Exponential Weights (EW)

EW with *learning rate* $\eta > 0$ plays weights in round t :

$$w_t^k = \frac{e^{-\eta \sum_{s=1}^{t-1} \ell_s^k}}{\sum_{j=1}^K e^{-\eta \sum_{s=1}^{t-1} \ell_s^j}}. \quad (\text{EW})$$

or, equivalently, $w_1^k = \frac{1}{K}$ and

$$w_{t+1}^k = \frac{w_t^k e^{-\eta \ell_t^k}}{\sum_{j=1}^K w_t^j e^{-\eta \ell_t^j}} \quad (\text{EW, incremental})$$

Exponential Weights / Hedge Algorithm

Algorithm: Exponential Weights (EW)

EW with learning rate $\eta > 0$ plays weights in round t :

$$w_t^k = \frac{e^{-\eta \sum_{s=1}^{t-1} \ell_s^k}}{\sum_{j=1}^K e^{-\eta \sum_{s=1}^{t-1} \ell_s^j}}. \quad (\text{EW})$$

or, equivalently, $w_1^k = \frac{1}{K}$ and

$$w_{t+1}^k = \frac{w_t^k e^{-\eta \ell_t^k}}{\sum_{j=1}^K w_t^j e^{-\eta \ell_t^j}} \quad (\text{EW, incremental})$$

Theorem (EW regret bd, Freund and Schapire 1997)

The regret of EW is bounded by $\mathcal{R}_T \leq \frac{\ln K}{\eta} + T \frac{\eta}{8}$.

Corollary

Tuning $\eta = \sqrt{\frac{8 \ln K}{T}}$ yields $\mathcal{R}_T \leq \sqrt{T/2 \ln K}$.

EW Analysis

Applying *Hoeffding's Lemma* to the loss of each round gives

$$\sum_{t=1}^T w_t^\top \ell_t \leq \underbrace{\sum_{t=1}^T \left(\frac{-1}{\eta} \ln \left(\sum_{k=1}^K w_t^k e^{-\eta \ell_t^k} \right) \right)}_{\text{"mix loss"}} + \underbrace{\eta/8}_{\text{overhead}}$$

Crucial observation is that cumulative mix loss *telescopes*

$$\begin{aligned} \sum_{t=1}^T \frac{-1}{\eta} \ln \left(\sum_{k=1}^K w_t^k e^{-\eta \ell_t^k} \right) &= \sum_{t=1}^T \frac{-1}{\eta} \ln \left(\sum_{k=1}^K \frac{e^{-\eta \sum_{s=1}^{t-1} \ell_s^k}}{\sum_{j=1}^K e^{-\eta \sum_{s=1}^{t-1} \ell_s^j}} e^{-\eta \ell_t^k} \right) \\ &= \sum_{t=1}^T \frac{-1}{\eta} \ln \left(\frac{\sum_{k=1}^K e^{-\eta \sum_{s=1}^t \ell_s^k}}{\sum_{j=1}^K e^{-\eta \sum_{s=1}^{t-1} \ell_s^j}} \right) \\ &\stackrel{\text{telescopes}}{=} \frac{-1}{\eta} \ln \left(\sum_{k=1}^K e^{-\eta \sum_{t=1}^T \ell_t^k} \right) + \frac{\ln K}{\eta} \\ &\leq \min_{k \in [K]} \sum_{t=1}^T \ell_t^k + \frac{\ln K}{\eta}. \end{aligned}$$

Summary so far

Balancing act: “model complexity” vs “overfitting”

Theorem (OGD)

$$\mathcal{R}_T \leq \frac{D^2}{2\eta} + \frac{\eta}{2} G^2 T$$

Theorem (EW)

$$\mathcal{R}_T \leq \frac{\ln K}{\eta} + \frac{\eta}{8} T$$

Summary so far

Balancing act: “model complexity” vs “overfitting”

Theorem (OGD)

$$\mathcal{R}_T \leq \frac{D^2}{2\eta} + \frac{\eta}{2} G^2 T$$

Theorem (EW)

$$\mathcal{R}_T \leq \frac{\ln K}{\eta} + \frac{\eta}{8} T$$

Generates many follow-up questions:

- What if horizon T is not fixed? Anytime guarantees?
- What if gradient bound G is not known a priori?
- Can we have the actual gradient norms?
- What if model complexity (D) is not known? Not uniformly bounded? See Orabona and Cutkosky ICML'20 tutorial.

Summary so far

Balancing act: “model complexity” vs “overfitting”

Theorem (OGD)

$$\mathcal{R}_T \leq \frac{D^2}{2\eta} + \frac{\eta}{2}G^2T$$

Theorem (EW)

$$\mathcal{R}_T \leq \frac{\ln K}{\eta} + \frac{\eta}{8}T$$

Generates many follow-up questions:

- What if horizon T is not fixed? Anytime guarantees?
- What if gradient bound G is not known a priori?
- Can we have the actual gradient norms?
- What if model complexity (D) is not known? Not uniformly bounded? See Orabona and Cutkosky ICML'20 tutorial.

Need refined analyses \Rightarrow Restarts (doubling trick), decreasing η_t (AdaGrad/AdaHedge), learning the learning rate η (MetaGrad), ...

Summary so far

Balancing act: “model complexity” vs “overfitting”

Theorem (OGD)

$$\mathcal{R}_T \leq \frac{D^2}{2\eta} + \frac{\eta}{2} G^2 T$$

Theorem (EW)

$$\mathcal{R}_T \leq \frac{\ln K}{\eta} + \frac{\eta}{8} T$$

Generates many follow-up questions:

- What if horizon T is not fixed? Anytime guarantees?
- What if gradient bound G is not known a priori?
- Can we have the actual gradient norms?
- What if model complexity (D) is not known? Not uniformly bounded? See Orabona and Cutkosky ICML'20 tutorial.

Need refined analyses \Rightarrow Restarts (doubling trick), decreasing η_t (AdaGrad/AdaHedge), learning the learning rate η (MetaGrad), ...
Active research area!

FTRL/MD “sneak peek”

Q: What if my **domain** does not look like either ball or simplex?

FTRL/MD “sneak peek”

Q: What if my domain does not look like either ball or simplex?

Algorithm: Follow the Regularized Leader (FTRL) (with linearized losses)

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{u} \in \mathcal{U}} \sum_{s=1}^t f_s(\mathbf{u}) + \frac{1}{\eta} R(\mathbf{u})$$

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{u} \in \mathcal{U}} \sum_{s=1}^t \langle \mathbf{u}, \nabla f_s(\mathbf{w}_s) \rangle + \frac{1}{\eta} R(\mathbf{u})$$

Algorithm: Mirror Descent (MD)

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{u} \in \mathcal{U}} \langle \mathbf{u}, \nabla f_t(\mathbf{w}_t) \rangle + \frac{1}{\eta} B(\mathbf{u} \| \mathbf{w}_t)$$

FTRL/MD “sneak peek”

Q: What if my domain does not look like either ball or simplex?

Algorithm: Follow the Regularized Leader (FTRL) (with linearized losses)

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{u} \in \mathcal{U}} \sum_{s=1}^t f_s(\mathbf{u}) + \frac{1}{\eta} R(\mathbf{u})$$

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{u} \in \mathcal{U}} \sum_{s=1}^t \langle \mathbf{u}, \nabla f_s(\mathbf{w}_s) \rangle + \frac{1}{\eta} R(\mathbf{u})$$

Algorithm: Mirror Descent (MD)

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{u} \in \mathcal{U}} \langle \mathbf{u}, \nabla f_t(\mathbf{w}_t) \rangle + \frac{1}{\eta} B(\mathbf{u} \| \mathbf{w}_t)$$

	Regularizer R	Bregman Divergence B
Examples:	OGD	sq. Euclidean norm
	EW	Shannon entropy
		sq. Euclidean distance
		Kullback-Leibler divergence

FTRL/MD “sneak peek”

Q: What if my **domain** does not look like either ball or simplex?

Algorithm: Follow the Regularized Leader (FTRL) (with linearized losses)

$$w_{t+1} = \arg \min_{u \in \mathcal{U}} \sum_{s=1}^t f_s(u) + \frac{1}{\eta} R(u)$$

$$w_{t+1} = \arg \min_{u \in \mathcal{U}} \sum_{s=1}^t \langle u, \nabla f_s(w_s) \rangle + \frac{1}{\eta} R(u)$$

Algorithm: Mirror Descent (MD)

$$w_{t+1} = \arg \min_{u \in \mathcal{U}} \langle u, \nabla f_t(w_t) \rangle + \frac{1}{\eta} B(u \| w_t)$$

	Regularizer R	Bregman Divergence B
Examples: OGD	sq. Euclidean norm	sq. Euclidean distance
EW	Shannon entropy	Kullback-Leibler divergence

Other entropies: Burg, Tsallis, Von Neumann, ... Connections to continuous exponential weights [van der Hoeven et al., 2018].

FTRL/MD “sneak peak” performance

Algorithm: Follow the Regularized Leader (FTRL) (with linearized losses)

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{u} \in \mathcal{U}} \sum_{s=1}^t f_s(\mathbf{u}) + \frac{1}{\eta} R(\mathbf{u})$$

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{u} \in \mathcal{U}} \sum_{s=1}^t \langle \mathbf{u}, \nabla f_s(\mathbf{w}_s) \rangle + \frac{1}{\eta} R(\mathbf{u})$$

Algorithm: Mirror Descent

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{u} \in \mathcal{U}} \langle \mathbf{u}, \nabla f_t(\mathbf{w}_t) \rangle + \frac{1}{\eta} B(\mathbf{u} \| \mathbf{w}_t)$$

FTRL/MD “sneak peak” performance

Algorithm: Follow the Regularized Leader (FTRL) (with linearized losses)

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{u} \in \mathcal{U}} \sum_{s=1}^t f_s(\mathbf{u}) + \frac{1}{\eta} R(\mathbf{u})$$

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{u} \in \mathcal{U}} \sum_{s=1}^t \langle \mathbf{u}, \nabla f_s(\mathbf{w}_s) \rangle + \frac{1}{\eta} R(\mathbf{u})$$

Algorithm: Mirror Descent

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{u} \in \mathcal{U}} \langle \mathbf{u}, \nabla f_t(\mathbf{w}_t) \rangle + \frac{1}{\eta} B(\mathbf{u} \| \mathbf{w}_t)$$

Theorem (AdaFTRL, Orabona and Pál 2015)

Fix a norm $\|\cdot\|$ with associated dual norm $\|\cdot\|_*$. Let $R : \mathcal{U} \rightarrow [0, D^2]$ be strongly convex w.r.t. $\|\cdot\|$. AdaFTRL ensures

$$\mathcal{R}_T \leq 2D \sqrt{\sum_{t=1}^T \|\nabla f_t(\mathbf{w}_t)\|_*^2} + 2 \cdot \text{loss range}.$$

Quadratic Losses

So far we used convexity to “linearise”

$$f_t(\mathbf{u}) \geq f_t(\mathbf{w}_t) + \langle \mathbf{u} - \mathbf{w}_t, \nabla f_t(\mathbf{w}_t) \rangle,$$

and our methods essentially operated on linear losses. But what if we **know** there is curvature?

- How to **represent/quantify** curvature?
- How to **efficiently** manipulate curvature?
- How much can we reduce the regret?

Curvature assumptions

Assumption: Quadratic loss lower bound

There is a matrix $M_t \succeq 0$ such that

$$f_t(\mathbf{u}) \geq \underbrace{f_t(\mathbf{w}_t) + \langle \mathbf{u} - \mathbf{w}_t, \nabla f_t(\mathbf{w}_t) \rangle + \frac{1}{2}(\mathbf{u} - \mathbf{w}_t)^\top M_t (\mathbf{u} - \mathbf{w}_t)}_{=: q_t(\mathbf{u})}$$

for each $\mathbf{u} \in \mathcal{U}$.

Curvature assumptions

Assumption: Quadratic loss lower bound

There is a matrix $M_t \succeq \mathbf{0}$ such that

$$f_t(\mathbf{u}) \geq \underbrace{f_t(\mathbf{w}_t) + \langle \mathbf{u} - \mathbf{w}_t, \nabla f_t(\mathbf{w}_t) \rangle + \frac{1}{2}(\mathbf{u} - \mathbf{w}_t)^\top M_t (\mathbf{u} - \mathbf{w}_t)}_{=: q_t(\mathbf{u})}$$

for each $\mathbf{u} \in \mathcal{U}$.

Two main classes of instances

- squared Euclidean distance: $f_t(\mathbf{u}) = \frac{1}{2} \|\mathbf{u} - \mathbf{x}_t\|^2$ satisfies the assumption with $M_t = I$. More generally, **strongly convex** functions have $M_t \propto I$.
- linear regression: $f_t(\mathbf{u}) = (y_t - \langle \mathbf{u}, \mathbf{x}_t \rangle)^2$ satisfies the assumption with $M_t = \mathbf{x}_t \mathbf{x}_t^\top$. More generally, **exp-concave** functions have $M_t \propto \nabla_t f_t(\mathbf{w}_t) \nabla_t f_t(\mathbf{w}_t)^\top$.

ONS Algorithm

Algorithm: Online Newton Step (FTRL variant)

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{u} \in \mathcal{U}} \sum_{s=1}^t q_s(\mathbf{u}) + \frac{1}{2} \|\mathbf{u}\|^2$$

Computing the iterate \mathbf{w}_{t+1} amounts to minimising a convex quadratic. Often (depending on \mathcal{U}) **closed-form solution** or **1d line search**.

- For $M_t \propto \mathbf{I}$, takes $O(d)$ per round.
- For rank-one M_t , can do update in $O(d^2)$ per round.
- In both cases, need to take care of projection onto \mathcal{U} .

ONS Performance

Algorithm: Online Newton Step (FTRL version)

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{u} \in \mathcal{U}} \sum_{s=1}^t q_s(\mathbf{u}) + \frac{1}{2} \|\mathbf{u}\|^2$$

Theorem (ONS strcvx bd, Hazan et al. 2006)

For the strongly convex case $\mathbf{M}_t \propto \mathbf{I}$, ONS guarantees

$$\mathcal{R}_T = O(\ln T)$$

Algorithm reduces to OGD with specific decreasing step-size η_t

Theorem (ONS expcv bd, Hazan et al. 2006)

For the exp-concave case $\mathbf{M}_t \propto \mathbf{g}_t \mathbf{g}_t^T$, ONS guarantees

$$\mathcal{R}_T = O(d \ln T)$$

ONS Discussion

- Convex quadratics closed under taking sums. Run-time independent of T .
- Curvature gives huge reduction in regret: \sqrt{T} to $\ln T$.
- Matrix **sketching** techniques allow trading off run-time $O(d^2)$ vs $O(d)$ with regret $O(\ln T)$ vs $O(\sqrt{T})$ [Luo et al., 2016].

Conclusion

- Online Learning a powerful and versatile tool
- Environment-as-black-box. Adversarial.
- Foundation for optimization, statistical learning, games, . . .
- Techniques we saw here will reappear when we discuss adversarial bandits and adversarial MDPs

Conclusion

- Online Learning a powerful and versatile tool
- Environment-as-black-box. Adversarial.
- Foundation for optimization, statistical learning, games, ...
- Techniques we saw here will reappear when we discuss adversarial bandits and adversarial MDPs

Some (of many) cool things we left out:

- First-order (small loss) and second-order (small variance) bounds
- Adaptivity to friendly stochastic environments (best of both worlds, interpolation)
- Optimistic MD (predicting the upcoming gradient)
- Non-stationarity (tracking, adaptive/dynamic regret, path length)
- Beyond convexity (star-convex, geometrically convex, ...)
- Supervised Learning and (stochastic) complexities (VC, Littlestone, Rademacher, ...)

- Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *J. Comput. Syst. Sci.*, 55(1): 119–139, August 1997.
- Elad Hazan, Adam Kalai, Satyen Kale, and Amit Agarwal. Logarithmic regret algorithms for online convex optimization. In *Learning Theory*, pages 499–513, 2006.
- Haipeng Luo, Alekh Agarwal, Nicolò Cesa-Bianchi, and John Langford. Efficient second order online learning by sketching. In *Advances in Neural Information Processing Systems 29*, pages 902–910. 2016.
- Francesco Orabona and Dávid Pál. Scale-free algorithms for online linear optimization. In *Algorithmic Learning Theory*, pages 287–301, 2015.
- Dirk van der Hoeven, Tim van Erven, and Wojciech Kotłowski. The many faces of exponential weights in online learning. volume 75 of *Proceedings of Machine Learning Research*, pages 2067–2092, 06–09 Jul 2018.