# 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

# Working Definitions

**Context:** interactive decision making in unknown environment
**Aim**: Design systems to amass reward in many environments.

**Context:** interactive decision making in unknown environment
**Aim**: Design systems to amass reward in many environments.

**Main distinction**: model of environment

- **Reinforcement Learning** action affects future state
- **Bandits** action affects observation
- **Full Inf. Online Learning** action affects reward

Coming up:

(1) Full Information Online Learning

(2) Bandit Problems (or just "Bandits")

(3) Regret analysis in RL

1. Two Basic Problems
   - Online Convex Optimization; Online Gradient Descent
   - 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

# Setup

- Focus on losses (negative rewards)
- Model Environment as Adversary
- Online Convex Optimization (OCO) abstraction.

# OCO Problem

**Protocol: Online Convex Optimization**

Given: game length $T$, convex action space $\mathcal{U} \subseteq \mathbb{R}^d$

For $t = 1, 2, \ldots, T$,
- The learner picks action $\boldsymbol{w}_t \in \mathcal{U}$
- The adversary picks convex loss $f_t : \mathcal{U} \to \mathbb{R}$
- The learner observes $f_t$    ◁ full information
- The learner incurs loss $f_t(\boldsymbol{w}_t)$

# OCO Problem

## Protocol: Online Convex Optimization

Given: game length $T$, convex action space $\mathcal{U} \subseteq \mathbb{R}^d$

For $t = 1, 2, \ldots, T$,
- The learner picks action $\boldsymbol{w}_t \in \mathcal{U}$
- The adversary picks convex loss $f_t : \mathcal{U} \to \mathbb{R}$
- The learner observes $f_t$    $\triangleleft$ full information
- The learner incurs loss $f_t(\boldsymbol{w}_t)$

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

$$\mathcal{R}_T = \sum_{t=1}^{T} f_t(\boldsymbol{w}_t) - \min_{\boldsymbol{u} \in \mathcal{U}} \sum_{t=1}^{T} f_t(\boldsymbol{u})$$

using a computationally efficient algorithm for learner.

## Design Principle

Learner needs to "chase" the best point $\arg\min_{\boldsymbol{u} \in \mathcal{U}} \sum_{t=1}^{T} f_t(\boldsymbol{w}_t)$. But doing so naively overfits.

Idea: add regularization. Two manifestations:

- Penalization "FTRL style"
- Update iterates, but only slowly "MD style"

Will see examples of both. For our purposes, these are roughly equivalent

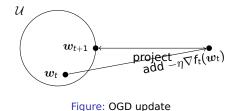# Online Gradient Descent (OGD) Algorithm

Let $\mathcal{U}$ be a convex set containing $\mathbf{0}$. Fix a learning rate $\eta > 0$.

### Algorithm: Online Gradient Descent (OGD)

OGD with learning rate $\eta > 0$ plays

$$\boldsymbol{w}_1 = \mathbf{0} \qquad \text{and} \qquad \boldsymbol{w}_{t+1} = \Pi_{\mathcal{U}}\left(\boldsymbol{w}_t - \eta \nabla f_t(\boldsymbol{w}_t)\right)$$

where $\Pi_{\mathcal{U}}(\boldsymbol{w}) = \arg\min_{\boldsymbol{u} \in \mathcal{U}} \|\boldsymbol{u} - \boldsymbol{w}\|$ is the projection onto $\mathcal{U}$.



Figure: OGD update

# Online Gradient Descent Result

### Algorithm: OGD

$$\boldsymbol{w_1} = \boldsymbol{0} \qquad \text{and} \qquad \boldsymbol{w_{t+1}} \,=\, \Pi_{\mathcal{U}}\left(\boldsymbol{w_t} - \eta \nabla f_t(\boldsymbol{w_t})\right)$$

### Assumption: Boundedness

Bounded domain $\max_{\boldsymbol{u} \in \mathcal{U}} \|\boldsymbol{u}\| \leq D$ and gradients $\|\nabla f_t(\boldsymbol{w_t})\| \leq G$.

# Online Gradient Descent Result

**Algorithm: OGD**

$$\boldsymbol{w_1} = \boldsymbol{0} \qquad \text{and} \qquad \boldsymbol{w}_{t+1} = \Pi_{\mathcal{U}}\left(\boldsymbol{w}_t - \eta \nabla f_t(\boldsymbol{w}_t)\right)$$

**Assumption: Boundedness**

Bounded domain $\max_{\boldsymbol{u} \in \mathcal{U}} \|\boldsymbol{u}\| \leq D$ and gradients $\|\nabla f_t(\boldsymbol{w}_t)\| \leq G$.

**Theorem (OGD regret bd, Zinkevich 2003)**

$$\mathcal{R}_T = \sum_{t=1}^{T} f_t(\boldsymbol{w}_t) - \min_{\boldsymbol{u} \in \mathcal{U}} \sum_{t=1}^{T} f_t(\boldsymbol{u}) \leq \frac{1}{2\eta}D^2 + \frac{\eta}{2}TG^2$$

# Online Gradient Descent Result

**Algorithm: OGD**

$$\boldsymbol{w_1} = \boldsymbol{0} \qquad \text{and} \qquad \boldsymbol{w_{t+1}} = \Pi_{\mathcal{U}}\left(\boldsymbol{w_t} - \eta \nabla f_t(\boldsymbol{w_t})\right)$$

**Assumption: Boundedness**

Bounded domain $\max_{\boldsymbol{u}\in\mathcal{U}}\|\boldsymbol{u}\| \leq D$ and gradients $\|\nabla f_t(\boldsymbol{w_t})\| \leq G$.

**Theorem (OGD regret bd, Zinkevich 2003)**

$$\mathcal{R}_T = \sum_{t=1}^{T} f_t(\boldsymbol{w_t}) - \min_{\boldsymbol{u}\in\mathcal{U}} \sum_{t=1}^{T} f_t(\boldsymbol{u}) \leq \frac{1}{2\eta}D^2 + \frac{\eta}{2}TG^2$$

**Corollary**

*Tuning $\eta = \frac{D}{G\sqrt{T}}$ results in $\mathcal{R}_T \leq DG\sqrt{T}$ .*

# Online Gradient Descent Result

## Algorithm: OGD

$$\boldsymbol{w_1} = \boldsymbol{0} \qquad \text{and} \qquad \boldsymbol{w_{t+1}} = \Pi_{\mathcal{U}}\left(\boldsymbol{w_t} - \eta \nabla f_t(\boldsymbol{w_t})\right)$$

## Assumption: Boundedness

Bounded domain $\max_{\boldsymbol{u} \in \mathcal{U}} \|\boldsymbol{u}\| \leq D$ and gradients $\|\nabla f_t(\boldsymbol{w_t})\| \leq G$.

## Theorem (OGD regret bd, Zinkevich 2003)

$$\mathcal{R}_T = \sum_{t=1}^{T} f_t(\boldsymbol{w_t}) - \min_{\boldsymbol{u} \in \mathcal{U}} \sum_{t=1}^{T} f_t(\boldsymbol{u}) \leq \frac{1}{2\eta} D^2 + \frac{\eta}{2} T G^2$$

## Corollary

*Tuning* $\eta = \frac{D}{G\sqrt{T}}$ *results in* $\mathcal{R}_T \leq DG\sqrt{T}$ .

Sublinear regret: learning overhead per round $\rightarrow 0$.

# Proof of OGD regret bound

Using convexity, we may analyse the tangent upper bound

$$f_t(\boldsymbol{w}_t) - f_t(\boldsymbol{u}) \leq \langle \boldsymbol{w}_t - \boldsymbol{u}, \nabla f_t(\boldsymbol{w}_t) \rangle$$

Moreover,

$$
\begin{aligned}
\|\boldsymbol{w}_{t+1} - \boldsymbol{u}\|^2 &= \|\Pi_{\mathcal{U}}\left(\boldsymbol{w}_t - \eta \nabla f_t(\boldsymbol{w}_t)\right) - \boldsymbol{u}\|^2 \\
&\leq \|\boldsymbol{w}_t - \eta \nabla f_t(\boldsymbol{w}_t) - \boldsymbol{u}\|^2 \\
&= \|\boldsymbol{w}_t - \boldsymbol{u}\|^2 - 2\eta\langle \boldsymbol{w}_t - \boldsymbol{u}, \nabla f_t(\boldsymbol{w}_t)\rangle + \eta^2 \|\nabla f_t(\boldsymbol{w}_t)\|^2
\end{aligned}
$$

Hence

$$\langle \boldsymbol{w}_t - \boldsymbol{u}, \nabla f_t(\boldsymbol{w}_t) \rangle \leq \frac{\|\boldsymbol{w}_t - \boldsymbol{u}\|^2 - \|\boldsymbol{w}_{t+1} - \boldsymbol{u}\|^2}{2\eta} + \frac{\eta}{2}\|\nabla f_t(\boldsymbol{w}_t)\|^2$$

# Proof of OGD regret bound (ctd)

Summing over $T$ rounds, we find

$$\mathcal{R}_T^{\boldsymbol{u}} \leq \sum_{t=1}^{T} \langle \boldsymbol{w}_t - \boldsymbol{u}, \nabla f_t(\boldsymbol{w}_t) \rangle$$

$$\leq \underbrace{\sum_{t=1}^{T} \frac{\|\boldsymbol{w}_t - \boldsymbol{u}\|^2 - \|\boldsymbol{w}_{t+1} - \boldsymbol{u}\|^2}{2\eta}}_{\text{telescopes}} + \frac{\eta}{2} \sum_{t=1}^{T} \|\nabla f_t(\boldsymbol{w}_t)\|^2$$

$$\leq \frac{\|\boldsymbol{u}\|^2 - \|\boldsymbol{w}_{T+1} - \boldsymbol{u}\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^{T} \|\nabla f_t(\boldsymbol{w}_t)\|^2$$

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

# OCO Lower Bound

Is OGD regret bound of $\mathcal{R}_T \leq GD\sqrt{T}$ any good?

## OCO Lower Bound

Is OGD regret bound of $\mathcal{R}_T \leq GD\sqrt{T}$ any good?
Scaling with $G$ and $D$ is natural. What about $\sqrt{T}$?

# OCO Lower Bound

Is OGD regret bound of $\mathcal{R}_T \leq GD\sqrt{T}$ any good?
Scaling with $G$ and $D$ is natural. What about $\sqrt{T}$?

### Theorem

*Any OCO algorithm can be made to incur $\mathcal{R}_T = \Omega(\sqrt{T})$.*

# OCO Lower Bound

Is OGD regret bound of $\mathcal{R}_T \leq GD\sqrt{T}$ any good?
Scaling with $G$ and $D$ is natural. What about $\sqrt{T}$?

### Theorem

*Any OCO algorithm can be made to incur $\mathcal{R}_T = \Omega(\sqrt{T})$.*

### Proof (by probabilistic argument).

Consider interval $\mathcal{U} = [-1, 1]$ and linear losses $f_t(u) = x_t \cdot u$ with i.i.d.
Rademacher coefficients $x_t \in \{\pm 1\}$. *Any* algorithm has expected loss zero. The
expected loss of the best action ($\pm 1$) is $-\mathbb{E}[|\sum_{t=1}^{T} x_t|] = -\Omega(\sqrt{T})$. Then as the
expected regret is $\mathbb{E}[\mathcal{R}_T] = \Omega(\sqrt{T})$, there is a deterministic witness. □

Here, the regret arises from *overfitting* of the best point.

# OGD Discussion

- Adversarial result, super strong!
- Proof reveals it is really about linear losses.
- Matching lower bounds

Successful in practise:

- Practically all deep learning uses versions of online gradient descent (e.g. TensorFlow has AdaGrad [Duchi et al., 2011]) even though objective not convex.

# 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 $t = 1, 2, \ldots, T$,
- Learner chooses a distribution $w_t \in \triangle_K$ on $K$ "experts".
- Adversary reveals loss vector $\ell_t \in [0, 1]^K$.
- Learner's loss is the **dot loss** $w_t^\mathsf{T} \ell_t = \sum_{k=1}^K w_t^k \ell_t^k$

# 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 $t = 1, 2, \ldots, T$,
- Learner chooses a distribution $w_t \in \triangle_K$ on $K$ "experts".
- Adversary reveals loss vector $\ell_t \in [0, 1]^K$.
- Learner's loss is the **dot loss** $w_t^\mathsf{T} \ell_t = \sum_{k=1}^K w_t^k \ell_t^k$

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

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

using a computationally efficient algorithm for learner.

# Let's apply what we know

Observations:

- Dot loss $\boldsymbol{u} \mapsto \boldsymbol{u}^\mathsf{T} \boldsymbol{\ell}_t$ is *linear* (hence convex).
- Gradient $\boldsymbol{\ell}_t \in [0, 1]^K$ bounded by $\|\boldsymbol{\ell}_t\| \leq \sqrt{K}$.
- Probability simplex $\triangle_K$ 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 $\boldsymbol{u} \mapsto \boldsymbol{u}^{\mathsf{T}} \boldsymbol{\ell}_t$ is *linear* (hence convex).
- Gradient $\boldsymbol{\ell}_t \in [0,1]^K$ bounded by $\|\boldsymbol{\ell}_t\| \leq \sqrt{K}$.
- Probability simplex $\triangle_K$ 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 $\boldsymbol{u} \mapsto \boldsymbol{u}^\mathsf{T} \boldsymbol{\ell}_t$ is *linear* (hence convex).
- Gradient $\boldsymbol{\ell}_t \in [0, 1]^K$ bounded by $\|\boldsymbol{\ell}_t\| \leq \sqrt{K}$.
- Probability simplex $\triangle_K$ 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 Weigths / Hedge Algorithm

### Algorithm: Exponential Weights (EW)

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

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

# Exponential Weigths / Hedge Algorithm

## Algorithm: Exponential Weights (EW)

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

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

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

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

# Exponential Weigths / Hedge Algorithm

## Algorithm: Exponential Weights (EW)

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

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

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

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

## 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

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

Crucial observation is that cumulative mix loss *telescopes*

$$\sum_{t=1}^{T} \frac{-1}{\eta} \ln \left( \sum_{k=1}^{K} w_t^k e^{-\eta \ell_t^k} \right) \ = \ \sum_{t=1}^{T} \frac{-1}{\eta} \ln \left( \sum_{k=1}^{K} \frac{e^{-\eta \sum_{s=1}^{t-1} \ell_s^k}}{\sum_{j=1}^{K} e^{-\eta \sum_{s=1}^{t-1} \ell_s^j}} e^{-\eta \ell_t^k} \right)$$

$$= \ \sum_{t=1}^{T} \frac{-1}{\eta} \ln \left( \frac{\sum_{k=1}^{K} e^{-\eta \sum_{s=1}^{t} \ell_s^k}}{\sum_{j=1}^{K} e^{-\eta \sum_{s=1}^{t-1} \ell_s^j}} \right)$$

$$\stackrel{\text{telescopes}}{=} \ \frac{-1}{\eta} \ln \left( \sum_{k=1}^{K} e^{-\eta \sum_{t=1}^{T} \ell_t^k} \right) + \frac{\ln K}{\eta}$$

$$\leq \ \min_{k \in [K]} \sum_{t=1}^{T} \ell_t^k + \frac{\ln K}{\eta}.$$

# 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)

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

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 $\eta_t$ (AdaGrad/AdaHedge), learning the learning rate $\eta$ (MetaGrad), ...

# Summary so far

Balancing act: "model complexity" vs "overfitting"

Theorem (OGD)

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

Theorem (EW)

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

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 $\eta_t$ (AdaGrad/AdaHedge), learning the learning rate $\eta$ (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)

$$w_{t+1} = \operatorname*{arg\,min}_{u \in \mathcal{U}} \sum_{s=1}^{t} f_s(u) + \frac{1}{\eta} R(u)$$

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

Algorithm: Mirror Descent (MD)

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

# 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)**

$$w_{t+1} = \operatorname*{arg\,min}_{u \in \mathcal{U}} \sum_{s=1}^{t} f_s(u) + \frac{1}{\eta} R(u)$$

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

**Algorithm: Mirror Descent (MD)**

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

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

# 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)**

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

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

**Algorithm: Mirror Descent (MD)**

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

| Examples: | | **Regularizer** $R$ | **Bregman Divergence** $B$ |
|---|---|---|---|
| | 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)**

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

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

**Algorithm: Mirror Descent**

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

# FTRL/MD "sneak peak" performance

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

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

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

**Algorithm: Mirror Descent**

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

**Theorem (AdaFTRL, Orabona and Pál 2015)**

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

$$\mathcal{R}_T \;\le\; 2D\sqrt{\sum_{t=1}^{T} \|\nabla f_t(\boldsymbol{w}_t)\|_\star^2 + 2 \cdot \textit{loss range}}.$$

## Quadratic Losses

So far we used convexity to "linearise"

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

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 $M_t \succeq \mathbf{0}$ such that

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

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

# Curvature assumptions

## Assumption: Quadratic loss lower bound

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

$$f_t(u) \geq \underbrace{f_t(w_t) + \langle u - w_t, \nabla f_t(w_t) \rangle + \frac{1}{2}(u - w_t)^\mathsf{T} M_t (u - w_t)}_{=:q_t(u)}$$

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

Two main classes of instances

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

# ONS Algorithm

**Algorithm: Online Newton Step (FTRL variant)**

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

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

- For $\boldsymbol{M}_t \propto \boldsymbol{I}$, 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)**

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

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

*For the strongly convex case $M_t \propto \boldsymbol{I}$, ONS guarantees*

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

*Algorithm reduces to OGD with specific decreasing step-size $\eta_t$*

**Theorem (ONS expccv bd, Hazan et al. 2006)**

*For the exp-concave case $M_t \propto \boldsymbol{g}_t \boldsymbol{g}_t^\mathsf{T}$, 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, ... )

John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011.

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):119âĂŞ139, August 1997.

Yoav Freund and Robert E Schapire. Adaptive game playing using multiplicative weights. *Games and Economic Behavior*, 29(1-2):79–103, 1999.

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.

Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the Twentieth International Conference on International Conference on Machine Learning*, ICML'03, page 928âĂŞ935, 2003.