

# Lecture 18 Online Learning (Part II) and Intro to Reinforcement Learning

Lei Li, Yu-Xiang Wang

# Recap: Online Learning

- Learning with expert advice
  - A summary of regret bound: # mistakes - Oracle # of mistakes

	Consistency	Halving	Weighted Majority	Randomized WM
Realizable setting	$\min(T,  \mathcal{H} )$	$\min(T, \log \mathcal{H} )$	$\min(T, \log \mathcal{H} )$	$\min(T, \log \mathcal{H} )$
Agnostic setting	n.a.	n.a.	$(1 + \epsilon)m + \log \mathcal{H}  / \epsilon$	$\sqrt{m \log \mathcal{H} } = O(\sqrt{T \log \mathcal{H} })$

# Recap: Hedge (aka Exponential weighted average) algorithm

---

EXPONENTIAL-WEIGHTED-AVERAGE ( $N$ )

```
1 for  $i \leftarrow 1$  to  $N$  do
2      $w_{1,i} \leftarrow 1$ 
3 for  $t \leftarrow 1$  to  $T$  do
4     RECEIVE( $x_t$ )
5      $\hat{y}_t \leftarrow \frac{\sum_{i=1}^N w_{t,i} y_{t,i}}{\sum_{i=1}^N w_{t,i}}$ 
6     RECEIVE( $y_t$ )
7     for  $i \leftarrow 1$  to  $N$  do
8          $w_{t+1,i} \leftarrow w_{t,i} e^{-\eta L(\hat{y}_{t,i}, y_t)}$ 
9 return  $\mathbf{w}_{T+1}$ 
```

- Works for linear loss function in its first argument
  - That is bounded
- Also works for any convex loss function in its first argument
  - Why?

# This lecture

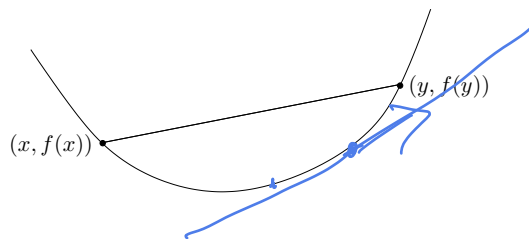
- Online Learning (Part II)
  - Online Gradient Descent
- Reinforcement Learning
  - Problem setup
  - Markov Decision Processes

# Recap: Convex functions / sets and subgradient

**Convex function:**  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  such that  $\text{dom}(f) \subseteq \mathbb{R}^n$  convex, and

$$f(tx + (1-t)y) \leq tf(x) + (1-t)f(y) \quad \text{for all } 0 \leq t \leq 1$$

and all  $x, y \in \text{dom}(f)$



- First order definition / subgradient

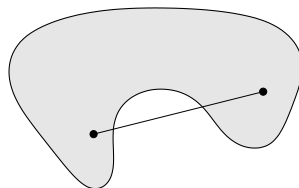
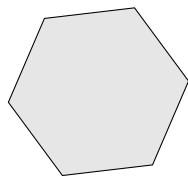
$$f(y) \geq f(x) + g_x^T (y-x) \quad \forall x, y$$

$g$  subgradient of  $f$  at  $x$

**Convex set:**  $C \subseteq \mathbb{R}^n$  such that

$$x, y \in C \implies tx + (1-t)y \in C \quad \text{for all } 0 \leq t \leq 1$$

$\Leftrightarrow f$ 's convex



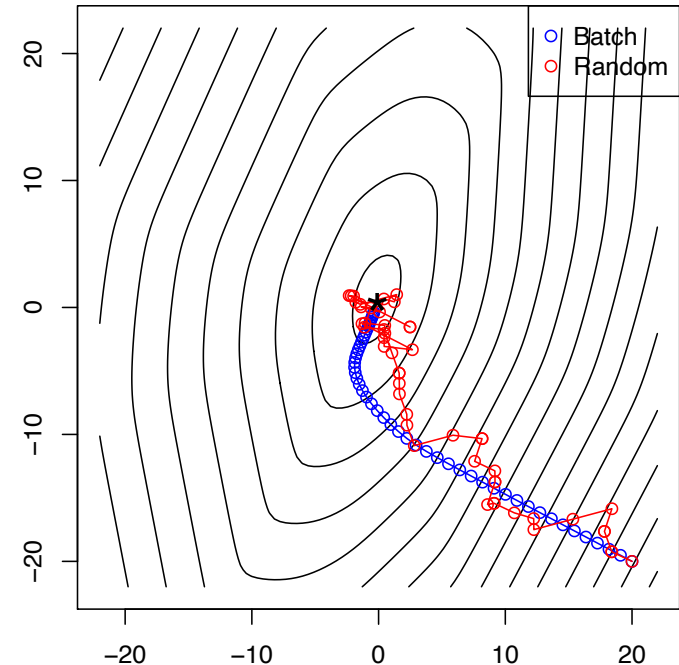
# Recap: Gradient Descent and SGD (from Lecture 4)

- Problem:  $\min_{\theta} f(\theta)$
- GD alg.:  $\theta_{t+1} = \theta_t - \eta_t \nabla f(\theta_t)$
- SGD alg.:  $\theta_{t+1} = \theta_t - \eta_t \hat{\nabla} f(\theta_t)$
  
- Example when solving ERM:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \ell(\theta, (x_i, y_i))$$

- Pick a **single** data point  $i$  uniformly at random

$$\nabla_{\theta} \ell(\theta, (x_i, y_i))$$



# (Projected) Subgradient "Descent"

- When we have constraints and non-differentiable convex functions, we can use

- Projected Subgradient method

$$\min_{x \in K} f(x)$$

$$\tilde{x}_{t+1} = x_t - \eta_t \cdot g_t$$

$g_t$  is a subgradient of  $f$  at  $x_t$

$$x_{t+1} = \Pi_K(\tilde{x}_{t+1}) = \underset{x \in K}{\operatorname{argmin}} \|x - \tilde{x}_{t+1}\|_2$$

- Projected Stochastic subgradient method

convex  $f$   
 Convergence  
 $O\left(\frac{1}{\sqrt{T}}\right)$

$$x_{t+1} = \Pi_K(x_t - \eta_t \hat{g}_t)$$

$$E[\hat{g}_t] = g_t$$

$$E[\|\hat{g}_t - g_t\|_2^2] = o(1)$$



# The problem of Online Convex Optimization

- Problem setup:

for  $t = 1, 2, \dots, T$   
learner plays  $x_t \in K$   
nature chooses  $f_t$   
learner incurs a loss  $f_t(x_t)$

- Performance metric --- regret

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

$$u = \underset{u \in K}{\text{argmin}} \sum f_t(u)$$



# Examples of OCO problems

- Example 1: Prediction with Expert Advice

$$f_t = \langle l_t, x_t \rangle \quad \|l_t\|_\infty \leq 1$$

$$K = \Delta^{n-1}$$

- Example 2: Online Linear models

$$f_t(x) = \ell_t(x^T \phi_t)$$

feature vector  $\in \mathbb{R}^d$

$$\frac{(x^T \phi_t - y_t)^2}{\log(1 + \exp(-y_t \cdot x^T \phi_t))}$$

$y_t \in \{-1, 1\}$

- Example 3: Portfolio Selection

$$f_t(x) = -\log(v_t^T x)$$

$v_t = \begin{pmatrix} 1 \\ 0.4 \\ 0.4 \\ 0.6 \end{pmatrix}$

$$-\log(v_t^T x) + \underbrace{(\log(\mathbb{1}^T x))}_{\text{reference}}$$

# Algorithm: OGD, i.e., Online (projected) (sub)Gradient Descent

- Standard projected subgradient updates

$$x_{t+1} = \text{Proj}_{\mathcal{K}} \left( x_t - \eta_t \underbrace{\nabla f_t(x_t)}_{g_{t,x_t}} \right)$$

- Assumptions needed:

- Bounded domain

$$\|x - x'\|_2 \leq D \quad \forall x, x' \in \mathcal{K}$$

- Lipschitz loss functions

$$f_t \text{ is } G\text{-Lipschitz} \quad \forall x, x' \quad |f_t(x) - f_t(x')| \leq G \|x - x'\|_2$$

$$\|g_t\|_2 \leq G \quad \forall x$$

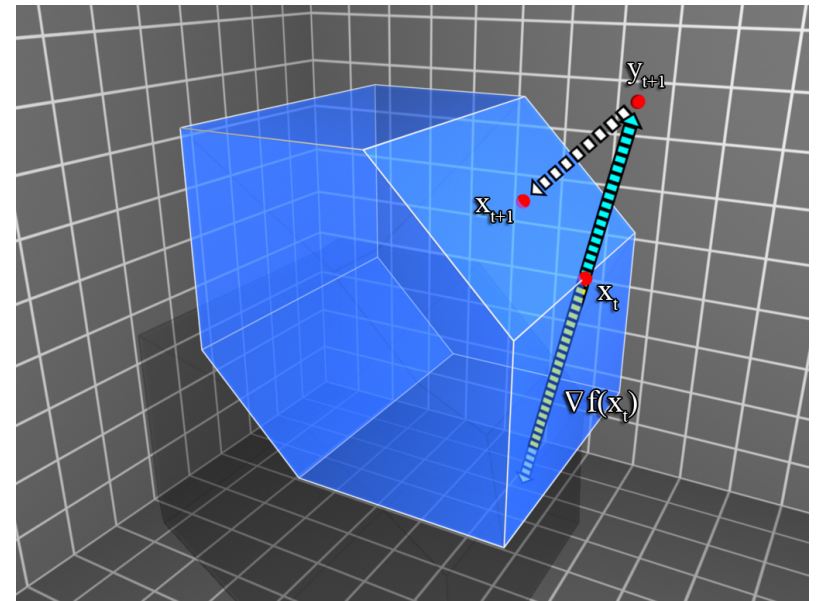
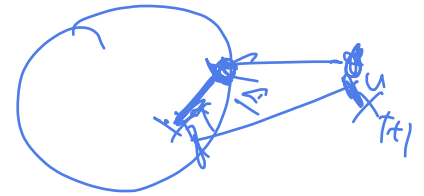


Figure 3.1: Online gradient descent: the iterate  $x_{t+1}$  is derived by advancing  $x_t$  in the direction of the current gradient  $\nabla_t$ , and projecting back into  $\mathcal{K}$ .

# Analysis of OGD



- By convex functions

$$f_t(x^*) \geq f_t(x_t) + \underbrace{g_t^T(x^* - x_t)}_{\text{subdifferential property}}$$

$x^*$

- By the update rule (and property of projection)

$$\|x_{t+1} - x^*\|^2 = \|\Pi_C(x_t - \eta_t g_t) - x^*\|^2 \leq \|x_t - \eta_t g_t - x^*\|^2 = \|x_t - x^*\|^2 + \eta_t^2 \|g_t\|^2 - 2\eta_t \underbrace{g_t^T(x_t - x^*)}_{\geq 0}$$

- Put them together!

$$f_t(x_t) - f_t(x^*) \leq g_t^T(x_t - x^*) \leq \frac{1}{2\eta_t} (\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2) + \frac{\eta_t \|g_t\|^2}{2}$$

# Analysis of OGD (continues)

- Telescoping

$$\begin{aligned}
 \sum_{t=1}^T f_t(x_t) - f_T(x^*) &\leq \frac{1}{2\eta} \left( \underbrace{\|x_1 - x^*\|^2}_{\leq D^2} - \underbrace{\|x_{T+1} - x^*\|^2}_{\leq 0} \right) + \sum_{t=1}^T \frac{\eta}{2} \underbrace{\|g_t\|^2}_{\leq G^2} \\
 &\leq \frac{1}{2\eta} D^2 + \frac{T\eta}{2} G^2 \\
 &= \underbrace{D \cdot G \sqrt{T}}_{\text{opt}} \quad \left( \frac{x^2 + y^2}{2} \geq xy \right)
 \end{aligned}$$

# Regret bound for OGD

**Theorem 3.1.** Online gradient descent with step sizes  $\{\eta_t = \frac{D}{G\sqrt{t}}, t \in [T]\}$  guarantees the following for all  $T \geq 1$ :

$$\text{regret}_T = \sum_{t=1}^T f_t(\mathbf{x}_t) - \min_{\mathbf{x}^* \in \mathcal{K}} \sum_{t=1}^T f_t(\mathbf{x}^*) \leq \frac{3}{2}GD\sqrt{T}$$

- “Any-time” algorithm with a decreasing learning rate schedule
- Learning rate depends on  $t$ . (exercise to prove that this works.)

# Online to Batch conversion: How do I use OGD to solve ERM?

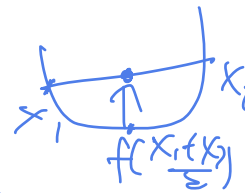
$f_t \stackrel{iid}{\sim} P$

$$\bar{X} = \frac{1}{T} \sum_{t=1}^T x_t$$

$$\mathbb{E}[f_t] = f$$

$$\min_x f(x) \leftarrow R_{\text{RHC}}$$

$$\frac{1}{T} \sum_{t=1}^T f_t(x) \quad \text{Empirical Risk}$$



$$\mathbb{E}[f(\bar{X})] - f(x^*) \leq \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T f(x_t)\right] - f^*$$

$$\mathbb{E}[f] = f \quad \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T f(x_t)\right] - f(x^*)$$

$$\left( \begin{array}{l} f(x^*) \geq f(x_t) \\ T \nabla f(x_t)^T (x_t - x^*) \end{array} \right)$$

$$\mathbb{E} g_t = \mathbb{E} \nabla f = \mathbb{E} \nabla f$$

$$\leq \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \nabla f(x_t)^T (x_t - x^*)\right]$$

$$= \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \underline{g_t^T (x_t - x^*)}\right]$$

$$\mathbb{E}[f] = f \leq \frac{1}{T} \mathbb{E}\left[\underbrace{\sum_{t=1}^T f_t(x_t) - f_t(x^*)}_{\text{Regret}}\right] = \frac{1}{T} \text{Regret} = O\left(\frac{GD}{\sqrt{T}}\right)$$

# Checkpoint

- Online learning
  - Operates in an adversarial environment
  - Almost no assumptions. Do not even use probability theory
- The idea of regret and no-regret learning algorithms
- Useful algorithmic ideas:
  - Hedge / Exponential Weighted Averages
  - Online gradient descent

# What we did not cover

- The strongly convex case
- Adapting to the geometry
  - AdaGrad / ADAM
- Adaptive regret / dynamic regret
  
- Modern applications to Ensemble learning, AutoML

[\[PDF\] Hyperband: A novel bandit-based approach to hyperparameter optimization](#)

[L Li, K Jamieson, G DeSalvo, A Rostamizadeh...](#) - The Journal of Machine ..., 2017 - jmlr.org

Performance of machine learning algorithms depends critically on identifying a good set of hyperparameters. While recent approaches use Bayesian optimization to adaptively select ...

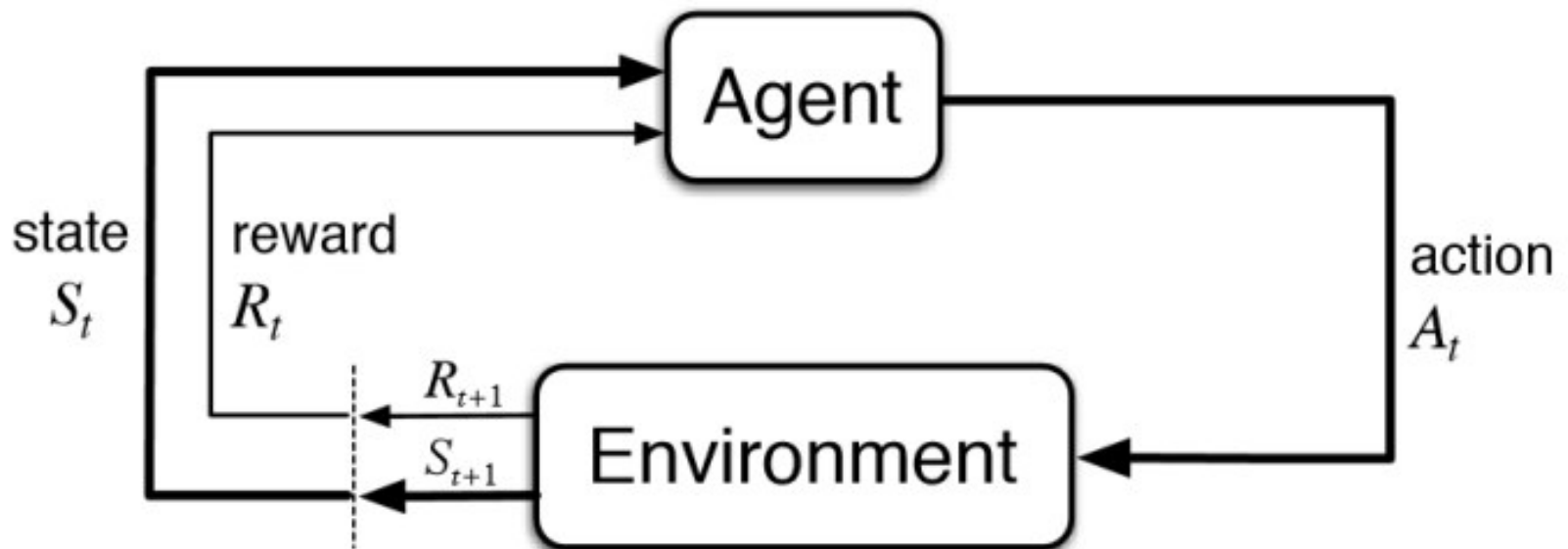
☆ Save  Cite Cited by 1612 Related articles Import into BibTeX 



# This lecture

- Online Learning (Part II)
  - Online Gradient Descent
- Reinforcement Learning
  - Problem setup
  - Markov Decision Processes

An RL agent learns **interactively** through the **feedbacks** of an environment.



- Learning how the world works (dynamics) and how to maximize the long-term reward (control) at the same time.

# Reinforcement learning is among the hottest area of research in ML!



“RL” is Top 1 Keyword at NeurIPS’2021, appearing 199 times  
“Deep Learning” only 129 times [[source](#)]

# Applications of RL in the real life

- RL for robotics.
- RL for dialogue systems.
- RL for personalized medicine.
- RL for self-driving cars.
- RL for new material discovery.
- RL for sustainable energy.
- RL for feature-based dynamic pricing.
- RL for maximizing user satisfaction.
- RL for QoE optimization in networking
- ...

# Reinforcement learning problem setup

- State, Action, Reward and Observation

$$S_t \in \mathcal{S} \quad A_t \in \mathcal{A} \quad R_t \in \mathbb{R} \quad O_t \in \mathcal{O}$$

- Policy:

$$\pi : \mathcal{S} \rightarrow \mathcal{A}$$

- When the state is observable:
- Or when the state is not observable

$$\pi_t : (\mathcal{O} \times \mathcal{A} \times \mathbb{R})^{t-1} \rightarrow \mathcal{A}$$

- Learn the best policy that maximizes the expected reward

- Finite horizon (episodic) RL:

$$\pi^* = \arg \max_{\pi \in \Pi} \mathbb{E} \left[ \sum_{t=1}^H R_t \right]$$

H: horizon

- Infinite horizon RL:

$$\pi^* = \arg \max_{\pi \in \Pi} \mathbb{E} \left[ \sum_{t=1}^{\infty} \gamma^{t-1} R_t \right]$$

$\gamma$ : discount factor

# RL for robot control



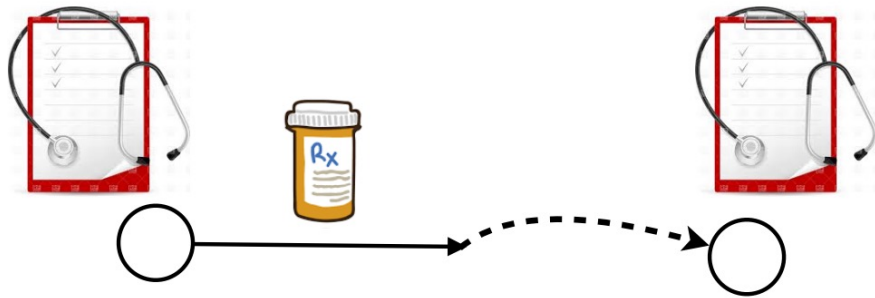
- States: The physical world, e.g., location/speed/acceleration and so on.
- Observations: camera images, joint angles
- Actions: joint torques
- Rewards: stay balanced, navigate to target locations, serve and protect humans, etc.

# RL for Inventory Management

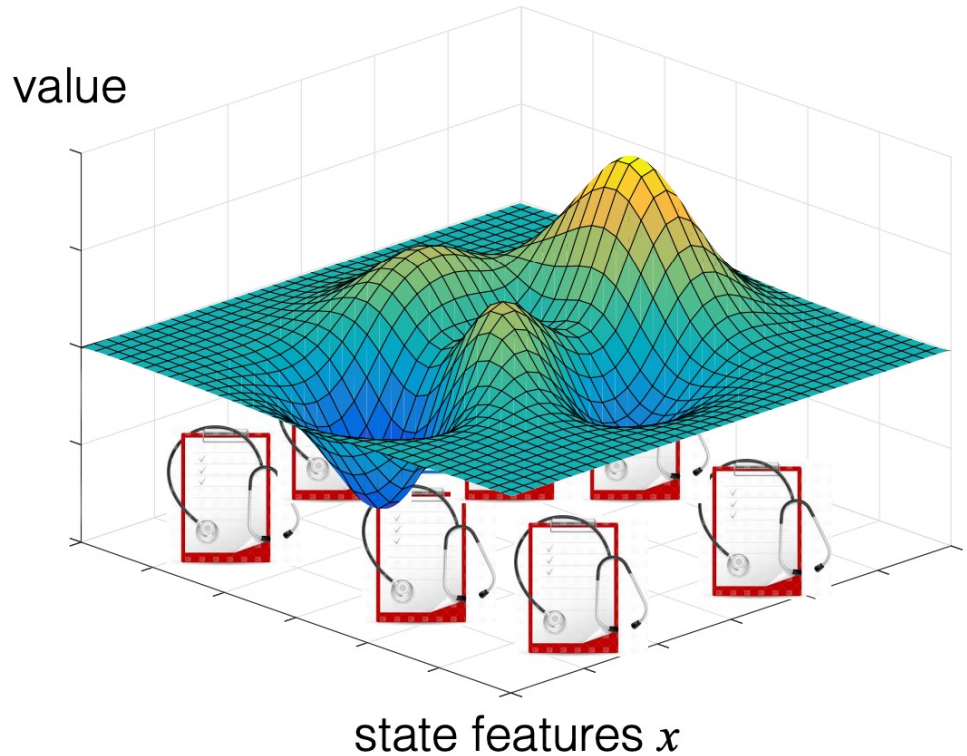


- State: Inventory level, customer demand, competitor's inventory
- Observations: current inventory levels and sales history
- Actions: amount of each item to purchase
- Rewards: profit

# RL for Adaptive medical treatment



- State: diagnosis
- Action: treatment
- Reward: progress in recovery



(example / illustration due to Nan Jiang)



# Example: Supervised learning vs RL in movie recommendation

- Bob is described by a feature vector
  - $s = [\text{Previous movies watched} / \text{Rating} / \text{Written reviews}]$
- Supervised learning predicts how likely Bob will click on “aliens vs predators”
- Reinforcement learning aims at controlling Bob
  - So in the future, Bob will develop a taste for “aliens vs predators” (e.g., from having watched “aliens” and “predators” both).