Algorithms for Easy and Hard Online Learning Environments

Yevgeny Seldin

University of Copenhagen

"Classical" (Batch) Machine Learning



How Online is different from "batch"?



Examples

- Investment in the stock market
- Online advertising/personalization
- Online routing
- Games
- Robotics

When do we need Online Learning?

Until recently, statistical theory has been restricted to the design and analysis of sampling experiments in which the size and composition of the samples are completely determined before the experimentation begins. The reasons for this are partly historical, dating back to the time when the statistician was consulted, if at all, only after the experiment was over, and partly intrinsic in the mathematical difficulty of working with anything but a fixed number of independent random variables. A major advance now appears to be in the making with the creation of a theory of the sequential design of experiments, in which the size and composition of the samples are not fixed in advance but are functions of the observations themselves.

(Robbins, 1952)



When do we need Online Learning?

- Interactive learning
- "Adversarial" game-theoretic settings
 - No i.i.d. assumption
- Intelligent data collection / experiment design
- Large-scale data analysis









Exploration-Exploitation Trade-off (1) 0/2 6/10 Never tried What drug to give to a new patient

When there are more patients to come...





We are building the dataset for ourselves





Learning in Adversarial Environments

• Game theoretic setting



• Cannot be treated in batch learning



- Evaluation measure: *regret*
 - Difference in performance compared to the best choice in hindsight (out of a limited set)
 - E.g. investment revenue vs. the best stock in hindsight





Delayed Feedback















Some "standard" algorithms



Typical performance scaling



"Standard" algorithms





Prediction with limited advice / Bandits with paid observations

[Seldin, Bartlett, Cramer, Abbasi-Yadkori, ICML, 2014]



Environmental resistance in full info

[Koolen & van Erven, COLT, 2015, Luo & Schapire, COLT, 2015, Wintenberger, 2015, van Erven, Kotłowski & Warmuth, COLT, 2014, Gaillard, Stoltz & van Erven, COLT 2014, Cesa-Bianchi, Mansour & Stoltz, MLJ, 2007, ...]



Environmental resistance in bandits

[Bubeck & Slivkins, COLT, 2012, Seldin & Slivkins, ICML, 2014, Auer & Chiang, COLT, 2016, Seldin & Lugosi, COLT, 2017, Wei & Luo, COLT, 2018, **Zimmert & Seldin, AISTATS, 2019**]



Prediction with Limited Advice

[Thune & Seldin, NIPS, 2018]



Prediction with Limited Advice Problem Setting





- Adversarial
 - ℓ_t^i arbitrary in [0,1]
- I.I.D.

•
$$\mathbb{E}[\ell_t^i] = \mu_i$$

- Gaps: $\Delta_i = \mu_i \min_j \mu_j$
- Effective Loss Range ε

•
$$\forall i, j, t: \left| \ell_t^i - \ell_t^j \right| \leq \varepsilon$$

• Evaluation measure: pseudo-regret

•
$$R_T = \mathbb{E}\left[\sum_{t=1}^T \ell_t^{A_t}\right] - \min_{i \in \{1, \dots, K\}} \mathbb{E}\left[\sum_{t=1}^T \ell_t^i\right]$$

SODA: Second Order Difference Adjustment

[Thune & Seldin, NIPS, 2018]

- The primary action A_t sampled according to
 - $p_t(i) \propto e^{-\eta_t D_{t-1}^i \eta_t^2 S_{t-1}^i}$
- The secondary observation B_t sampled
 - uniformly from $\{1, ..., K\} \setminus \{A_t\}$
- Unbiased loss difference estimates

•
$$\widetilde{\Delta \ell_s^i} = (K-1)\mathbb{I}(B_t = i)\left(\ell_t^{B_t} - \ell_t^{A_t}\right)$$

- Loss difference estimator and its second moment
 - $D_t^i = \sum_{s=1}^t \widetilde{\Delta \ell_s^i}$ $S_t^i = \sum_{s=1}^t \left(\widetilde{\Delta \ell_s^i} \right)^2$
- The learning rate

•
$$\eta_t \approx \sqrt{\frac{\ln K}{\max_i S_{t-1}^i}}$$

SODA: Regret Guarantees

[Thune & Seldin, NIPS, 2018]

- Adversarial
 - $R_T = O(\varepsilon \sqrt{KT \ln K})$
- Stochastic

•
$$R_T = O\left(\frac{K\epsilon^2}{\ln K}\sum_{i:\Delta_i>0}\frac{1}{\Delta_i}\right)$$

- Knowledge of i.i.d./adversarial not required!
- Knowledge of ε not required!
- Simultaneous adaptation to two types of easiness!

An optimal algorithm for i.i.d. and adversarial bandits

[Zimmert & Seldin, AISTATS, 2019]



I.I.D. and adversarial bandits Problem Setting

time



- Adversarial as before
 - Stochastically Constrained Adversary [Wei & Luo, 2018]
 - *i*^{*} best action

•
$$\mathbb{E}[\ell_t^i - \ell_t^{i^*}] = \Delta_i \ge 0$$

• I.I.D.: special case when

•
$$\mathbb{E}[\ell_t^i] = \mu_i$$

• Evaluation measure: pseudo-regret (as before)

•
$$R_T = \mathbb{E}\left[\sum_{t=1}^T \ell_t^{A_t}\right] - \min_{i \in \{1, \dots, K\}} \mathbb{E}\left[\sum_{t=1}^T \ell_t^i\right]$$

α -Tsallis-INF

[Zimmert & Seldin, AISTATS, 2019]

• Tsallis entropy

•
$$H_{\alpha}(x) = \frac{1}{1-\alpha} (1 - \sum_{i} x_{i}^{\alpha})$$

- INF Implicitly Normalized Forecaster
 - [Audibert & Bubeck, 2009]



α -Tsallis-INF

[Zimmert & Seldin, AISTATS, 2019]



• Sample $A_t \sim w_t$

• Update
$$\tilde{L}_t^i = \tilde{L}_{t-1}^i + \frac{\ell_t^i \mathbb{I}(A_t=i)}{w_t^i}$$

- α = 1 corresponds to entropic regularization
 EXP3 algorithm
- $\alpha = 0$ corresponds to log-barrier



1 2 [Zimmert & Seldin, AISTATS, 2019]

- $\eta_t = \sqrt{1/t}$
- Adversarial
 - $R_T \le 4\sqrt{KT} + 1$
- Stochastically Constrained Adversary

•
$$R_T \leq \sum_{i \neq i^*} \frac{4 \ln T + 20}{\Delta_i} + 4\sqrt{K}$$

- i.i.d. is a special case
- Both results match the lower bounds (up to constants)!
- No knowledge of i.i.d./adversarial is required!

Tsallis-INF in relation to prior work

[Zimmert & Seldin, AISTATS, 2019]

	Regime	Upper Bound Lower Bound
BROAD [Wei&Luo,2018] Log-barrier + doubling $(\alpha = 0)$	i.i.d.	$O(K^2)$
	adversarial	$O(\sqrt{\ln T})$
$\alpha = \frac{1}{2}$	i.i.d. & adversarial	0 (1)
EXP3++ [Seldin&Lugosi,2017] Entropic regularization + mix in extra exploration $(\alpha = 1)$	i.i.d.	$O(\ln T)$
	adversarial	$O(\sqrt{\ln K})$

$\frac{1}{2}$ -Tsallis-INF: Experiments [Zimmert & Seldin, AISTATS, 2019]



2500

5000

7500

10000





Stochastically Constrained **Adversary**

1 2 [Zimmert & Seldin, AISTATS, 2019]



- Optimality in the moderately contaminated stochastic regime
- I.I.D. and adversarial optimality in utility-based dueling bandits













If you want to learn more

- Check my homepage
 - https://sites.google.com/site/yevgenyseldin/
 - (or Google me)
 - Papers, tutorials, lecture notes, etc...
- Join our Advanced Topics in Machine Learning course
 - Co-taught by Christian Igel
 - September-October 2019
 - WARNING: Math-heavy course

α -Tsallis-INF: Some Considerations Γ

[Zimmert & Seldin, AISTATS, 2019]

- Exploration rate $\left(\eta_t^i (\tilde{L}_t^i \tilde{L}_t^{i^*})\right)^{-\frac{1}{1-\alpha}}$
- For i.i.d. regret of $\frac{\ln T}{\Delta_i}$ the exploration rate must be $\frac{1}{\Delta_i^2 t}$
- $\mathbb{E}\left[\tilde{L}_t^i \tilde{L}_t^{i^*}\right] = \Delta_i t$
- $\left(\eta_t^i \Delta_i t\right)^{-\frac{1}{1-\alpha}} = \frac{1}{\Delta_i^2 t} \implies \eta_t^i = \Delta_i^{1-2\alpha} t^{-\alpha}$
- $\alpha = \frac{1}{2}$ is the only value for which η_t^i requires no tuning by (unknown) Δ_i

α -Tsallis-INF: Some Considerations

[Zimmert & Seldin, AISTATS, 2019]

- $\mathbb{E}\left[\tilde{L}_t^i \tilde{L}_t^{i^*}\right] = \Delta_i t$
- The variance of $\tilde{L}_t^i \tilde{L}_t^{i^*}$ is of order $\Delta_i^2 t^2$
- $\tilde{L}_t^i \tilde{L}_t^{i^*}$ cannot be efficiently controlled by Bernstein's inequality
- Our analysis is based on a self-bounding property of the regret