Invited talk for INFORMS Applied Probability Society

An Adversarial Analysis of Thompson Sampling for Full-information Online Learning: from Finite to Infinite Action Spaces

Joint work with Jeff Negrea



Cornell University

Alexander Terenin HTTPS://AVT.IM/ · ♥ @AVT_IM

Probabilistic Decision-making



This talk: adversarial analogs of problems like this

The Online Learning Game

At each time t = 1, .., T:

- 1. Learner picks a random *action* $x_t \sim p_t \in \mathcal{M}_1(X)$.
- 2. Adversary responds with a *reward function* $y_t : X \to \mathbb{R}$, chosen adaptively from a given function class.

Regret:



Non-discrete Online Learning



Adversarial problem: seems to have nothing to do with Bayes' Rule?

No-regret Online Learning Algorithms

Given various no-regret online learning oracles, one can obtain:

Online-to-PAC Conversions: Generalization Bounds via Regret Analysis No-Regret Learning Dynamics for Extensive-Form The Statistical Complexity of Interactive Decision Making **Correlated Equilibrium** Confidence Sequences for Generalized Gábor Lugosi ICREA, Universita Dylan J. Foster Sham M. Kakade Jian Qian Alexander Rakhlin t Panneu Fabra, and Barcelana School of Economics, Barcelona, Snah Linear Models via Regret Analysis Gergely Neu Bounces Fabra, Barcelona, Spain Andrea Celli* olitecnico di Mila Alberto Marchesi* Politecnico di Milano Eugenio Clerico¹, Hamish Flynn¹, Wojciech Kotłowski², and Gergely Neu ¹Universitet Pompeu Falme, Barcelona, Spain. Abstract 1 Introduction study the standard model of statistical learning. We are given a training sample of π i.i.d. da Z_n) drawn from a distribution μ over a measurable hg sample to an output $W_n = A(S_n)$ taking values in a α trially randomized way. In other words, a randomized le Contents em 2 a probability distril tion over W and draws a sample from that distri ent is denoted by W_n . More pro-gorithm can thus be formally v 1 Introductio mance of the learning algorithm measured by a loss function $\ell : W \times Z \rightarrow \mathbb{R}_+$ est are the risk (or test error) $\mathbb{E}[\ell(w, Z')]$ and the compirical risk (or training error $\kappa_i Z_i)$ of a hypothesis $w \in W_i$ where the random element Z' has the same distribution t of S_n . The ultimate goal in statistical learning is to find algorithms with small excess Jence sets for parameters of statistical models is one or use a parsitoms of statistics. In this paper, we consider this problem (generalized linear models (GLMs), where one has access t servations $(X^n, Y^n) = (X_n, Y_{1,n})_{n-1}$. Here, $X_n \in \mathbb{R}^d$ is a vecto fortures), $Y_n \in \mathbb{R}$ is a reak-valued lakel, and the likelihood of en by the exponential-family model 2 Preliminaries 2.1 Minimax Re 3 A Theory of Learnability for Interactive Decision Making 3.1 Lower Bound: The Decision-Estimation Coefficient is a Fund 3.2 E2D: A Unified Meto-Algorithm for Interactive Decision Mak 1 Introduction $\mathcal{E}(W_u) = \mathbb{E}\left[\ell(W_u, Z')|W_u\right] - \inf_{u \in W} \mathbb{E}\left[\ell(w, Z')\right],$ The Nash canilibrium (NE) [37] is the more ecomposed into the task of minimizing the empirical risk $L(\cdot, S_n)$, and showing that the risk and the empirical risk is small. This gap is commonly called the generalization error d is defined formally as $p(y|X_t, \theta^*) = \exp(\langle \theta^*, X_t \rangle y - \psi(\langle \theta^*, X_t \rangle))h(y)$, interplay between computer science and game theory (see, e.g., to limit poker by Brown and Sandholm [5] and Moravčík et al. with $\theta^* \in \mathbb{R}^d$ the unknown parameter, $\psi : \mathbb{R} \to \mathbb{R}$ a convex function (often callee the log-partition function), and $h : \mathbb{R} \to \mathbb{R}$ the reference distribution (independent of X_1 or θ^*). The model can be alternatively written as $Y_1 = \mu(|\theta^*, X_1) + \epsilon_1$ $gen(W_n, S_n) = \mathbb{E} [\ell(W_n, Z') | W_n] - L(W_n, S_n),$ The E2D Meta-Algorithm: General Toolkit 4.2 Dual Perspective and 4.3 General Divergences a eralization error measures the extent of overfitting occurring during training, re specifically, upper bounding) the generalization error has been in the center heory ever since its inception. Over the past half century, numerous approach 5 Illustrative Examples on Neural Information Processing Systems (NeurIPS 2020), Vanc

Generalization

bounds

Confidence sets

Sample-efficient reinforcement learning Equilibrium computation algorithms **Online Learning: Discrete Action Spaces**

Typical algorithm: mirror descent or follow-the-regularized-leader

$$p_t = rgmax_{p \in \mathcal{M}_1(X)} \mathbb{E} \sum_{x \sim p} \sum_{t=1}^T y_t(x) - \Gamma(p)$$

Parameterized by convex regularizer $\Gamma:X
ightarrow\mathbb{R}$

- Not obvious how to pick Γ in non-discrete settings
- Standard choice of KL ignores adversary's smoothness class
- Unclear how to perform numerics in general

Adversarial problem: seemingly nothing to do with Bayesian learning?

Bayesian Algorithms for Adversarial Learning

Idea: think of this game in a Bayesian way

- 1. Place a prior $q^{(\gamma)}$ on the *adversary's future reward functions* $\gamma_1, ..., \gamma_T \in \mathcal{M}_1(Y)$
- 2. Condition on observed rewards $\gamma_1 = y_1, ..., \gamma_{t-1} = y_{t-1}$
- 3. Draw a sample from the posterior, and play

$$x_t = rgmax_{x\in X} \sum_{ au=1}^{t-1} y_ au(x) + \sum_{ au=t}^T \gamma_ au(x)$$

Algorithm: Thompson sampling

Bayesian Algorithms for Online Learning

Thompson sampling: FTPL with very specific learning rates

• Why should this be a good idea?

Result (Gravin, Peres, and Sivan, 2016). If $X = \{1, 2, 3\}$, the (infinite-horizon discounted) online learning game's Nash equilibrium is Thompson sampling with respect to a certain optimal prior.

Conjecture: Thompson sampling with strong priors is minimax-strong **This work**: prove the finite and simplest infinite-dimensional case

Idea: Analyze Bayesian Algorithm in a Bayesian Way

Regret decomposition:

$$\underbrace{R(p,y)}_{ ext{regret}} = \underbrace{R(p,q^{(\gamma)})}_{ ext{prior regret}} + \underbrace{E_{q^{(\gamma)}}(p,y)}_{ ext{excess regret}}$$

where

$$E_{q^{(\gamma)}}(p,y) = \sum_{t=1}^T \Gamma^*_{t+1}(y_{1:t}) - \Gamma^*_t(y_{1:t-1}) - \langle y_t | p_t
angle + \mathbb{E} \left\langle \gamma_t | p_t^{(\gamma)}
ight
angle$$

and $\Gamma^*_t(f) = \sup_{x \in X} f(x) + \gamma_{t:T}(x).$

Strong Priors for Adversarial Feedback

Strong prior over adversary:

- 1. Equalizing: $\mathbb{E}\,\gamma(x)=\mathbb{E}\,\gamma(x')$ for all x,x'
- 2. Certifies a sharp regret lower bound

$$R(\cdot,q^{(\gamma)}) \leq \min_p \max_q R(p,q)$$

Practical approximation: Gaussian process with matching smoothness

- Matérn priors: widely-used in Bayesian optimization today
- Straightforward and well-understood numerics

A Bregman Divergence Bound on Excess Regret

For a strong prior:

$$E_{q^{(\gamma)}}(p,y) \leq \sum_{t=1}^{T} \underbrace{D_{\Gamma_t^*}(y_{1:t} \mid\mid y_{1:t-1})}_{ ext{Bregman divergence} \ ext{induced by } \Gamma_t^*}$$

Gaussian adversary: rates

- Finite X with ℓ^∞ adversary: $R(p,q) \leq 2\sqrt{T\log N}$
- X = [0,1], BL adversary: $R(p,q) \leq \mathcal{O}\left(eta\sqrt{Td\log(1+\sqrt{d}rac{\lambda}{eta})}
 ight)$

Hessian Bounds for Gaussian Adversaries

Taylor form of the Bregman divergence:

$$egin{aligned} D_{\Gamma^*_t}(y_{1:t} \mid\mid y_{1:t-1}) &= rac{1}{2} \int_0^1 \underbrace{\partial^2_{y_t,y_t} \Gamma^*_t(y_{1:t-1} + lpha y_t)}_{ ext{Gâteaux Hessian}} ext{d}lpha \ \partial^2_{u,v} \Gamma^*_t(f) &= rac{1}{\sqrt{T-t+1}} \, \mathbb{E} \, u(\underbrace{x^*_{f+\gamma_{t:T}}}_{ ext{perturbed}}) ig\langle \gamma_{t:T} ert \mathcal{K}^{-1} v ig
angle \ & ext{covariance} \ & ext{operator} \end{aligned}$$

Classical finite-dimensional argument: $\|\Gamma^*_t(f)\|_{L_{\infty,1}} \leq 2\operatorname{tr} \Gamma^*_t(f)$

- Only works if ${\cal K}$ is the identity matrix
- Essentially all obvious workarounds give vacuous rates

Probabilistic Hessian Bounds

Idea: condition on maximizer and work with truncated normals

• Algebraic properties from discrete case have probabilistic analogs

Theorem. Suppose adversary satisfies $\|y\|_{\infty} \leq \beta$ and prior is IID over time with constant variance $k(x,x) = \sigma^2$. Suppose that

$$\sup_{y\in Y}y(x)-y(x')rac{k(x,x')}{k(x',x')}\leq C_{Y,k}\left(1-rac{k(x,x')}{k(x',x')}
ight).$$

Then we have

$$D_{\Gamma^*_t}(y_{1:t}\mid\mid y_{1:t-1}) \leq rac{eta^2+eta C_{Y,k}}{2\sigma^2\sqrt{T-t+1}}\,\mathbb{E}\sup_{x\in X}\gamma_{t:T}(x).$$

Bounded Lipschitz Adversary

Bounded Lipschitz (eta, λ) adversary: Matérn kernel with length scale κ

$$\sup_{y\in Y}y(x)-y(x')rac{k(x,x')}{k(x',x')}\leq rac{eta(\lambda+rac{1}{\kappa})}{rac{1}{\kappa}\left(1-e^{-rac{2}{\lambda\kappa+1}}
ight)}\left(1-rac{k(x,x')}{k(x',x')}
ight)$$

Thompson sampling regret:

$$R(p,q) \leq eta \left(32 + rac{32}{1-rac{1}{e}}
ight) \sqrt{Td\log\left(1+\sqrt{d}rac{\lambda}{eta}
ight)}$$

Algorithmic design principle:

To explore by random actions, don't be too predictable, and match smoothness

- Non-discrete online learning:
 - Thompson sampling: perturbation based algorithm
 - Bayesian viewpoint: match smoothness using Gaussian priors
 - Analysis: probabilistic argument for non-discrete Hessian bounds

Future work:

- More general smoothness classes
- Learning in games
- Bandit feedback

Thank you! https://avt.im/· 🎔 @avt_im

A. Terenin and J. Negrea. An Adversarial Analysis of Thompson Sampling for Full-information Online Learning: from Finite to Infinite Action Spaces. *arXiv:2502.14790*, 2025.





Cornell University