

# Distributionally Robust Counterfactual Risk Minimization

Louis Faury<sup>1,2</sup>, Ugo Tanielian<sup>1,3</sup>, Elena Smirnova<sup>1</sup>  
Flavian Vasile<sup>1</sup>, Elvis Dohmatob<sup>1</sup>

<sup>1</sup> Criteo AI Labs

<sup>2</sup> LTCI, Telecom ParisTech

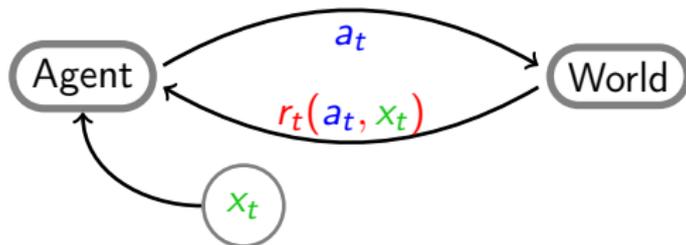
<sup>3</sup> LPSM, Paris 6

# Outline

- We are interested in **off-line policy evaluation and improvement** in a contextual bandit setting.
- We propose to use tools from **Distributionally Robust Optimization** (DRO) for this task, motivated by asymptotic guarantees.
- We introduce a **new algorithm** for off-line policy improvement, based on the DRO framework, that outperforms the state-of-the-art on classical datasets.

# Contextual Bandits (CB)

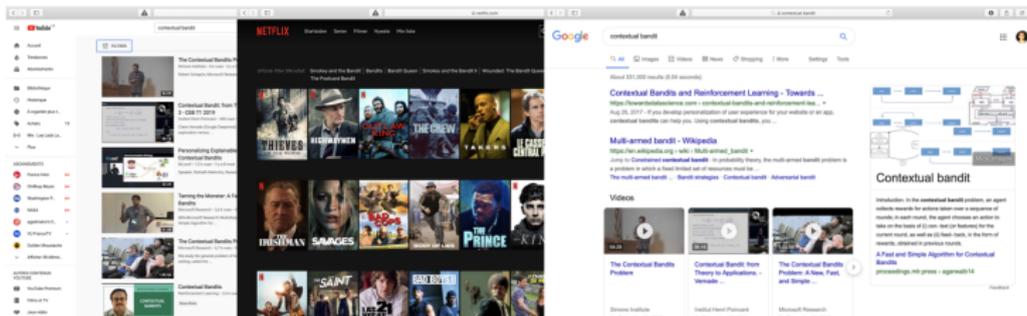
The contextual bandits (CB) is an extension to the classical multi-arm bandit setting.



In CB, an agent is presented with a context  $x_t$  (exogenous) and plays an action  $a_t$ . The environment then generates a reward  $r_t$ .

The agent's goal is to maximize its expected reward.

# Contextual Bandit (CB, cntd')



- Recommender system:

$x_t$ =user embedding,  $a_t$ =recommandation,  $r$ =click

- Clinical trials:

$x_t$ =patient information,  $a_t$ =medication,  $r$ =remission

# Contextual Bandit (CB, cntd')

**Goal:** Maximize the expected reward, under two settings:

- **Online setting:** at every round  $t$ , the agents interacts with the world to minimize its cumulative regret. The challenge is the **exploration-exploitation** trade-off.
- **Offline setting:** the agent only has access to past interactions and must find a way to improve its performance. The challenge is off-line policy **evaluation** and **improvement**.

We will consider the offline setting.

## Some notations

- Let contexts  $x \in \mathcal{X}$  and action  $a \in \mathcal{A}$
- Let the cost  $c(x, a) := -r(x, a)$ .
- The contexts are drawn under  $\nu$  (unknown).

An agent is characterized by its **policy**: a function that maps contexts to a distribution on the actions.

The goal is to find the policy  $\pi$  with minimal **risk**:

$$R(\pi) := \mathbb{E}_{x \sim \nu, y \sim \pi(\cdot|x)} [c(x, y)]$$

which is the expected cost suffered when playing the policy  $\pi$ .

# Offline Contextual Bandits (OCB)

In OCB, the agent cannot interact with the environment. The only available data are **interaction logs** from a logging policy  $\pi_0$ :

$$\mathcal{H}_0 = \left( x_i, a_i, p_i = \pi_0(x_i|a_i), c_i = c(x_i, a_i) \right)_{1 \leq i \leq n}$$

A standard estimator for  $R(\pi)$  involves **inverse propensity scores**:

$$R(\pi) = \mathbb{E}_{x \sim \nu, a \sim \pi_0} \left[ c(x, a) \frac{\pi(a|x)}{\pi_0(a|x)} \right]$$

usually estimated with capping:

$$\hat{R}_n(\pi) = \frac{1}{n} \sum_{\mathcal{H}_0} c_i \min \left( M, \frac{\pi(a_i|x_i)}{p_i} \right)$$

sometimes called the IPS estimator.

# Counterfactual Risk Minimization (CRM)

**Problem:** the estimator  $\hat{R}_n(\pi)$  can have very large variance for some  $\pi$  and may be **over-confident** (optimizer's curse).

**Solution:** [Swaminathan et al, 2015]<sup>1</sup> suggest looking at an variance-sensitive upper-bound on the true risk:

$$R(\pi) \leq \hat{R}_n(\pi) + \lambda \sqrt{\widehat{\text{Var}}_n(\pi)/n} \quad \text{w.h.p}$$

leading to the CRM principle for policy improvement:

$$\operatorname{argmin}_{\pi} \hat{R}_n(\pi) + \lambda \sqrt{\widehat{\text{Var}}_n(\pi)/n}$$

which gave rise to the **POEM** algorithm (state-of-the-art).

Can be augmented with **variance-reduction** techniques (Self-Normalized estimator, Doubly Robust).

---

<sup>1</sup>Counterfactual Risk Minimization: Learning from Logged Bandit Feedback

# Our contribution

We show that Distributionally Robust Optimization (DRO) tools can be applied to OCB in order to:

- provide a unified framework to build a collection of (asymptotic) **variance-sensitive upper-bounds on the risk**
- derive existing CRM algorithms
- derive **new** CRM algorithms outperforming state-of-the-art

⇒ DRO provides principled tools for the OCB problem. It is a **general** framework that generalizes existing CRM solutions.

# Distributionally Robust Optimization (DRO)

Denote  $\xi := (x, a)$  with distribution  $P := \nu \times \pi_0$ . We write the empirical risk as follows

$$\hat{R}_n(\pi) = \mathbb{E}_{\xi \sim \hat{P}_n} [\ell_\pi(\xi)] = \frac{1}{n} \sum_{i=1}^n \ell_\pi(\xi_i)$$

where  $\ell_\pi(\xi_i) = c_i \min(M, \pi(a_i|x_i)/p_i)$  (capped propensity-costs).

In DRO, we treat  $\hat{P}_n$  with **skepticism** and introduce a **robust risk**:

$$\tilde{R}_n^{\mathcal{U}}(\pi, \varepsilon) := \sup_{Q \in \mathcal{U}_\varepsilon} \mathbb{E}_{\xi \sim Q} [\ell_\pi(\xi)]$$

where  $\mathcal{U}_\varepsilon$  is an **ambiguity set**: a «ball» of radius  $\varepsilon$  around  $\hat{P}_n$ .

# DRO (cnt'd)

We define  $\mathcal{U}_\epsilon$  using **coherent**  $\varphi$ -divergences

$$\mathcal{U}_\epsilon = \left\{ Q \text{ s.t. } D_\varphi(Q \parallel \hat{P}_n) \leq \epsilon \right\}$$

where for  $Q \ll P$ :

$$D_\varphi(Q \parallel P) = \int \varphi \left( \frac{dQ}{dP} \right) dP$$

and **1)**  $\varphi$  is a convex function, **2)**  $\varphi(t) \geq \varphi(1) = 0$ , **3)**  $\varphi'(1) = 0$ , **4)**  $\varphi''(1) > 0$  (**coherent conditions**).

We will consider robust risk defined through coherent  $\varphi$ -divergences:

$$\tilde{R}_n^\varphi(\pi, \epsilon) = \sup_{D_\varphi(Q \parallel \hat{P}_n) \leq \epsilon} \mathbb{E}_{\xi \sim Q} \left[ \ell_\pi(\xi) \right]$$

# DRO for CRM: guarantees

**Guarantee 1** The robust risk  $\tilde{R}_n^\varphi(\pi, \varepsilon)$  provides an **asymptotic performance certificate** for the true risk.

Lemma 1: Risk upper-bound

For any  $\delta > 0$ :

$$\lim_{n \rightarrow \infty} \mathbb{P} \left[ R(\pi) \leq \tilde{R}_n^\varphi(\pi, \varepsilon_n) \right] \leq 1 - \delta$$

where  $\varepsilon_n = \varphi''(1) \chi_{1, 1-\delta}^2 / (2n)$ .

This result can be derived from Proposition 1 in [Duchi 2016]<sup>2</sup>.

---

<sup>2</sup>Statistics of Robust Optimization: A Generalized Empirical Likelihood Approach

## DRO for CRM: guarantees (2)

**Guarantee 2** The robust risk **penalizes high-variance** estimates:

Lemma 2: Asymptotic variance decomposition

$$\tilde{R}_n^\varphi(\pi, \varepsilon/n) = \hat{R}_n(\pi) + \sqrt{\frac{\varepsilon}{n} \widehat{\text{Var}}_n(\pi)} + o\left(\frac{1}{\sqrt{n}}\right)$$

This result can be obtained as a Corollary of Theorem 2 of [Duchi 2016].

⇒ Lemma 1 and Lemma 2 imply that the upper-bounds provide variance-sensitive performance certificate, making it a reliable tool for **off-line policy evaluation**.

## DRO for CRM: guarantees (3)

**Guarantee 3** With exponentially parametrized policy, minimizing the robust risk with  $\chi^2$  ambiguity sets is **exactly** the POEM algorithm.

Lemma 3: Exact variance decomposition

For  $\varepsilon$  small enough:

$$\tilde{R}_n^{\chi^2}(\pi, \varepsilon) = \hat{R}_n(\pi) + \sqrt{\varepsilon \widehat{\text{Var}}_n(\pi)}$$

⇒ Existing CRM algorithms are already instances of DRO estimators!

## DRO for CRM: guarantees (4)

Sketch of proof By strong duality we have

$$\sup_{D_\varphi(Q||\hat{P}_n)} \mathbb{E}_Q[l_\pi(\xi)] = \inf_{\gamma \geq 0} \gamma \varepsilon + \inf_Q \left\{ \mathbb{E}_Q[l_\pi(\xi)] - \gamma D_\varphi(Q||\hat{P}_n) \right\} \quad (1)$$

Using the Envelope Theorem of [Rockafellar18]<sup>3</sup> one gets:

$$\sup_{D_\varphi(Q||\hat{P}_n)} \mathbb{E}_Q[l_\pi(\xi)] = \inf_{\gamma \geq 0} \gamma \varepsilon + \inf_c \left\{ c + \gamma \mathbb{E}_{\hat{P}_n} [\varphi^*((l_\pi(\xi) - c)/\gamma)] \right\} \quad (2)$$

For the  $\chi^2$ -divergence,  $\varphi(z) = (z - 1)^2$  and  $\varphi^*(s) = s^2/4 + s$  for  $s \geq -2$ . Solving leads to the result.

---

<sup>3</sup>Risk and utility in the duality framework of convex analysis

# DRO for CRM: guarantees (5)

## Sum-up so far

- Guarantees 1 and 2: DRO is a general tool for building variance-sensitive upper-bounds on the risk
- Guarantee 3: POEM is actually DRO with  $\chi^2$  divergences.

## In what follows

We introduce a **new** CRM algorithm inspired from DRO, and derived from Kullback-Leibler divergence ambiguity sets.

# New KL-based CRM algorithm

We now consider **KL-based** ambiguity sets:

$$\tilde{R}_n^{\text{KL}}(\pi, \varepsilon) = \min_{\text{KL}(Q \parallel \hat{P}_n)} \mathbb{E}_{\xi \sim Q} [\ell_\pi(\xi)]$$

There is a **tractable computation** for the worst-case distribution.

Lemma 4: KL robustified risk

It exists  $\gamma > 0$  such that

$$\tilde{R}_n^{\text{KL}}(\pi, \varepsilon) = \sum_{i=1}^n \frac{\ell_\pi(\xi_i) e^{\ell_\pi(\xi_i)/\gamma}}{\sum_{j=1}^n e^{\ell_\pi(\xi_j)/\gamma}} \quad (3)$$

The line of proof follows the one of Lemma 3, and uses the convex conjugate of  $\varphi_{\text{KL}}(z) = z \log(z) - z + 1$ .

# New CRM algorithms

**Policy optimization:** Minimize the upper-bound given by the robust risk! This gives rise to the **KL-CRM** algorithm:

$$\text{minimize}_{\pi} \left[ \tilde{R}_n^{\text{KL}}(\pi, \varepsilon) = \sum_{i=1}^n \frac{\ell_{\pi}(\xi_i) e^{\ell_{\pi}(\xi_i)/\gamma}}{\sum_{j=1}^n e^{\ell_{\pi}(\xi_j)/\gamma}} \right] \quad (\text{KL-CRM})$$

where  $\gamma$  is treated as a hyper-parameter (cross-validation).

Temperature  $\gamma$  dictates the level of pessimism:

- $\gamma \rightarrow \infty$  reduces to the IPS estimator
- $\gamma \rightarrow 0$  only consider the worst case propensity cost.

# New CRM algorithms

A finer analysis reveals a good approximation  $\gamma$ .

Lemma 5: aKL-CRM

$$\gamma_* \simeq \sqrt{\frac{\widehat{\text{Var}}_n(\pi)}{2\varepsilon}}$$

The proof relies on a second-order Taylor approximation of the log-m.g.f of the loss in the dual objective.

This gives rise to **aKL-CRM** which minimizes the KL-CRM objective and concurrently updates  $\gamma_*$ .

# Experimental Results

We evaluate on **standard datasets** (supervised→bandit) and compare **KL-CRM** and **aKL-CRM** with the basic **IPS** approach and the **POEM** algorithm.

Hyper-parameters are determined through cross-validation.  
Experiments are average over 20 different random initialization.

The performance of a policy is reported by its **expected** instant regret or by the **instant regret** of its **greedy** policy.

## Experimental Results (ctn'd)

Expected instant regret:

	Scene	Yeast	RCV1-Topics	TMC2009
$\pi_0$	1.529	5.542	1.462	3.435
CIPS	1.163	4.658	0.930	2.776
POEM	1.157	<b>4.535</b>	0.918	2.191
KL-CRM	1.146	4.604	0.922	2.136
aKL-CRM	<b>1.128</b>	<b>4.553</b>	<b>0.783</b>	<b>2.126</b>
CRF	0.646	2.817	0.341	1.187

**Table:** Expected Hamming loss on  $\mathcal{D}_{\text{test}}^*$  for the different algorithms, averaged over 20 independent runs. Bold font indicate that one or several algorithms are statistically better than the rest, according to a one-tailed paired difference t-test at significance level of 0.05.

Rq: CRF is a **skyline** that has access to full supervised feedback.

## Experimental Results (ctn'd)

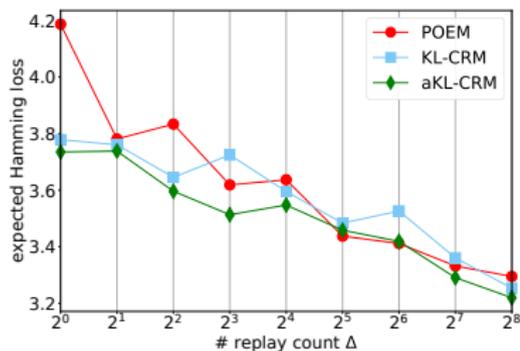
Greedy instant regret:

	Scene	Yeast	RCV1-Topics	TMC2009
CIPS	1.163	4.369	0.929	2.774
POEM	1.157	<b>4.261</b>	0.918	2.190
KL-CRM	1.146	4.316	0.922	2.134
aKL-CRM	<b>1.128</b>	<b>4.271</b>	<b>0.779</b>	<b>2.034</b>

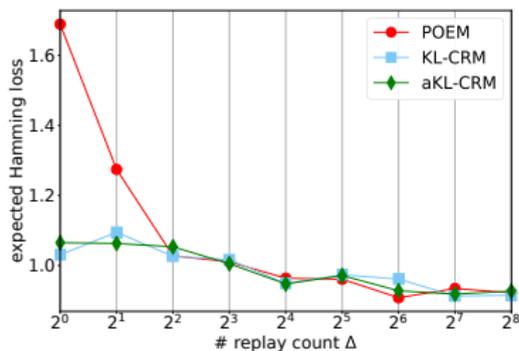
**Table:** Hamming loss on  $\mathcal{D}_{\text{test}}^*$  for the different greedy policies, averaged over 20 independent runs. Bold font indicates that one or several algorithms are statistically better than the rest, according to a one-tailed paired difference t-test at significance level of 0.05.

# Experimental Results (ctn'd)

Influence of the size of the logged history:



(a) Yeast dataset



(b) Scene dataset

**Figure:** Impact of the replay count  $\Delta$  on the expected Hamming loss. Results are average over 10 independent runs, that is 10 independent train/test split and bandit dataset creation. KL-CRM and aKL-CRM outperform POEM in the small data regime.

# Conclusion and future work

DRO is a principled tool for OCB and lead to **competitive** CRM algorithms.

Future work:

- further experimental evaluations (SNIPS, DR)
- solving the primal problem can be easy! we can use performance certificate given by many  $\varphi$  divergences.
- can we derive **finite samples** guarantees?

Thank you!