

# Large-scale Validation of Counterfactual Learning Methods: A Test-Bed

Damien Lefortier<sup>1,2</sup>, Adith Swaminathan<sup>3</sup>, Xiaotao Gu<sup>4</sup>, Thorsten Joachims<sup>3</sup>, and Maarten de Rijke<sup>2</sup>

<sup>1</sup> Facebook  
<sup>2</sup> University of Amsterdam  
<sup>3</sup> Cornell University, Ithaca, NY  
<sup>4</sup> Tsinghua University, Beijing, China

## Contributions

- ▶ We provide the **first public dataset with accurately logged propensities** from a production interactive system with recorded user feedback:
  - ▷ The dataset was collected at Criteo;
  - ▷ The dataset enables research into the problem of Batch Learning from Bandit Feedback (BLBF).
- ▶ We propose new sanity checks and evaluation methodologies when running BLBF experiments.
  - ▷ We provide a standardized test-bed that implements our workflow and benchmark several counterfactual learning algorithms in a sample BLBF task.

## Motivation

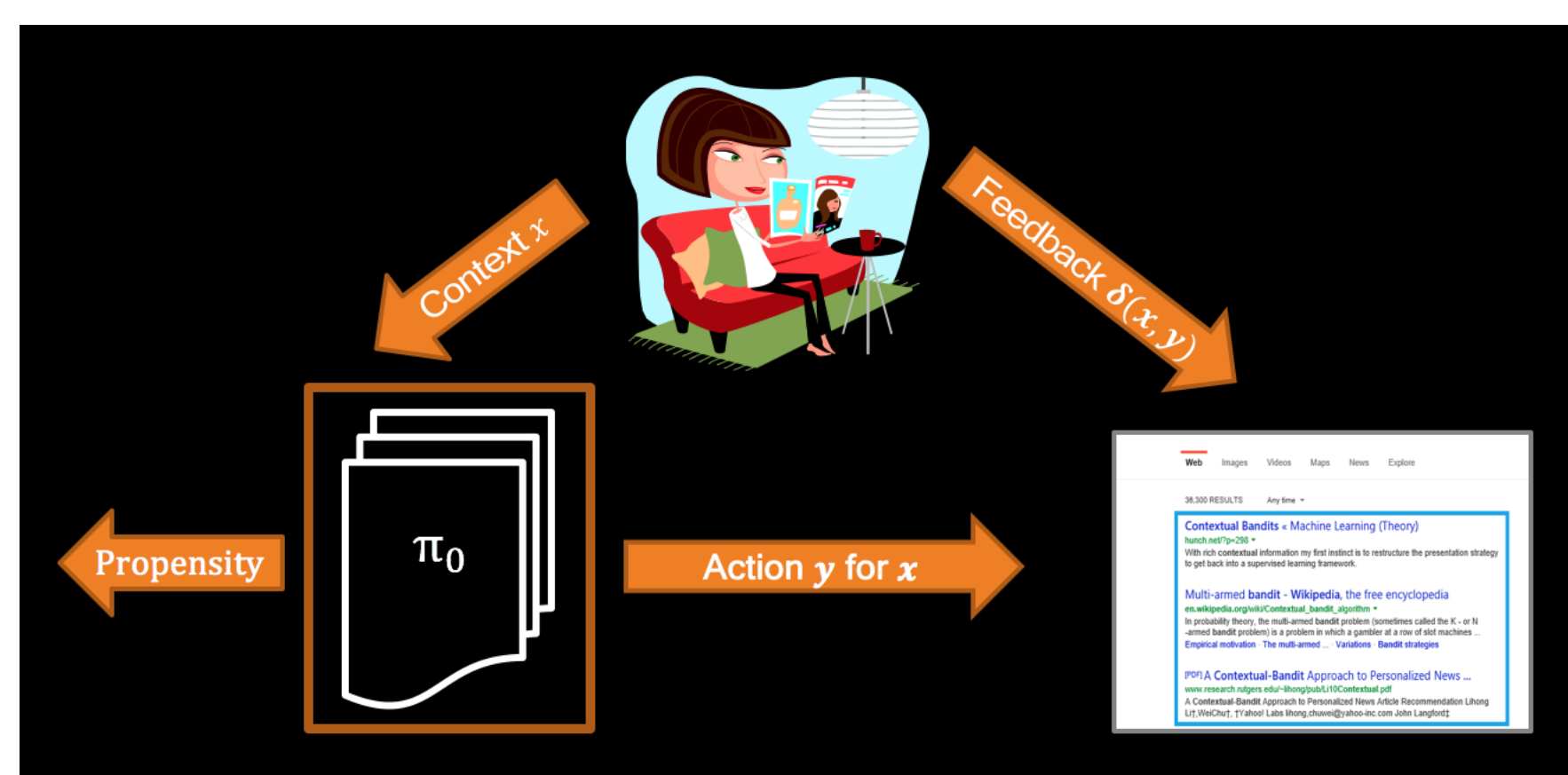


Figure: BLBF algorithm.

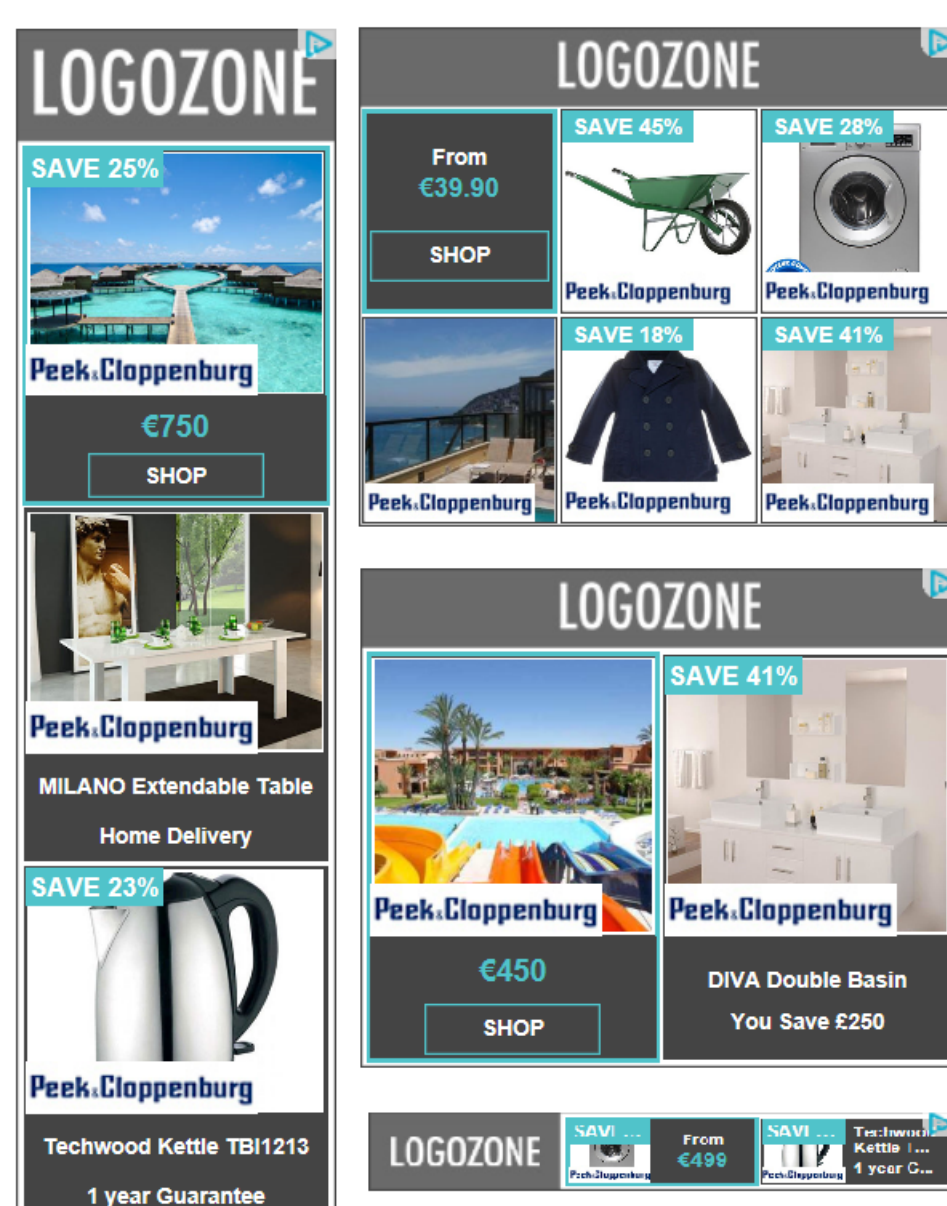


Figure: Concrete example: banner-filling task at Criteo.

This dataset and test-bed will hopefully enable research into:

- ▶ New training objectives, learning algorithms, and regularization mechanisms;
- ▶ Improved model selection procedures (analogous to cross-validation);
- ▶ Effective and tractable policy classes  $\pi \in \Pi$  for the specified task  $x \mapsto y$ ; and
- ▶ Algorithms that can scale to massive amounts of data.

## Dataset

The logging policy  $\pi_0$  stochastically selects products to construct a banner by first computing non-negative scores  $f_p$  for all candidate products  $p \in P_c$ , and using:

$$P(\text{slot1} = p) = \frac{f_p}{\sum_{p' \in P_c} f_{p'}} \quad P(\text{slot2} = p' \mid \text{slot1} = p) = \frac{f_{p'}}{\sum_{p'' \in P_c \wedge p'' \neq p} f_{p''}}, \quad \dots$$

The propensity of a chosen banner ad  $\langle p_1, p_2, \dots \rangle$  is  $P(\text{slot1} = p_1) * P(\text{slot2} = p_2 \mid \text{slot1} = p_1) * \dots$  and our dataset was logged as follows:

```
example {exID}: {hashID} {wasAdClicked} {propensity} {nbSlots}
{nbCandidates} {displayFeat1}:{v_1} ...
{wasProduct1Clicked} exid:{exID} {productFeat1_1}:{v1_1} ...
...
{wasProductMClicked} exid:{exID} {productFeatM_1}:{vM_1} ...
```

## Download our dataset at:

▶ <http://www.cs.cornell.edu/~adith/Criteo/index.html>

## Statistics

Sub-sampling to limit dataset size. Accounted for in the statistics and subsequent evaluation in our code.

#Slots	1	2	3	4	5	6
#Impressions	$2.13e+07$	$3.55e+07$	$2.27e+07$	$6.92e+06$	$2.95e+06$	$1.40e+07$
$\hat{N}$	$2.03e+08$	$3.39e+08$	$2.15e+08$	$6.14e+07$	$2.65e+07$	$1.30e+08$
Avg(InvPropensity)	11.96	$3.29e+02$	$1.87e+04$	$2.29e+06$	$2.62e+07$	$3.51e+09$
Max(InvPropensity)	$5.36e+05$	$3.38e+08$	$3.23e+10$	$9.78e+12$	$2.03e+12$	$2.34e+15$

Table: Number of impressions and propensity statistics for slices of traffic with k-slot banners with  $1 \leq k \leq 6$ . Estimated sample size ( $\hat{N}$ ) corrects for 10% sub-sampling of non-clicked impressions.

Consequences:

- ▶ Don't rely on a single point estimate (like IPS), but report multiple estimates.
- ▶ Confidence intervals can mislead (esp. when  $k \geq 4$ ).

## Benchmark Learning Algorithms

- ▶ Slice of traffic can enable logged contextual bandit learning: 1-slot filling task.
  - ▷ Regression to predict CTR of candidates. Pick best estimated CTR;
  - ▷ Off-policy learning method like DRO or POEM.

## Results for 1-slot task

Approach	Test set estimates		
	$\hat{R}(\pi_\epsilon) \times 10^4$	$\hat{R}(\pi_\epsilon) \times 10^4 / \hat{C}(\pi_\epsilon)$	$\hat{C}(\pi_\epsilon)$
Random	$44.676 \pm 2.112$	$45.446 \pm 0.001$	$0.983 \pm 0.021$
$\pi_0$	$53.540 \pm 0.224$	$53.540 \pm 0.000$	$1.000 \pm 0.000$
Regression	$48.353 \pm 3.253$	$48.162 \pm 0.001$	$1.004 \pm 0.041$
IPS	$54.125 \pm 2.517$	$53.672 \pm 0.001$	$1.008 \pm 0.016$
DRO	$57.356 \pm 14.008$	$57.086 \pm 0.005$	$1.005 \pm 0.025$
POEM	$58.040 \pm 3.407$	$57.480 \pm 0.001$	$1.010 \pm 0.018$

Table: Test set performance of policies learnt using different counterfactual learning baselines. Errors bars are 99% confidence intervals under a normal distribution. Confidence interval for SNIPS is constructed using the delta method.

Where  $\hat{C}(\pi) = \frac{1}{N} \sum_{i=1}^N \frac{\pi(y_i|x_i) \mathbb{1}\{o_i=1\}}{q_i \Pr(O=1|\delta_i)}$  and  $\hat{R}(\pi) = \frac{1}{N} \sum_{i=1}^N \delta_i \frac{\pi(y_i|x_i) \mathbb{1}\{o_i=1\}}{q_i \Pr(O=1|\delta_i)}$ .

## Grand BLBF challenges

- ▶ **Size of the action space:** Increase the size of the action space.
- ▶ **Feedback granularity:** Use per item feedback.
- ▶ **Contextualization:** We can learn a separate model for each banner type or learn a contextualized model across multiple banner types.

We hope you find this first public user impressions dataset with logged propensities useful for your research.

## References

- [1] Counterfactual reasoning and learning systems: the example of computational advertising. L. Bottou et al. JMLR 2013.
- [2] Batch learning from logged bandit feedback through counterfactual risk minimization. A. Swaminathan et al. JMLR 2015.
- [3] Doubly Robust Policy Evaluation and Learning. M. Dudik et al. ICML 2011.
- [4] Unbiased Offline Evaluation of Contextual-bandit-based News Article Recommendation Algorithms. L. Li et al. WSDM 2011.
- [5] The self-normalized estimator for counterfactual learning. A. Swaminathan et al. NIPS 2015.