

Nonstationary Policy Evaluation

John Langford @Microsoft Research

{ With help from many }

NYU Large Scale Learning Class, April 16, 2013

A Reminder: The Contextual Bandit Setting

For $t = 1, \dots, T$:

- 1 The world produces some context $x_t \in X$
- 2 The learner chooses an action $a_t \in \{1, \dots, K\}$
- 3 The world reacts with reward $r_t(a_t) \in [0, 1]$

Goal: Efficiently competing with a large reference class of possible policies $\Pi = \{\pi : X \rightarrow \{1, \dots, K\}\}$:

$$\text{Regret} = \max_{\pi \in \Pi} \sum_{t=1}^T r_t(\pi(x_t)) - \sum_{t=1}^T r_t(a_t)$$

A Reminder: The Contextual Bandit Setting

For $t = 1, \dots, T$:

- 1 The world produces some context $x_t \in X$
- 2 The learner chooses an action $a_t \in \{1, \dots, K\}$
- 3 The world reacts with reward $r_t(a_t) \in [0, 1]$

Goal: Efficiently competing with a large reference class of possible policies $\Pi = \{\pi : X \rightarrow \{1, \dots, K\}\}$:

$$\text{Regret} = \max_{\pi \in \Pi} \sum_{t=1}^T r_t(\pi(x_t)) - \sum_{t=1}^T r_t(a_t)$$

Which combined explore/exploit algorithm is best for your setting?

A Rejection Sampling approach

Rejection_Sampler(policy π , events $(\vec{x}, a, r, p)^T$)

Let $h = \emptyset$ a history, $R = 0$

For each event (\vec{x}, a, r, p)

- 1 If $\pi(h, \vec{x}) = a$
- 2 then with probability $\frac{p_{\min}}{p}$
 - 1 $h \leftarrow h \cup (\vec{x}, a, r)$
 - 2 $R \leftarrow R + r$

Return $R/|h|$

A Rejection Sampling approach

Rejection_Sampler(policy π , events $(\vec{x}, a, r, p)^T$)

Let $h = \emptyset$ a history, $R = 0$

For each event (\vec{x}, a, r, p)

- 1 If $\pi(h, \vec{x}) = a$
- 2 then with probability $\frac{p_{\min}}{p}$
 - 1 $h \leftarrow h \cup (\vec{x}, a, r)$
 - 2 $R \leftarrow R + r$

Return $R/|h|$

Theorem: For all history lengths T , For all **nonstationary** policy π , and all IID worlds D , the probability of a simulated history of length T = the probability of the same history of length T in the real world.

The Master Evaluator

Eval(policy π , events $(\vec{x}, a, r, \rho)^T$, quantile ρ , bound b)

Let $h = \emptyset$, $R = 0$, $C = 0$, $Q = \emptyset$, $c = b$

For each event (\vec{x}, a, r, ρ)

$$\textcircled{1} \quad R \leftarrow R + c \left(\frac{\pi(a|x, h)}{\rho} (r - \hat{r}(x, a)) + \sum_{a'} \pi(a'|x, h) \hat{r}(x, a') \right)$$

$$\textcircled{2} \quad C \leftarrow C + c$$

$$\textcircled{3} \quad Q \leftarrow Q \cup \left\{ \frac{\rho}{\pi(a|x, h)} \right\}$$

$\textcircled{4}$ With probability $\frac{c\pi(a|x, h)}{\rho}$:

$$\textcircled{1} \quad h \leftarrow h + (x, a, r)$$

$$\textcircled{2} \quad c \leftarrow \min\{b, \rho\text{-th quantile of } Q\}$$

Return R/C

The Master Evaluator

Eval(policy π , events $(\vec{x}, a, r, p)^T$, quantile ρ , bound b)

Let $h = \emptyset$, $R = 0$, $C = 0$, $Q = \emptyset$, $c = b$

For each event (\vec{x}, a, r, p)

① $R \leftarrow R + c \left(\frac{\pi(a|x, h)}{p} (r - \hat{r}(x, a)) + \sum_{a'} \pi(a'|x, h) \hat{r}(x, a') \right)$

② $C \leftarrow C + c$

③ $Q \leftarrow Q \cup \left\{ \frac{p}{\pi(a|x, h)} \right\}$

④ With probability $\frac{c\pi(a|x, h)}{p}$:

① $h \leftarrow h + (x, a, r)$

② $c \leftarrow \min\{b, \rho\text{-th quantile of } Q\}$

Return R/C

Incorporates Double Robust + Nonstationary evaluation.

Theorem: Introduces bounded bias + much more efficient.

Empirically, an order of magnitude better for nonstationary eval.

Reject L. Li, W. Chu, J. Langford, and RE Schapire, “A Contextual-Bandit Approach to Personalized News Article Recommendation”, WWW 2010.

Improved L. Li, W. Chu, J. Langford, and X. Huang, “Unbiased Offline Evaluation of Contextual-bandit-based News Article Recommendation Algorithms”, WSDM 2011.

Master M. Dudik, D. Erhan, J. Langford, and L. Li, “Sample-efficient Nonstationary Policy Evaluation for Contextual Bandits”, UAI 2012.