

Intervention Bandits

Blah blah

January 15, 2015

Abstract

An abstract.

1 Introduction

Useful references are: ?.

2 Notation

Assume we have a known causal model with binary variables $\mathbf{X} = \{X_1 \dots X_K\}$ that independently cause a target variable of interest Y . We can run sequential experiments on the system, where at each timestep t we can select a variable on which to intervene and then we observe the complete result, (\mathbf{X}_t, Y_t) . This problem can be viewed as a variant of the multi-armed bandit problem.

Let $p \in [0, 1]^K$ be a fixed and known vector. In each time-step t :

1. The learner chooses an $I_t \in \{1, \dots, K\}$ and $J_t \in \{0, 1\}$.
2. Then $X_t \in \{0, 1\}^K$ is sampled from a product of Bernoulli distributions, $X_{t,i} \sim \text{Bernoulli}(p_i)$
3. The learner observes $\tilde{X}_t \in \{0, 1\}^K$, which is defined by

$$\tilde{X}_{t,i} = \begin{cases} X_{t,i} & \text{if } i \neq I_t \\ J_t & \text{otherwise.} \end{cases}$$

4. The learner receives reward $Y_t \sim \text{Bernoulli}(q(\tilde{X}))$ where $q : \{0, 1\}^K \rightarrow [0, 1]$ is unknown and arbitrary.

The expected reward of taking action i, j is $\mu_{i,j} = \mathbb{E}[q(X) | do(X_i = j)]$. The optimal reward and action are μ^* and (i^*, j^*) respectively, where $(i^*, j^*) = \arg \max_{i,j} \mu_{i,j}$ and $\mu^* = \mu(i^*, j^*)$. The n -step cumulative expected regret is

$$R_n = \mathbb{E} \sum_{t=1}^n (\mu^* - \mu_{I_t, J_t}).$$

3 Estimating $\mu_{i,j}$

The most natural way to estimate $\mu_{i,j}$ is to compute an empirical estimate based on samples when that action was taken. This approach would lead directly to the UCB algorithm with $2K$ actions and a regret bound that depended linearly on K .

In this instance we can significantly outperform the standard approach by exploiting the known causal structure of the problem.

$$\begin{aligned}
P(Y|do(X_i = j)) &= P(Y|X_i = j) \\
&= \sum_b P(Y|X_i = j, X_a = b)P(X_a = b|X_i = j) \\
&= \sum_b P(Y|X_i = j, X_a = b)P(X_a = b), \forall a \in \{1 \dots K\}/i \text{ as } X_a \perp\!\!\!\perp X_i \\
&= \sum_b P(Y|X_i = j, do(X_a = b))P(X_a = b)
\end{aligned}$$

3.1 Estimators

Fix some time-step t and $i \in \{1, \dots, K\}$ and $j \in \{0, 1\}$.

Let $\hat{\mu}_a$ be an empirical estimator for $P(Y|do(X_i = j))$ obtained via marginalization over X_a .

$$\hat{\mu}_a = \begin{cases} \frac{m_{a,1}}{n_{a,1}}p_a + \frac{m_{a,0}}{n_{a,0}}(1 - p_a) & \text{if } a \neq i \\ \frac{m_{i,j}}{n_{i,j}} & \text{if } a = i \end{cases}$$

where:

$$\begin{aligned}
m_{a,b} &= \sum_{s=1}^t \mathbb{1}\{X_i = j, I = a, J = b, Y = 1\}_s \\
n_{a,b} &= \sum_{s=1}^t \mathbb{1}\{X_i = j, I = a, J = b\}_s
\end{aligned}$$

This gives K estimators $\{\hat{\mu}_1 \dots \hat{\mu}_K\}$ to be pooled into a single estimator $\hat{\mu}$.

$$\hat{\mu} = \sum_{a=1}^K w_a \hat{\mu}_a = w_i \frac{m_{i,j}}{n_{i,j}} + \sum_{a \neq i} w_a \left[p_a \frac{m_{a,1}}{n_{a,1}} + (1 - p_a) \frac{m_{a,0}}{n_{a,0}} \right]$$

If p is not known, these expression are unchanged except that p_a is replaced with \hat{p}_a

$$\hat{p}_a = \frac{\sum_{s=1}^t \mathbb{1}\{X_a = 1, I \neq a\}_s}{\sum_{s=1}^t \mathbb{1}\{I \neq a\}_s}$$

3.2 Observe then exploit

Basic idea - observe until the uncertainty on all actions is smaller than ϵ , then exploit (or switch to standard UCB). This is only really going to work if $p_a \sim 0.5$.

Bounds on regret: when observing the regret for each timestep is at most 1 (since the reward is $\in [0, 1]$). When exploiting regret it at most ϵ . If the horizon is n , the cost of exploiting $< \epsilon n$ and the cost of observing is $O(1/\epsilon^2)$ (this comes from how quickly we converge to within ϵ for all arms).

$$R_n = O(n\epsilon + \frac{1}{\epsilon^2}) \quad (1)$$

Differentiating and selecting $\epsilon = (\frac{2}{n})^{1/3}$ to minimize the regret yields:

$$R_n = O(n^{2/3}) \quad (2)$$

For each arm $X_i = j$, we have:

$$\hat{\mu}_{ij} = \frac{m_{ij}}{n_{ij}} \quad (3)$$

where

$$m_{i,j} = \sum_{s=1}^t \mathbb{1}\{X_i = j, Y = 1\} \quad (4)$$

$$n_{i,j} = \sum_{s=1}^t \mathbb{1}\{X_i = j\} \quad (5)$$

$$(6)$$

and Hoeffding's Inequality gives:

$$P\left(|\hat{\mu}_{ij} - \mu_{ij}| > \sqrt{\frac{1}{2n_{ij}} \log \frac{2}{\delta}}\right) \leq \delta \quad (7)$$

4 Observe until number of plausibly optional arms $< \alpha$

Define the set of plausibly optimal arms to be those with an upper confidence bound that is higher than the largest lower confidence bound.

$$A_t = \{k : UB(k) \geq \max_{k'}(LB(k'))\} \quad (8)$$

We then observe until the size of this set is less than some value α before switching to UCB. This takes into account that we don't need narrow confidence bounds on the arms with very low expected reward. Once we start doing UCB, we keep track of the ucb bound separately and for each arm use whichever bounds are narrower.

4.1 Regret during observe phase

$$R_n = \max_{i=1 \dots K} E \left[\sum_{t=1}^n Y_{it} - \sum_{t=1}^n Y_{I_t, t} \right] \quad (9)$$

$$= E[n] (P(Y|X_{i^*} = j^*) - P(Y)) \text{ where } X_{i^*} = j^* \text{ is the best arm} \quad (10)$$

4.1.1 A very simple model

Lets try to analyse a very simple model in which the reward depends only on the value of X_1 .

$$P(Y|X_1 = 0, X_2, \dots, X_N) = \frac{1}{2} - a \quad (11)$$

$$P(Y|X_1 = 1, X_2, \dots, X_N) = \frac{1}{2} + a \quad (12)$$

$$P(X_i = j) = \frac{1}{2} \quad \forall (i, j) \quad (13)$$

$$\implies P(Y) = \frac{1}{2} \quad (14)$$

$$\implies R_n = aE[n] \quad (15)$$

The goal is now to find $E[n]$ (which will be a function of a)

Lets start by considering how long on average we need to observe until the lower bound on the optimal arm is higher than the upper bound on a single other arm, i , (with expected reward $\frac{1}{2}$)

If the following three statements hold, then there is no overlap $LB^* = \hat{\mu}^* - \epsilon > UB_i = \hat{\mu}_i + \epsilon$.

$$\mu^* - \hat{\mu}^* \leq \epsilon \quad (16)$$

$$\hat{\mu}_i - \mu_i \leq \epsilon \quad (17)$$

$$\epsilon \leq a/4 \quad (18)$$

If 18 is true, we can bound the probability of overlap at some timestep t with the union bound of the probability that either 16 or 17 are false.

$$P(\text{overlap}) = P(n \geq t) \leq P(\mu^* - \hat{\mu}^* \geq \epsilon) + P(\hat{\mu}_i - \mu_i \geq \epsilon) \quad \text{if } \epsilon \leq a/4 \quad (19)$$

Let $\epsilon = a/4$

$$P(n > t) \leq e^{-nD(p_{ij} + p(a/4) || p_{ij})} + e^{-nD(p_{ij} - p(a/4) || p_{ij})} \quad (20)$$

$$\leq ? \quad (21)$$

Where $p_{ij} = P(Y = 1, X_i = j) \leq P(X_i = j) = p \quad \forall (i, j)$

Now

$$E[n] = \sum_{t=1}^{\infty} P(n \geq t) \quad (22)$$

But this is just the expected number of steps until one arm does not overlap. With more arms its more complicated - as we need the probability that α of them overlap (another union bound?).

4.2 Bounding a weighted sum of unbiased estimators with McDiarmid's Inequality

McDiarmid's Inequality states: If $X_i \perp\!\!\!\perp X_j$ and

$$|\phi(X_1 \dots X_i \dots X_N) - \phi(X_1 \dots X'_i \dots X_N)| < c_i \quad \forall i \quad (23)$$

$$P(|\phi(\mathbf{X}) - E[\phi(\mathbf{X})]| \geq \epsilon) \leq 2 \exp\left(-\frac{2\epsilon^2}{\sum_i c_i^2}\right) \quad (24)$$

$$P\left(|\phi(\mathbf{X}) - E[\phi(\mathbf{X})]| \geq \sqrt{\frac{\sum_i c_i^2}{2} \log \frac{2}{\delta}}\right) \leq \delta \quad (25)$$

In our problem, for each variable V_i and value j we have:

$$\hat{\mu}_a = \tilde{P}(Y|V_i = j) = \begin{cases} \frac{p_a}{n_{a1}} \sum_{s=1}^{n_{a1}} X_{a1,s} + \frac{1-p_a}{n_{a0}} \sum_{s=1}^{n_{a0}} X_{a0,s} & \text{if } a \neq i \\ \sum_{s=1}^{n_i} X_s & \text{if } a = i \end{cases} \quad (26)$$

where $X_{a1,s} \sim P(Y|V_i = j, do(V_a = 1))$ and $X_{a0,s} \sim P(Y|V_i = j, do(V_a = 0))$, both $\in [0, 1]$. We then pool $\{\hat{\mu}_1 \dots \hat{\mu}_K\}$ to get a single estimate.

$$\hat{\mu} = \tilde{P}(Y|V_i = j) = \sum_a w_a \hat{\mu}_a, \text{ where } \sum_a w_a = 1 \quad (27)$$

Let

$$\hat{\mu} = \phi(\mathbf{X}) = \sum_{a \neq i} w_a \left[\frac{p_a}{n_{a1}} \sum_{s=1}^{n_{a1}} X_{a1,s} + \frac{1-p_a}{n_{a0}} \sum_{s=1}^{n_{a0}} X_{a0,s} \right] + \frac{w_i}{n_i} \sum_{s=1}^{n_i} X_{i,s} \quad (28)$$

Note: this treats the number of samples n_i, n_{a1} , and n_{a0} as fixed. Even for fixed actions, the latter two are still random variables that depend on \mathbf{p} Will this bound still hold?

$$\phi(\dots X_{\alpha 1, s} \dots) - \phi(\dots X'_{\alpha 1, s} \dots) \leq w_\alpha \frac{p_\alpha}{n_{\alpha 1}} \quad (29)$$

$$\phi(\dots X_{\alpha 0, s} \dots) - \phi(\dots X'_{\alpha 0, s} \dots) \leq w_\alpha \frac{1-p_\alpha}{n_{\alpha 0}} \quad (30)$$

$$\phi(\dots X_i \dots) - \phi(\dots X'_i \dots) \leq \frac{w_i}{n_i} \quad (31)$$

$$\sum_i c_i^2 = n_i \left(\frac{w_i}{n_i} \right)^2 + \sum_a \left[n_{a1} \left(w_a \frac{p_a}{n_{a1}} \right)^2 + n_{a0} \left(w_a \frac{1-p_a}{n_{a0}} \right)^2 \right] \quad (32)$$

$$= \frac{w_i^2}{n_i} + \sum_a w_a^2 \left(\frac{p_a^2}{n_{a1}} + \frac{(1-p_a)^2}{n_{a0}} \right) \quad (33)$$

$$= \sum_a w_a^2 f(a), \text{ where} \quad (34)$$

$$f(a) = \begin{cases} \frac{p_a^2}{n_{a1}} + \frac{(1-p_a)^2}{n_{a0}} & a \neq i \\ \frac{1}{n_i} & a = i \end{cases} \quad (35)$$

We want to choose weights w so as to minimize 34 subject to the constraint $\sum_a w_a = 1$.

The minimum (assuming it exists) should occur at a critical point of:

$$L(w_1 \dots w_k, \lambda) = \sum_a w_a^2 f(a) + \lambda \left(\sum_a w_a - 1 \right) \quad (36)$$

$$\frac{\partial L}{\partial w_a} = 2w_a f(a) + \lambda = 0 \quad (37)$$

$$\implies w_a = \frac{-\lambda}{2f(a)} \quad (38)$$

$$\frac{\partial L}{\partial \lambda} = \left(\sum_a w_a \right) - 1 = 0 \quad (39)$$

$$\implies -\frac{\lambda}{2} \sum_a \frac{1}{f(a)} = 1 \implies \lambda = -\frac{2}{\sum_a \frac{1}{f(a)}} \quad (40)$$

$$\implies w_a = \frac{1}{f(a) \sum_a \frac{1}{f(a)}} \quad (41)$$

$$\implies \sum_i c_i^2 = \frac{1}{\sum_a \frac{1}{f(a)}} \quad (42)$$

Substituting 42 into the McDiarmid inequality 25:

$$P \left(|\hat{\mu}^{ij} - \mu^{ij}| > \sqrt{\frac{1}{2 \sum_a \eta_a^{ij}} \log \frac{2}{\delta}} \right) \leq \delta \quad (43)$$

where:

$$\mu^{ij} = P(Y|V_i = j) \quad (44)$$

$$\eta_a^{ij} = \begin{cases} \frac{n_{a1}^{ij} n_{a0}^{ij}}{n_{a1}^{ij} (1-p_a)^2 + n_{a0}^{ij} p_a^2} & a \neq i \\ n_{ij} & a = i \end{cases} \quad (45)$$

$$n_{al}^{ij} = \sum_{s=1}^t \mathbb{1}\{do(V_a = l), V_i = j\} \quad (46)$$

$$n_{ij} = \sum_{s=1}^t \mathbb{1}\{do(V_i = j)\} \quad (47)$$

All of the above is to get estimates for $P(Y|V_i = j)$ for some fixed i, j . Suppose we explore for some fixed total number of rounds h , the goal is to minimize the worst confidence bound - that doesn't quite make sense as a goal - we just need to be sure that all actions are worse than our estimate of the best with probability δ anyway

We have control of are the number of times we do each action, n_{al} , which obviously also influences the random variables n_{al}^{ij} , subject to the constraint $\sum n_{al} = h$.

Objective: minimize the expectation of the largest confidence bound after a fixed number of rounds n .

Let $\mathbf{n} = [n_{11}, n_{10}, n_{21}, n_{20} \dots n_{k1}, n_{k0}]$

$$\tilde{\mathbf{n}} = \arg \min_{\{\mathbf{n}: \|\mathbf{n}\|=h\}} \max_{i,j} \left(E \left[\sqrt{\frac{1}{2 \sum_a \eta_a^{ij}} \log \frac{2}{\delta}} \right] \right) \quad (48)$$

We can bound the first part of η_a^{ij}

$$\eta_a^{ij} = \frac{n_{a1}^{ij} n_{a0}^{ij}}{n_{a1}^{ij} (1-p_a)^2 + n_{a0}^{ij} p_a^2} \leq \frac{n_{a1}^{ij} n_{a0}^{ij}}{\max(n_{a1}^{ij} (1-p_a)^2 + n_{a0}^{ij} p_a^2)} \quad (49)$$

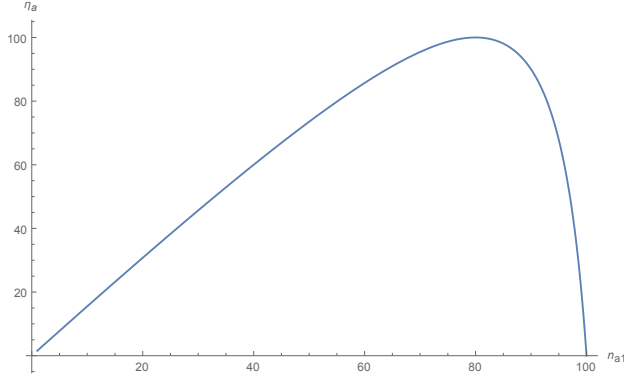
$$= \min \left(\frac{n_{a0}^{ij}}{(1-p_a)^2}, \frac{n_{a1}^{ij}}{p_a^2} \right) \quad (50)$$

$$E \left[\min \left(\frac{n_{a0}^{ij}}{(1-p_a)^2}, \frac{n_{a1}^{ij}}{p_a^2} \right) \right] \leq \min \left(E \left[\frac{n_{a0}^{ij}}{(1-p_a)^2} \right], E \left[\frac{n_{a1}^{ij}}{p_a^2} \right] \right) \quad (51)$$

$$= P(V_i = j) * \min \left(\frac{n_{a0}}{(1-p_a)^2}, \frac{n_{a1}}{p_a^2} \right) \quad (52)$$

$$(53)$$

Figure 1: The effective number of samples η_a versus n_{a1} , where $p_a = .8$ and the total number of samples, $n_{a1} + n_{a0} = 100$. The effective number of samples is maximized (and equals the total) if we sample each side according to its probability.



4.3 Targeted Sampling

Sample a equal number of points from each arm - until the error on all arms is smaller than some threshold. Then exploit. This should do better than observe-exploit, particularly where the probability of some events occurring without intervention is low.

A more sophisticated version:

4.4 UCB variant

Combine estimators according to:

$$w_a = \frac{n_a}{\sum_{a=1}^K n_a} \text{ and } n_{i,j} = \begin{cases} n_{i,j} & \text{if } a = i \\ \frac{1}{2} \min \left\{ \frac{n_{a,1}}{p_a}, \frac{n_{a,0}}{1-p_a} \right\} & \text{otherwise} \end{cases}$$

In this case the estimators from section 3.1 are biased - so its a lot harder to prove a regret bound.

Theorem 1. (Probably False) With probability at least $1 - \delta$ we have that: $|\hat{\mu}_t - \mu| \leq \sqrt{\frac{\beta}{\sum_a n_a} \log \frac{1}{\delta}}$, where $\beta > 0$ is some constant.

Proof. First note that $n_{a,b}$ is a random variable that is bounded by t for all a, b . We use the short-hand $\mu_{i,j}^{a,b} = \mathbb{E}[q(X)|X_i = j, X_a = b]$. Then

$$\mu_{i,j} = p_a \mu_{i,j}^{a,1} + (1 - p_a) \mu_{i,j}^{a,0}.$$

Now we can apply Hoeffding's bound and the union bound to show that

$$\mathbb{P} \left\{ \left| \frac{m_{a,b}}{n_{a,b}} - \mu_{i,j}^{a,b} \right| \geq \sqrt{\frac{1}{2n_{a,b}} \log \frac{4t}{\delta}} \right\} \leq \frac{\delta}{2}.$$

Therefore by the union bound

$$\mathbb{P} \left\{ \left| p_a \frac{m_{a,1}}{n_{a,1}} + (1-p_a) \frac{m_{a,0}}{n_{a,0}} - \mu_{i,j} \right| \geq p_a \sqrt{\frac{1}{2n_{a,1}} \log \frac{4t}{\delta}} + (1-p_a) \sqrt{\frac{1}{2n_{a,0}} \log \frac{4t}{\delta}} \right\} \leq \delta$$

Now by Jensen's inequality

$$\begin{aligned} p_a \sqrt{\frac{1}{2n_{a,1}} \log \frac{4t}{\delta}} + (1-p_a) \sqrt{\frac{1}{2n_{a,0}} \log \frac{4t}{\delta}} &\leq \sqrt{\left(\frac{p_a}{2n_{a,1}} + \frac{1-p_a}{2n_{a,0}} \right) \log \frac{4t}{\delta}} \\ &\leq \sqrt{\max \left\{ \frac{p_a}{n_{a,1}}, \frac{1-p_a}{n_{a,0}} \right\} \log \frac{4t}{\delta}} \\ &= \sqrt{\frac{1}{2n_a} \log \frac{4t}{\delta}}. \end{aligned}$$

Similarly,

$$\mathbb{P} \left\{ \left| \frac{m_{i,j}}{n_{i,j}} - \mu_{i,j} \right| \geq \sqrt{\frac{1}{2n_a} \log \frac{4t}{\delta}} \right\} \leq \mathbb{P} \left\{ \left| \frac{m_{i,j}}{n_{i,j}} - \mu_{i,j} \right| \geq \sqrt{\frac{1}{2n_a} \log \frac{2t}{\delta}} \right\} \leq \delta.$$

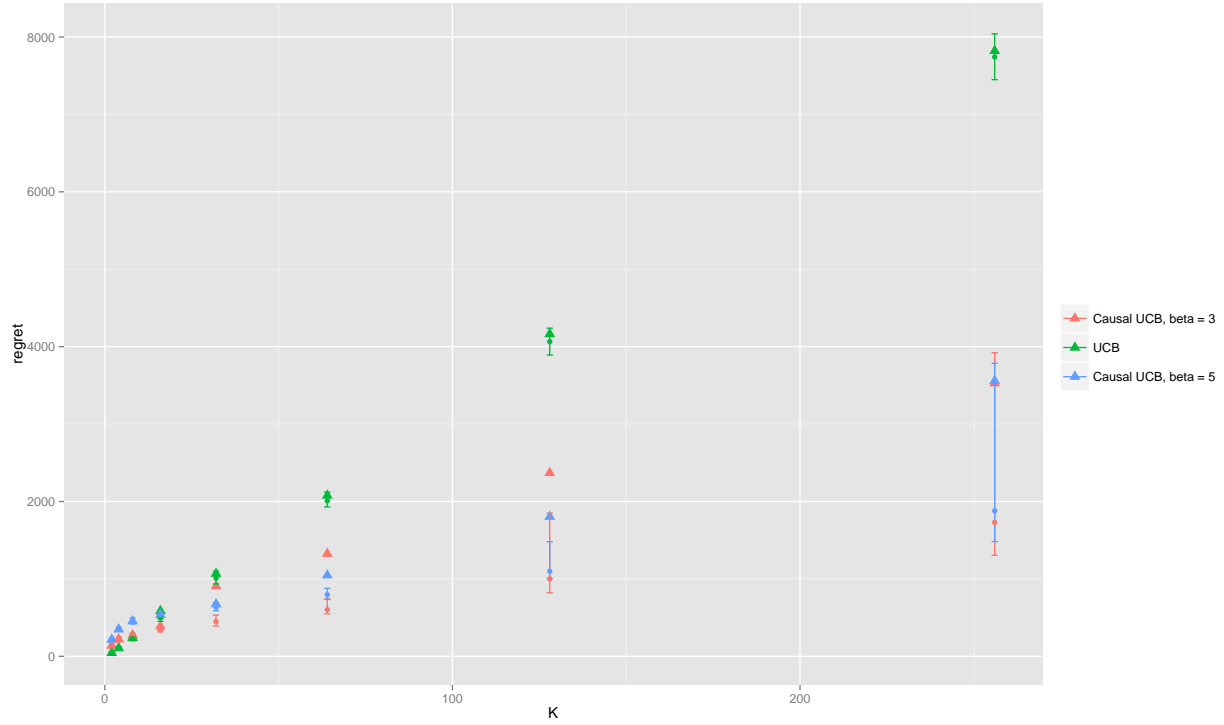
□

5 Algorithm

Algorithm 1 UCB

- 1: **Input:** Number of variables K , vector $p \in [0, 1]^K$, horizon n
 - 2: **for** $t \in 1, \dots, n$ **do**
 - 3: **for** $i \in 1, \dots, K$ **do**
 - 4: **for** $j \in \{0, 1\}$ **do**
 - 5: Compute $\tilde{\mu}_{i,j} = \hat{\mu}_{i,j} + \sqrt{\frac{\alpha}{\sum_a n_a} \log n}$
 - 6: **end for**
 - 7: **end for**
 - 8: Choose $I_t, J_t = \arg \max_{i,j} \tilde{\mu}_{i,j}$
 - 9: **end for**
-

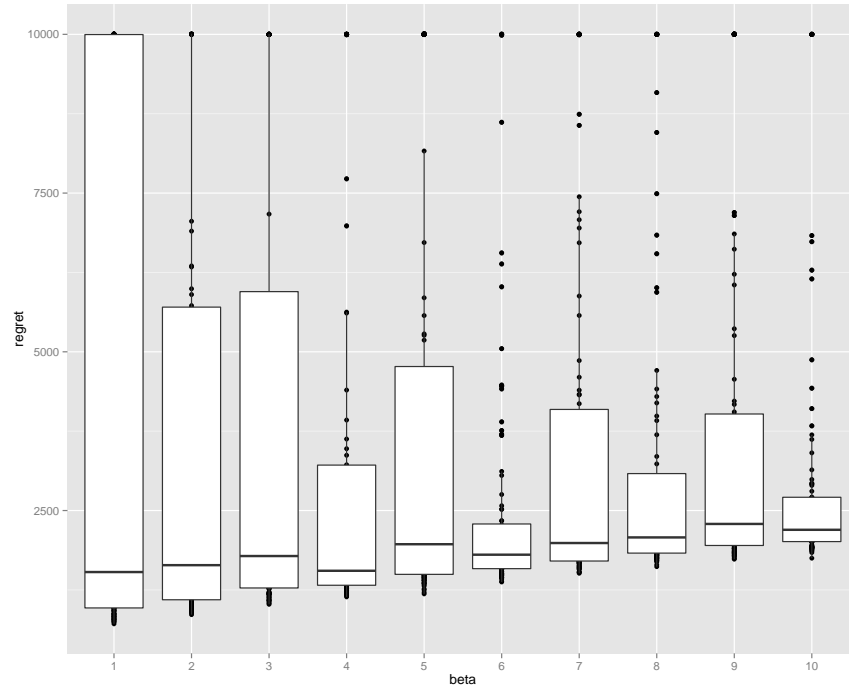
Figure 2: Comparison of the performance of standard UCB versus causal UCB with $\beta = 3$ and $\beta = 5$. 100 simulations were run for each algorithm up to a horizon of 10^5 per value of K . Error bars span the 1st to 3rd quantile of the regret, round points mark the median and triangular points show the mean. For standard UCB the regret increases linearly with the number of arms K . For causal UCB the increase is sub-linear. Increasing β leads to slower convergence but lower variance.



6 Theorems

7 Experiments

Figure 3: The distribution of regret varies with the β parameter in the bound in the estimator. As beta increases, the mean regret increases but the variance decreases. The plot shows the results of running 100 independent bandits, with $K = 256$ and $\epsilon = 0.1$, up to a horizon $h = 10^5$ for each value of β .



Simulations to compare the performance of standard UCB with our modified algorithm. For each number of arms, 100 bandits of each type were created and run upto to a horizon of 1000 timesteps. The mean regret and its standard error from these simulations is plotted in figure ?? The true data was generated from a model where:

$$p = [0.5]^K$$

$$q(\mathbf{X}) = \begin{cases} 0.5 & \text{if } X_1 = 0 \\ 0.6 & \text{otherwise} \end{cases}$$

8 Conclusion

9 Notes

9.1 Why Hoeffdings bound doesn't hold if n is a random variable

Let $\{Z_1, Z_2 \dots Z_n\} \sim \text{Bernoulli}(\frac{1}{2})$ and $X_i = 2Z_i - 1 \implies X_i \in \{-1, 1\}$ and $E[X_i] = 0$.

For a fixed n , Hoeffdings inequality says:

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i\right| > \epsilon\right) \leq 2e^{-n\epsilon^2} \quad (54)$$

$$\implies P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i\right| > \sqrt{\frac{\log 4}{n}}\right) \leq \frac{1}{2} \quad (55)$$

$$\implies P\left(\left|\sum_{i=1}^n X_i\right| > \sqrt{n \log 4}\right) \leq \frac{1}{2} \quad (56)$$

If n is allowed to be dependent on the sequence of values sampled, this inequality no longer holds.

Proof. Choose n based on the sequence of samples seen so far such that:

$$n = \min\{n : n > 4 \text{ and } \sum_{i=1}^n X_i > \sqrt{n \log \log n}\}$$

By the law of iterated logarithms this quantity is finite with probability ~ 1 .

$$\begin{aligned} \implies P\left(\left|\sum_{i=1}^n X_i\right| > \sqrt{n \log \log n}\right) &\sim 1 \\ \implies P\left(\left|\sum_{i=1}^n X_i\right| > \sqrt{n \log 4}\right) &\sim 1, \text{ Since } \log \log n > \log 4 \quad \forall n > 4 \end{aligned}$$

Thus Hoeffdings inequality (equation 54) does not hold in general if n is not independent of the samples $\{X_i\}$ (The bound does not work if you decide to stop sampling as soon as you reach a point where random walk fluctuations take you outside it)

□