

UCB Formulation: Establishes Some Confidence

2024年5月20日 星期一 17:01

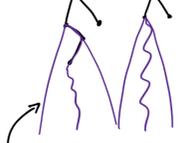
More expansion, more correct

Essentially putting the tree to 2 parts

1. Rational tree growing
2. Data that help to understand the quality better

more you grow the tree, less random Monte Carlo sample

more rational part on the top



Enough coverage, correct values (when growing, distribution is dynamic)

What is the optimal in 2 degree of stochastic



$\mu_1 = 0.8$ TRUE
 $\mu_2 = 0.2$ IE[X₁]
IE[X₂]

$X_1 \begin{cases} 1 & 0.8 \\ 0 & 0.2 \end{cases}$ r.v.
 $X_2 \begin{cases} 1 & 0.2 \\ 0 & 0.8 \end{cases}$ r.v.

If you know this, $\mu_1(x_1)$ is better but when don't know, essentially same with a learning

option \rightarrow some probabilistic distribution \rightarrow each time only see 1 trajectory of such distribution

both stochastic outcome

What to actually choose, what is a mathematical strategy

In ϵ -greedy or 50%/50% choose: There would be certain % of option that is wrong, in the long run essentially the same in expectation

Regret

Mathematical construct

$$IE[\sum_{i=1}^T X_{C_i}]$$

Variable on coin also random 2
value of coin is random

Coin $C_i \in \{C_{blue}, C_{red}\}$ random 1
choosing coin is random

2 stochastic flows taken into it

Optimistically:

Always play the better coin μ^*

Best value: $\mu^* \cdot T$

Lowest value: $\mu_{low} \cdot T$

Regret is the gap between getting and optimal

$$R_0(T) = IE[\sum_{i=0}^T X_{C_i^*} - \sum_{i=0}^T X_{C_i}]$$

some are C^* , some are not
Zero regret when $\sum X_{C_i} = \sum X_{C_i^*} \rightarrow (\mu^* - \mu) = \Delta$
Make one given coin value not random

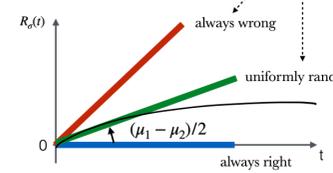
50% 50% Policy:

$$\frac{1}{2}T(\mu^* - \mu)$$

"Linear Regret"

All naive approach is linear regret (No learning, look at the same with fresh eye)
Come up with something not linear?

- Non learning approach boils to the same idea



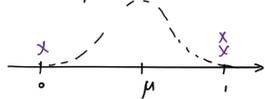
Monotonic increase, can't decrease, but can regret less and less, which is the definition of where learning actually is

$IE[R(N)] \sim O(N^{2/3} (K \log N)^{1/3})$ better than linear a bit for ϵ greedy

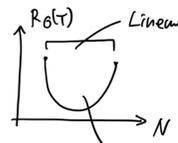
$IE[R(N)] \sim O(\log N)$ for UCB (The best)

Give it a confidence

Essentially building confidence interval, statistics is useful



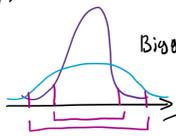
$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \text{ as } n \rightarrow \infty, \bar{x} \rightarrow \mu \text{ by CLT and LLN}$$



Something in the middle

Compare distributions with Concentration Inequality and Pin Point expectation!
(Chebyshev inequality is a Concentration inequality)

Assuming the sample is somewhat consistent with statistical beliefs



Bigger Tail when not concentrated

Give ϵ range and ask how likely to have data point there

$$IP(|\bar{x} - \mu| > \epsilon) \leq 2e^{-2N\epsilon^2}$$

Probability of deviation from ϵ neighborhood

Sample tolerance for deviation from the expectation

Big ϵ : Not concentrated, less convergence

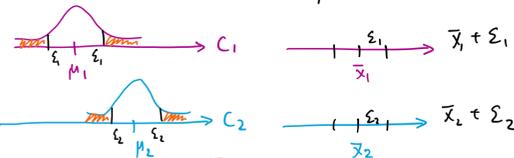
Small ϵ : concentrated, more convergence, more sure about actual value

1. $\epsilon \uparrow$, P \downarrow (Greater confidence, less accuracy to match μ)
2. $N \uparrow$, Tail decrease exponentially fast \rightarrow more concentrated \rightarrow less ϵ needed for same accuracy of knowing μ

Adaptive ϵ choice: Bonifert of doubt

Assume that we set a starting point where set ϵ as $\sqrt{\frac{2 \log T}{N_i}} \rightarrow 2e^{-2N_i \frac{2 \log T}{N_i}} = 2e^{-4 \log T} \sim O(\frac{1}{T^2})$

Adaptive ϵ , band the tail in the same way



Optimistically: Choose with sample t doubt accordingly who has higher value

Quickly decrease what's outside of ϵ
As T (total sample) \uparrow , probability of deviation shrinks really quickly

Optimism in face of uncertainty: Since μ is in there range, and we don't know the true μ to choose the better coin, the best true can be do for μ_1 is $\bar{x}_1 \in \epsilon_1$ (bonifert of doubt)

of choosing less optimal coin (Accumulation of Regret)

Assume C_1 bad and C_2 good and the sample happens to say C_1 is bad

$$\bar{x}_1 + \epsilon_1 \geq \bar{x}_2 + \epsilon_2$$

$$\left\{ \begin{array}{l} \bar{x}_1 \in [\mu_1 - \epsilon_1, \mu_1 + \epsilon_1] \\ \bar{x}_2 \in [\mu_2 - \epsilon_2, \mu_2 + \epsilon_2] \end{array} \right\}$$

Assume that our sample is within true value \pm dashes
Assume suboptimal higher and optimal lower

$$\mu_1 + \epsilon_1 \geq \mu_2 - \epsilon_2$$

Relax to upper band
Constrict to lower band

With particular ϵ , Prob of μ not in $\bar{x} \pm \epsilon$ or \bar{x} not in $\mu \pm \epsilon$ decrease exponentially

$$\mu_1 + \epsilon_1 + \epsilon_1 \geq \bar{x}_1 + \epsilon_1 \geq \bar{x}_2 + \epsilon_2 \geq \mu_2 + \epsilon_2 - \epsilon_2$$

Connection between sample \bar{x} and expectation μ
Sample \pm dashes contain μ , $\mu \pm$ dashes contain sample (assume)

$$= \mu_1 + 2\epsilon_1 \geq \mu_2$$

$$= 2\epsilon_1 \geq \mu_1 - \mu_2 = 2\sqrt{\frac{2 \log T}{N_1}} \geq \Delta$$

(dist between μ_2 and μ_1)

$$N_1 \leq \frac{8 \log T}{\Delta^2}$$

NO way for the number of choosing less optimal coin (N_1) to be greater than log growth for the total number of coin chosen.

Exponential less play with less optimal coin

\rightarrow only small enough sample misguide you

Intuition 水落石出 (ϵ 是水, μ 是石)

If bonifert of doubt choose the bad coin, well, choose it, and then the band constricts really quickly, small ϵ and the next round the true bad coin's high value would be better shown

ONLY the true appears

More tries, less advantage of ϵ , more constrict, more confident about the value, give less doubt bonifert when comparison

UCB is everywhere in Decision making and Algorithm

Practical: UCT

Assume the 2 option fixed distribution: UCB proved to be deriving from i.i.d. samples

\rightarrow AS you grow the tree, the distribution changes, favour the better one

\rightarrow Not i.i.d. any more

$\bar{x} + C \sqrt{\frac{2 \log T}{N}}$ Progress driven to reduce hyperparameters, when mean is not damn near (clear), there are more hyperparameters attached

Hence the name UCB on Tree with out former proof \rightarrow UCT

