The two opposing forces

    # 1   collect data ( exploration )
    # 2   select choice   with highest win rate (exploitation )


Algorithm

- Epsilon - greedy
- optimistic initial values
- UCB 1
- Thompson sampling


CTR の 13°l

- 2つの値、どっちにする？
  → 100 たーうし , 1000 たーうし でのCTRをCて採.
  100 , 1000ぴ要ずが不明 ( infinite number のたーうしが対象 )


Epsilon- Greedy

- greedy - method

  - picking the bandit with highest MLE win rate with no regard to confidence in prediction or amount of data.

  ex ) bandits with 90% and 80%
  $$E(R) = (1-\varepsilon) 0.9 + \varepsilon \left( \frac{0.8 + 0.9}{2} \right)$$


Sample mean

- $x$ 2) は傾にするので , $\bar{x}_n = \frac{1}{N} \sum_{c=1}^{N} x_i$ から 2分に p゛/分裂
  → $O(N)$ と $O(1)$ になてない

$$\overline{X}_N = \frac{1}{N}\left(\sum_{i=1}^{N-1} X_i + X_N\right)$$

$$\overline{X}_{N-1} = \frac{1}{N-1}\sum_{i=1}^{N-1} X_i \quad\longleftrightarrow\quad (N-1)\overline{X}_{N-1} = \sum_{i=1}^{N-1} X_i$$

$$\overline{X}_N = \frac{1}{N}\left((N-1)\overline{X}_{N-1} + X_N\right)$$

$$= \frac{N-1}{N}\overline{X}_{N-1} + \frac{1}{N}X_N$$

$$= \left(1 - \frac{1}{N}\right)\overline{X}_{N-1} + \frac{1}{N}X_N$$

$$= X_{N-1} + \frac{1}{N}\left(X_N - \overline{X}_{N-1}\right)$$

## Optimistic Initial Values 　　　　実(価値?)

estimation の初期値を大きく設定する

→ arithmetic average はデータを集めれば集めるだけ小さくなる

→ もしあるbandit の estimated mean が max estimated mean より $d\sim1$ くらつたら, 探索してない

## Pole of Initial Value

→ High initial value : 推定値が真の値に下がるまで時間がより必要

$\updownarrow$

Low initial value : 推定値が真の値に下かっで時間はすぐ上がすくかからない

# UCB 1

- upper confidence bound
- so far

  Epsilon greedy : small prob of random exploration

  Optimistic : naturally start at large value (each bandit will be chosen often)

- applying probability

  $p(\text{sample mean} - \text{true mean} \geq \text{error}) <= f(\text{error})$

  et) $p(\text{sample mean} - \text{true mean} \geq t) \leq 1/t$

  - 大きな誤差 より 大きくなる prob は 小さく
  - $N \sim I$ なら ..  大きくなる

  Markov inequality / Chebyshev inequality / Hoeffding

  $$p(\bar{X}_n - E(X) \geq t) \leq e^{-2nt^2}$$

  <span style="color:red">↓<br>n が大きくなると、右辺は<br>小さくなる</span>

  $$\hat{j} = \underset{\hat{j}}{\text{argmax}} \left( \bar{X}_{nj} + \sqrt{2\,\frac{\log N}{n\hat{j}}} \right)$$

  ↓
  sample mean

  $N$ : total plays
  $n\hat{j}$ : plays made on bandit $\hat{j}$

# Thompson Sampling

CI を思い出す。

Small Dataset : Large confidence interval

Fat → Explore more , skinny → explore less

CLT : sum of RVs     RVs ~= normal Distribution

what we want?

    — the distribution of mean (win rate)

    → ∴ a distribution の parameter を

    (何)か a distribution に従う

$P(X|\theta)$ : 尤度

ある θ が与えられた時の, データの確率

なぜ evidence を考えずに, にて例 としていいか?

    → evidence は 年表側であり, 計算が難しい

    ( Monte Carlo は 計算がへで — )

事後分布 $\propto$ 尤度 ・ 事前分布

共役事前分布 を使えよ

ex) Gaussian $\propto$ Gaussian × Gaussian

/s. 尤度は ベルヌーイ分布

    → この 共役事前分布 ( conjugate prior ) は

Beta 分布。

$$p(\theta|X) \propto \left( \prod_{i=0}^{N} \theta^{x_i} (1-\theta)^{(1-x_i)} \right) \left( \frac{1}{B(\alpha,\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \right)$$

$$\propto \left( \prod_{i=1}^{N} \theta^{x_i} (1-\theta)^{(1-x_i)} \right) \left( \theta^{\alpha-1} (1-\theta)^{\beta-1} \right)$$

$$= \left( \theta^{\sum_{i=1}^{N} x_i} (1-\theta)^{\sum_{i=1}^{N}(1-x_i)} \right) \left( \theta^{\alpha-1} (1-\theta)^{\beta-1} \right)$$

$$= \left( \theta^{\sum_{i=1}^{N} x_i + \alpha - 1} (1-\theta)^{\sum_{i=1}^{N}(1-x_i) + \beta - 1} \right)$$

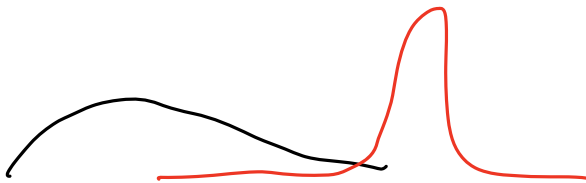どのように prior を推定する？

→ Beta (1, 1) は一様分布
もしドメイン知識があれば それを使用する.

$X$ を得る度に分布を更新する

ex) prior = Beta(1.1) , collect 1, post (|+1, |×|-1)=
                  Beta(2,1)
          (2.1)               1             Beta(3,1)
          (3.1)               0             Beta(3,2)



→ このような分布では, 黒はあまり選択されない
suboptimal な分布を使用してよい

虎盤



1 が出てないことがダメ L,
コブにてする

Rewards coming from normal distribution

- precision = 1 / variance

$$p(X \mid \mu, \tau) = \prod_{i=1}^{N} \sqrt{\frac{\tau}{2\pi}} \, e^{-\frac{\tau}{2}(x_i - \mu)^2}$$

$$X \sim N(\mu, \tau^{-1}) \quad, \quad \mu \mid X \sim N(m, \partial^{-1})$$

$$p(\mu \mid X) = \frac{p(\mu, X)}{p(X)} = \frac{p(X \mid \mu) \, p(\mu)}{p(X)}$$

$$X \sim N(\mu, \tau^{-1}) \quad, \quad \mu \sim N(m_0, \partial_0^{-1}) \quad, \quad \mu \mid X \sim N(m, \partial^{-1})$$

$$p(\mu \mid X) \propto p(X \mid \mu) \, p(\mu)$$

$$= \left( \prod_{i=1}^{N} \sqrt{\frac{\tau}{2\pi}} \, e^{-\frac{\tau}{2}(x_i - \mu)^2} \right) \left( \sqrt{\frac{\partial_0}{2\pi}} \, e^{-\frac{\partial_0}{2}(\mu - m_0)^2} \right)$$

$$= \left( \left[ \frac{\tau}{2\pi} \right]^N e^{-\frac{\tau}{2} \sum_{i=1}^{N}(x_i - \mu)^2} \right) \left( \quad \cdots \quad \right)$$

<span style="color:red">const $\leftarrow \sqrt{\frac{\partial_0}{2\pi}}$</span>

$$\propto \left( e^{-\frac{\tau}{2} \sum_{i=1}^{N}(x_i - \mu)^2} \right) \left( e^{-\frac{\partial_0}{2}(\mu - m_0)^2} \right)$$

$$= e^{-\frac{\tau}{2} \sum_{i=1}^{N}(x_i - \mu) - \frac{\partial_0}{2}(\mu - m_0)^2}$$

$$= e^{-\frac{\tau}{2} \sum_{i}^{N}(\mu^2 - 2\mu x_i + x_i)^2 - \frac{\partial_0}{2}(\mu^2 - 2\mu m_0 + m_0^2)}$$

$$= \exp\left( -\frac{\tau}{2} \sum_{i=1}^{N}(\mu^2 - 2\mu x_i + x_i^2) - \frac{\partial_0}{2}(\mu^2 - 2\mu m_0 + m_0^2) \right)$$

$$= \exp\left( -\frac{\tau}{2}\left( N\mu^2 - 2\mu \sum_{i=1}^{N} x_i + \sum_{i=1}^{N} x_i^2 \right) - \frac{\partial_0}{2}\mu^2 - \frac{\partial_0}{2} 2\mu m_0 - \frac{\partial_0}{2} m_0^2 \right)$$

$$\propto \exp\left( -\frac{\tau}{2}\left( N\mu^2 - 2\mu \sum_{i=1}^{N} x_i \right) - \frac{\partial_0}{2}(\mu^2 - 2\mu m_0) \right)$$

$$= \exp\left( -\frac{\tau N + \partial_0}{2}\mu^2 + \left( \tau \sum_{i=1}^{N} x_i + \partial_0 m_0 \right)\mu \right)$$

$$p(\mu|x) = \sqrt{\frac{\lambda}{2\pi}} \exp\left(-\frac{\lambda}{2}(\mu - m)^2\right)$$

$$= \sqrt{\frac{\lambda}{2\pi}} \exp\left(-\frac{\lambda}{2}(\mu^2 - 2m\mu + m^2)\right)$$

$$\propto \exp\left(-\frac{\lambda}{2}(\mu^2 - 2m\mu)\right)$$

$$= \exp\left(-\frac{\lambda}{2}\mu^2 + m\lambda\mu\right)$$

$$\exp\left(-\frac{\tau N + \lambda_0}{2}\mu^2 + \left(\tau \sum_{i=1}^{N} x_i + \lambda_0 m_0\right)\mu\right)$$

$$||$$

$$\exp\left(-\frac{\lambda}{2}\mu^2 + m\lambda\mu\right)$$

$$\lambda = \tau N + \lambda_0$$

$$m\lambda = \tau \sum_{i=1}^{N} x_i + \lambda_0 m_0 \iff m = \frac{1}{\lambda}\left(\tau \sum_{i=1}^{N} x_i + \lambda_0 m_0\right)$$

$$m = \frac{1}{\tau N + \lambda_0}\left(\tau \sum_{i=1}^{N} x_i + \lambda_0 m_0\right)$$

$$\frac{x - \mu}{\sigma} = z$$

$$x = \sigma z + \mu$$