

Lowerbound:

$$\mathbb{E}[\tau] \geq C^*(v) \log\left(\frac{1}{\delta}\right)$$

$$C^*(v)^{-1} = \sup_{\alpha \in P_K} \inf_{v' \in \mathcal{E}_{\text{alt}}(v)} \sum \alpha_i D(v_i, v'_i)$$

Insights:

$$\alpha^* \approx \frac{\mathbb{E}[N_i(\tau)]}{\mathbb{E}[\tau]} \approx \text{ratio of expected number of pulls of arm } i \text{ up to stoppage to expected stoppage} \quad (1)$$

$$\inf_{v' \in \mathcal{E}_{\text{alt}}(v)} \sum \mathbb{E}_{v'}[N_i(\tau)] D(v_i, v'_i) \approx \log\left(\frac{1}{\delta}\right) \quad (2)$$

The condition (1) guides "sampling"
The condition (2) guides "stopping"

TRACK & STOP

Sampling: Pull $\arg\max_i (\alpha_i^*(\hat{v})t - N_i(t))$

1) Here \hat{v} is an estimate of the instance v . Further, this is generally possible when the family is expressed as a single-parameter family or so

2) $\alpha_i^*(\hat{v})t$ is like a target pull fraction
 $N_i(t)$ is the actual value

This is the "Track" part of the algorithm

3) α^* is the value of α which solves the optimization problem

$$\sup_{\alpha \in P_K} \inf_{v' \in \mathcal{E}_{\text{alt}}(v)} \sum \alpha_i D(v_i, v'_i)$$

$\alpha^*(\hat{v})$ means that we feed the estimate to the optimization problem, instead of the true value.

Stopping: $\inf_{v' \in \mathcal{E}_{\text{alt}}(\hat{v})} \sum_{i=1}^k N_i(t) D(\hat{v}_i, v'_i) \geq \log\left(\frac{1}{\delta}\right)$

1) The space over which we are infimizing is the space of alternate solutions of the instance \hat{v} .

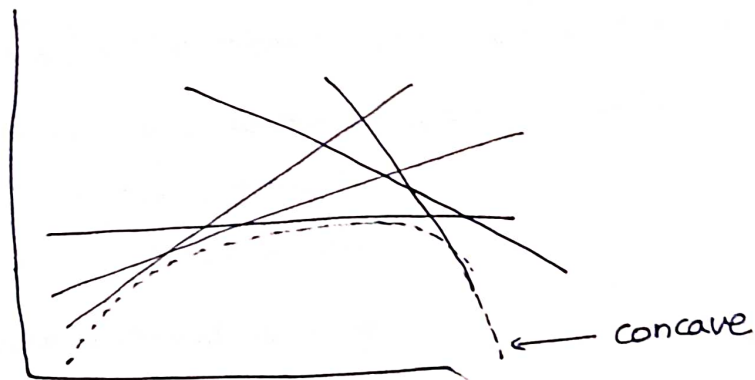
When can this go wrong?

1) This will work only when $\hat{v} \rightarrow v$. This is possible only when all arms get sampling enough so that convergence is possible asymptotically over time.
Hence, this algorithm requires FORCED EXPLORATION

2) If $\delta \rightarrow 0$, then $\mathbb{E}[T] \rightarrow \infty$, $\hat{v} \rightarrow v$

Track and Stop was proposed by Garivier & Kaufmann (2016).

Note
1) The infimum of a family of function is concave



2) Smoothness & Differentiability is not guaranteed

3) Computation of α^* is a problem in convex optimization

Assume all arms to be Bernoulli

For arms a or b ($a \neq b$)

$$\bar{Z}_{a,b}(t) = \log \left(\frac{\max_{\mu_a' \geq \mu_b'} L_{\mu_a'}(X_{t,a}) L_{\mu_b'}(X_{t,b})}{\max_{\mu_b' \geq \mu_a'} L_{\mu_a'}(X_{t,a}) L_{\mu_b'}(X_{t,b})} \right)$$

μ_a' is $P(1$ on instance μ' for arm a)

$X_{t,a}$ is observation of arm a

$L_{\mu_a'}$ is the likelihood under observations of arm a

Numerator: a is better than b

Denominator: b is better than a

The statistic $\bar{Z}_{a,b}(t)$ assigns the ratio of the hypotheses "arm a is superior to arm b " and compares it to "arm b is better than arm a "

$$L_{\mu_a'}(X_{t,a}) = (\mu_a')^{N_a^1} (1 - \mu_a')^{N_a^0}$$

$\bar{Z}_{a,b}(b)$ is termed a "GENERALIZED LOG-LIKELIHOOD STATISTIC"

Claim 1: For $\hat{\mu}_a \geq \hat{\mu}_b$

$$Z_{a,b}(t) = N_a(t) d(\hat{\mu}_a, \hat{\mu}_{a,b}) + N_b(t) d(\hat{\mu}_b, \hat{\mu}_{a,b})$$

$$\text{where } \hat{\mu}_{a,b} = \frac{N_a(t)}{N_a(t) + N_b(t)} \hat{\mu}_a + \frac{N_b(t)}{N_a(t) + N_b(t)} \hat{\mu}_b$$

$d(\hat{\mu}_a, \hat{\mu}_{a,b})$ is the relative entropy between the two Bernoulli dist.

$\hat{\mu}_{a,b}$ is the average number (fraction) of heads

$$= \frac{\text{total \# heads}}{\text{total \# pulls}}$$

* $Z_{a,b}(t) = -Z_{b,a}(t)$, which is trivial to prove

$$Z(t) = \max_a \min_{b \neq a} Z_{a,b}(t)$$

$Z(t)$ is our stoppage statistic

The outer maximization will be achieved by ONLY the empirically best arm.
 $Z(t)$ hence tells us how well separated the best arm is from the closest challenging arm.

Claim 2

$$Z(t) = \min_{\lambda \in \mathcal{E}_{\text{alt}}(\hat{\mu})} \sum \frac{N_a(t)}{t} d(\hat{\mu}_a, \lambda_a)$$

The claim is that stopping via $Z(t)$ is the same as achieving the information-theoretic bound which we showed earlier.

It is absolutely identical to $\inf_{\lambda' \in \mathcal{E}_{\text{alt}}(\hat{\mu})} \sum_{i=1}^k N_i(t) D(\hat{\nu}_i, \lambda'_i) \geq \log\left(\frac{1}{\delta}\right)$

We STOP when $Z(t) \geq \beta(t, \delta)$, a certain threshold

FORMAL DEFINITION OF TRACK & STOP

• Pull each arm once

• While $Z(t) \leq \beta(t, \delta) = \log\left(\frac{2t(k-1)}{\delta}\right)$,

- if $\arg\max_i N_i(t) < \sqrt{t}$, pull $\arg\min_i N_i(t)$ ①

- else pull $\arg\max_i (td_i^*(\hat{\mu}) - N_i(t))$ ②

① IS FORCED EXPLORATION ② IS TRACKING

• Output $i^*(\hat{\mu})$

Thm 1: Track and Stop is δ -sound

Thm 2: $\limsup_{\delta \downarrow 0} \frac{\mathbb{E}[\tau]}{\log(1/\delta)} \leq C^*(\nu) \leftarrow \text{Asymptotic Optimality}$

Note that we cannot guarantee anything about a fixed δ .

Eventually, we can also argue equality to hold in the Thm 2.

Thm 2 talks about the expected stopping time, we can also have an almost sure inequality

Thm 3: $\mathbb{P}_\nu \left(\left(\limsup_{\delta \downarrow 0} \frac{\tau}{\log(1/\delta)} \right) \leq C^*(\nu) \right) = 1$