

On the Complexity of Best Arm Identification with Fixed Confidence

Discrete Optimization in the Presence of Noise

Aurélien Garivier[†], joint work with Emilie Kaufmann^{*}

Rencontres UT1-UT3, 20 septembre 2017

[†] Institut de Mathématiques de Toulouse
LabeX CIMI
Université Paul Sabatier, France

^{*} Université Lille, CNRS UMR 9189
Laboratoire CRISTAL
F-59000 Lille, France

Table of contents

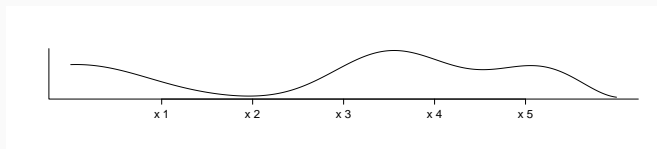
1. The Problem
2. Prolegomenon: Large Deviation Bounds for Bandits
3. Lower Bound
4. The Track-and-Stop Strategy

The Problem

Best-Arm Identification with Fixed Confidence

K options = probability distributions $\nu = (\nu_a)_{1 \leq a \leq K}$

$\nu_a \in \mathcal{F}$ exponential family parameterized by its expectation μ_a



At round t , you may:

- choose an **option** $A_t = \phi_t(A_1, X_1, \dots, A_{t-1}, X_{t-1}) \in \{1, \dots, K\}$
- observe a new **independent sample** $X_t \sim \nu_{A_t}$

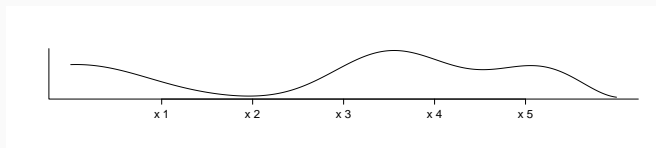
so as to **identify the best option** $a^* = \operatorname{argmax}_a \mu_a$ and $\mu^* = \max_a \mu_a$
as fast as possible: **stopping time** τ .

| Fixed-budget setting | Fixed-confidence setting |
|---|--|
| given $\tau = T$ | minimize $\mathbb{E}[\tau]$ |
| minimize $\mathbb{P}(\hat{a}_\tau \neq a^*)$ | under constraint $\mathbb{P}(\hat{a}_\tau \neq a^*) \leq \delta$ |

Best-Arm Identification with Fixed Confidence

K options = probability distributions $\nu = (\nu_a)_{1 \leq a \leq K}$

$\nu_a \in \mathcal{F}$ exponential family parameterized by its expectation μ_a



At round t , you may:

- choose an option $A_t = \phi_t(A_1, X_1, \dots, A_{t-1}, X_{t-1}) \in \{1, \dots, K\}$
- observe a new independent sample $X_t \sim \nu_{A_t}$

so as to identify the best option $a^* = \operatorname{argmax}_a \mu_a$ and $\mu^* = \max_a \mu_a$
as fast as possible: stopping time τ_δ .

| Fixed-budget setting | Fixed-confidence setting |
|--|--|
| given $\tau = T$ | minimize $\mathbb{E}[\tau_\delta]$ |
| minimize $\mathbb{P}(\hat{a}_\tau \neq a^*)$ | under constraint $\mathbb{P}(\hat{a}_\tau \neq a^*) \leq \delta$ |

Intuition: a Simple Example

Most simple setting: for all $a \in \{1, \dots, K\}$,

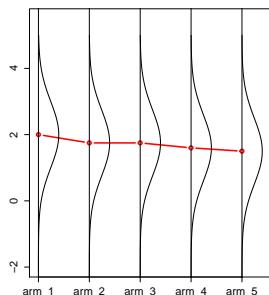
$$\nu_a = \mathcal{N}(\mu_a, 1)$$

For example: $\mu = [2, 1.75, 1.75, 1.6, 1.5]$.

At time t :

→ you have sampled n_a times the option a

→ your empirical average is \bar{X}_{a,n_a} .



→ if you stop at time t , your **probability of preferring arm $a \geq 2$ to arm $a^* = 1$** is:

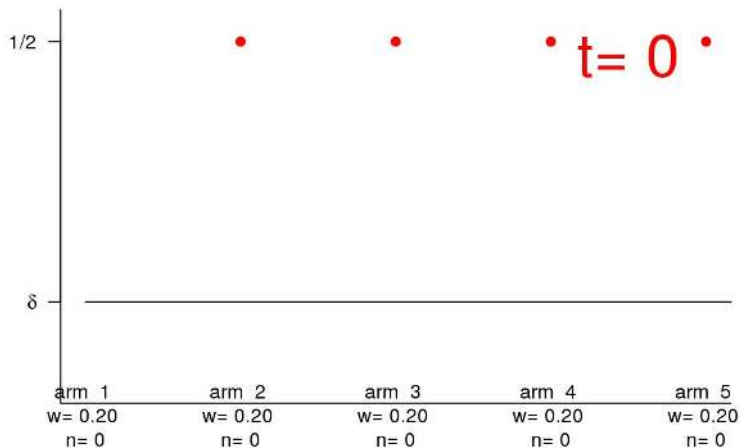
$$\begin{aligned} \mathbb{P}(\bar{X}_{a,n_a} > \bar{X}_{1,n_1}) &= \mathbb{P}\left(\frac{\bar{X}_{a,n_a} - \mu_a - (\bar{X}_{1,n_1} - \mu_1)}{\sqrt{1/n_1 + 1/n_a}} > \frac{\mu_1 - \mu_a}{\sqrt{1/n_1 + 1/n_a}}\right) \\ &= \bar{\Phi}\left(\frac{\mu_1 - \mu_a}{\sqrt{1/n_1 + 1/n_a}}\right) \end{aligned}$$

where $\bar{\Phi}(u) = \int_u^\infty \frac{e^{-u^2/2}}{\sqrt{2\pi}} du$

Uniform Sampling



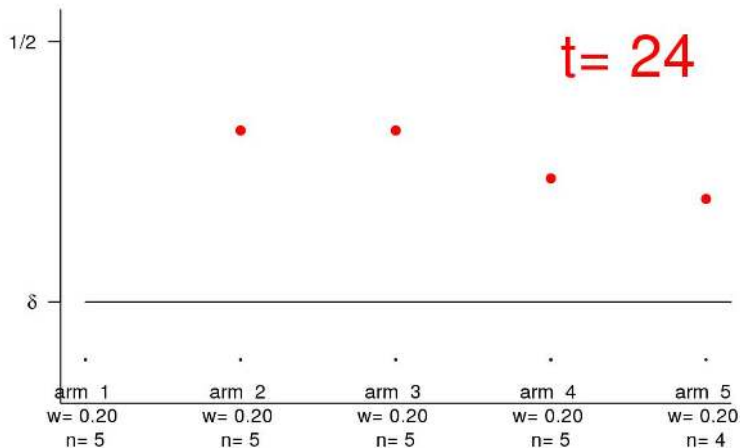
P(confusion)



Uniform Sampling



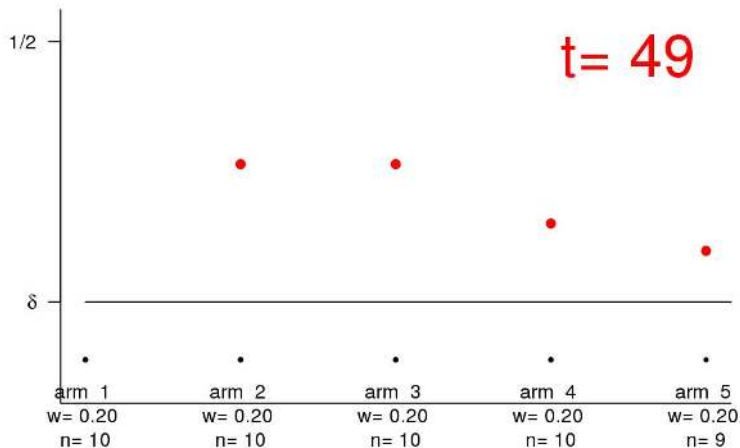
P(confusion)



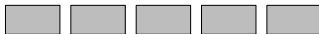
Uniform Sampling



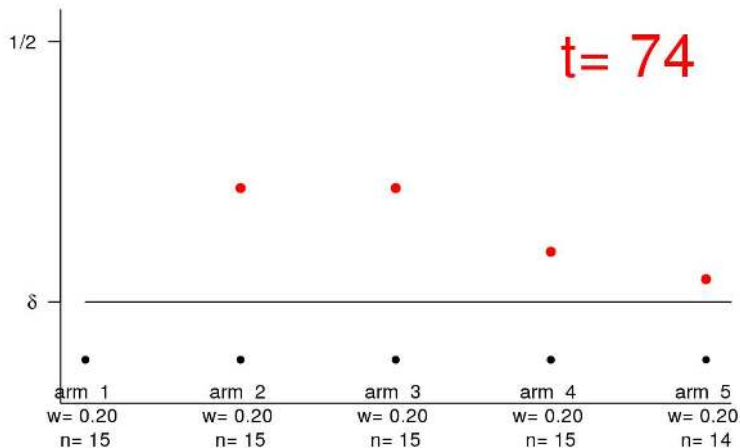
P(confusion)



Uniform Sampling



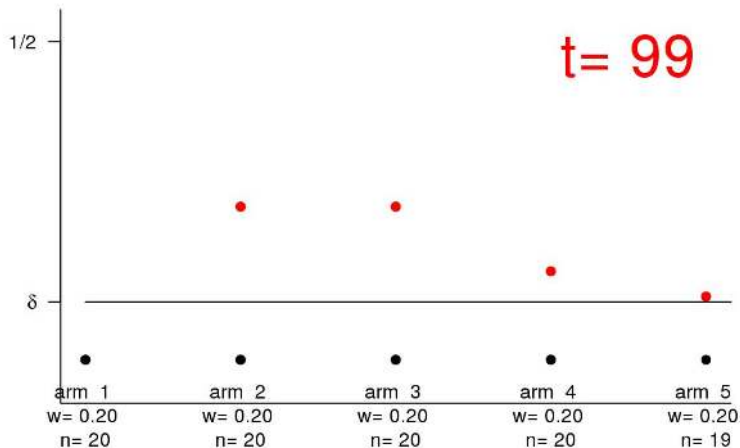
P(confusion)



Uniform Sampling



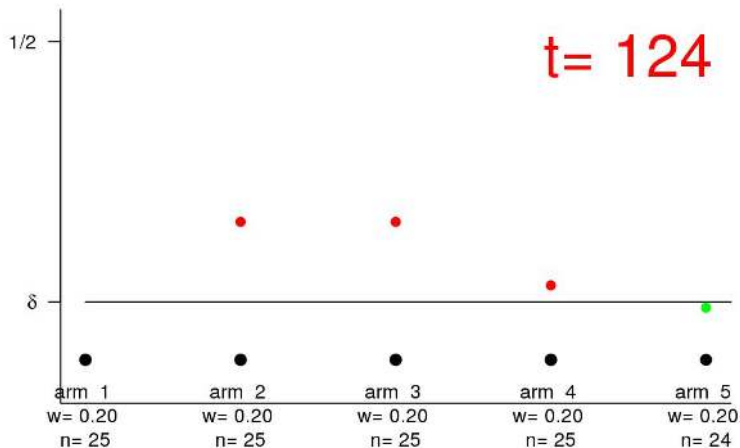
P(confusion)



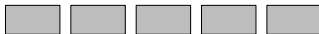
Uniform Sampling



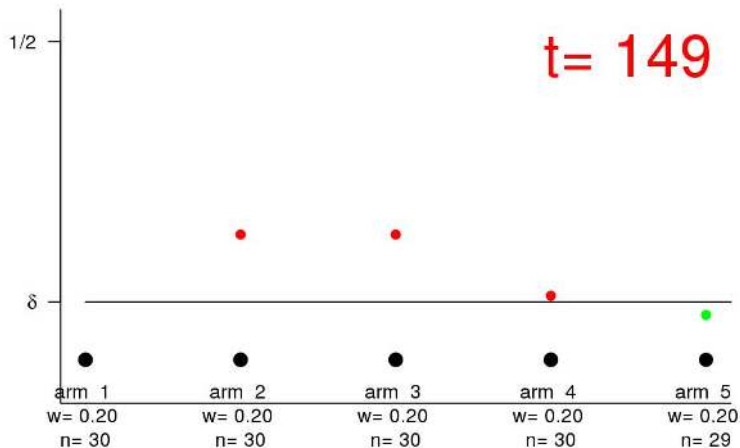
P(confusion)



Uniform Sampling



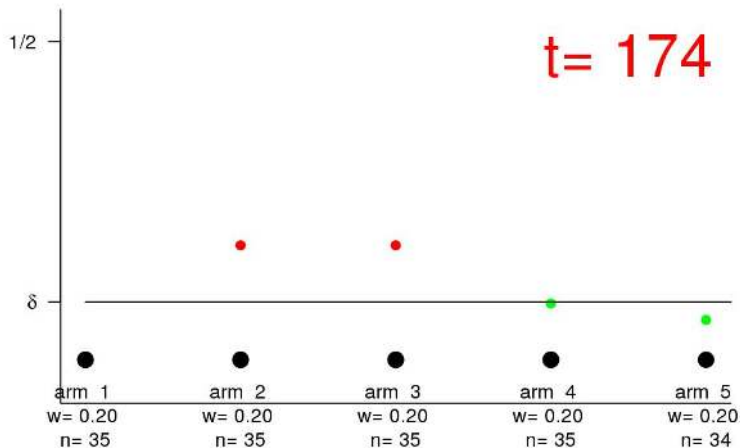
P(confusion)



Uniform Sampling



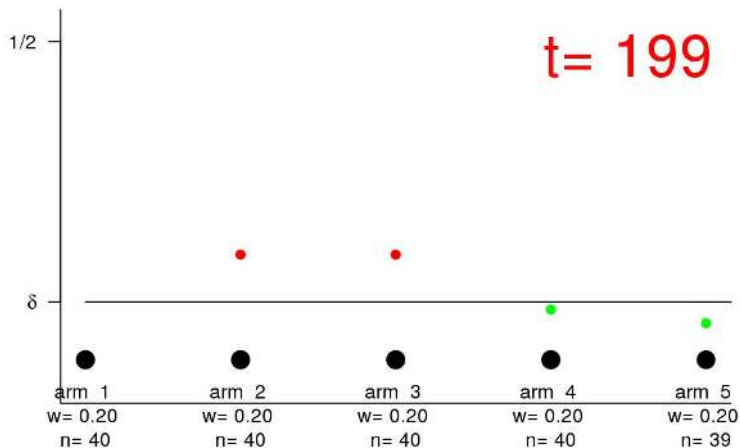
P(confusion)



Uniform Sampling



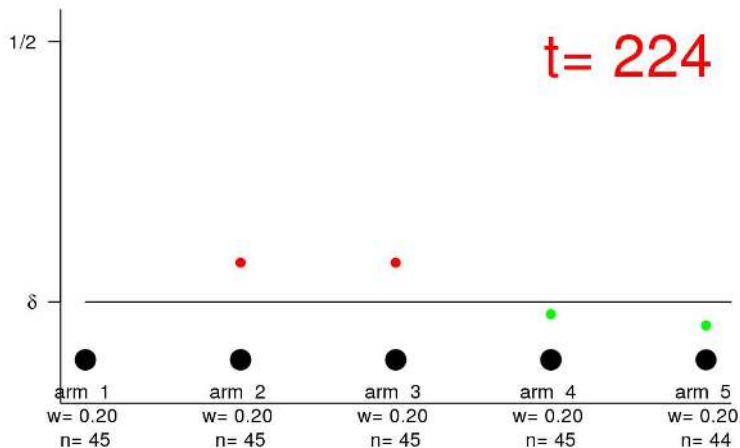
P(confusion)



Uniform Sampling



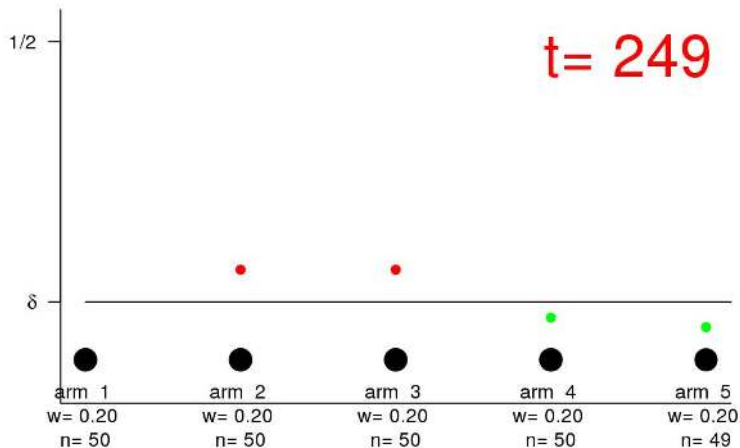
P(confusion)



Uniform Sampling



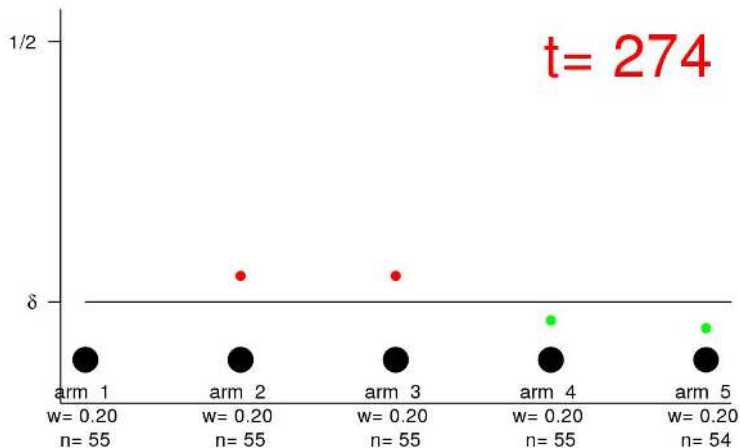
P(confusion)



Uniform Sampling



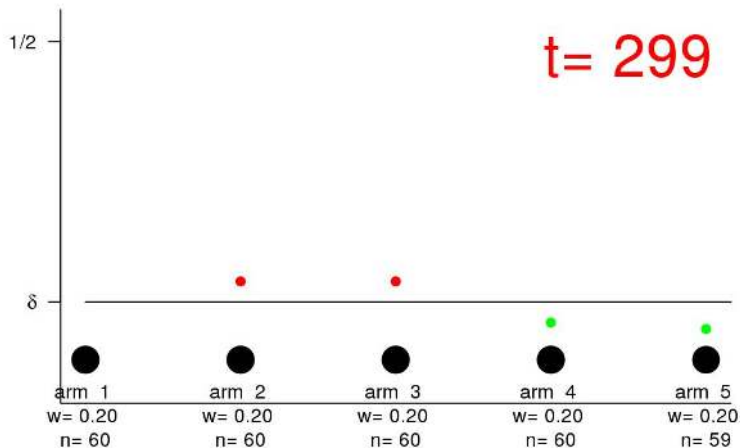
P(confusion)



Uniform Sampling



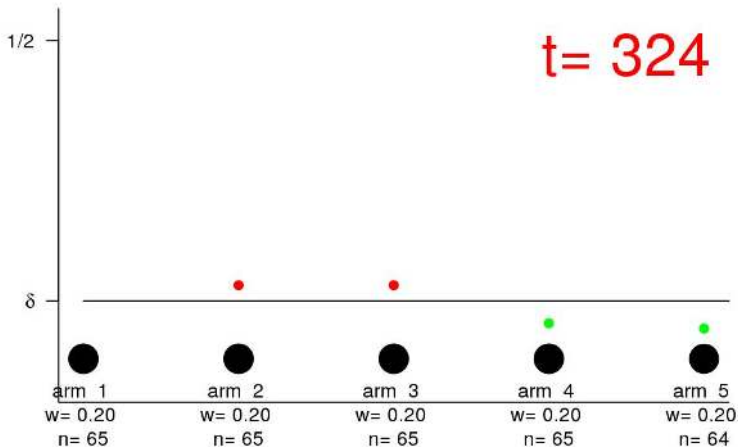
P(confusion)



Uniform Sampling



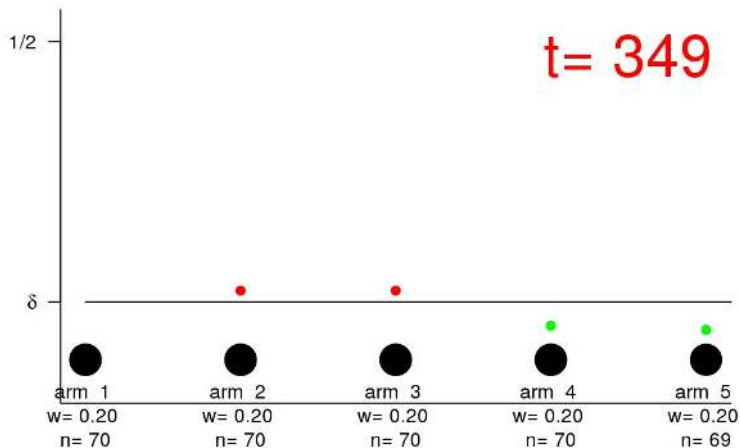
P(confusion)



Uniform Sampling



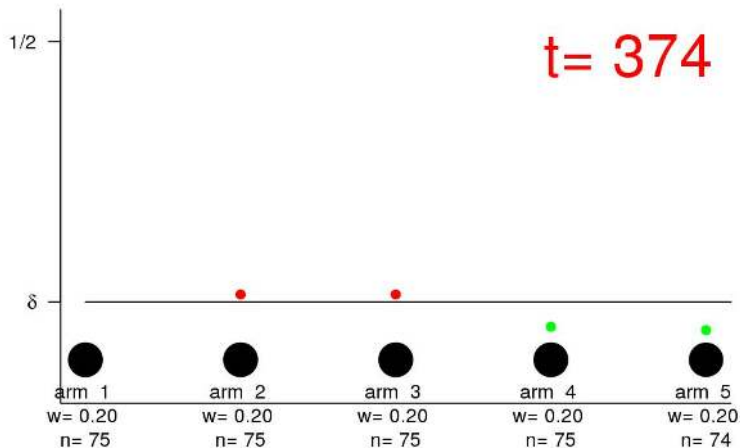
P(confusion)



Uniform Sampling



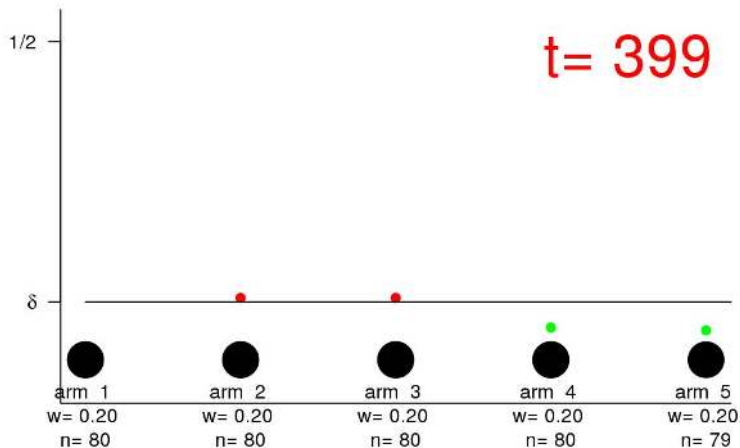
P(confusion)



Uniform Sampling



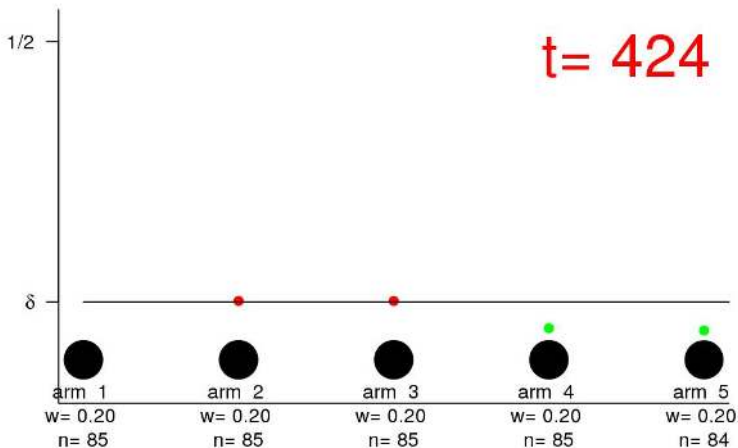
P(confusion)



Uniform Sampling



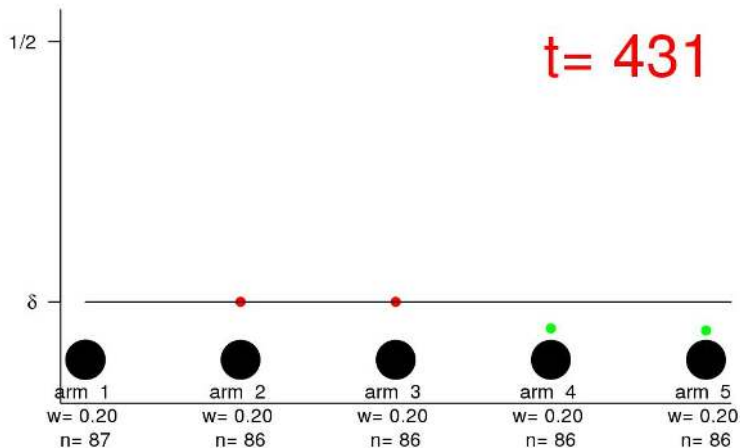
P(confusion)



Uniform Sampling



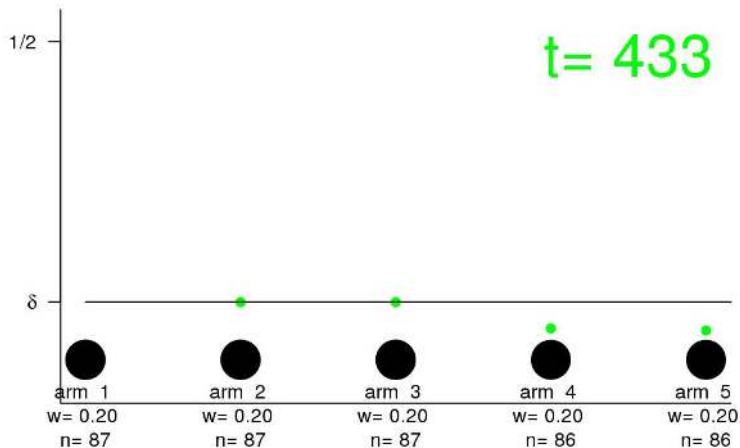
P(confusion)



Uniform Sampling



P(confusion)



Intuition: Equalizing the Probabilities of Confusion

Most simple setting: for all $a \in \{1, \dots, K\}$,

$$\nu_a = \mathcal{N}(\mu_a, 1)$$

For example: $\mu = [2, 1.75, 1.75, 1.6, 1.5]$.

Active Learning

→ You allocate a **relative budget** w_a to option a , with $w_1 + \dots + w_K = 1$.

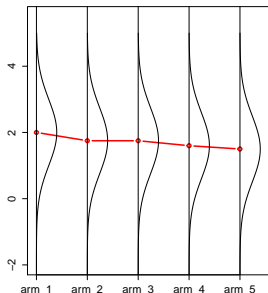
At time t :

→ you have sampled $\mathbf{n}_a \approx \mathbf{w}_a \mathbf{t}$ times the option a

→ your empirical average is \bar{X}_{a, n_a} .

→ if you stop at time t , your **probability of preferring arm $a \geq 2$ to arm $a^* = 1$** is:

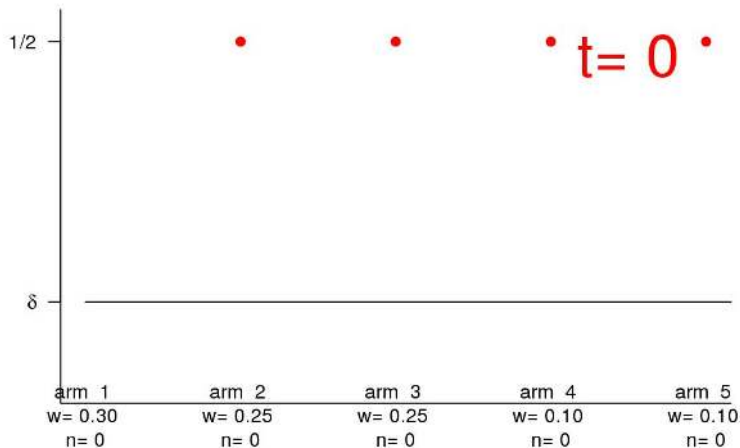
$$\begin{aligned} \mathbb{P}(\bar{X}_{a, n_a} > \bar{X}_{1, n_1}) &= \mathbb{P}\left(\frac{\bar{X}_{a, n_a} - \mu_a - (\bar{X}_{1, n_1} - \mu_1)}{\sqrt{1/n_1 + 1/n_a}} > \frac{\mu_1 - \mu_a}{\sqrt{1/n_1 + 1/n_a}}\right) \\ &= \bar{\Phi}\left(\frac{\mu_1 - \mu_a}{\sqrt{1/n_1 + 1/n_a}}\right) \end{aligned}$$



Improving: trial 1



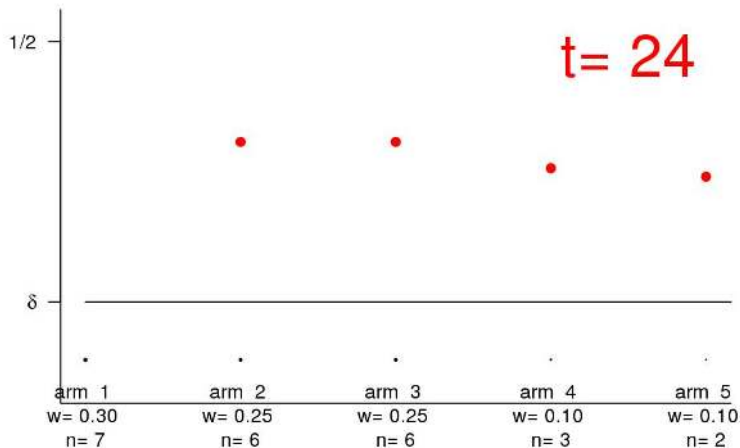
P(confusion)



Improving: trial 1



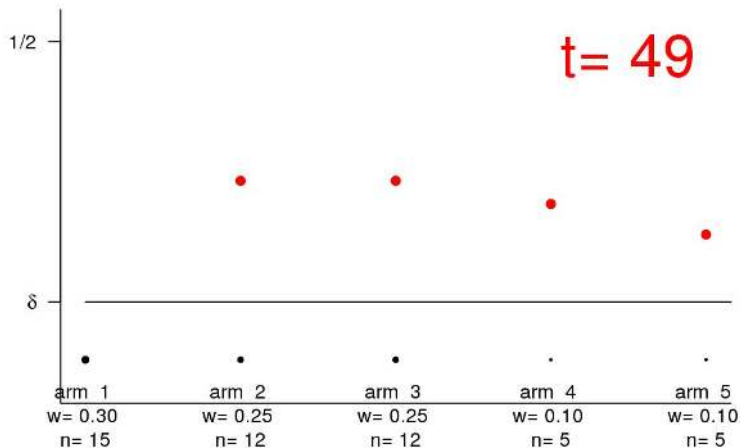
P(confusion)



Improving: trial 1



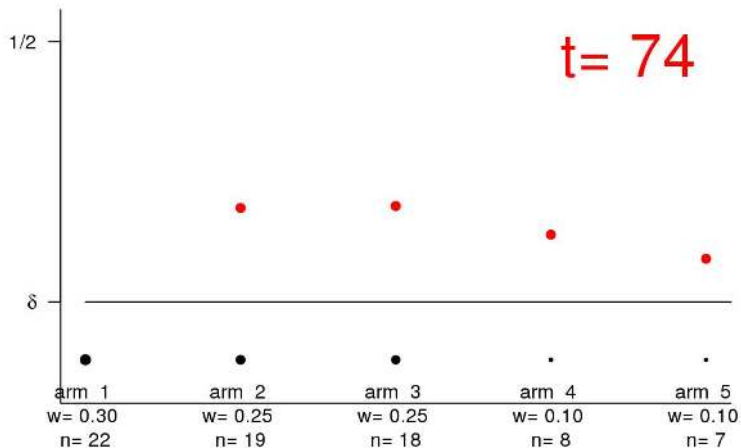
P(confusion)



Improving: trial 1



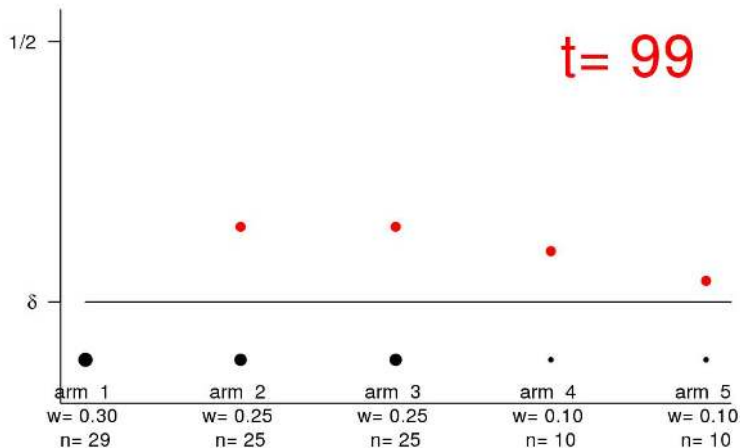
P(confusion)



Improving: trial 1



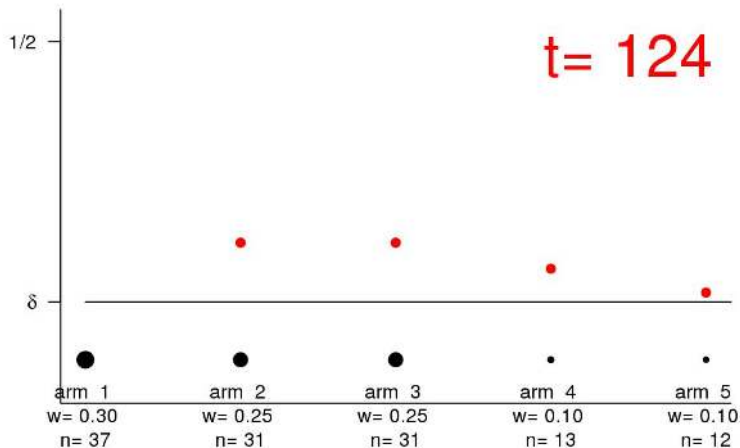
P(confusion)



Improving: trial 1



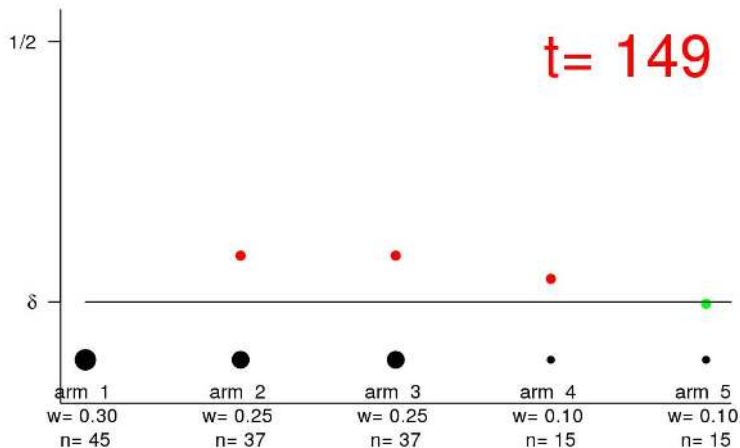
P(confusion)



Improving: trial 1



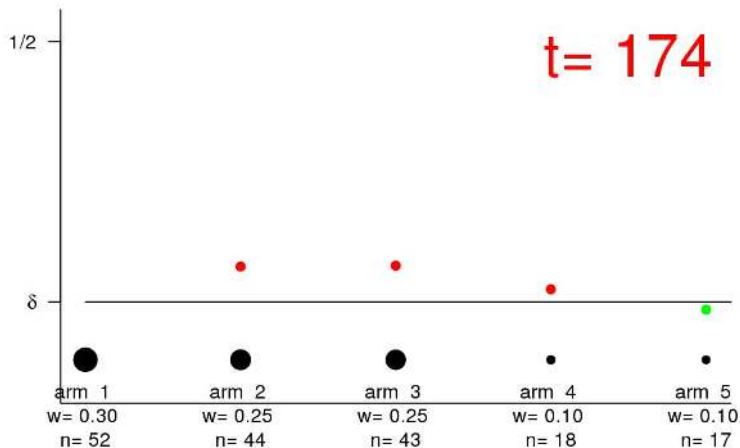
P(confusion)



Improving: trial 1



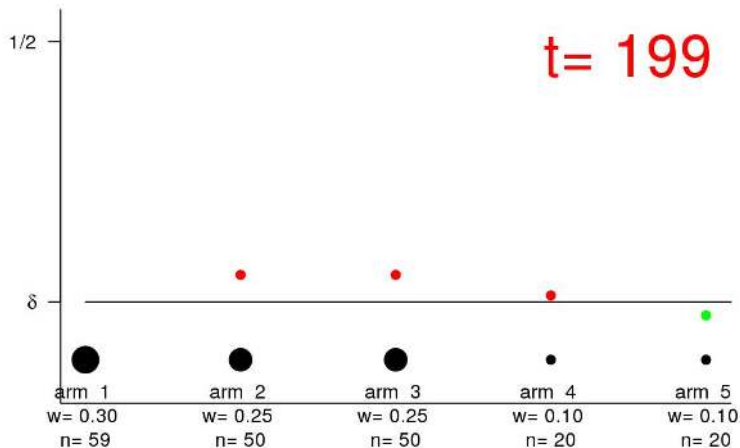
P(confusion)



Improving: trial 1



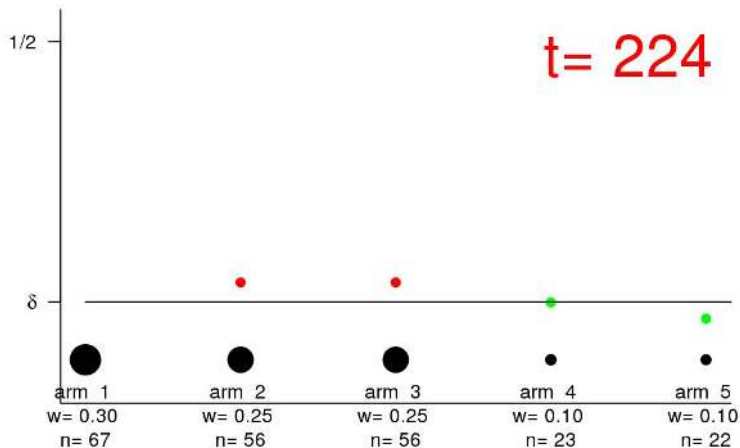
P(confusion)



Improving: trial 1



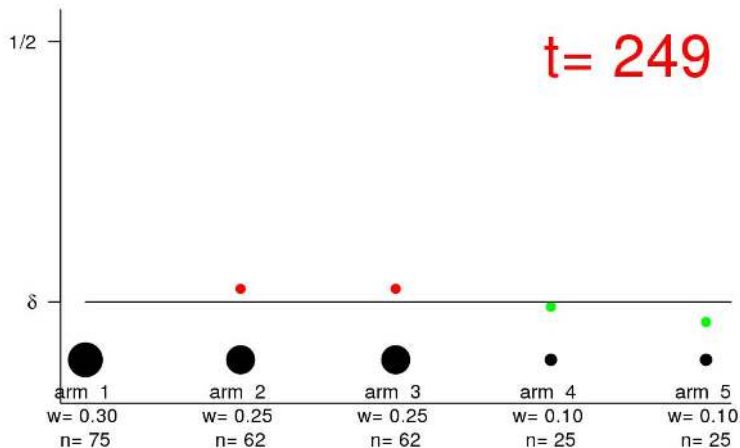
P(confusion)



Improving: trial 1



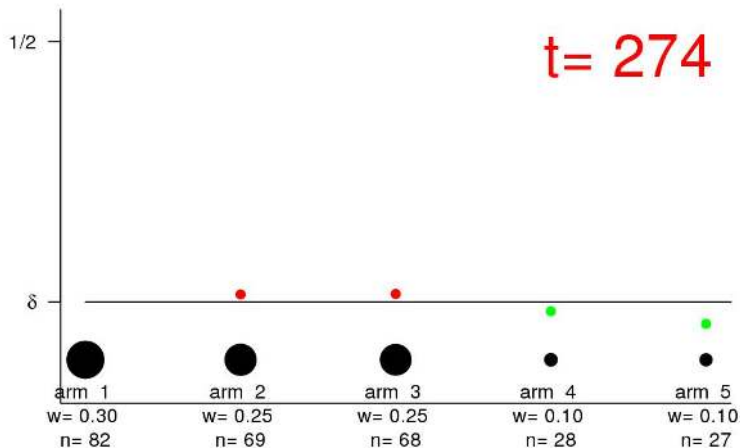
P(confusion)



Improving: trial 1



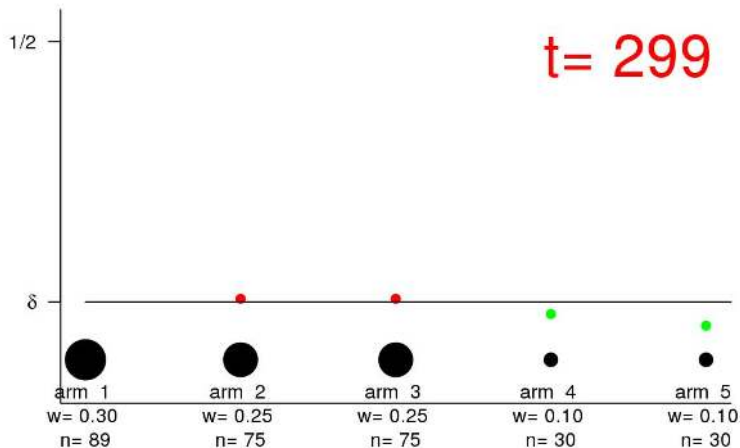
P(confusion)



Improving: trial 1



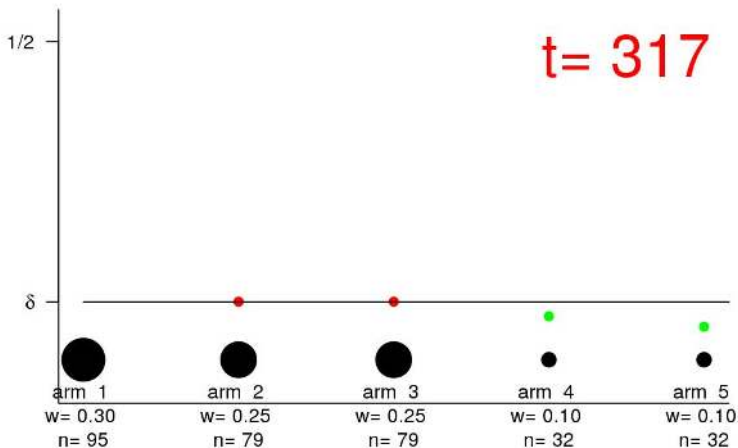
P(confusion)



Improving: trial 1



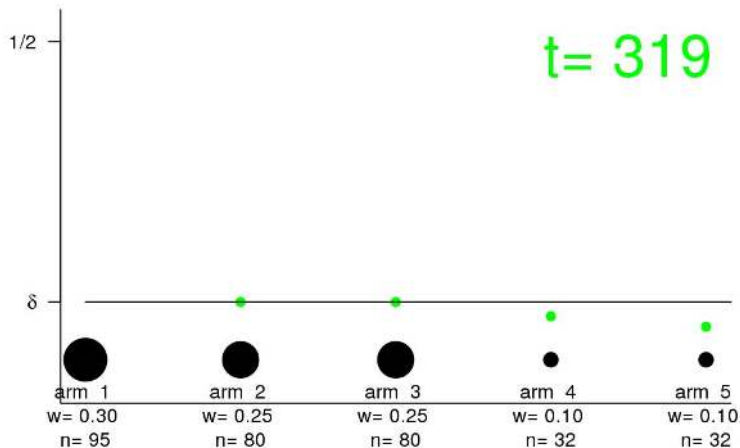
P(confusion)



Improving: trial 1



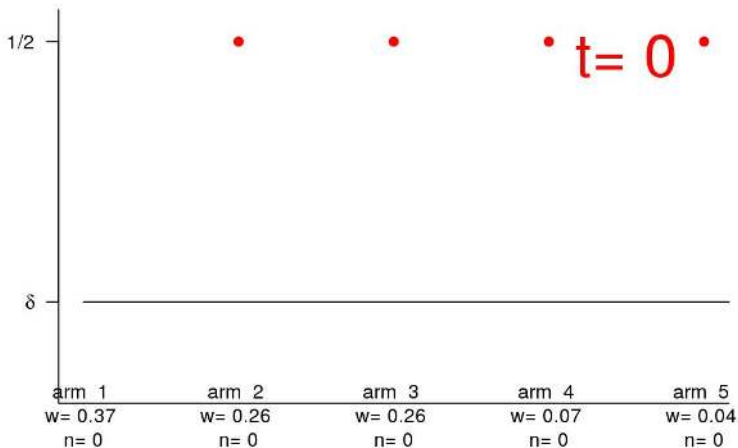
P(confusion)



Optimal Proportions



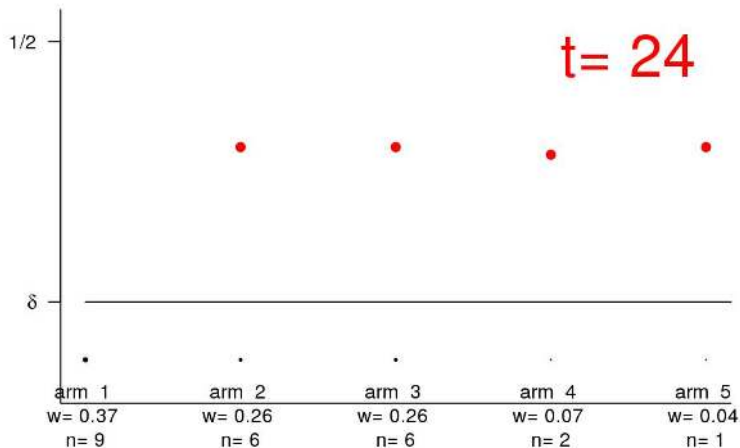
P(confusion)



Optimal Proportions



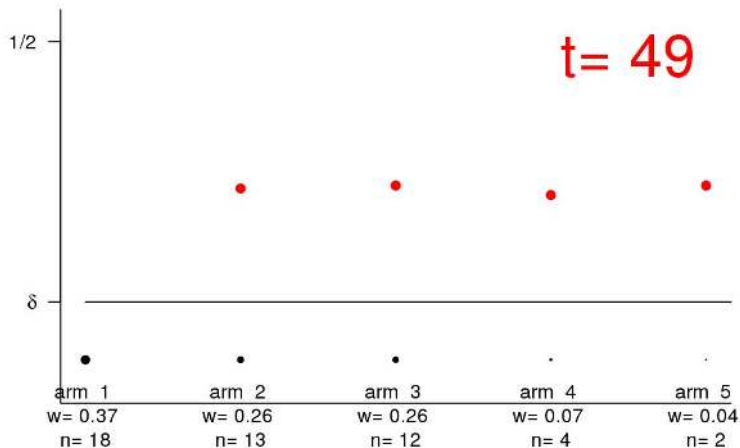
P(confusion)



Optimal Proportions



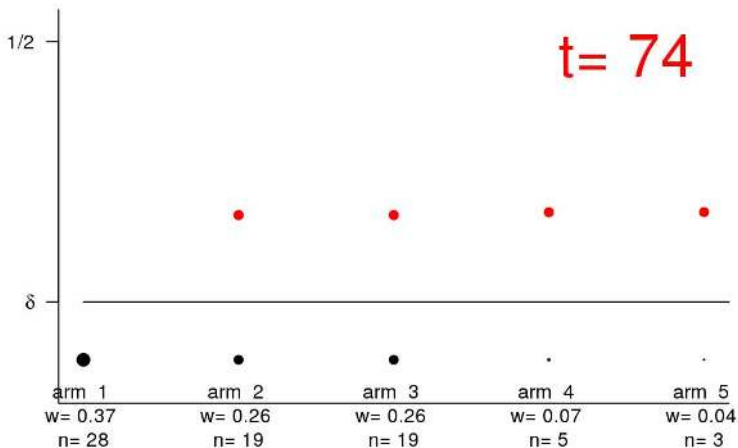
P(confusion)



Optimal Proportions



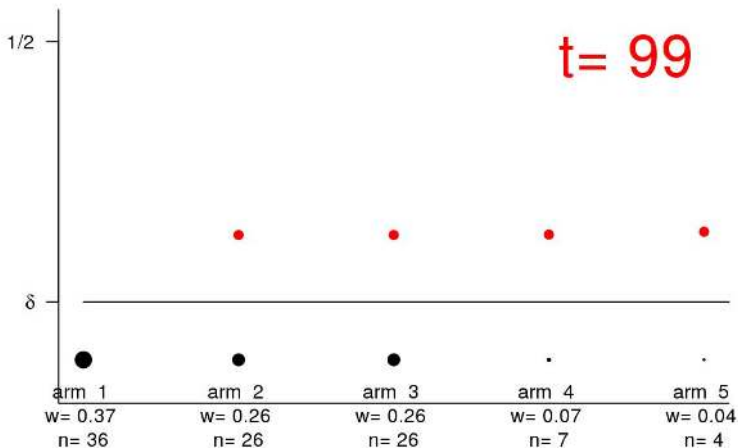
P(confusion)



Optimal Proportions



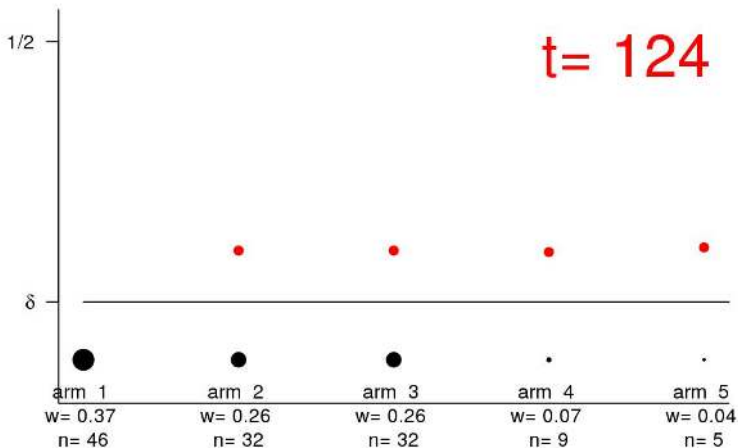
P(confusion)



Optimal Proportions



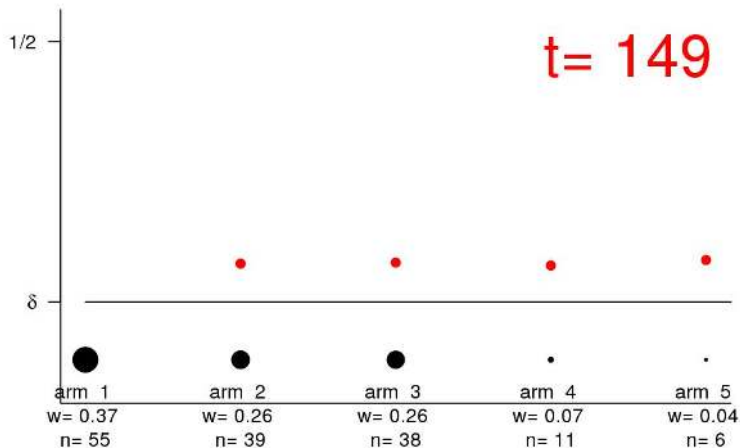
P(confusion)



Optimal Proportions



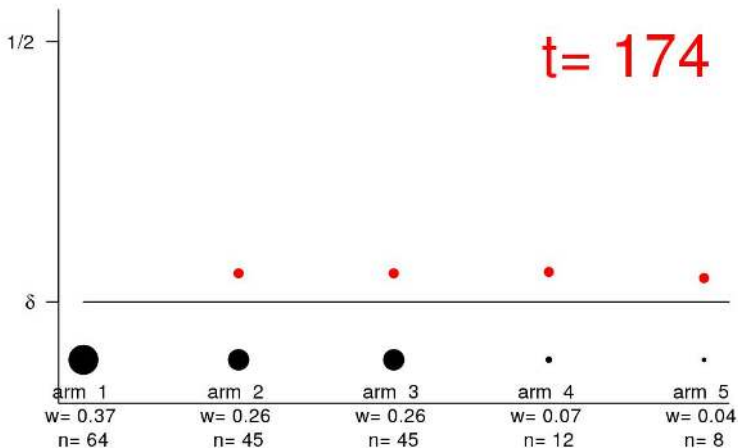
P(confusion)



Optimal Proportions



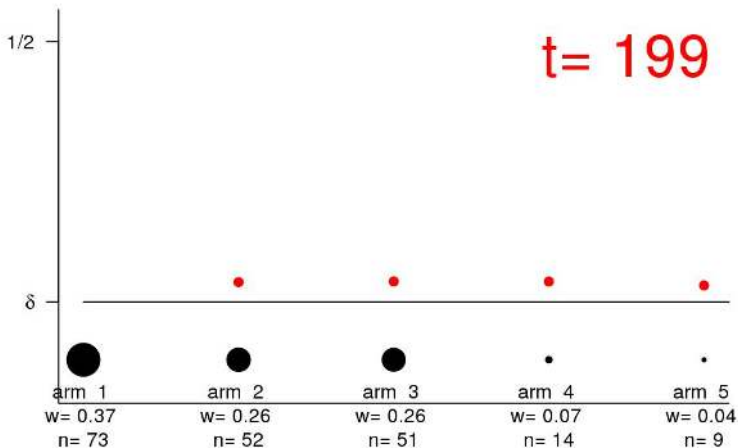
P(confusion)



Optimal Proportions



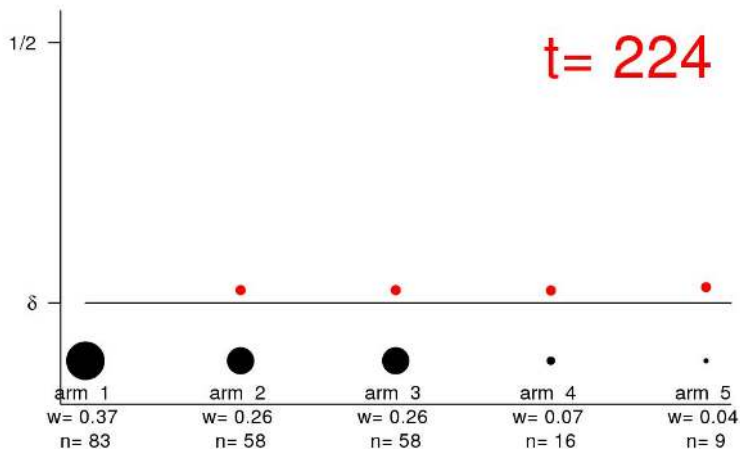
P(confusion)



Optimal Proportions



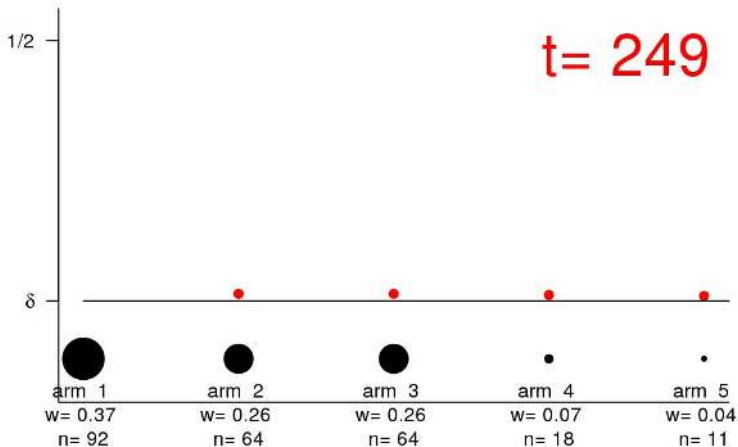
P(confusion)



Optimal Proportions



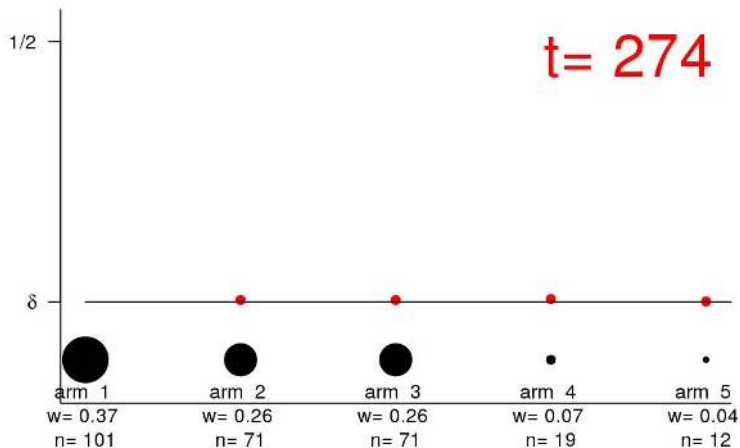
P(confusion)



Optimal Proportions



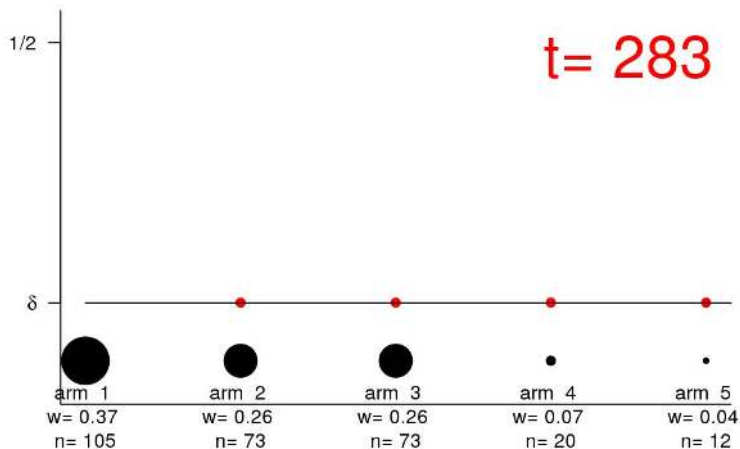
P(confusion)



Optimal Proportions



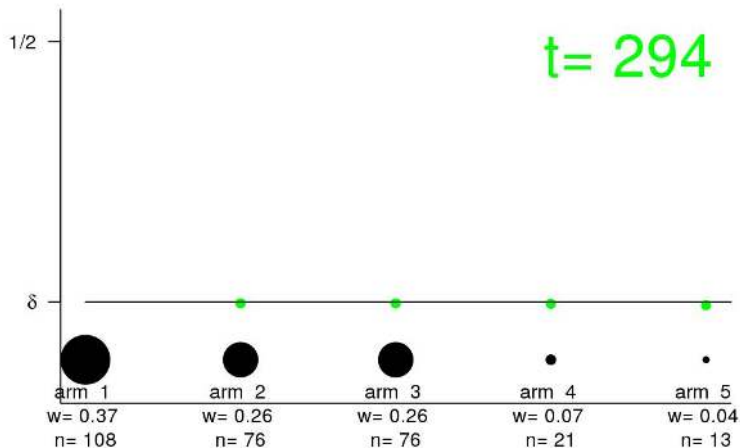
P(confusion)



Optimal Proportions



P(confusion)



How to Turn this Intuition into a Theorem?

- The arms are **not Gaussian** (no formula for probability of confusion)
 - large deviations (Sanov, KL)
- You do not allocate a relative budget at first, but you use **sequential sampling**
 - no fixed-size samples: *sequential experiment*
 - tracking lemma
- How to **compute the optimal proportions**?
 - lower bound, game
- The **parameters** of the distribution are **unknown**
 - (sequential) estimation
- **When** should you **stop**?
 - Chernoff's stopping rule

Exponential Families

ν_1, \dots, ν_K belong to a **one-dimensional exponential family**

$$\mathbb{P}_{\lambda, \Theta, b} = \{ \nu_\theta, \theta \in \Theta : \nu_\theta \text{ has density } f_\theta(x) = \exp(\theta x - b(\theta)) \text{ w.r.t. } \lambda \}$$

Example: Gaussian, Bernoulli, Poisson distributions...

- ν_θ can be parametrized by its mean $\mu = \dot{b}(\theta) : \nu^\mu := \nu_{\dot{b}^{-1}(\mu)}$

Notation: Kullback-Leibler divergence

For a given exponential family,

$$d(\mu, \mu') := \text{KL}(\nu^\mu, \nu^{\mu'}) = \mathbb{E}_{X \sim \nu^\mu} \left[\log \frac{d\nu^\mu}{d\nu^{\mu'}}(X) \right]$$

is the **KL-divergence between the distributions of mean μ and μ'** .

We identify $\nu = (\nu^{\mu_1}, \dots, \nu^{\mu_K})$ and $\mu = (\mu_1, \dots, \mu_K)$ and consider

$$\mathcal{S} = \left\{ \mu \in (\dot{b}(\Theta))^K : \exists a \in \{1, \dots, K\} : \mu_a > \max_{i \neq a} \mu_i \right\}$$

Prolegomenon: Large Deviation Bounds for Bandits

Back to: Equalizing the Probabilities of Confusion

Most simple setting: for all $a \in \{1, \dots, K\}$, $\nu_a = \mathcal{N}(\mu_a, 1)$.

For example: $\mu = [2, 1.75, 1.75, 1.6, 1.5]$.

You allocate a **relative budget** w_a to option a , with $w_1 + \dots + w_K = 1$.

At time t , you have sampled $n_a \approx w_a t$ times option a and the empirical average is \bar{X}_{a, n_a} .

→ if you stop at time t , your probability of preferring arm $a \geq 2$ to arm $a^* = 1$ is:

$$\begin{aligned}\mathbb{P}(\bar{X}_{a, n_a} > \bar{X}_{1, n_1}) &= \mathbb{P}\left(\frac{\bar{X}_{a, n_a} - \mu_a - (\bar{X}_{1, n_1} - \mu_1)}{\sqrt{1/n_1 + 1/n_a}} > \frac{\mu_1 - \mu_a}{\sqrt{1/n_1 + 1/n_a}}\right) \\ &= \bar{\Phi}\left(\frac{\mu_1 - \mu_a}{\sqrt{1/n_1 + 1/n_a}}\right) \\ &\leq e^{-\frac{(\mu_1 - \mu_a)^2}{2(1/n_1 + 1/n_a)}}\end{aligned}$$

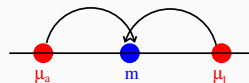
Chernoff's bound

Back to: Equalizing the Probabilities of Confusion

Most simple setting: for all $a \in \{1, \dots, K\}$, $\nu_a = \mathcal{N}(\mu_a, 1)$.

→ if you stop at time t , your probability of preferring arm $a \geq 2$ to arm $a^* = 1$ is:

$$\begin{aligned}\mathbb{P}(\bar{X}_{a,n_a} > \bar{X}_{1,n_1}) &= \mathbb{P}\left(\frac{\bar{X}_{a,n_a} - \mu_a - (\bar{X}_{1,n_1} - \mu_1)}{\sqrt{1/n_1 + 1/n_a}} > \frac{\mu_1 - \mu_a}{\sqrt{1/n_1 + 1/n_a}}\right) \\ &= \bar{\Phi}\left(\frac{\mu_1 - \mu_a}{\sqrt{1/n_1 + 1/n_a}}\right) \\ &\approx e^{-\frac{(\mu_1 - \mu_a)^2}{2(1/n_1 + 1/n_a)}} \quad \text{Large Deviation Principle} \\ &= e^{-\frac{n_1(\mu_1 - m)^2}{2}} \times e^{-\frac{n_a(\mu_a - m)^2}{2}} \quad \text{with } m = \frac{n_1\mu_1 + n_a\mu_a}{n_1 + n_a} \\ &\approx \mathbb{P}(\bar{X}_{1,n_1} < m) \times \mathbb{P}(\bar{X}_{a,n_a} \geq m) \\ &= \max_{\mu_a < m < \mu_1} \mathbb{P}(\bar{X}_{1,n_1} < m) \times \mathbb{P}(\bar{X}_{a,n_a} \geq m) \quad \text{cf. Sanov}\end{aligned}$$

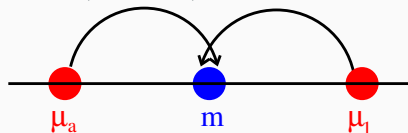


Entropic Method for Large Deviations Lower Bounds

Let $d(\mu, \mu') = \text{KL}(\mathcal{N}(\mu, 1), \mathcal{N}(\mu', 1)) = \frac{(x-y)^2}{2}$,

$\mathcal{KL}(\mathcal{L}(Y), \mathcal{L}(Z)) = \text{KL}(\mathcal{L}(Y), \mathcal{L}(Z))$, $\epsilon > 0$, $\mu_a \leq m \leq \mu_1$ and

- $X_{1,1}, \dots, X_{1,n_1} \stackrel{iid}{\sim} \mathcal{N}(\mu_1, 1)$
- $X'_{1,1}, \dots, X'_{1,n_1} \stackrel{iid}{\sim} \mathcal{N}(m - \epsilon, 1)$
- $X_{a,1}, \dots, X_{a,n_a} \stackrel{iid}{\sim} \mathcal{N}(\mu_a, 1)$
- $X'_{a,1}, \dots, X'_{a,n_a} \stackrel{iid}{\sim} \mathcal{N}(m + \epsilon, 1)$



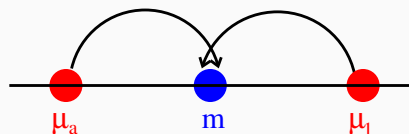
$$n_1 d(m - \epsilon, \mu_1) + n_a d(m + \epsilon, \mu_a) = \mathcal{KL}((X'_{a,i})_{a,i}, (X_{a,i})_{a,i}) = \text{KL}(P \otimes P', Q \otimes Q') = \text{KL}(P, Q) + \text{KL}(P', Q')$$

$$\geq \mathcal{KL}(\mathbb{1}\{\bar{X}'_{a,n_a} > \bar{X}'_{1,n_1}\}, \mathbb{1}\{\bar{X}_{a,n_a} > \bar{X}_{1,n_1}\}) \quad \begin{array}{l} \text{contraction of entropy} \\ = \text{data-processing inequality} \end{array}$$

$$= \text{kl}\left(\mathbb{P}(\bar{X}'_{a,n_a} > \bar{X}'_{1,n_1}), \mathbb{P}(\bar{X}_{a,n_a} > \bar{X}_{1,n_1})\right) \quad \text{kl}(p, q) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$$

$$\geq \mathbb{P}(\bar{X}'_{a,n_a} > \bar{X}'_{1,n_1}) \log \frac{1}{\mathbb{P}(\bar{X}_{a,n_a} > \bar{X}_{1,n_1})} - \log(2) \quad \text{kl}(p, q) \geq p \log \frac{1}{q} - \log 2$$

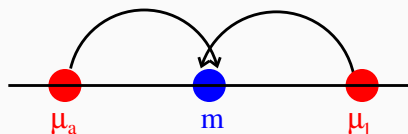
Entropic Method for Large Deviations Lower Bounds



$$\begin{aligned}
 n_1 d(m - \epsilon, \mu_1) + n_a d(m + \epsilon, \mu_a) &= \mathcal{KL}((X'_{a,i})_{a,i}, (X_{a,i})_{a,i}) = \mathcal{KL}(P \otimes P', Q \otimes Q') \\
 &= \mathcal{KL}(P, Q) + \mathcal{KL}(P', Q') \\
 &\geq \mathcal{KL}(\mathbb{1}\{\bar{X}'_{a,n_a} > \bar{X}'_{1,n_1}\}, \mathbb{1}\{\bar{X}_{a,n_a} > \bar{X}_{1,n_1}\}) \quad \begin{array}{l} \text{contraction of entropy} \\ = \text{data-processing inequality} \end{array} \\
 &= \text{kl}(\mathbb{P}(\bar{X}'_{a,n_a} > \bar{X}'_{1,n_1}), \mathbb{P}(\bar{X}_{a,n_a} > \bar{X}_{1,n_1})) \quad \text{kl}(p, q) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q} \\
 &\geq \mathbb{P}(\bar{X}'_{a,n_a} > \bar{X}'_{1,n_1}) \log \frac{1}{\mathbb{P}(\bar{X}_{a,n_a} > \bar{X}_{1,n_1})} - \log(2) \quad \text{kl}(p, q) \geq p \log \frac{1}{q} - \log 2
 \end{aligned}$$

$$\begin{aligned}
 \mathbb{P}(\bar{X}_{a,n_a} > \bar{X}_{1,n_1}) &\geq \max_{\mu_1 \leq m \leq \mu_a} \exp\left(-\frac{n_1 d(m - \epsilon, \mu_1) + n_a d(m + \epsilon, \mu_a) + \log(2)}{1 - e^{-(n_1 + n_a)\epsilon^2/2}}\right) \\
 &= \exp\left(-\frac{\frac{(\mu_1 - \mu_a + \epsilon)^2}{1/n_1 + 1/n_a} + \log(2)}{2(1 - e^{-(n_1 + n_a)\epsilon^2/2})}\right) \quad m = \frac{n_1 \mu_1 + n_a \mu_a}{n_1 + n_a}
 \end{aligned}$$

Entropic Method for Large Deviations Lower Bounds



$$n_1 d(m - \epsilon, \mu_1) + n_a d(m + \epsilon, \mu_a) = \mathcal{KL}((X'_{a,i})_{a,i}, (X_{a,i})_{a,i}) = \mathcal{KL}(P \otimes P', Q \otimes Q')$$

$$\geq \mathcal{KL}(\mathbb{1}\{\bar{X}'_{a,n_a} > \bar{X}'_{1,n_1}\}, \mathbb{1}\{\bar{X}_{a,n_a} > \bar{X}_{1,n_1}\})$$

contraction of entropy
= data-processing inequality

$$= \text{kl}\left(\mathbb{P}(\bar{X}'_{a,n_a} > \bar{X}'_{1,n_1}), \mathbb{P}(\bar{X}_{a,n_a} > \bar{X}_{1,n_1})\right)$$

$$\text{kl}(p, q) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$$

$$\geq \mathbb{P}(\bar{X}'_{a,n_a} > \bar{X}'_{1,n_1}) \log \frac{1}{\mathbb{P}(\bar{X}_{a,n_a} > \bar{X}_{1,n_1})} - \log(2)$$

$$\text{kl}(p, q) \geq p \log \frac{1}{q} - \log 2$$

$$\Rightarrow T(w_1 d(m - \epsilon, \mu_1) + w_a d(m + \epsilon, \mu_a)) \gtrsim \log \frac{1}{\mathbb{P}(\bar{X}_{a,n_a} > \bar{X}_{1,n_1})}$$

→ if you want to have $\mathbb{P}(\bar{X}_{a,n_a} > \bar{X}_{1,n_1}) \leq \delta$ then you need

$$T \gtrsim \frac{\log(1/\delta)}{w_1 d(m, \mu_1) + w_a d(m, \mu_a)}$$

It Works in the Bandit Setting!

Theorem (see Garivier, Ménard and Stoltz, M.O.R. to appear)

For all bandit problems μ and λ , all stopping time τ and $\sigma(\mathcal{F}_\tau)$ -measurable random variables Z with values in $[0, 1]$,

$$\sum_{a=1}^K \mathbb{E}_\mu[N_a(\tau)] d(\mu_a, \lambda_a) \geq \text{kl}(\mathbb{E}_\mu[Z], \mathbb{E}_\lambda[Z]).$$

Proof: if $I_\tau = (A_1, X_{A_1,1}, \dots, A_\tau, X_{A_\tau, N_{A_\tau}(\tau)})$,

$$\sum_{a=1}^K \mathbb{E}_\mu[N_a(\tau)] d(\mu_a, \lambda_a) = \text{KL}(\mathbb{P}_\mu^{I_\tau}, \mathbb{P}_\lambda^{I_\tau}) \geq \text{KL}(\mathbb{P}_\mu^Z, \mathbb{P}_\lambda^Z) \geq \text{kl}(\mathbb{E}_\mu[Z], \mathbb{E}_\lambda[Z])$$

by *tensorization* and *contraction* of entropy (and small lemma).

Lower Bound

Lower-Bounding the Sample Complexity

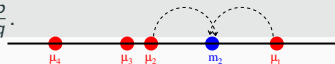
Let $\mu = (\mu_1, \dots, \mu_K)$ and $\lambda = (\lambda_1, \dots, \lambda_K)$ be two elements of \mathcal{S} .

Uniform δ -correct Constraint [Kaufmann, Cappé, G. '15]

If $a^*(\mu) \neq a^*(\lambda)$, any δ -correct algorithm satisfies

$$\sum_{a=1}^K \mathbb{E}_{\mu} [N_a(\tau_{\delta})] d(\mu_a, \lambda_a) \geq \text{kl}(\delta, 1 - \delta)$$

where $\text{kl}(p, q) = p \log \frac{p}{q} + (1 - p) \log \frac{1-p}{1-q}$.



Let $\text{Alt}(\mu) = \{\lambda : a^*(\lambda) \neq a^*(\mu)\}$. Take: $\lambda_1 = m_2 - \epsilon$ $\lambda_2 = m_2 + \epsilon$

$$\mathbb{E}_{\mu} [N_1(\tau_{\delta})] d(\mu_1, m_2 - \epsilon) + \mathbb{E}_{\mu} [N_2(\tau_{\delta})] d(\mu_2, m_2 + \epsilon) \geq \text{kl}(\delta, 1 - \delta)$$

Lower-Bounding the Sample Complexity

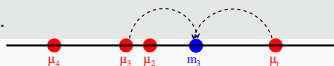
Let $\mu = (\mu_1, \dots, \mu_K)$ and $\lambda = (\lambda_1, \dots, \lambda_K)$ be two elements of \mathcal{S} .

Uniform δ -correct Constraint [Kaufmann, Cappé, G. '15]

If $a^*(\mu) \neq a^*(\lambda)$, any δ -correct algorithm satisfies

$$\sum_{a=1}^K \mathbb{E}_{\mu} [N_a(\tau_{\delta})] d(\mu_a, \lambda_a) \geq \text{kl}(\delta, 1 - \delta)$$

where $\text{kl}(p, q) = p \log \frac{p}{q} + (1 - p) \log \frac{1-p}{1-q}$.



Let $\text{Alt}(\mu) = \{\lambda : a^*(\lambda) \neq a^*(\mu)\}$. Take: $\lambda_1 = m_3 - \epsilon$ $\lambda_3 = m_3 + \epsilon$

$$\mathbb{E}_{\mu} [N_1(\tau_{\delta})] d(\mu_1, m_2 - \epsilon) + \mathbb{E}_{\mu} [N_2(\tau_{\delta})] d(\mu_2, m_2 + \epsilon) \geq \text{kl}(\delta, 1 - \delta)$$

$$\mathbb{E}_{\mu} [N_1(\tau_{\delta})] d(\mu_1, m_3 - \epsilon) + \mathbb{E}_{\mu} [N_3(\tau_{\delta})] d(\mu_3, m_3 + \epsilon) \geq \text{kl}(\delta, 1 - \delta)$$

Lower-Bounding the Sample Complexity

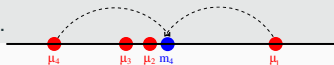
Let $\mu = (\mu_1, \dots, \mu_K)$ and $\lambda = (\lambda_1, \dots, \lambda_K)$ be two elements of \mathcal{S} .

Uniform δ -correct Constraint [Kaufmann, Cappé, G. '15]

If $a^*(\mu) \neq a^*(\lambda)$, any δ -correct algorithm satisfies

$$\sum_{a=1}^K \mathbb{E}_{\mu} [N_a(\tau_{\delta})] d(\mu_a, \lambda_a) \geq \text{kl}(\delta, 1 - \delta)$$

where $\text{kl}(p, q) = p \log \frac{p}{q} + (1 - p) \log \frac{1-p}{1-q}$.



Let $\text{Alt}(\mu) = \{\lambda : a^*(\lambda) \neq a^*(\mu)\}$. Take: $\lambda_1 = m_4 - \epsilon$ $\lambda_4 = m_4 + \epsilon$

$$\mathbb{E}_{\mu} [N_1(\tau_{\delta})] d(\mu_1, m_2 - \epsilon) + \mathbb{E}_{\mu} [N_2(\tau_{\delta})] d(\mu_2, m_2 + \epsilon) \geq \text{kl}(\delta, 1 - \delta)$$

$$\mathbb{E}_{\mu} [N_1(\tau_{\delta})] d(\mu_1, m_3 - \epsilon) + \mathbb{E}_{\mu} [N_3(\tau_{\delta})] d(\mu_3, m_3 + \epsilon) \geq \text{kl}(\delta, 1 - \delta)$$

$$\mathbb{E}_{\mu} [N_1(\tau_{\delta})] d(\mu_1, m_4 - \epsilon) + \mathbb{E}_{\mu} [N_4(\tau_{\delta})] d(\mu_4, m_4 + \epsilon) \geq \text{kl}(\delta, 1 - \delta)$$

Lower-Bounding the Sample Complexity

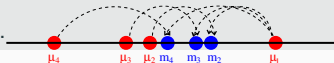
Let $\mu = (\mu_1, \dots, \mu_K)$ and $\lambda = (\lambda_1, \dots, \lambda_K)$ be two elements of \mathcal{S} .

Uniform δ -correct Constraint [Kaufmann, Cappé, G. '15]

If $a^*(\mu) \neq a^*(\lambda)$, any δ -correct algorithm satisfies

$$\sum_{a=1}^K \mathbb{E}_{\mu} [N_a(\tau_{\delta})] d(\mu_a, \lambda_a) \geq \text{kl}(\delta, 1 - \delta)$$

where $\text{kl}(p, q) = p \log \frac{p}{q} + (1 - p) \log \frac{1-p}{1-q}$.



Let $\text{Alt}(\mu) = \{\lambda : a^*(\lambda) \neq a^*(\mu)\}$.

$$\inf_{\lambda \in \text{Alt}(\mu)} \sum_{a=1}^K \mathbb{E}_{\mu} [N_a(\tau_{\delta})] d(\mu_a, \lambda_a) \geq \text{kl}(\delta, 1 - \delta)$$

$$\mathbb{E}_{\mu} [\tau_{\delta}] \times \inf_{\lambda \in \text{Alt}(\mu)} \sum_{a=1}^K \frac{\mathbb{E}_{\mu} [N_a(\tau_{\delta})]}{\mathbb{E}_{\mu} [\tau_{\delta}]} d(\mu_a, \lambda_a) \geq \text{kl}(\delta, 1 - \delta)$$

$$\mathbb{E}_{\mu} [\tau_{\delta}] \times \left(\sup_{w \in \Sigma_K} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{a=1}^K w_a d(\mu_a, \lambda_a) \right) \geq \text{kl}(\delta, 1 - \delta)$$

Lower Bound: the Complexity of BAI

Theorem [G. and Kaufmann 2016]

For any δ -correct algorithm,

$$\mathbb{E}_{\mu}[\tau_{\delta}] \geq T^*(\mu) \text{kl}(\delta, 1 - \delta),$$

where

$$T^*(\mu)^{-1} = \sup_{\mathbf{w} \in \Sigma_K} \inf_{\lambda \in \text{Alt}(\mu)} \left(\sum_{a=1}^K w_a d(\mu_a, \lambda_a) \right).$$

- $\text{kl}(\delta, 1 - \delta) \sim \log(1/\delta)$ when $\delta \rightarrow 0$, $\text{kl}(\delta, 1 - \delta) \geq \log(1/(2.4\delta))$
 - cf. [Graves and Lai 1997, Vaidhyan and Sundaresan, 2015]
- the **optimal proportions of arm draws** are

$$\mathbf{w}^*(\mu) = \operatorname{argmax}_{\mathbf{w} \in \Sigma_K} \inf_{\lambda \in \text{Alt}(\mu)} \left(\sum_{a=1}^K w_a d(\mu_a, \lambda_a) \right)$$

→ they **do not depend on δ**

Given a parameter $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K)$:

- the statistician chooses proportions of arm draws $\mathbf{w} = (w_a)_a$
- the opponent chooses an alternative model $\boldsymbol{\lambda}$
- the payoff is the minimal number $T = T(\mathbf{w}, \boldsymbol{\lambda})$ of draws necessary to ensure that he does not violate the δ -PAC constraint

$$\sum_{a=1}^K T w_a d(\mu_a, \lambda_a) \geq \text{kl}(\delta, 1 - \delta)$$

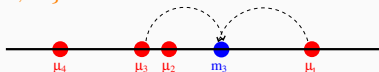
- $T^*(\boldsymbol{\mu}) \text{kl}(\delta, 1 - \delta) = \text{value of the game}$
 $\mathbf{w}^* = \text{optimal action for the statistician}$

PAC-BAI as a Game

Given a parameter $\mu = (\mu_1, \dots, \mu_K)$ such that $\mu_1 > \mu_2 \geq \dots \geq \mu_K$:

- the statistician chooses proportions of arm draws $\mathbf{w} = (w_a)_a$
- the opponent chooses an arm $a \in \{2, \dots, K\}$ and

$$\lambda_a = \arg \min_{\lambda} w_1 d(\mu_1, \lambda) + w_a d(\mu_a, \lambda)$$



- the payoff is the minimal number $T = T(\mathbf{w}, a, \delta)$ of draws necessary to ensure that

$$T w_1 d(\mu_1, \lambda_a - \epsilon) + T w_a d(\mu_a, \lambda_a + \epsilon) \geq \text{kl}(\delta, 1 - \delta)$$

$$\text{that is } T(\mathbf{w}, a, \delta) = \frac{\text{kl}(\delta, 1 - \delta)}{w_1 d(\mu_1, \lambda_a - \epsilon) + w_a d(\mu_a, \lambda_a + \epsilon)}$$

- $T^*(\mu) \text{kl}(\delta, 1 - \delta) = \text{value of the game}$
 $\mathbf{w}^* = \text{optimal action for the statistician}$

Properties of $T^*(\mu)$ and $w^*(\mu)$

1. **Unique** solution, solution of **scalar equations** only
2. For all $\mu \in \mathcal{S}$, for all a , $w_a^*(\mu) > 0$
3. w^* is **continuous** in every $\mu \in \mathcal{S}$
4. If $\mu_1 > \mu_2 \geq \dots \geq \mu_K$, one has $w_2^*(\mu) \geq \dots \geq w_K^*(\mu)$
(one may have $w_1^*(\mu) < w_2^*(\mu)$)
5. Case of **two arms** [Kaufmann, Cappé, G. '14]:

$$\mathbb{E}_\mu[\tau_\delta] \geq \frac{\text{kl}(\delta, 1 - \delta)}{d_*(\mu_1, \mu_2)} .$$

where d_* is the 'reversed' Chernoff information

$$d_*(\mu_1, \mu_2) := d(\mu_1, \mu_*) = d(\mu_2, \mu_*) .$$

6. **Gaussian arms** : algebraic equation but no simple formula for $K \geq 3$.

$$\sum_{a=1}^K \frac{2\sigma^2}{\Delta_a^2} \leq T^*(\mu) \leq 2 \sum_{a=1}^K \frac{2\sigma^2}{\Delta_a^2} .$$

The Track-and-Stop Strategy

The Problem

Prolegomenon: Large Deviation Bounds for Bandits

Lower Bound

The Track-and-Stop Strategy

Sampling Rule

Stopping Rule

Optimality

Sampling rule: Tracking the optimal proportions

$\hat{\mu}(t) = (\hat{\mu}_1(t), \dots, \hat{\mu}_K(t))$: vector of empirical means

Introducing

$$U_t = \{a : N_a(t) < \sqrt{t}\},$$

the arm sampled at round $t + 1$ is

$$A_{t+1} \in \begin{cases} \operatorname{argmin}_{a \in U_t} N_a(t) & \text{if } U_t \neq \emptyset & (\text{forced exploration}) \\ \operatorname{argmax}_{1 \leq a \leq K} t w_a^*(\hat{\mu}(t)) - N_a(t) & (\text{tracking}) \end{cases}$$

Lemma

Under the Tracking sampling rule,

$$\mathbb{P}_{\mu} \left(\lim_{t \rightarrow \infty} \frac{N_a(t)}{t} = w_a^*(\mu) \right) = 1.$$

The Problem

Prolegomenon: Large Deviation Bounds for Bandits

Lower Bound

The Track-and-Stop Strategy

Sampling Rule

Stopping Rule

Optimality

Sequential Generalized Likelihood Test

High values of the Generalized Likelihood Ratio

$$Z_{a,b}(t) := \log \frac{\max_{\{\lambda: \lambda_a \geq \lambda_b\}} dP_\lambda(X_1, \dots, X_t)}{\max_{\{\lambda: \lambda_a \leq \lambda_b\}} dP_\lambda(X_1, \dots, X_t)}$$
$$= N_a(t) d(\hat{\mu}_a(t), \hat{\mu}_{a,b}(t)) + N_b(t) d(\hat{\mu}_b(t), \hat{\mu}_{a,b}(t)) \quad \begin{array}{l} \text{if } \hat{\mu}_a(t) > \hat{\mu}_b(t) \\ -Z_{b,a}(t) \text{ otherwise} \end{array}$$

reject the hypothesis that $(\mu_a \leq \mu_b)$.

We stop when **one arm is assessed to be significantly larger than all other arms**, according to a GLR test:

$$\tau_\delta = \inf \left\{ t \in \mathbb{N} : \exists a \in \{1, \dots, K\}, \forall b \neq a, Z_{a,b}(t) > \beta(t, \delta) \right\}$$
$$= \inf \left\{ t \in \mathbb{N} : Z(t) := \max_{a \in \{1, \dots, K\}} \min_{b \neq a} Z_{a,b}(t) > \beta(t, \delta) \right\}$$

Chernoff stopping rule [Chernoff '59]

Two other possible interpretations of the stopping rule:

→ MDL:

$$Z_{a,b}(t) = (N_a(t) + N_b(t)) H(\hat{\mu}_{a,b}(t)) - \left[N_a(t) H(\hat{\mu}_a(t)) + N_b(t) H(\hat{\mu}_b(t)) \right]$$

Sequential Generalized Likelihood Test

High values of the Generalized Likelihood Ratio

$$Z_{a,b}(t) := \log \frac{\max_{\{\lambda: \lambda_a \geq \lambda_b\}} dP_{\lambda}(X_1, \dots, X_t)}{\max_{\{\lambda: \lambda_a \leq \lambda_b\}} dP_{\lambda}(X_1, \dots, X_t)}$$

reject the hypothesis that $(\mu_a \leq \mu_b)$.

We stop when **one arm is assessed to be significantly larger than all other arms**, according to a GLR test:

$$\tau_{\delta} = \inf \left\{ t \in \mathbb{N} : Z(t) := \max_{a \in \{1, \dots, K\}} \min_{b \neq a} Z_{a,b}(t) > \beta(t, \delta) \right\}$$

Chernoff stopping rule [Chernoff '59]

Two other possible interpretations of the stopping rule:

→ **plug-in complexity estimate**: if $F(w, \mu) := \inf_{\lambda \in \text{Alt}(\mu)} \sum_{a=1}^K w_a d(\mu_a, \lambda_a)$,

stop when $Z(t) = t F\left(\frac{N_a(t)}{t}, \hat{\mu}(t)\right) \geq \beta(t, \delta)$ instead of the lower bound

$$\frac{t}{T^*(\mu)} = t F(w^*, \mu) \geq \text{kl}(\delta, 1 - \delta).$$

Theorem

The Chernoff rule is δ -PAC for $\beta(t, \delta) = \log \left(\frac{2(K-1)t}{\delta} \right)$

Lemma

If $\mu_a < \mu_b$, whatever the sampling rule,

$$\mathbb{P}_{\mu} \left(\exists t \in \mathbb{N} : Z_{a,b}(t) > \log \left(\frac{2t}{\delta} \right) \right) \leq \delta$$

The proof uses:

- Barron's lemma (change of distribution)
- and Krichevsky-Trofimov's universal distribution
(very information-theoretic ideas)

The Problem

Prolegomenon: Large Deviation Bounds for Bandits

Lower Bound

The Track-and-Stop Strategy

Sampling Rule

Stopping Rule

Optimality

Theorem

The Track-and-Stop strategy, that uses

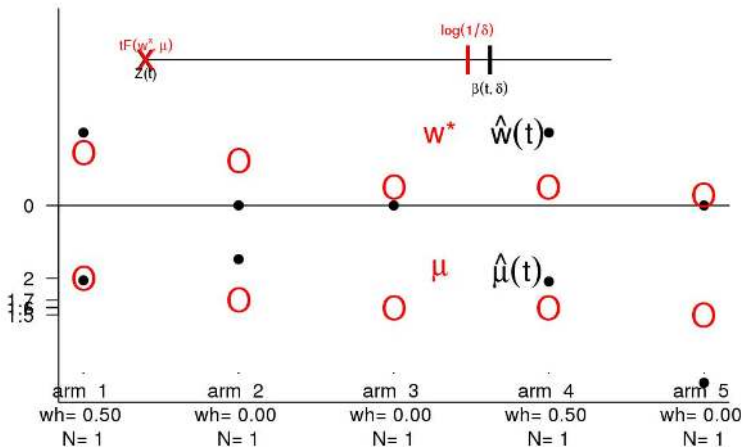
- the Tracking sampling rule
- the Chernoff stopping rule with $\beta(t, \delta) = \log\left(\frac{2(K-1)t}{\delta}\right)$
- and recommends $\hat{a}_{\tau_\delta} = \operatorname{argmax}_{a=1\dots K} \hat{\mu}_a(\tau_\delta)$

is δ -PAC for every $\delta \in (0, 1)$ and satisfies

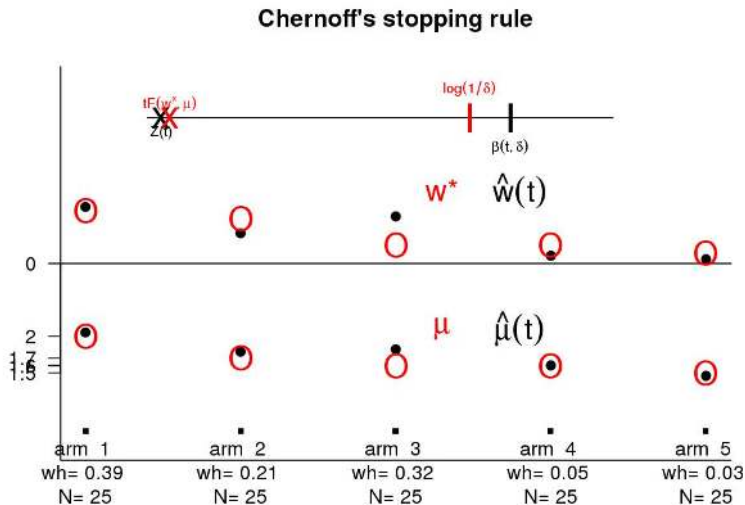
$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_\mu[\tau_\delta]}{\log(1/\delta)} = T^*(\mu).$$

Why is the T&S Strategy asymptotically Optimal?

Chernoff's stopping rule

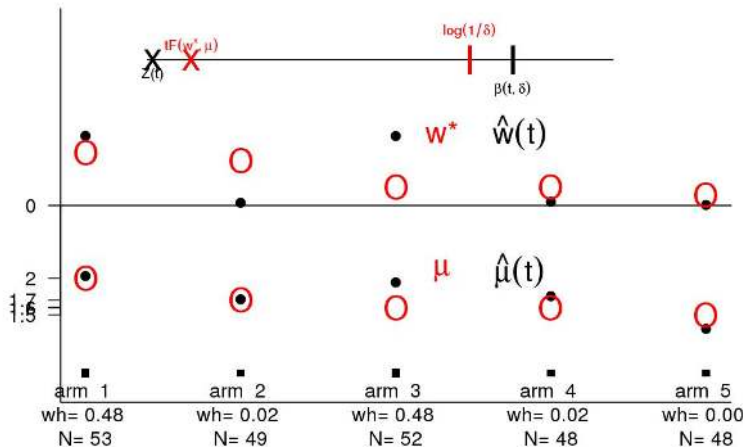


Why is the T&S Strategy asymptotically Optimal?



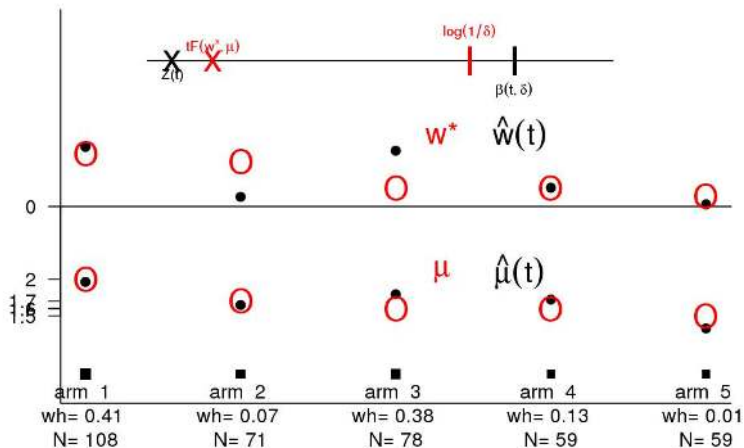
Why is the T&S Strategy asymptotically Optimal?

Chernoff's stopping rule

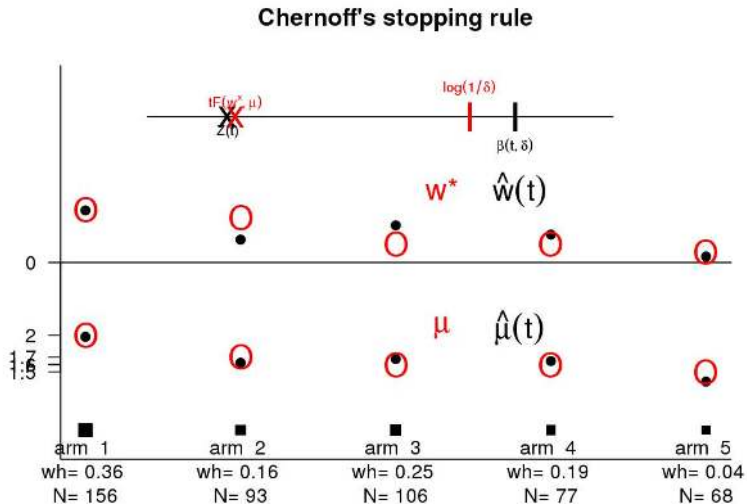


Why is the T&S Strategy asymptotically Optimal?

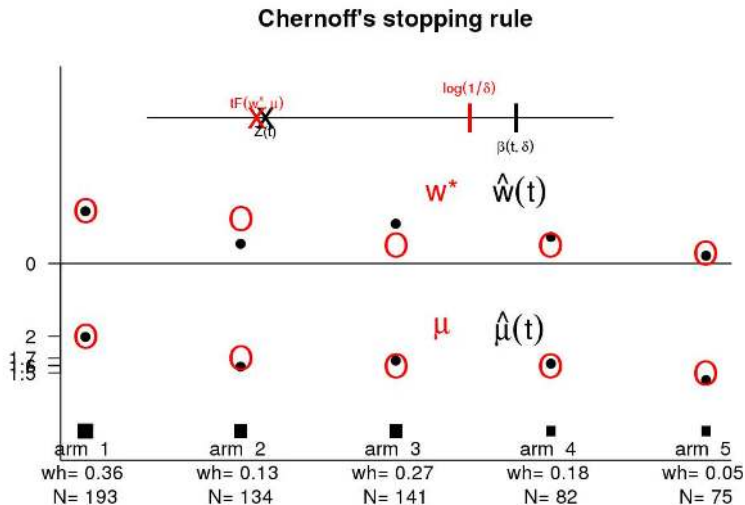
Chernoff's stopping rule



Why is the T&S Strategy asymptotically Optimal?

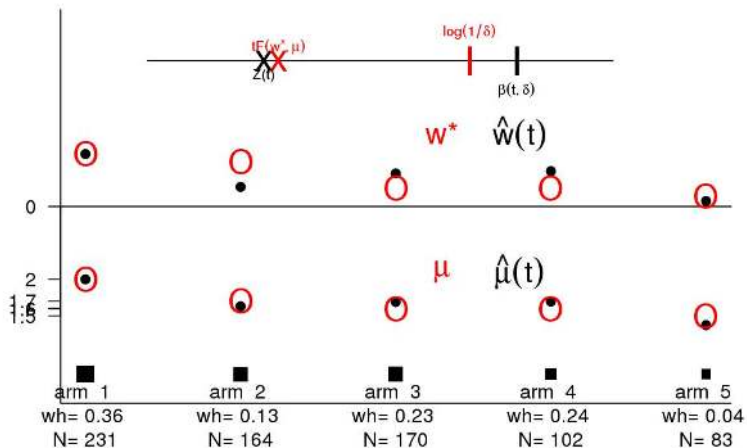


Why is the T&S Strategy asymptotically Optimal?

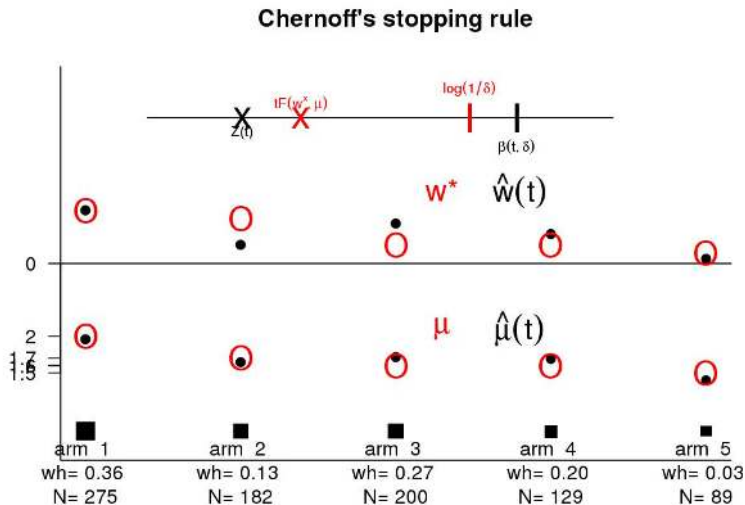


Why is the T&S Strategy asymptotically Optimal?

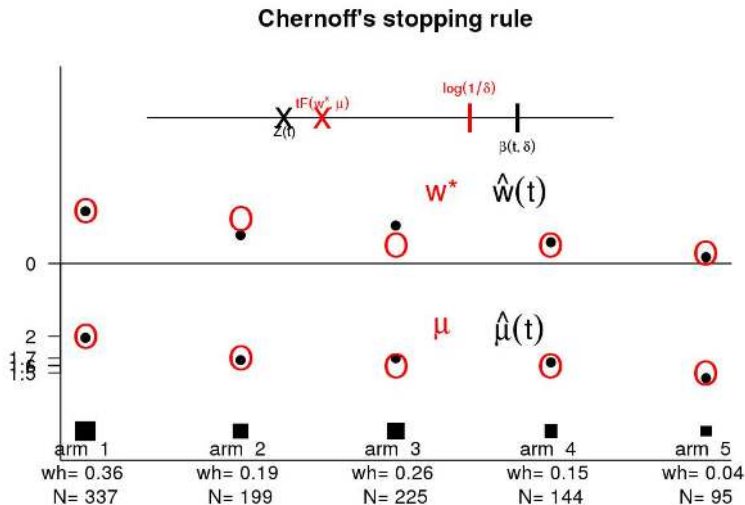
Chernoff's stopping rule



Why is the T&S Strategy asymptotically Optimal?

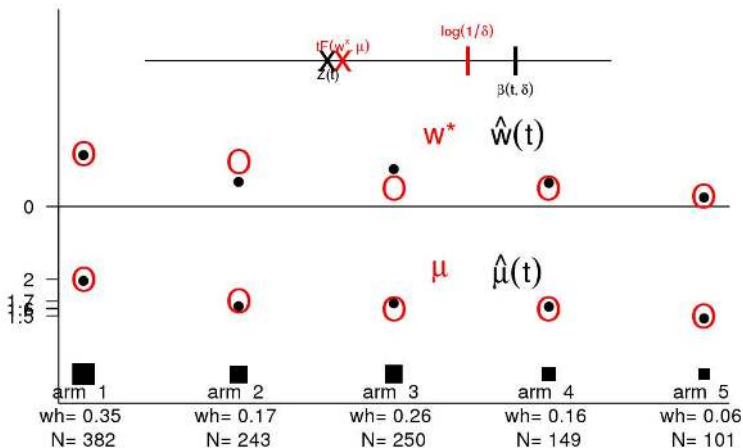


Why is the T&S Strategy asymptotically Optimal?



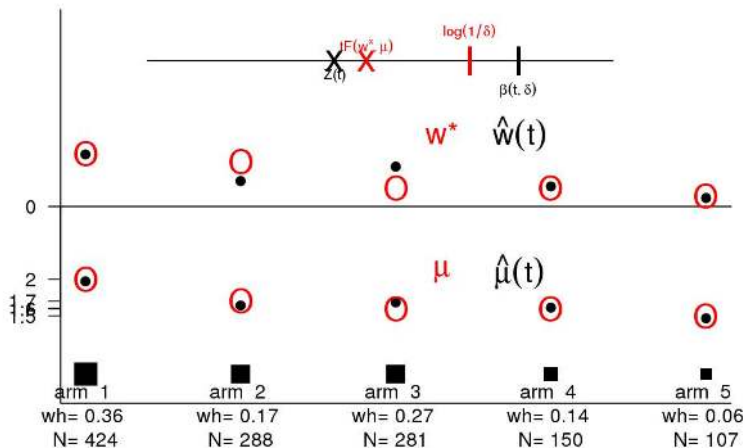
Why is the T&S Strategy asymptotically Optimal?

Chernoff's stopping rule



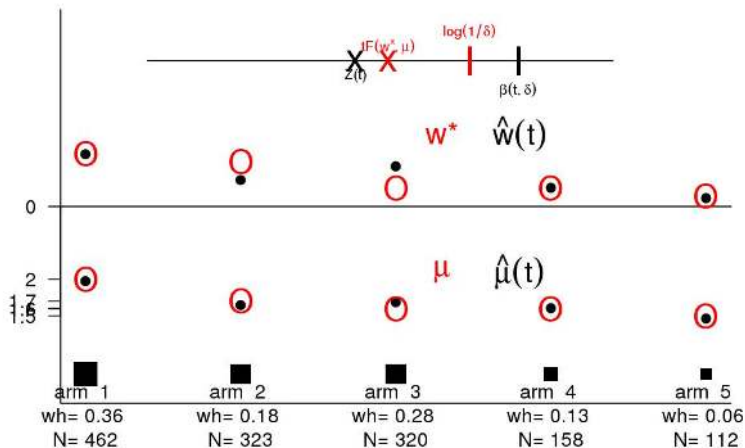
Why is the T&S Strategy asymptotically Optimal?

Chernoff's stopping rule



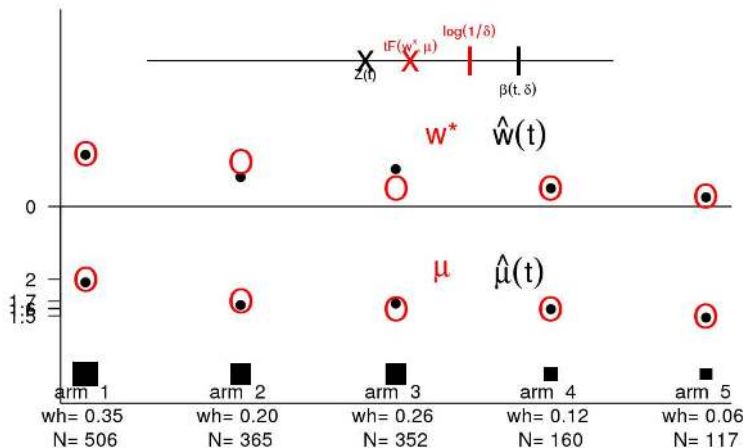
Why is the T&S Strategy asymptotically Optimal?

Chernoff's stopping rule



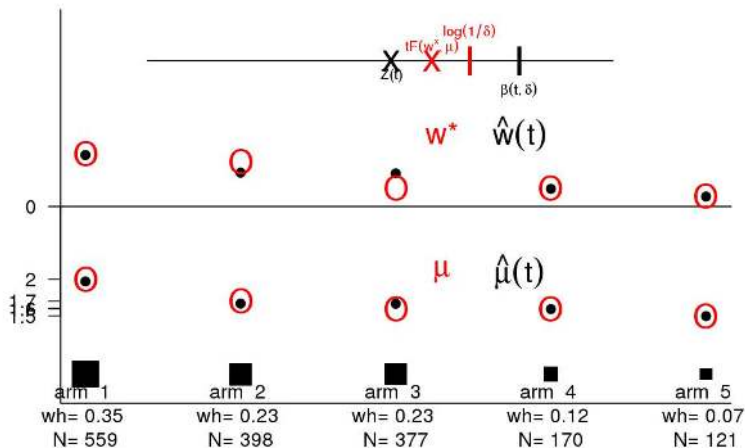
Why is the T&S Strategy asymptotically Optimal?

Chernoff's stopping rule



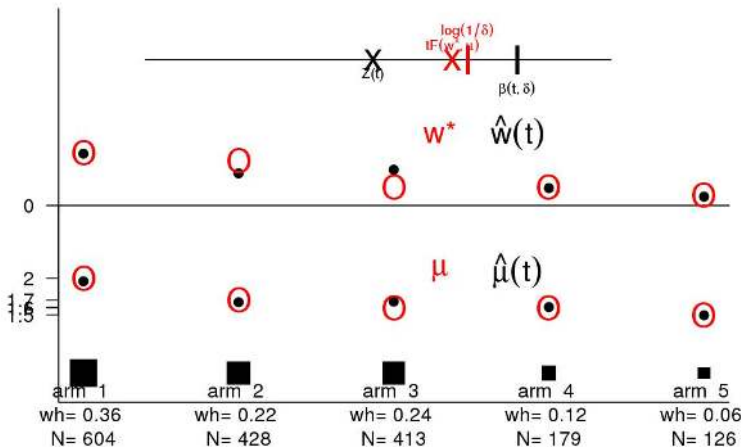
Why is the T&S Strategy asymptotically Optimal?

Chernoff's stopping rule



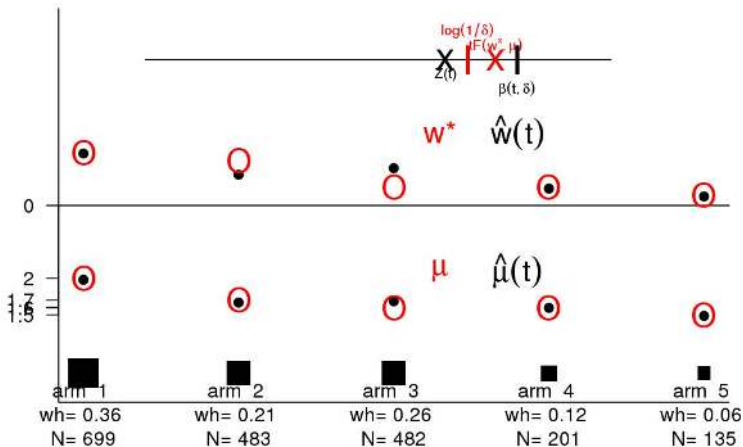
Why is the T&S Strategy asymptotically Optimal?

Chernoff's stopping rule

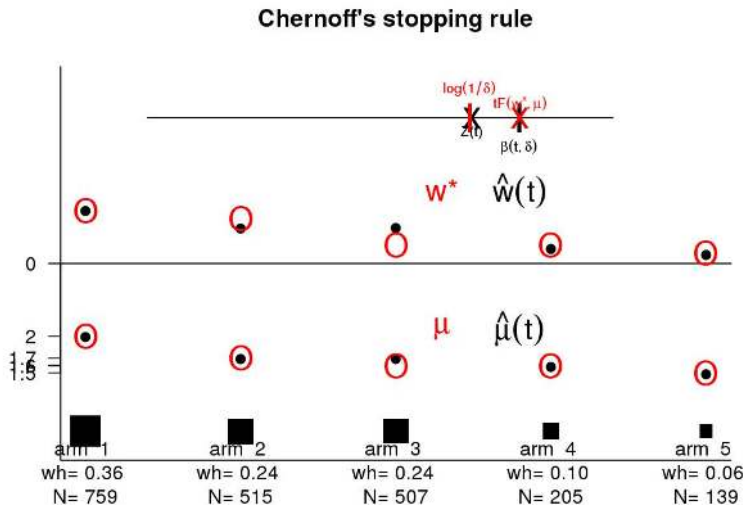


Why is the T&S Strategy asymptotically Optimal?

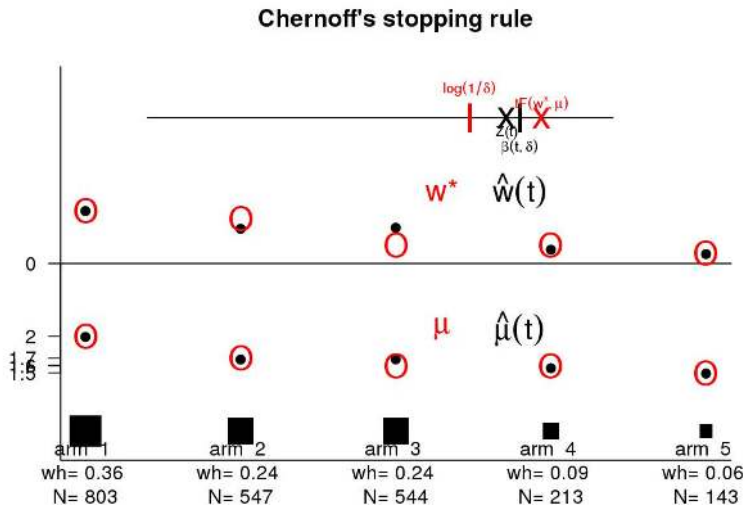
Chernoff's stopping rule



Why is the T&S Strategy asymptotically Optimal?

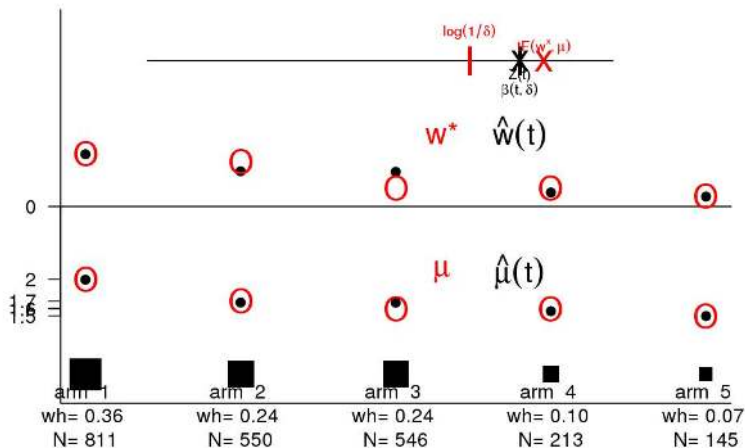


Why is the T&S Strategy asymptotically Optimal?



Why is the T&S Strategy asymptotically Optimal?

Chernoff's stopping rule



Sketch of proof (almost-sure convergence only)

- forced exploration $\implies N_a(t) \rightarrow \infty$ a.s. for all $a \in \{1, \dots, K\}$
- $\rightarrow \hat{\mu}(t) \rightarrow \mu$ a.s.
- $\rightarrow \mathbf{w}^*(\hat{\mu}(t)) \rightarrow \mathbf{w}^*$ a.s.
- \rightarrow tracking rule: $\frac{N_a(t)}{t} \xrightarrow{t \rightarrow \infty} w_a^*$ a.s.

- but the mapping $F : (\mu', \mathbf{w}) \mapsto \inf_{\lambda \in \text{Alt}(\mu')} \sum_{a=1}^K w_a d(\mu'_a, \lambda_a)$ is continuous at $(\mu, \mathbf{w}^*(\mu))$:

- $\rightarrow Z(t) = t \times F(\hat{\mu}(t), (N_a(t)/t)_{a=1}^K) \sim t \times F(\mu, \mathbf{w}^*) = t \times T^*(\mu)^{-1}$
and for every $\epsilon > 0$ there exists t_0 such that

$$t \geq t_0 \implies Z(t) \geq t \times (1 + \epsilon)^{-1} T^*(\mu)^{-1}$$

$$\implies \text{Thus } \tau_\delta \leq t_0 \wedge \inf \left\{ t \in \mathbb{N} : (1 + \epsilon)^{-1} T^*(\mu)^{-1} t \geq \log(2(K-1)t/\delta) \right\}$$

and $\limsup_{\delta \rightarrow 0} \frac{\tau_\delta}{\log(1/\delta)} \leq (1 + \epsilon) T^*(\mu) \quad \text{a.s.}$

Numerical Experiments

- $\mu_1 = [0.5 \ 0.45 \ 0.43 \ 0.4] \rightarrow w^*(\mu_1) = [0.42 \ 0.39 \ 0.14 \ 0.06]$
- $\mu_2 = [0.3 \ 0.21 \ 0.2 \ 0.19 \ 0.18] \rightarrow w^*(\mu_2) = [0.34 \ 0.25 \ 0.18 \ 0.13 \ 0.10]$

In practice, set the threshold to $\beta(t, \delta) = \log\left(\frac{\log(t)+1}{\delta}\right)$ (δ -PAC OK)

| | Track-and-Stop | Chernoff-Racing | KL-LUCB | KL-Racing |
|---------|----------------|-----------------|---------|-----------|
| μ_1 | 4052 | 4516 | 8437 | 9590 |
| μ_2 | 1406 | 3078 | 2716 | 3334 |

Table 1: Expected number of draws $\mathbb{E}_\mu[\tau_\delta]$ for $\delta = 0.1$, averaged over $N = 3000$ experiments.

- Empirically good even for 'large' values of the risk δ
- Racing is sub-optimal in general, because it plays $w_1 = w_2$
- LUCB is sub-optimal in general, because it plays $w_1 = 1/2$

For best arm identification, we showed that

$$\limsup_{\delta \rightarrow 0} \inf_{\delta\text{-correct strategy}} \frac{\mathbb{E}_{\mu}[\tau_{\delta}]}{\log(1/\delta)} = \left(\sup_{w \in \Sigma_K} \inf_{\lambda \in \text{Alt}(\mu)} \left(\sum_{a=1}^K w_a d(\mu_a, \lambda_a) \right) \right)^{-1}$$

and provided an efficient strategy asymptotically matching this bound.

Future work:

- * anytime stopping \rightarrow gives a confidence level
- ** find an ϵ -optimal arm (PAC-setting)
- * find the m -best arms
- *** design and analyze more stable algorithm (hint: optimism)
- *** give a simple algorithm with a finite-time analysis
candidate: play action maximizing the expected increase of $Z(t)$
- *** extend to structured and continuous settings



References

- O. Cappé, A. Garivier, O-A. Maillard, R. Munos, and G. Stoltz. Kullback-Leibler upper confidence bounds for optimal sequential allocation. *Annals of Statistics*, 2013.
- H. Chernoff. Sequential design of Experiments. *The Annals of Mathematical Statistics*, 1959.
- E. Even-Dar, S. Mannor, Y. Mansour, Action Elimination and Stopping Conditions for the Multi-Armed Bandit and Reinforcement Learning Problems. *JMLR*, 2006.
- T.L. Graves and T.L. Lai. Asymptotically Efficient adaptive choice of control laws in controlled markov chains. *SIAM Journal on Control and Optimization*, 35(3):715743, 1997.
- S. Kalyanakrishnan, A. Tewari, P. Auer, and P. Stone. PAC subset selection in stochastic multi- armed bandits. *ICML*, 2012.
- E. Kaufmann, O. Cappé, A. Garivier. On the Complexity of Best Arm Identification in Multi-Armed Bandit Models. *JMLR*, 2015
- A. Garivier, E. Kaufmann. Optimal Best Arm Identification with Fixed Confidence, COLT'16, New York, arXiv:1602.04589
- A. Garivier, P. Ménard, G. Stoltz. Explore First, Exploit Next: The True Shape of Regret in Bandit Problems.
- E. Kaufmann, S. Kalyanakrishnan. The information complexity of best arm identification, COLT 2013
- T.L. Lai and H. Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 1985.
- D. Russo. Simple Bayesian Algorithms for Best Arm Identification, COLT 2016
- N.K. Vaidhyan and R. Sundaresan. Learning to detect an oddball target. arXiv:1508.05572, 2015.