

Lecture 3: Bandit models

Initial question:

Two probability distributions: v_1, v_2
At time t , you choose $a_t \in \{1, 2\}$
and you observe y_t
if $a_t = j$, y_t is an independent sample of v_j

Example: $v_j = \mathcal{N}(\mu_j, \sigma^2=1)$, $j=1, 2$

Question: how do you find $j^* = \arg \max_{j \in \{1, 2\}} \mu_j$?

First approach:

total budget T
policy: assign m_1 times v_1
 m_2 times v_2 with $m_1 + m_2 = T$

pick $\hat{j} = \arg \max_{j \in \{1, 2\}} \bar{X}_{m_1}^j, \bar{X}_{m_2}^j$

where X_k^j = the k -th outcome when using distribution v_j
and $\bar{X}_m^j = \frac{1}{m} (X_1^j + \dots + X_m^j)$

We do the computation in the gaussian case.

$$\bar{X}_{m_1}^1 = \mathcal{N}(\mu_1, \frac{1}{m_1}), \quad \bar{X}_{m_2}^2 = \mathcal{N}(\mu_2, \frac{1}{m_2})$$

Probability of mistake (assume $j^* = 1$ wlog).

$$\begin{aligned} P(\hat{j}_T \neq j^*) &= P(\bar{X}_{m_2}^2 < \bar{X}_{m_1}^1) \\ &= P\left(\frac{\bar{X}_{m_2}^2 - \bar{X}_{m_1}^1}{\sqrt{\frac{1}{m_1} + \frac{1}{m_2}}} < -\frac{\mu_1 - \mu_2}{\sqrt{\frac{1}{m_1} + \frac{1}{m_2}}}\right) = \Phi\left(-\frac{\mu_1 - \mu_2}{\sqrt{\frac{1}{m_1} + \frac{1}{m_2}}}\right) \end{aligned}$$

where $\Phi(x) = \int_{-\infty}^x \frac{e^{-\frac{v^2}{2}}}{\sqrt{2\pi}} dv$

Let $p_1 = \frac{m_1}{T}$ and $p_2 = \frac{m_2}{T}$

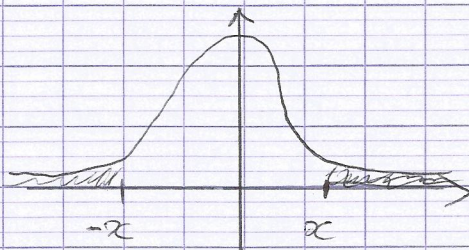
$$\frac{1}{m_1} + \frac{1}{m_2} = \frac{m_1 + m_2}{m_1 m_2} = \frac{1}{T p_1 p_2}$$

$$P(\hat{j}_T \neq j^*) = \Phi\left(-(\mu_1 - \mu_2) \sqrt{T p_1 p_2}\right)$$

decreasing with T

minimale if $p_1 = p_2 = \frac{1}{2}$ i.e. $m_1 = m_2 = \frac{T}{2}$

I want a more explicit formula.



$x > 0$

$$\begin{aligned} \bar{\Phi}(-x) &= \int_x^{+\infty} \frac{e^{-\frac{v^2}{2}}}{\sqrt{2\pi}} dv \\ &\stackrel{z}{=} \int_x^{+\infty} \frac{v}{x} \frac{e^{-\frac{v^2}{2}}}{\sqrt{2\pi}} dv \\ &= \frac{1}{x\sqrt{2\pi}} \left[-v e^{-\frac{v^2}{2}} \right]_x^{+\infty} \\ &= \frac{e^{-\frac{x^2}{2}}}{x\sqrt{2\pi}} = \frac{\phi(x)}{x} \end{aligned}$$

where $\phi(x) = \frac{e^{-\frac{x^2}{2}}}{\sqrt{2\pi}}$ ($= \bar{\Phi}'(x)$)

$$\left(\bar{\Phi}(x) \approx \left(\frac{1}{x} - \frac{1}{3x^3} \right) \phi(x) \right), \text{ in fact } \bar{\Phi}(x) \sim \frac{\phi(x)}{x} \text{ as } x \rightarrow +\infty$$

$$Z \sim \mathcal{N}(0, 1) \quad x > 0$$

$$P(Z > x) = P(e^{tZ} > e^{tx})$$

$$\text{Markov} \leq \frac{\mathbb{E}[e^{tZ}]}{e^{tx}}$$

$$= \frac{e^{t^2/2}}{e^{tx}} = e^{-(tx - t^2/2)}$$

$$\text{Thus } P(Z > x) \leq e^{-\sup_{t>0} tx - t^2/2}$$

$$= e^{-x^2/2}$$

$$\text{Thus } P(\hat{j}_T \neq j^*) \leq e^{-\frac{T(\mu_1 - \mu_2)^2 p_1 p_2}{2}} \quad (1)$$

→ sample complexity of the pb: if we fix a probability of failure δ , what is the total sample size T that's we need to ensure that

$$P(\hat{j}_T \neq j^*) \leq \delta ?$$

From (1) we have the upper-bound/approx

$$e^{-\frac{T(\mu_1 - \mu_2)^2 p_1 p_2}{2}}$$

$$\Leftrightarrow T = \frac{2}{p_1 p_2} \times \frac{\log(\frac{1}{\delta})}{(\mu_1 - \mu_2)^2}$$

In particular, the optimal fixed-design strategie has $p_1 = p_2 = \frac{1}{2}$ and a sample complexity

$$K_{FD}(\delta) = \frac{8 \log(\frac{1}{\delta})}{(\mu_1 - \mu_2)^2}$$

→ useful only if the gap $\Delta = |\mu_1 - \mu_2|$ is known in advance

Sollec 2: sequential statistics

At time t , \mathcal{B} pick $A_t \in \{1, 2\}$ according to a sampling strategy

$$(\psi_t)_{t \geq 1}, \quad \psi_t = (\{1, 2\} \times \mathbb{R})^{t-1} \rightarrow \{1, 2\}$$

$$A_t = \psi(A_1, Y_1, A_2, Y_2, \dots, A_{t-1}, Y_{t-1})$$

Stopping^v: a stopping^v τ with respect to $(\mathcal{F}_t)_{t \geq 1}$ where

$$\mathcal{F}_t = \sigma(A_1, Y_1, \dots, A_t, Y_t)$$

• A decision rule: \hat{A}_τ measurable w.r.t. (\mathcal{F}_τ)

• Rules: for a failure probability δ fixed in advance, we are obliged to choose sampling strategy, a stopping rule and a decision rule such that

$$\mathbb{P}(\hat{A}_\tau \neq j^*) \leq \delta$$

• goal: choose the strategy (= stopping rule, stopping time, decision rule) to minimize $\mathbb{E}_v[\tau]$, $v = (v_1, v_2)$.

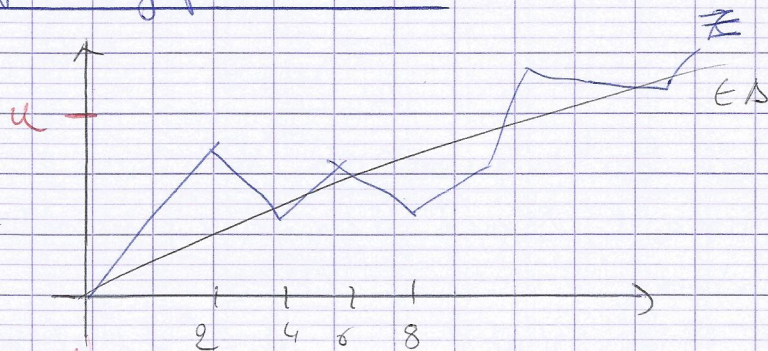
• Intuitively: $\hat{A}_\tau = \arg \max_{j \in \{1, 2\}} \overline{X}_{N_j}^j(\tau)$

where $N_j(t) = \sum_{s \leq t} \mathbb{1}\{A_s = j\}$

$$A_t = 1 + (t \bmod 2)$$

Decision rule?

Ⓐ If the gap Δ is known.



$$S_k = Z_k = x_1^1 + \dots + x_k^1 - x_1^2 - \dots - x_k^2$$

$$= (x_1^1 - x_1^2) + \dots + (x_k^1 - x_k^2)$$

Dr $\mu_1 > \mu_2$, wlog $\mu_1 > \mu_2$

Stopping rule $\tau \leq 2k \iff |S_k| > u$

We need to choose u so that

$$P(\exists k \geq 1, S_k < -u) < \delta$$

~~and $\forall k, S_k < u$~~

Lemma: $P(\exists n \geq 1, S_n < -u) \leq e^{-\frac{\Delta u}{\sigma^2}}$

Proof: $M_k = e^{-\Delta(x_1^1 + \dots + x_k^1)}$

$$\mathbb{E}[M_k] = \mathbb{E}[e^{-\Delta(x_1^1 + \dots + x_k^1)}] \text{ if } z \sim N(0, 1)$$

$$= e^{-\Delta^2 \frac{\sigma^2 k}{2}} = 1$$

Then M_k is a martingale. Let $\tau'_m = \tau \wedge m$ defined after

But for every $m \geq 1$, τ'_m is a bounded stopping

\rightarrow From the Doob's stopping theorem, we conclude that

$$\mathbb{E}[M_{\tau'_m}] = 1$$

for the stopping time $\tau' = \inf\{k, D_k \leq -u\}$

$$\text{If } m > \tau', \quad S_{\tau' \wedge m} \leq -u$$

$$\text{and } M_{\tau' \wedge m} = e$$

$$1 = \mathbb{E}[M_{\tau' \wedge m}] \geq \mathbb{E}\left[M_{\tau'} \mathbb{1}(\tau' \leq m)\right]$$

$$\geq e^{-\Delta u} P(\tau' \leq m)$$

$$= e^{-\Delta u} P(\tau' \leq m)$$

and $P(\tau' \leq m) \leq e^{-\Delta u}$ for all $m \geq 1$.

Thus

$$P(\tau' < \infty) \leq e^{-\Delta u}.$$

□

Consequence: choosing u st $e^{-\Delta u} \leq \delta$
ie $u \geq \frac{\ln(1/\delta)}{\Delta}$

yields an admissible stopping rule:

$$P(\hat{A}_\tau \neq j^*) \leq \delta$$

Summarize:

$$A_\tau = 1 + (S + 1 \bmod 2)$$

$$\tau = \inf\{k: |x_1^2 + \dots + x_k^2 - x_1^2 - \dots - x_k^2| > \frac{\ln(1/\delta)}{\Delta}\}$$

$$\hat{A}_\tau = \arg \max_j \bar{X}_{\tau/2}$$

We need to upper-bound its sample complexity

$$K_S(\delta) = \mathbb{E}_r [\tau_S]$$

Lemma: Let M be a \mathbb{N} -valued r.v.

Then $\mathbb{E}[M] = \sum_{k=1}^{+\infty} P(M \geq k)$

Proof:

$$\begin{aligned} \sum_{k=1}^{+\infty} P(M \geq k) &= \sum_{k=1}^{+\infty} \sum_{j=k}^{+\infty} P(M=j) \\ &= \sum_{j=1}^{+\infty} \sum_{k=1}^j P(M=j) \\ &= \sum_{j=1}^{+\infty} j P(M=j) = \mathbb{E}[M] \end{aligned}$$

□

$$\mathbb{E}[\tau_S] \leq \sum_{k=1}^{+\infty} P(\tau_S \geq k)$$

But $\{\tau_S \geq k\} \subset \{S_k < \alpha\}$ and $S_k \sim d(k\Delta, \sigma^2)$

So $\frac{S_k - k\Delta}{\sqrt{2k}} \sim d(0, 1)$

Thus, if $k \geq \frac{\ln \frac{1}{\delta}}{\Delta^2}$

$$\begin{aligned} P(S_k < \alpha) &= P\left(\frac{S_k - k\Delta}{\sqrt{2k}} < \frac{\alpha - k\Delta}{\sqrt{2k}}\right) \\ &= P\left(\frac{S_k - k\Delta}{\sqrt{2k}} < \frac{(\ln \frac{1}{\delta})/\Delta - k\Delta}{\sqrt{2k}}\right) \\ &\leq \exp\left(-\frac{1}{2} \left(\frac{(\ln \frac{1}{\delta})/\Delta - k\Delta}{\sqrt{2k}}\right)^2\right) \\ &= \exp\left(-\frac{1}{4k\Delta^2} \left(k\Delta^2 - \ln \frac{1}{\delta}\right)^2\right) \end{aligned}$$

We have seen that $\mathbb{P}(S_n < \frac{cn}{\delta}) \leq \exp\left(-\frac{(cn - cn/\delta)^2}{4n\delta^2}\right)$

Thus, $\mathbb{E}[\tilde{\zeta}_\delta] = \sum_{k=1}^{\infty} \mathbb{P}(\tilde{\zeta}_\delta \geq k)$

$$\leq 2 + 2 \sum_{k=1}^{\infty} \mathbb{P}(\tilde{\zeta}_\delta \geq 2k)$$

$$\leq 2 + 2 \sum_{k=\lceil \frac{cn}{\delta^2} + \frac{\sqrt{cn}}{\delta} \rceil}^{\infty} \mathbb{P}(\tilde{\zeta}_\delta \geq 2k)$$

and $\sum_{k=\lceil \frac{cn}{\delta^2} + \frac{\sqrt{cn}}{\delta} \rceil}^{\infty} \mathbb{P}(\tilde{\zeta}_\delta \geq 2k) \leq \sum_{k=\lceil \frac{cn}{\delta^2} + \frac{\sqrt{cn}}{\delta} \rceil}^{\infty} \mathbb{P}(S_n < \frac{cn}{\delta}) \leq \sum_{k=\lceil \frac{cn}{\delta^2} + \frac{\sqrt{cn}}{\delta} \rceil}^{\infty} \exp\left(-\frac{(kn - cn/\delta)^2}{4n\delta^2}\right)$

$$\leq \sum_{u=\lceil \frac{\sqrt{cn}}{\delta} \rceil}^{\infty} \exp\left(-\frac{(u\delta^2)^2}{4\delta^2(\frac{cn}{\delta^2} + u)}\right) = \sum_{u=\lceil \frac{\sqrt{cn}}{\delta} \rceil}^{\infty} \exp\left(-\frac{u\delta^2}{4(\frac{cn}{\delta^2} + 1)}\right)$$

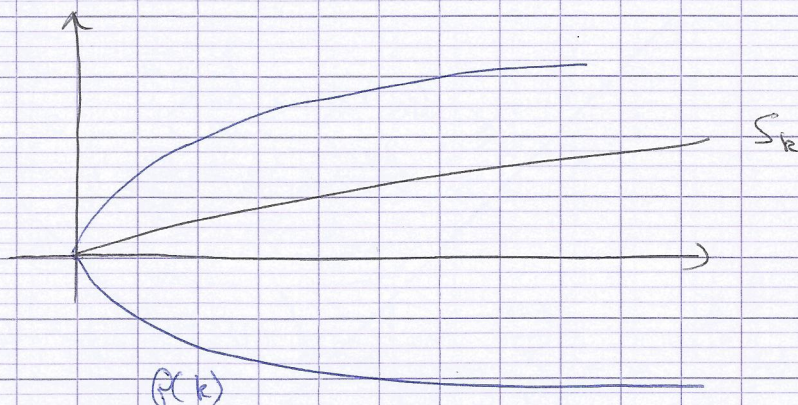
But $u\delta^2 > \sqrt{cn} \Rightarrow \frac{\sqrt{cn}}{u\delta^2} \times \sqrt{cn} \leq \sqrt{cn}$ and thus

$$\sum_{k=\lceil \frac{cn}{\delta^2} + \frac{\sqrt{cn}}{\delta} \rceil}^{\infty} \mathbb{P}(\tilde{\zeta}_\delta \geq 2k) \leq \sum_{u=\lceil \frac{\sqrt{cn}}{\delta} \rceil}^{\infty} \exp\left(-\frac{u\delta^2}{4(\sqrt{cn} + 1)}\right)$$

$$\leq 1 + \int_0^{+\infty} \exp\left(-\frac{u\delta^2}{4(\sqrt{cn} + 1)}\right) du = 1 + \frac{4(\sqrt{cn} + 1)}{\delta^2}$$

Thus, $\mathbb{E}[\tilde{\zeta}_\delta] \leq \frac{2cn}{\delta^2} + 4 + \frac{4(\sqrt{cn} + 1)}{\delta^2} \sim \frac{2cn}{\delta^2}$ as $\delta \rightarrow 0$

(B) Δ is unknown



Remark: $\frac{S_k}{\sqrt{k}} \sim d(0, 1)$

$\limsup_{k \rightarrow \infty} \frac{S_k}{\sqrt{k}} = +\infty$; $\liminf_{k \rightarrow \infty} \frac{S_k}{\sqrt{k \log k}} = -1$

Law of iterated logarithm

II Lower bounds in bandit model

Let k be a positive integer

Let $\nu_1, \nu_2, \dots, \nu_k$ be probability distributions on \mathbb{R}

We say that $\nu = (\nu_1, \dots, \nu_k)$ is a bandit model

For every $j \in \{1, \dots, k\}$, let $x_1^j, x_2^j, x_3^j, \dots$ be indep. samples of ν_j

A bandit strategy is a sampling strategy $(\psi_i)_{i \geq 1}$ where $\psi_i : (\{1, \dots, k\} \times \mathbb{R})^{i-1} \rightarrow \{1, \dots, k\}$

It defines a sequence of r.v.

- (2)
- $A_1 = \psi_1(\emptyset)$, $y_1 = X_1^{A_1}$, $N_j(1) = \mathbb{1}\{A_1 = j\}$
 - $A_2 = \psi_2(A_1, y_1)$, $y_2 = X_2^{A_2}$, $N_j(2) = \sum_{s=1}^2 \mathbb{1}\{A_s = j\}$
 - \vdots
 - $A_t = \psi_t(A_1, y_1, \dots, A_{t-1}, y_{t-1})$, $y_t = X_t^{A_t}$

$$N_j(H) = \sum_{s \leq t} \mathbb{1}\{A_s = j\}$$

2 possible goals:

- Best arm identification

\rightarrow find $j^* \in \arg \max_{1 \leq j \leq K} \mathbb{E}[X_1^j]$ (previous example)

- Regret minimization

An oracle is someone who knows j^* in advance, on average, the oracle gathers

$$\mathbb{E}\left[\sum_{t=1}^T X_t^{j^*}\right] = T\mu^*, \text{ where } \mu^* = \max_j \mathbb{E}[X_1^j]$$

from the time 1 to T. A strategy ψ gathers

$$\mathbb{E}\left[\sum_{t=1}^T X_t^{A_t}\right] \text{ during the same time.}$$

(This is the expected sum of rewards)

The regret is

$$R_T = T\mu^* - \mathbb{E}\left[\sum_{t=1}^T X_t^{A_t}\right]$$

Regret minimization aims to choosing the strategy ψ so as to minimize R_T .

• To fix ideas, we will assume that the $(v_j)_j$ belongs to an exponential family.

$$v^j = \mathcal{D}(\mu_j, \lambda)$$

$$\text{or } v^j = \mathcal{B}(\mu_j)$$

$$\text{or } v^j = \mathcal{P}(\mu_j)$$

In particular, we will assume that we can compute the KL divergence between distribut^o $KL(v^j, v^k)$

Regret minimization strategies

(greedy, ϵ -greedy, randomized)

Lower bound.

Let v and v' be two bandit models.

Then, for all $T \geq 1$

$$\sum_{j=1}^k \mathbb{E}_v [N_j(T)] KL(v_j, v'_j) \geq KL(\mathbb{E}_v[z], \mathbb{E}_{v'}[z])$$

for every z w. $0 \leq z \leq 1$

Proof: Let $H_\epsilon = (A_1, g_1, A_2, g_2, \dots, A_\epsilon, g_\epsilon)$ and let v^{H_ϵ} be the distribut^o of the H_ϵ under the model v .

Step 1: $KL(v^{H_\epsilon}, v'^{H_\epsilon}) \stackrel{\text{claim}}{=} \sum_{j=1}^k \mathbb{E}[N_j(T)] KL(v_j, v'_j)$

Proof by induct^o. $\epsilon = 0$ obvious

$$\epsilon = 1 \quad H^1 = \begin{matrix} (A_1, g_1) \\ \parallel \\ a_1 \quad x_{a_1}^1 \end{matrix}, \quad \begin{matrix} v^{H^1} = (g_{a_1}, v_{a_1}) \\ v'^{H^1} = (g'_{a_1}, v'_{a_1}) \end{matrix}$$

$$\begin{aligned}
 KL(v^{H_1}, v^{H_2}) &= \int dv^{H_1}(h) \ln \frac{dv^{H_1}(h)}{dv^{H_2}(h)} \\
 &= \int dv_{a_1}(g) \ln \frac{dv_{a_1}(h)}{dv'_{a_1}(h)} = KL(v_{a_1}, v'_{a_1})
 \end{aligned}$$

$$N_j(h) = \mathbb{1}\{j = a_1\} = \mathbb{E}[N_j^o(h)]$$

$$\sum_{j=1}^K \mathbb{E}[N_j^o(h)] K(v_j^o, v_j^{\prime o}) = KL(v_{a_1}, v'_{a_1})$$

Induction:

$$\begin{aligned}
 &KL(v^{H_{e+1}}, v^{H_{e+2}}) \\
 &= KL(v^{H_e}, v^{H_e}) + \mathbb{E}[KL(v^{(H_{e+1}, g_{e+1})}, v^{(H_{e+1}, g_{e+1})})] \\
 &= \sum_{j=1}^K \mathbb{E}[N_j^o(h) KL(v_j, v_j^{\prime o})] + \mathbb{E}\left[\sum_{j=1}^K \mathbb{1}\{A_{e+1} = j\} \times \right. \\
 &\quad \left. \times KL(v^{(A_{e+1}, g_{e+1})}, v^{(A_{e+1}, g_{e+1})})\right] \\
 &= \sum_{j=1}^K \mathbb{E}\left[\sum_{s=1}^e \mathbb{1}\{A_s = j\}\right] KL(v_j, v_j^{\prime o}) \\
 &\quad + \mathbb{E}\left[\sum_{s=1}^K \mathbb{1}\{A_{e+1} = j\}\right] K(v_j, v_j^{\prime o}) \\
 &= \sum_{j=1}^K \mathbb{E}\left[\sum_{s=1}^{e+1} \mathbb{1}\{A_s = j\}\right] KL(v_j, v_j^{\prime o})
 \end{aligned}$$

Step 2: Contract^o of entropy (data processing ineq.)
 If Z is $\sigma(H_e)$ -measurable
 $KL(v^Z, v'^Z) \leq KL(v^{H_e}, v'^{H_e})$

Step 3: if $0 \leq Z \leq 1$, then $KL(v^Z, v'^Z) \geq K(\mathbb{E}_g(Z), \mathbb{E}_{v'}(Z))$

proof: $\mathbb{E}_v(Z) = \mathbb{E}[U \leq Z]$ where U is $\mathcal{U}([0, 1])$ and $\mathbb{1}\{Z =$

$$\begin{aligned}
 \mathbb{P}(U \leq Z) \int dP(z) \mathbb{1}\{U \leq Z\} &= \int dP(z) \int \mathbb{1}\{U \leq Z\} dU \\
 &= \int Z dP(z) = \mathbb{E}[Z]
 \end{aligned}$$

$$\begin{aligned}
& \text{Then } KL((v_0)_t^{(z,u)}, (v_0)_t^{(z,u)}) \\
&= KL(v^z, v'^z) + K(\cancel{L^u}, \cancel{L^u}) \\
&\geq KL((v_0)_t^{\mathbb{1}\{U \leq z\}}, (v_0)_t^{\mathbb{1}\{U \leq z\}}) \\
&= KL(P_{v_0}(U \leq z), P_{v_0}(U \leq z)) \\
&= KL(\mathbb{E}_v(z), \mathbb{E}_{v'}(z)) \quad \square
\end{aligned}$$

To put everything together:

$$\begin{aligned}
\sum_{j=2}^K \mathbb{E}[N_j(t)] KL(v_j, v_j') &= KL(v^{\mathbb{H}_t}, v'^{\mathbb{H}_t}) \\
&\geq KL(v^z, v'^z) \\
&\geq KL(\mathbb{E}_v(z), \mathbb{E}_{v'}(z))
\end{aligned}$$

□

Example 2 arms gaussian, gap Δ unknown
 $v = (v_1, v_2)$ $U_j \sim \mathcal{N}(\mu_j, 1)$

wlog $\mu_1 > \mu_2$, $j^* = 1$

$v' = (v_2, v_2)$, $Z = \mathbb{1}\{\hat{A}_t = 1\}$

$$\begin{aligned}
& \mathbb{E}[N_1(t)] KL(v_1, v_2) + \mathbb{E}[N_2(t)] KL(v_2, v_2) \\
&\geq KL(\mathbb{E}_v(\mathbb{1}\{\hat{A}_t = 1\}), \mathbb{E}_{v'}(\mathbb{1}\{\hat{A}_t = 1\}))
\end{aligned}$$

(\Rightarrow)

$$\frac{(\mu_1 - \mu_2)^2}{2} \mathbb{E}[Z] \geq KL(P_v(\hat{A}_t = 1), P_{v'}(\hat{A}_t = 1))$$

$\geq 1 - \delta$

$\leq \delta$

error almost with prob δ

δ admissible strategy

$$\geq KL(1 - \delta, \delta) = (1 - \delta) \ln\left(\frac{1 - \delta}{\delta}\right) + \delta \ln\left(\frac{\delta}{1 - \delta}\right)$$

$$\geq \ln\left(\frac{1}{2.4\delta}\right)$$

$$\sim \ln\left(\frac{1}{\delta}\right)$$

13)

$$\frac{(\mu_1 - \mu_2)^2}{2} E[\tau] \geq \ln \frac{1}{2.48}$$

$$E[\tau] \geq \frac{2 \ln \left(\frac{1}{2.48} \right)}{\Delta^2}$$