



Ch6. General Decision Making

- Introduce a unified framework for decision making, including structured bandit problem, contextual bandit problem etc.
- Show that the Decision-Estimation Coefficient (DEC) and associated meta-algorithm (E2D) extend to this general framework and also show that boundedness of the DEC is sufficient and necessary for low regret (constitute a fundamental limit)

Setting: Focus on a framework called Decision Making with Structured Observations (DMSO)

- for $t=1, 2, \dots, T$ round, the learner select a decision $\pi^t \in \Pi$, the decision space
- Nature selects reward $r^t \in \mathcal{R}$ and observation $o^t \in \mathcal{O}$ based on π^t
observation space reward space

- Both r^t and o^t are observed by the learner

Two key Assumptions:

Assumption 1 (Stochastic Rewards and Observations):

r^t and o^t are independently generated via

$$(r^t, o^t) \sim M^*(\cdot | \pi^t)$$

Where $M^*: \Pi \rightarrow \Delta(R \times O)$ is the underlying true model.

Assumption 2 (Realizability): The learner has access to a flexible enough model class \mathcal{M} , where $M^* \in \mathcal{M}$

(Here, \mathcal{M} could be of linear models, neural networks, random forest, and other function approximation)

Objectives: For a model $M \in \mathcal{M}$, $\mathbb{E}^{M, \pi}(\cdot)$ is

the expectation under $(r, o) \sim M(\pi)$, let

$$f^M(\pi) = \mathbb{E}^{M, \pi}(r)$$

be the mean reward function; and let

$$\pi_M = \arg \max_{\pi \in \Pi} f^M(\pi)$$

be the optimal decision with maximal expected reward;
and also let

$$\mathcal{F}_M = \{f^M \mid M \in \mathcal{M}\}$$

be the induced class of mean reward functions.

Finally, for evaluation of the learner's performance, we consider regret to the optimal decision for M^*

$$\text{Reg} = \sum_{t=1}^T \mathbb{E}_{\pi^t \sim p^t} [f^{M^*}(\pi_{M^*}) - f^{M^*}(\pi^t)]$$

where $p^t \in \Delta(\Pi)$ is the learner's distribution over decisions at round t . (Shorthand $f^* = f^{M^*}$, $\pi^* = \pi_{M^*}$)

Rank: Compared to basic bandit problem, $y^t \sim M^*(\cdot \mid \pi^t)$

without the observations, and the mean reward is $\mathbb{E}(y \mid \pi)$
only, while regret is $\sum_{t=1}^T f^*(\pi^t) - \underbrace{\mathbb{E}(y \mid \pi)}_{f^*(\pi)}$
 $\sum_{t=1}^T \mathbb{E}_{\pi^t \sim p^t} [f^*(\pi^t)]$

That is, observation provided extra information gain.

Sec 1. Some Examples

(Example 1) Structural Bandits. When $\mathcal{O} = \{\emptyset\}$, i.e. no observations, DMSD reduces to Structural Bandits problem. We can start with a set of models \mathcal{M} and then define induced class $\mathcal{F}_{\mathcal{M}}$, serving as the class \mathcal{F} of mean reward functions before. Different $\mathcal{F}_{\mathcal{M}}$ will include different structural bandit problems such as linear, nonparametric, etc.

(Example 2) Contextual Bandits. In Context Bandit, reward $r^t \sim M^*(\cdot | \pi^t, x^t)$ for some covariate x^t , and $f^*(x, \pi) = \mathbb{E}(r | x, \pi)$ for $r \sim M^*(\cdot | \pi, x)$.

Think of π^t as functions mapping x^t to an action in $\Pi = [A]$. On round t , the decision-maker selects a mapping $\pi^t: \mathcal{X} \rightarrow [A]$, and the context $\mathcal{O}^t = x^t$ is observed for each round. This is basically observing x^t , then

select $\pi^t(x^t) \in [A]$, (which is called behavioral decision rule in decision theory). Let $\mathcal{O} = \mathcal{X}$ be the space of context, $\mathcal{I} = [A]$ be the set of actions, and $\Pi: \mathcal{X} \rightarrow [A]$ be the space of decisions. Then, $(r, x) \sim M(\cdot | \pi)$ has: $x \sim D^M$ for some context dist, $r \sim R^M(\cdot | x, \pi(x))$ for some reward dist R^M . Here, D^M for x is part of M (treating x as \mathcal{O} , observation).

(Example 3) Online reinforcement learning. For online reinforcement learning, learner selects a randomized, non-stationary policy $\Pi = (\pi_1, \dots, \pi_H)$, $\pi_h: \mathcal{S} \rightarrow \Delta(A)$
 \uparrow
 $\pi_h \in \Pi_{\text{rns}}$

Beginning from state $s_1 \sim d_1 \in \Delta(\mathcal{S})$, for $h = 1, \dots, H$

$$a_h \sim \pi_h(s_h)$$

$$r_h \sim R_h^M(s_h, a_h), \quad R_h^M: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{R})$$

(reward dist)

and

$$s_{h+1} \sim P_h^M(s_h, a_h), \quad P_h^M: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$$

(transition kernel)

then $f^M(\pi) = \mathbb{E}^{M, \pi} \left(\sum_{h=1}^H r_h \right)$
 stands for MDP $\{S, A, P_h^M, R_h^M\}$

Here, let's take $\Pi = \Pi_{\text{rns}}$, $\gamma^t = \sum_{h=1}^H \gamma_h^t$ and $\mathcal{O}^t = \tau^t = (S_1^t, a_1^t, r_1^t), \dots, (S_H^t, a_H^t, r_H^t)$, the trajectory of learning. Then

$$(\gamma^t, \mathcal{O}^t) = \left(\sum_{h=1}^H \gamma_h^t, \tau^t \right) \sim (R_h^M, P_h^M | A) \quad \uparrow \pi$$

can be considered as $M(\cdot | \pi^t)$ for some dens M and $\pi^t \in \Pi_{\text{rns}}$.

Sec 2. DEC for General Decision Making

Optimally explore and make decisions for M is central to understanding the optimal statistical complexity for M .

As seen before, any notion of complexity needs to capture

- (i) Simple problems like multi-armed bandit
- (ii) problems with structural feedback where observations or structures in the noise can provide extra information

(Decision Estimation Coefficient)

For \mathcal{M} , reference model $\hat{M} \in \mathcal{M}$ and $\gamma > 0$ (scalar). DEC for general decision making is

regret of decision

$$\text{dec}_\gamma(\mathcal{M}, \hat{M}) = \inf_{P \in \Delta(\mathcal{I})} \sup_{M \in \mathcal{M}} \mathbb{E}_{\pi \sim P} \left[f^M(\pi_M) - f^M(\pi) - \gamma \cdot \mathcal{D}_H^2(M(\pi), \hat{M}(\pi)) \right]$$

where $\mathcal{D}_H^2(P, Q) = \int (\sqrt{p} - \sqrt{q})^2 d\nu$ information gain from obs.

for $P, Q \ll \nu$. Also, define

$$\text{dec}_\gamma(\mathcal{M}) = \sup_{\hat{M} \in \text{co}(\mathcal{M})} \text{dec}_\gamma(\mathcal{M}, \hat{M})$$

where $\text{co}(\mathcal{M})$ is the convex hull of \mathcal{M} , as did for $\text{co}(\mathcal{F})$

Remark: Compared to structural bandit problem, the major difference is that instead of "max" and considering $f(\pi_f)$

$- f(\pi)$, i.e. restricts on a class of reward functions, here the general DEC is defined over \mathcal{M} (model class for both reward and observations).

Also, rather than measuring the information gain via

$f(\pi) - \hat{f}(\pi)$, here we consider information gain from matching the dsns over rewards and observations of M and \hat{M} (for learner's decision π), $\mathbb{E}_{\pi \sim p} [D_H^2(M(\pi), \hat{M}(\pi))]$

(i) incorporates observations O^t or τ^t in reinforcement learning

(ii) even for bandit problems, it measures the distance between dsns rather than the means.

Sec 2.1 E2D Algorithm for General Decision Making

Estimation to Decisions (E2D) for general decision making is readily extended from structured bandit problems.

Parameter $\gamma > 0$ (given)

For $t = 1, \dots, T$ do

Obtain \hat{M}^t from online estimation oracle with $(\pi^1, r^1, O^1), \dots, (\pi^{t-1}, r^{t-1}, O^{t-1})$ by minimizing $\text{dss}_\gamma(M, \hat{M}^t)$. Compute

$$p^t = \arg\min_{p \in \Delta(\Pi)} \sup_{M \in \mathcal{M}} \mathbb{E}_{\pi \sim p} [f^{M(\pi)} - f^{M(\pi)} - \gamma \cdot D_H^2(M(\pi), \hat{M}(\pi))]$$

(*)

Simple decision $\pi^t \sim p^t$ and update estimation with (π^t, r^t, o^t) .

Link: i) Rather attempting to estimate the reward function f^* , one estimates the underlying model M^* (details later)

$$\text{i) } \text{DEO} \xrightarrow{\uparrow \mathcal{H}^{t+1} = (\pi^1, r^1, o^1), \dots, (\pi^t, r^t, o^t)} \hat{M}^t \xrightarrow{(\pi^t)} p^t \rightarrow \pi^t \rightarrow \text{DEO}$$

Proposition 1 (26 in the note): Running E2D, the regret is bounded by DEC and estimation error, which is defined here as

$$\text{as } \text{Est}_H = \sum_{t=1}^T \mathbb{E}_{\pi^t \sim p^t} [D_H^2(M^*(\pi^t), \hat{M}^t(\pi^t))]$$

That is, for $\delta > 0$, E2D admits

$$\text{Reg} \leq \sup_{\hat{M} \in \hat{\mathcal{M}}} \text{dec}_\gamma(M, \hat{M}) \cdot T + \delta \cdot \text{Est}_H$$

almost surely, where $\hat{\mathcal{M}}$ is any set s.t. $\hat{M}^t \in \hat{\mathcal{M}}$ for all $t \in [T]$

proof: $\text{Reg} = \sum_{t=1}^T \mathbb{E}_{\pi^t \sim p^t} [f^*(\pi_{M^*}) - f^{M^*}(\pi^t)]$

$$= \sum_{t=1}^T \mathbb{E}_{\pi^t \sim p^t} [f^*(\pi_{M^*}) - f^{M^*}(\pi^t)] \\ - \gamma \mathbb{E}_{\pi^t \sim p^t} [D_H^2(M^*(\pi^t), \hat{M}^t(\pi^t))] + \gamma \mathbb{E}_{St} H$$

For each t , as $M^* \in \mathcal{M}$ (Assumption 2)

$$\mathbb{E}_{\pi^t \sim p^t} [f^*(\pi_{M^*}) - f^{M^*}(\pi^t)] \\ - \gamma \mathbb{E}_{\pi^t \sim p^t} [D_H^2(M^*(\pi^t), \hat{M}^t(\pi^t))]$$

$$\leq \sup_{M \in \mathcal{M}} \mathbb{E}_{\pi^t \sim p^t} [f^M(\pi_M) - f^M(\pi^t)] \\ - \gamma \mathbb{E}_{\pi^t \sim p^t} [D_H^2(M(\pi^t), \hat{M}^t(\pi^t))]$$

by def of p^t in E2D

$$= \inf_{p \in \Delta(\Pi)} \sup_{M \in \mathcal{M}} \mathbb{E}_{\pi \sim p} [f^M(\pi_M) - f^M(\pi) - \gamma \cdot D_H^2(M(\pi), \hat{M}^t(\pi))]$$

$$= \text{dec}_\gamma(M, \hat{M}^t). \text{ Summing over } t,$$

$$\text{Reg} \leq \sup_{\hat{M} \in \hat{\mathcal{M}}} \text{dec}_\gamma(M, \hat{M}) \cdot T + \gamma \mathbb{E}_{St} H \quad \#$$

Remark: One can optimize over γ above, to yield

$$\text{Reg} \leq \inf_{\gamma} \left\{ \sup_{\hat{M} \in \hat{\mathcal{M}}} \text{dec}_\gamma(M, \hat{M}) \cdot T + \gamma \cdot \mathbb{E}_{St} H \right\}$$

$$\leq 2 \inf_{\gamma > 0} \left\{ \sup_{M \in \hat{\mathcal{M}}} \text{dec}_\gamma(M, \hat{M}) \cdot T, \gamma \cdot \text{Esc}_H \right\}$$

For any finite class M , the averaged exponential weights (#) algorithm with the log loss achieves $\text{Esc}_H \leq \log(|M|/\delta)$ w.p at least $1 - \delta$. We can take $\hat{M} = \mathcal{W}(M)$, $\mathcal{W} \leq O_\delta$.

Then, for any finite class, w.p. $1 - \delta$,

$$\text{Reg} \leq \text{dec}_\gamma(M) \cdot T + \gamma \log(|M|/\delta).$$

Review: Exponential weights is a main online learning algorithm, applicable to finite class. At each time, the algorithm compute a den $q^t \in \Delta(F)$, where $\Delta(F) \propto \exp\{-\eta \sum_{i=1}^{t-1} \ell(f(x^i), y^i)\}$ for $\eta > 0$. Specifically,

for $t = 1, \dots, T$ do

compute q^t above.

let $f^t = \mathbb{E}_{f \sim q^t}(f)$

observe (x^t, y^t) , incur $\ell(f^t(x^t), y^t)$

loss function

↑

Sec 2.2 Examples

(Example 4, Multi-armed bandit with Gaussian-rewards) Let $\mathcal{I} = [A]$, $\mathcal{R} = \mathbb{R}$, $\mathcal{O} = \{\phi\}$

(no observations). Define

$$\mathcal{M}_{\text{MAB-G}} = \left\{ M: M(\pi) = N(f(\pi), 1) \right. \\ \left. f: \mathcal{I} \rightarrow [0, 1] \right\}$$

Consider $\hat{M} \in \mathcal{M}$. From results of multi-armed bandit problem, $\text{dely}(\hat{f}, \hat{f}) \propto \frac{A}{\gamma}$, it is thus sufficient to argue that the squared Hellinger-distance for Gaussian reduces to squared difference between the means.

$$\begin{aligned} \text{In fact, } D_H^2(M(\pi), \hat{M}(\pi)) \\ \leq D_{KL}(M(\pi) \parallel \hat{M}(\pi)) \\ \stackrel{(\#)}{=} \frac{1}{2} (f^M(\pi) - f^{\hat{M}}(\pi))^2 \end{aligned}$$

$$\begin{aligned} \text{On the other hand, } D_H^2(M(\pi), \hat{M}(\pi)) \\ = 1 - \exp \left\{ -\frac{1}{8} (f^M(\pi) - f^{\hat{M}}(\pi))^2 \right\} \end{aligned}$$

$$\text{So that } D_H^2(M(\pi), \hat{M}(\pi)) \geq c \cdot (f^M(\pi) - f^{\hat{M}}(\pi))^2$$

as $1 - e^{-x} \geq (1 - e^{-1})x$ for $x \in [0, 1]$. Hence these two inequalities give

$$\text{dec}_\gamma(M_{\text{MAB-G}}) \propto \frac{A}{\gamma}$$

By Proposition 1, then for $M_{\text{MAB-G}}$.

$$\text{Reg} \lesssim \frac{AT}{\gamma} + \gamma \text{Est}_H$$

Review: (\S) As densities of any two dens \mathbb{P} and \mathbb{Q} ,

p.g. satisfies $|\mathbb{P} - \mathbb{Q}| = |(\sqrt{\mathbb{P}} + \sqrt{\mathbb{Q}})(\sqrt{\mathbb{P}} - \sqrt{\mathbb{Q}})| \geq (\sqrt{\mathbb{P}} - \sqrt{\mathbb{Q}})^2$

and $\text{TV}(\mathbb{P}, \mathbb{Q}) = \frac{1}{2} \int |\mathbb{P} - \mathbb{Q}| d\nu$, it follows that

$$H^2(\mathbb{P}, \mathbb{Q}) \leq \text{TV}(\mathbb{P}, \mathbb{Q})$$

so $D_H^2(\mathbb{P}, \mathbb{Q}) \leq D_{\text{KL}}(\mathbb{P} \parallel \mathbb{Q})$ by Pinsker's inequality.

Let $\gamma = \sqrt{AT / \text{Est}_H}$, then

$$\text{Reg} \lesssim \sqrt{AT \cdot \text{Est}_H}$$

(Example 5. Bandit problems with structured noise)

Let $\mathcal{I} = [A]$, $\mathcal{R} = \mathbb{R}$, $\mathcal{O} = \{\phi\}$. Define

$$\mathcal{M}_{\text{MAB-SN}} = \{M_1, \dots, M_A\} \cup \{\hat{M}\}$$

where $M_i(\pi) := \mathcal{N}(\frac{1}{2}, 1)$ for $\pi \neq i$ and $M_i(\pi) = \text{Ber}(3/4)$ for $\pi = i$. Define $\hat{M}(\pi) = \mathcal{N}(\frac{1}{2}, 1)$ for all $\pi \in \mathcal{I} = [A]$.

The valuable information contained in the reward obs. is reflected in the Hellinger divergence, which attains the maximum when comparing a continuous obs. to a discrete one.

$$D_H^2(M_i(\pi), \hat{M}(\pi)) = 2\mathbb{I}(\pi = i)$$

Notice that the maximum over M in the definition of $\text{deg}_\gamma(\mathcal{M}_{\text{MAB-SN}}, \hat{M})$ is NOT attained at $M = \hat{M}$, as both the divergence and " regret " term will be zero.

regardless of \mathcal{P} . Take $p = \text{unif}[A]$, for any
 $M \in \{M_1, \dots, M_A\}$

$$\mathbb{E}_{\pi \sim p} [f^M(\pi_M) - f^M(\pi)] = \left(1 - \frac{1}{A}\right) \left(\frac{3}{4} - \frac{1}{2}\right)$$

$$\text{as } f^M(\pi) = \mathbb{E}(\mathbb{E}(r|\pi)) = \left(1 - \frac{1}{A}\right) \cdot \frac{1}{2} + \frac{1}{A} \cdot \frac{3}{4}$$

$$\text{and } \pi_M = \arg \max_{\Pi \in [A]} f^M(\pi) = 3/4$$

$$\begin{aligned} \text{and } \text{dec}_\gamma(\mathcal{M}, \hat{M}) &\lesssim \left(1 - \frac{1}{A}\right) \left(\frac{3}{4} - \frac{1}{2}\right) - \gamma \frac{2}{A} \\ &\lesssim \mathbb{I}(\gamma \leq \frac{A}{4}) \end{aligned}$$

$$\text{as } \sup_{M \in \mathcal{M}} \mathbb{E}_{\pi \sim p} [D_H^2(M_i(\pi), \hat{M}(\pi))] = \frac{1}{A}.$$

$$\text{Hence } \text{dec}_\gamma(\mathcal{M}_{\text{MAB-SN}}, \hat{M}) \lesssim \mathbb{I}(\gamma \leq \frac{A}{4})$$

With Proposition 1,

$$\text{Reg} \lesssim \mathbb{I}(\gamma \leq \frac{A}{4}) \cdot T + \gamma \cdot \text{Est}_H$$

Let $\gamma = A$, we have

$$\text{Reg} \lesssim A \cdot \text{Est}_H.$$

Sec 23 Some structural Properties of DEC

Without proofs, we list a few structural properties of DEC, which are useful in practice for computing the DEC of specific model class.

Property 1 (Proposition 41 in note, Square loss is sufficient for structural bandit problem)

Consider any structural bandit problem with decision space Π , function class $\mathcal{F} \subseteq (\Pi \rightarrow [0, 1])$, and $\mathcal{O} = \{\phi\}$. ^(no obsers) Denote

$$\mathcal{M}_{\mathcal{F}} = \{M \mid f^M \in \mathcal{F}, M(\pi) \text{ is } 1\text{-sub-Gaussian} \forall \pi\}$$

$$\text{Let } \text{dec}_{\gamma}^{\text{sq}}(\mathcal{F}, \hat{f}) = \inf_{P \in \mathcal{O}(\Pi)} \sup_{\mathcal{F}} \mathbb{E}_{\pi \sim P} [f(\pi_{\mathcal{F}}) - f(\pi) - \gamma(f(\pi) - \hat{f}(\pi))^2]$$

$$\text{Then, } \text{dec}_{c_1 \gamma}^{\text{sq}}(\mathcal{F}) \leq \text{dec}_{\gamma}(\mathcal{M}_{\mathcal{F}}) \leq \text{dec}_{c_2 \gamma}^{\text{sq}}(\mathcal{F})$$

for $G, G' \geq 0$.

Property 2 (Proposition 42 in Note, Filtrating irrelevant information)

Adding observations that are irrelevant to the model does not change the DEC!

Consider \mathcal{M} with observation space \mathcal{O}_1 , and a class of conditional dist \mathcal{D} over another observation space \mathcal{O}_2 , where $\forall D \in \mathcal{D}$ has $D(\pi) \in \Delta(\mathcal{O}_2)$. For $M \in \mathcal{M}$, and $D \in \mathcal{D}$, let $(M \otimes D)(\pi)$ be the model that given $\pi \in \Pi$, $(r_1, o_1) \sim M(\pi)$, $o_2 \sim D(\pi)$. then $(r, (o_1, o_2))$ is obtained. Set $\mathcal{M} \otimes \mathcal{D} = \{M \otimes D : M \in \mathcal{M}, D \in \mathcal{D}\}$.
Then for $\forall \hat{M} \in \mathcal{M}, \hat{D} \in \mathcal{D}$
$$\text{dec}_g(\mathcal{M} \otimes \mathcal{D}, \hat{M} \otimes \hat{D}) = \text{dec}_g(\mathcal{M}, \hat{M})$$

Property 3 (Proposition 43 in Note, Data Processing)

Passing observations through a channel never reduce

DEC. Consider \mathcal{M} with \mathcal{O} . Let

$$\rho: \mathcal{O} \rightarrow \mathcal{O}'$$

be a given mapping. Define $\rho \circ M$ as the model that given decision π , sample $(r, o) \sim M(\pi)$ or $(r, o) \sim M(\cdot | \pi)$, then observer/provides $(r, \rho(o))$. Let

$$\rho \circ \mathcal{M} = \{ \rho \circ M \mid M \in \mathcal{M} \}$$

For all $\hat{M} \in \mathcal{M}$. we have

$$\text{dec}_g(\mathcal{M}, \hat{M}) \leq \text{dec}_g(\rho \circ \mathcal{M}, \rho \circ \hat{M})$$

This is an immediate consequence of the data processing inequality for Hellinger-distance that

$$\begin{aligned} D_H^2((\rho \circ M)(\pi), (\rho \circ \hat{M})(\pi)) \\ \leq D_H^2(M(\pi), \hat{M}(\pi)). \end{aligned}$$

Sec 2.4 Online Estimation with D_H^2

Estimation of model M is more challenging compared to regression problem, such as estimating the reward function. As we have seen before, estimating M^* w.r.t the Hellinger distance can be solved using online conditional density estimation with log loss.

Given (π^t, r^t, o^t) , log loss for M is

$$l_{\log}^t(M) = \log \left(\frac{1}{m^M(r^t, o^t | \pi^t)} \right)$$

where $m^M(\cdot | \pi)$ is the conditional density for (r.o) under model M . Define

$$\text{Reg}_{\text{KL}} = \sum_{t=1}^T l_{\log}^t(\hat{M}^t) - \inf_{M \in \mathcal{M}} \sum_{t=1}^T l_{\log}^t(M)$$

Then, a bound on the log-loss regret yields an immediate bound on the Hellinger estimation error, as follows.

Lemma 1 (21 in lecture): For any online estimation algorithm such as averaged weights, whenever Assumption 2 holds, $\mathbb{E}[\text{Reg}_{KL}] \geq \mathbb{E}\left[\sum_{t=1}^T D_{KL}(M^*(\pi^t) \parallel \hat{M}^t(\pi^t))\right]$

So that $\mathbb{E}[\text{Est}_H] \leq \mathbb{E}[\text{Reg}_{KL}] \quad (\#)$

Also, $\forall \delta \in (0, 1)$. W.p. at least $1-\delta$

$$\text{Est}_H \leq \text{Reg}_{KL} + \underbrace{2\log(1/\delta)}_{\text{red}} \quad (\#\#)$$

Proof: By assumption, Does not scale with T !
due to Mgf of Log-loss.

$$\sum_{t=1}^T \ell_{\log}^*(\hat{M}^t) - \sum_{t=1}^T \ell_{\log}^t(M^*) \leq \text{Reg}_{KL}$$

So by Assumption 2 that $M^* \in \mathcal{M}$.

$$\sum_{t=1}^T \mathbb{E}[D_{KL}(M^*(\pi^t) \parallel \hat{M}^t(\pi^t))] \leq \mathbb{E}[\text{Reg}_{KL}]$$

By Definition of Est_H , and fact that $D_H^2(\mathbb{P}, \mathbb{Q}) \leq D_{KL}(\mathbb{P} \parallel \mathbb{Q})$, $(\#)$ follows.

To prove $(\# \#)$, we employ the tail bound for martingales (Review: \forall real-valued r.v. $(X_t)_{t \leq T}$, adapted to Filtration $(\mathcal{F}_t)_{t \leq T}$, w.p. $1-\delta$, $\forall T' \leq T$

$$\sum_{t=1}^{T'} X_t \leq \sum_{t=1}^{T'} \log [E_{t-1}(e^{X_t})] + \log(1/\delta)$$

Define $Z_t = \frac{1}{2} (\ell_{\log}^t(\hat{M}^t) - \ell_{\log}^t(M^*))$.

Applying tail bound for martingales to $(-Z_t)_{t \leq T}$

w.p. at least $1-\delta$,

$$\begin{aligned} \sum_{t=1}^T -\log(E_{t-1}(e^{-Z_t})) &\leq \sum_{t=1}^T Z_t + \log(1/\delta) \\ &= \frac{1}{2} \sum_{t=1}^T (\ell_{\log}^t(\hat{M}^t) - \ell_{\log}^t(M^*)) + \log(1/\delta) \end{aligned}$$

Fix t , define

$$z^t = (r^t, o^t),$$

let $\nu(\cdot | \pi)$ be any conditional dominating measure for $m^{\hat{M}^t}$ and m^{M^*} . Notice

$$E_{t-1}(e^{-Z_t} | \pi^t) = E_{t-1} \left[\sqrt{\frac{m^{\hat{M}^t}(z^t | \pi^t)}{m^{M^*}(z^t | \pi^t)}} \mid \pi^t \right]$$

$$\begin{aligned}
&= \int m^{M^*}(z | \pi^t) \sqrt{\frac{m^{\hat{M}^t}(z | \pi^t)}{m^{M^*}(z | \pi^t)}} \nu(dz | \pi^t) \\
&= \int \sqrt{m^{M^*}(z | \pi^t) m^{\hat{M}^t}(z | \pi^t)} \nu(dz | \pi^t) \\
&= 1 - \frac{1}{2} D_H^2(M^*(\pi^t), \hat{M}^t(\pi^t))
\end{aligned}$$

Hence, $\mathbb{E}_{\mathcal{G}_1}(e^{-2\epsilon}) = 1 - \frac{1}{2} \mathbb{E}_{\mathcal{G}_1} [D_H^2(M^*(\pi^t), \hat{M}^t(\pi^t))]$

and by $-\log(1-x) \geq x$ for $x \in [0, 1]$,

$$\begin{aligned}
&\frac{1}{2} \sum_{t=1}^T \mathbb{E}_{\mathcal{G}_1} [D_H^2(M^*(\pi^t), \hat{M}^t(\pi^t))] \\
&\leq \frac{1}{2} \sum_{t=1}^T (\ell_{\log}^t(\hat{M}^t) - \ell_{\log}^t(M^*)) + \log(1/2)
\end{aligned}$$

#

Pink: Lemma 1 is useful as regret minimization.

w.r.t. log loss is well-studied, such as averaged.

Weights algorithm which admits

$$\text{Reg}_{KL} \leq \log |M| \text{ for finite class } M.$$

Also. for linear model where $m^M(r, o | \pi) = \langle \phi(r, o, \pi), \theta \rangle$,
for some feature map $\phi \in \mathbb{R}^d$.

$$\text{Reg}_{KL} = O(d \log T)$$

Sec 3. Optimality for General Decision Making — DEC: lower bound on Regret

Classical question: for a given class of models \mathcal{M} .
What is the best regret that can be achieved
by ANY algorithm?

Answer: Minimax optimality — for a model class
 \mathcal{M} . define minimax regret as

$$\mathcal{R}(\mathcal{M}, T) = \inf_{P^1, \dots, P^T} \sup_{M^* \in \mathcal{M}} \mathbb{E}^{M^*, P} [\text{Reg}(T)]$$

- Where:
- i) $p^t = p^t(\cdot | \mathcal{H}^{t-1})$ is the algorithm for step t
as a function of history \mathcal{H}^{t-1} .
 - ii) $\text{Reg}(T)$ makes its dependence on T explicitly.

An algorithm is minimax optimal if it achieves $M(M, T)$ up to a constant free from M and T .

Sec 3.1. Constrained DEC

How to lower bound the minimax regret for any model class M in terms of DEC for M ?

Working on "constrained DEC" instead of $\text{dec}_\varepsilon(M)$ in Proposition 1. which is called the offset DEC. Here for $\varepsilon > 0$, Constrained DEC is defined by

$$\text{dec}_\varepsilon^c(M, \hat{M}) = \inf_{P \in \Delta(\Pi)} \sup_{M \in M} \left\{ \mathbb{E}_{\pi \sim P} [f^M(\pi_M) - f^M(\pi)] \mid \mathbb{E}_{\pi \sim P} [D_H^2(M(\pi), \hat{M}(\pi))] \leq \varepsilon^2 \right\}$$

$$\text{where } \text{dec}_\varepsilon^c(M) = \sup_{\hat{M} \in \text{co}(M)} \text{dec}_\varepsilon^c(M \cup \hat{M}, \hat{M})$$

Remark: Similar to $\text{dec}_\varepsilon(M)$, instead of subtracting the information gain due to observations, $\text{dec}_\varepsilon^c(M)$ puts a

hard constraint on the information gain. Both of them bias the max learner/player towards model where the gain is small.

offset / traditional DEC can be viewed as a Lagrangian relaxation of DEC with constraints. and

$$\begin{aligned} \text{dec}_{\varepsilon}^c(\mathcal{M}, \hat{\mathcal{M}}) &= \inf_{P \in \Delta(\Pi)} \sup_{M \in \mathcal{M}} \left\{ \mathbb{E}_{\pi \sim P} [f^M(\pi_M) - f^M(\pi)] \right. \\ &\quad \left. \mid \mathbb{E}_{\pi \sim P} [\mathcal{D}_H^2(M(\pi), \hat{M}(\pi))] \leq \varepsilon^2 \right\} \\ &= \inf_{P \in \Delta(\Pi)} \sup_{M \in \mathcal{M}} \inf_{\gamma \geq 0} \left\{ \mathbb{E}_{\pi \sim P} [f^M(\pi_M) - f^M(\pi)] \right. \\ &\quad \left. - \gamma (\mathbb{E}_{\pi \sim P} [\mathcal{D}_H^2(M(\pi), \hat{M}(\pi))] - \varepsilon^2) \right\} \quad \text{V.O.} \end{aligned}$$

$$\leq \inf_{\gamma \geq 0} \inf_{P \in \Delta(\Pi)} \sup_{M \in \mathcal{M}} \left\{ \mathbb{E}_{\pi \sim P} [f^M(\pi_M) - f^M(\pi)] \right. \\ \left. - \gamma (\mathbb{E}_{\pi \sim P} [\mathcal{D}_H^2(M(\pi), \hat{M}(\pi))] - \varepsilon^2) \right\} \quad \text{V.O.}$$

$$= \inf_{\gamma \geq 0} \left\{ \text{dec}_{\gamma}(\mathcal{M}, \hat{\mathcal{M}}) + \gamma \varepsilon^2 \right\} \quad \text{V.O.}$$

It is easy to see that $\text{dec}_{\gamma}(\mathcal{M}) \leq \text{dec}_{\gamma^{-1/2}}^c(\mathcal{M})$

Pink: Some classes the constrained DEC is meaningfully smaller than the offset DEC. However, if we restrict to a "localized" sub-class of models that are not "too far" from \hat{M} , we may have

Proposition 2 (26 in Note): Given a model \hat{M} and parameter α , define

$$M_{\alpha}(\hat{M}) = \{M \in \mathcal{M} : f^{\hat{M}}(\pi_{\hat{M}}) \geq f^M(\pi_M) - \alpha\}$$

For all $\varepsilon > 0$ and $\gamma \geq C_1/\varepsilon$,

$$\text{dec}_{\varepsilon}^C(\mathcal{M}) \leq C_3 \sup_{\gamma \geq C_1/\varepsilon} \sup_{\hat{M} \in \mathcal{C}_0(\mathcal{M})} \text{dec}_{\gamma} (M_{\alpha(\varepsilon, \gamma)}(\hat{M}), \hat{M})$$

with $\alpha(\varepsilon, \gamma) = C_2 \cdot \gamma \varepsilon^2$.

(The lengthy proof is referred to Foster, Golowich & Han, 2023)

key message of Proposition 2 is that for well-behaved model class such as multi-armed bandits, linear bandits,

$$\text{dec}_Y(M_{\alpha(\varepsilon, r)}(\hat{M}), \hat{M}) \approx \text{dec}_Y(M, \hat{M})$$

Whenever $\text{dec}_Y(M, \hat{M}) \approx r\varepsilon^2$. i.e., localization does not change the complexity. So. lower bound in terms of constrained DEC immediately implies that in terms of the offset DEC. (Though refined F2D may lead to tighter upper bound for some cases).

Sec 3.2 Lower Bound

Proposition 3 (20 in Note, DEC lower bound): Let

$$\underline{\varepsilon}_T = \frac{c}{\sqrt{T}} \text{ for } c > 0 \text{ sufficiently small. For all } T$$

$$\text{s.t.} \quad \text{dec}_{\underline{\varepsilon}_T}^c(M) \geq \log \underline{\varepsilon}_T$$

for any algorithm. $\exists M \in \mathcal{M}$ that

$$\mathbb{E}[\text{Reg}(T)] \geq \text{dec}_{\underline{\varepsilon}_T}^c(M) \cdot T.$$

Remark: (1) For any algorithm & model class \mathcal{M} . the optimal.

regret must scale with the constrained DEC in the worst-case. For example, by Example 4. (multi-armed bandit with A actions).

$$\text{dec}_\gamma(\mathcal{M}, \hat{\mathbf{M}}) \lesssim \frac{A}{\gamma}$$

By results in Sec 3.1

$$\text{dec}_\varepsilon^c(\mathcal{M}, \hat{\mathbf{M}}) = \inf_{\gamma} \{ \text{dec}_\gamma(\mathcal{M}, \hat{\mathbf{M}}) + \gamma \varepsilon^2 \} \quad \forall \varepsilon > 0$$

so $\text{dec}_\varepsilon^c(\mathcal{M}) \propto \varepsilon \sqrt{A}$. Then:

$$\mathbb{E}[\text{Reg}(T)] \gtrsim \sqrt{AT} \quad \text{for } \varepsilon = \frac{c}{\sqrt{T}}.$$

② Combining Propositions 2 & 3, we have

Corollary 1 (1 in Moe, lower bound based on localized offset DEC): Fix $T \in \mathbb{N}$. for any algorithm.

there exists model $M \in \mathcal{M}$ for which

$$\mathbb{E}[\text{Reg}(T)] \gtrsim \sup_{\gamma \geq \sqrt{T}} \sup_{\hat{\mathbf{M}} \in \mathcal{CO}(\mathcal{M})} \text{dec}_\gamma(\mathcal{M}_{\alpha(T, \gamma)}^{(\hat{\mathbf{M}})}, \hat{\mathbf{M}})$$

with $\alpha(T, \gamma) = c \cdot \gamma / T$.

③ In Foster, Golowich & Han (2023), the authors design an algorithm based on a refined variant of E2D, s.t. the upper bound on regret is based on the constrained DEC

Proposition 4 (29 in Note): For a finite class \mathcal{M} .

$\bar{\epsilon}_T = C \cdot \sqrt{\frac{\log(|\mathcal{M}|/\delta)}{T}}$ with sufficiently small C . Under technical conditions, \exists an algorithm s.t.

$$\mathbb{E}[\text{Reg}(T)] \lesssim \text{dec}_{\bar{\epsilon}_T}^C(\mathcal{M}) \cdot T.$$

w. p. at least $1-\delta$.

(Though there exists a $\log |\mathcal{M}|$ gap, for class with finite $\log |\mathcal{M}|$, $\text{dec}_{\bar{\epsilon}}^C$ is necessary & sufficient to lower bound regret)

Proof of Proposition 3: Basic idea of establishing any lower bound is similar: finding a pair of models M and \hat{M} s.t.

i) any algorithm achieving low regret must be able to distinguish M and \hat{M}

i) M and \hat{M} are difficult to distinguish statistically
i.e., some information-theoretic difference between them
is small

\Rightarrow Algorithm must have large regret on either
 M or \hat{M} (Similar idea in Hajek 1973)

Some simplifications:

1) $\exists C$ s.t. $D_{KL}(M(\pi) \| M'(\pi)) \leq C D_H^2(M(\pi), M'(\pi))$
for all $M, M' \in \mathcal{M}$ and $\pi \in \Pi$

2) Rather than proving a lower bound scaling with
 $\text{dec}_\varepsilon^C(\mathcal{M}) = \sup_{\hat{M} \in \mathcal{C}(\mathcal{M})} \text{dec}_\varepsilon^C(\mathcal{M} \cup \{\hat{M}\})$, one

focuses on a weaker one that scales with $\sup_{\hat{M} \in \mathcal{M}} \text{dec}_\varepsilon^C(\mathcal{M}, \hat{M})$

Fix T and an algorithm, defined by a sequence of
mappings p^1, \dots, p^T where $p^t = p^t(\cdot | \mathcal{H}^{t-1})$. Let
 \mathbb{P}^M denote the dsn over \mathcal{H}^T for the algorithm when
 M is the true model, and denote \mathbb{E}^M the expectation.

Each p^t is a RV as a fcn of \mathcal{H}^{t-1} , we can consider its expected value under M . For any $M \in \mathcal{M}$,

$$\text{let } \bar{P}_M = \mathbb{E}^M \left[\frac{1}{T} \sum_{t=1}^T p^t \right] \in \Delta(\mathcal{I})$$

be the algorithm's average action dsu when M is the true model. Our goal is to show that we can find a model M for which the algorithm's regret is at least as large as the lower bound

$$\sup_{\gamma \geq \frac{1}{T}} \sup_{\hat{M} \in \mathcal{CO}(\mathcal{M})} \text{dec}_{\alpha(t, \gamma)}(M(\hat{M}), \hat{M})$$

Fix $\varepsilon > 0$ and arbitrary model $\hat{M} \in \mathcal{M}$, set

$$M = \arg \max_{\mathcal{M}} \left\{ \underbrace{\mathbb{E}_{\pi \sim P_{\hat{M}}} [f^M(\pi_M) - f^M(\pi)]}_{\text{dec}_{\varepsilon}^c} \mid \mathbb{E}_{\pi \sim P_{\hat{M}}} [D_{\mathcal{H}}^2(M(\pi), \hat{M}(\pi))] \leq \varepsilon^2 \right\} \quad (\#)$$

Model M should be considered as the "worst-case alternative" to \hat{M} , but only for the algorithm fixed now. Next, we'll show that the algorithm needs to

have large regret on either M or \hat{M} . To this end,
define $g^M(\pi) = f^M(\pi_M) - f^M(\pi)$, we will establish:

① for all models M , (as $\text{Reg}(T) = \sum_{t=1}^T \mathbb{E}_{\pi_t, p_t} \dots$)

$$(1) \quad \frac{1}{T} \mathbb{E}^M [\text{Reg}(T)] = \mathbb{E}_{\pi \sim P_M} [g^M(\pi)]$$

So, to prove lower bound on Reg, we need to show that either $\mathbb{E}_{\pi \sim P_M} [g^M(\pi)]$ or $\mathbb{E}_{\pi \sim P_{\hat{M}}} [g^{\hat{M}}(\pi)]$ is large.

$$(2) \quad \mathbb{E}_{\pi \sim P_{\hat{M}}} [g^M(\pi)] \geq \text{dec}_{\varepsilon}^c(M, \hat{M}) := \Delta \quad (2)$$

by the definition of $\text{dec}_{\varepsilon}^c$ and by the construction of M above in (1), M is the best response to a potentially sub-optimal choice $P_{\hat{M}}$. Then, it remains to fill in the gap that g^M is about M and $P_{\hat{M}}$ is about \hat{M} .

(3) Using the chain rule of KL-divergence. ($P_{\hat{M}}$ and P_M)

$$D_{\text{KL}}(P_{\hat{M}} \| P_M) \geq \mathbb{E}^{\hat{M}} \left[\sum_{t=1}^T \mathbb{E}_{\pi_t, p_t} D_{\text{KL}}(\hat{M}(\pi_t) \| M(\pi_t)) \right]$$

simplification ①

$$\leq C \cdot \mathbb{E}^{\hat{M}} \left[\sum_{t=1}^T \mathbb{E}_{\pi^t \sim p^t} D_H^2(\hat{M}(\pi^t), M(\pi^t)) \right]$$

$$= C \cdot \mathbb{E}_{\pi \sim p_M^{\hat{M}}} [D_H^2(\hat{M}(\pi), M(\pi))]$$

Here, the first equality follows from the chain rule.

Apply the chain rule to sequence $\pi^1, z^1, \dots, \pi^T, z^T$, with

$z^t = (r^t, o^t)$. Then

$$D_{KL}(\mathbb{P}^{\hat{M}} \| \mathbb{P}^M)$$

$$= \mathbb{E}^{\hat{M}} \left[\sum_{t=1}^T D_{KL}(\mathbb{P}^{\hat{M}}(z^t | \mathcal{H}^{t-1}, \pi^t) \| \mathbb{P}^M(z^t | \mathcal{H}^{t-1}, \pi^t)) \right. \\ \left. + D_{KL}(\mathbb{P}^{\hat{M}}(\pi^t | \mathcal{H}^{t-1}) \| \mathbb{P}^M(\pi^t | \mathcal{H}^{t-1})) \right]$$

$$= \mathbb{E}^{\hat{M}} \left[\sum_{t=1}^T D_{KL}(\hat{M}(\pi^t) \| M(\pi^t)) \right]$$

as π^t is free from the model if conditional on \mathcal{H}^{t-1} .

↓ [Review: Chain Rule for KL Divergence:

Let $(\mathcal{X}_1, \mathcal{F}_1), \dots, (\mathcal{X}_n, \mathcal{F}_n)$ be a sequence of measurable spaces. Let $\mathcal{X}^i = \prod_{t=1}^i \mathcal{X}_t$, $\mathcal{F}^i = \bigotimes_{t=1}^i \mathcal{F}_t$

For each i , let $P^i(\cdot|\cdot)$ and $Q^i(\cdot|\cdot)$ be probability kernels from $(\mathcal{X}^{i-1}, \mathcal{F}^{i-1})$ to $(\mathcal{X}^i, \mathcal{F}^i)$.

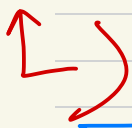
Let P and Q be the dsns of X_1, \dots, X_n under $X_i \sim P^i(\cdot | X_{1:i-1})$ and $X_i \sim Q^i(\cdot | X_{1:i-1})$, respectively.

Then, it holds that:

$$D_{KL}(P \parallel Q) = \mathbb{E}_P \left[\sum_{i=1}^n D_{KL}(P^i(\cdot | X_{1:i-1}) \parallel Q^i(\cdot | X_{1:i-1})) \right]$$

An easier way to understand chain rule of KL divergence is considering $P(X, Y)$, $Q(X, Y)$. then

$$\begin{aligned} D_{KL}(P(X, Y) \parallel Q(X, Y)) &= \mathbb{E}_P [D_{KL}(P(Y|X=x) \parallel Q(Y|X=x))] \\ &\quad + D_{KL}(P(X) \parallel Q(X)) \\ &= D_{KL}(P(Y|X) \parallel Q(Y|X)) + D_{KL}(P(X) \parallel Q(X)) \end{aligned}$$



Now, we can choose $\Sigma = c/\sqrt{n}$ for $c > 0$ sufficiently small s.t. (in dec_{Σ}^c)

$$(3) \quad TV^2(P^{\hat{M}}, P^M) \leq D_{KL}(\hat{P}^M \parallel P^M) \leq \frac{1}{100}.$$

That is, with non-trivial probability, the algorithm fixed here fails to separate M and \hat{M} .

④ Finally, as $r \in [0, 1]$,

$$\mathbb{E}_{\pi \sim P_{\hat{M}}} [f^M(\pi) - f^{\hat{M}}(\pi)]$$

$$\leq \mathbb{E}_{\pi \sim P_{\hat{M}}} [TV(M(\pi), \hat{M}(\pi))]$$

(4)

$$\leq \sqrt{\mathbb{E}_{\pi \sim P_{\hat{M}}} [D_H^2(M(\pi), \hat{M}(\pi))]}$$

$$\leq \epsilon.$$

Step 1. Define $G_M = \{\pi \in \Pi \mid g^M(\pi) \leq \Delta/10\}$, where $\Delta = \text{dec}_{\epsilon}^c(M, \hat{M})$. Notice that

$$\mathbb{E}_{\pi \sim P_M} [g^M(\pi)] \underset{\text{(Markov)}}{\geq} \frac{\Delta}{10} P_M(\pi \notin G_M)$$

$$\geq \frac{\Delta}{10} (P_{\hat{M}}(\pi \notin G_M) - TV(P_M, P_{\hat{M}}))$$

(5)

$$\geq \frac{\Delta}{10} (P_{\hat{M}}(\pi \notin G_M) - 1/10)$$

as $TV(P_M, P_{\hat{M}}) \leq TV(P^M, P^{\hat{M}}) \leq 1/10$ by

above choice of ϵ in (3) and data-processing inequality

(Review: Data processing inequality in Notes: $D_H^2(M(\pi), \hat{M}(\pi))$)

$\geq D_H^2((P \circ M)(\pi), (P \circ \hat{M})(\pi))$, fixed P ; also held for TV
e.g. Proposition 43)

Next, assume that

$$(b) \quad \mathbb{E}_{\pi \sim P_{\hat{M}}} [g^{\hat{M}}(\pi)] \leq \frac{\Delta}{10},$$

otherwise, we are done by (1). Our goal is to show that under (b), $P_{\hat{M}}(\pi \notin G_M) \geq 1/2$, which will imply that $\mathbb{E}_{\pi \sim P_M} [g^M(\pi)] \geq \Delta$ by (5).

Step 2: Adding (b) and (2)

$$f^M(\pi_M) - f^{\hat{M}}(\pi_{\hat{M}})$$

$$= \mathbb{E}_{\pi \sim P_{\hat{M}}} [f^M(\pi_M) - f^{\hat{M}}(\pi_{\hat{M}})]$$

$$= \mathbb{E}_{\pi \sim P_{\hat{M}}} [f^M(\pi_M) - f^M(\pi) + f^M(\pi) - f^{\hat{M}}(\pi) + f^{\hat{M}}(\pi) - f^{\hat{M}}(\pi_{\hat{M}})]$$

$$= \mathbb{E}_{\pi \sim P_{\hat{M}}} [\underbrace{f^M(\pi_M) - f^M(\pi)}_{g^M(\pi)} - \underbrace{f^{\hat{M}}(\pi_{\hat{M}}) - f^{\hat{M}}(\pi)}_{g^{\hat{M}}(\pi)}] + \mathbb{E}_{\pi \sim P_{\hat{M}}} [f^M(\pi) - f^{\hat{M}}(\pi)]$$

$$\begin{aligned}
&\geq \mathbb{E}_{\pi \sim p_{\hat{M}}} [g^M(\pi) - g^{\hat{M}}(\pi)] \\
&\quad - \mathbb{E}_{\pi \sim p_{\hat{M}}} [|f^M(\pi) - f^{\hat{M}}(\pi)|] \\
&\geq \frac{9}{10} \Delta - \mathbb{E}_{\pi \sim p_{\hat{M}}} [|f^M(\pi) - f^{\hat{M}}(\pi)|].
\end{aligned}$$

By (4), $\mathbb{E}_{\pi \sim p_{\hat{M}}} [|f^M(\pi) - f^{\hat{M}}(\pi)|] \leq \varepsilon$, so

$$f^M(\pi_m) - f^{\hat{M}}(\pi_{\hat{m}}) \geq \frac{9}{10} \Delta - \varepsilon$$

As long as $\varepsilon \leq \frac{1}{10} \Delta$, which is by the condition that $\text{dec}_{\varepsilon_T}^c(\mathcal{M}) \geq 10 \varepsilon_T$,

$$f^M(\pi_m) - f^{\hat{M}}(\pi_{\hat{m}}) \geq \frac{4}{5} \Delta.$$

Step 3: Observe that if $\pi \in G_m$,

$$\begin{aligned}
|f^M(\pi) - f^{\hat{M}}(\pi)|_+ &\geq |f^M(\pi_m) - f^{\hat{M}}(\pi) - \frac{\Delta}{10}|_+ \\
&\geq |f^M(\pi_m) - f^{\hat{M}}(\pi_{\hat{m}}) - \frac{\Delta}{10}|_+
\end{aligned}$$

$$\geq \frac{7}{10} \Delta$$

by step 2. Using (4) again,

$$\begin{aligned} \varepsilon &\geq \mathbb{E}_{\pi \sim P_M^{\hat{M}}} [|f^M(\pi) - f^{\hat{M}}(\pi)|_+] \\ &\geq \frac{7}{10} \Delta \cdot P_{\hat{M}}(\pi \in G_M) \quad (\text{Markov}) \end{aligned}$$

Since $\varepsilon \leq \Delta/10$ by assumption, we have

$$\frac{\Delta}{10} \geq \frac{7}{10} \Delta \cdot P_{\hat{M}}(\pi \in G_M).$$

$\therefore P_{\hat{M}}(\pi \in G_M) \leq 1/7$. Combining this with (5) gives

$$\begin{aligned} \frac{1}{T} \mathbb{E}^M [\text{Reg}(T)] &= \mathbb{E}_{\pi \sim P_M} [g^M(\pi)] \\ &\geq \frac{\Delta}{10} \left(1 - \frac{1}{7} - \frac{1}{10} \right) \\ &\geq \frac{\Delta}{20} \end{aligned}$$

Finally, notice that the choice of $\hat{M} \in \mathcal{M}$ is arbitrary, we are free to choose \hat{M} to maximize $\text{dec}_\varepsilon^c(\mathcal{M}, \hat{M})$ #.

Sec 3.3 Examples

Consider a few concrete model classes to demonstrate.

(Example 4 continued, Multi-armed bandit with Gaussian reward) What is the constrained DEC for this case?

Set $\hat{M}(\pi) = N(\frac{1}{\Sigma}, 1)$, let $\{M_1, \dots, M_A\} \subseteq \mathcal{M}$ be a sub-family of models that $M_i(\pi) = N(f^{M_i}(\pi), 1)$ where $f^{M_i}(\pi) = \frac{1}{\Sigma} + \Delta I(\pi = i)$ for parameter Δ .

For all i , $\mathbb{E}_{\pi \sim p} [D_H^2(M_i(\pi), \hat{M}(\pi))] \leq \frac{1}{\Sigma} \Delta^2 p_i$

(Hellinger distance for Gaussian is squared difference in means as seen before). and

$$\mathbb{E}_{\pi \sim p} [f^{M_i}(\pi_{M_i}) - f^{M_i}(\pi)] = (1 - p_i) \Delta$$

Then. $\text{dec}_\Sigma^c(\mathcal{M}, \hat{M})$

$$= \inf_{p \in \Delta(\Pi)} \sup_{\mathcal{M}} \left\{ \mathbb{E}_{\pi \sim p} [f^M(\pi_M) - f^M(\pi)] \right\} \left| \mathbb{E}_{\pi \sim p} [D_H^2(M(\pi), \hat{M}(\pi))] \leq \Sigma \right\}$$

$$\geq \inf_{p \in \Delta(\mathcal{I})} \max_i \{ (1 - p(i)) \Delta \mid p(i) \frac{\Delta^2}{2} \leq \varepsilon^2 \}$$

For any p , $\exists i$ s.t. $p(i) \leq 1/A$. Choose $\Delta = \varepsilon \sqrt{2A}$
 then this choice for i satisfies $p(i) \frac{\Delta^2}{2} \leq \varepsilon^2$.

$$\begin{aligned} \text{Hence, } \text{dec}_{\varepsilon}^c(M, \hat{M}) &\geq (1 - p(i)) \Delta \\ &\geq \varepsilon \sqrt{A/2} \end{aligned}$$

as $1 - p(i) \geq 1/2$. By Proposition 3,

$$\mathbb{E}[\text{Reg}] \geq \tilde{\Omega}(\sqrt{AT})$$

Remark: One can generalize this to any M that "embeds" the multi-armed bandit problem in a certain sense

Proposition 5: Given reference model \hat{M} , suppose that a class M contains a sub-class $\{M_1, \dots, M_N\}$, and a collection of decisions π_1, \dots, π_N that for each i

- (i) $D_H^2(M_i(\pi), \hat{M}(\pi)) \leq \beta^2 I(\pi = \pi_i)$
- (ii) $f^{M_i}(\pi_{M_i}) - f^{M_i}(\pi) \geq \alpha I(\pi \neq \pi_i)$

Then. $\text{dec}_\varepsilon^c(M, \hat{M}) \gtrsim \alpha \cdot I(\varepsilon \geq \beta/\sqrt{N})$

Notice conditions (i) & (ii) are exactly the two basic
techniques to derive any lower bound.

(Example 5 Continued, Bandits with Structured noise)

Recall that $M = \{M_1, \dots, M_A\}$ with $M_i(\pi) = N(\frac{1}{2}, 1)$
 $I(i \neq \pi) + \text{Ber}(3/4) I(i = \pi)$. If we consider reference
model $\hat{M}(\pi) = N(\frac{1}{2}, 1)$. then by Proposition 5
above, $\alpha = 1/4$, and $\beta^2 = 2$.

$$(f^{M_i}(\pi_{M_i}) - f^{M_i}(\pi) \geq 1/4 \text{ if } \pi \neq i)$$

Thus, $\text{dec}_\varepsilon^c(M_{\text{MAB-SN}}) \gtrsim I(\varepsilon \geq \sqrt{2/A})$, yielding

$$\mathbb{E}(\text{Reg}) \gtrsim O(A)$$

by Proposition 3.

(Example 6. Linear Bandit and Lipschitz bandit)

$$\text{Linear: } \mathcal{F} = \{ \pi \rightarrow \langle \theta, \phi(\pi) \rangle \mid \theta \in \Theta \}$$

$\mathcal{D} \subseteq \mathcal{B}_2^d(1)$. $\phi: \mathcal{I} \rightarrow \mathbb{R}^d$ is known.
feature map

\mathcal{M} is set of all reward dsns with $f^M \in \mathcal{F}$.
and 1-sub-Gaussian noise. Then

$$\text{dec}_{\varepsilon}^c(\mathcal{M}) \gtrsim \varepsilon \sqrt{d}$$

$$\text{and } \mathbb{E}(\text{Reg}) \gtrsim \sqrt{dT}$$

Levy: $\mathcal{F} = \{f: \mathcal{I} \rightarrow [0,1] \mid f \text{ is } 1\text{-Lip wrt } p\}$
 \mathcal{I} is a metric space with metric p .

\mathcal{M} is set of all reward dsns with $f^M \in \mathcal{F}$.
and 1-sub-Gaussian noise. Assume

covering s.t. $N_p(\mathcal{I}, \varepsilon) \geq 1/\varepsilon^d$ for $d > 0$

$$\text{Then, } \text{dec}_{\varepsilon}^c(\mathcal{M}) \gtrsim \varepsilon^{\frac{2}{d+2}}$$

$$\text{and } \mathbb{E}(\text{Reg}) \gtrsim T^{\frac{d+1}{d+2}}$$