# Regret Analysis for MEX

BY JACK YANSONG LI

University of Illinois Chicago

*Email:* `yli340@uic.edu`

## 1 Regret Analysis for Finite Hypothesis Set

Through this note, we denote $\pi^k$ as the guesses policy of player 2 in episode $k$ and $\pi^*$ as the true policy of player 2. For each player 2's policy $\pi$, the set of all best response policies is denoted as $\mathrm{BR}(\pi)$, i.e.,

$$\mathrm{BR}(\pi) = \operatorname*{argmax}_{\mu \in \mathcal{U}} V(\mu, \pi),$$

where $\mathcal{U}$ is the set of all possible policies for player 1. The hypothesis set of all possible policies of player 2 is denoted as $\mathcal{H}$.

### 1.1 Oracle and Type Number

For any player 2's policy $\pi$, we assume the existence of an oracle that can return a best response from $\mathrm{BR}(\pi)$.

**Definition 1.** *(Oracle) A best response oracle $\psi$ refers to a function that, upon receiving policies as input, yields a best response as its output, i.e., $\psi$ is a function $\psi \colon \mathcal{H} \to \mathcal{U}$ such that*

$$\psi(\pi) \in \mathrm{BR}(\pi).$$

With the definition of an oracle, we can categorize policies within a hypothesis set into various types.

**Definition 2.** *($\psi$-type) We call two policies $\pi$ and $\pi'$ to be of the same type under oracle $\psi$ if for any $\mu = \psi(\pi)$ and $\mu' = \psi(\pi')$ we have*

$$V(\mu, \pi) = V(\mu', \pi').$$

*The relationship is denoted as $\pi \overset{\psi}{\sim} \pi'$. On the contrary, two policies $\pi$ and $\pi'$ not of the same type under oracle $\psi$ is denoted as $\pi \overset{\psi}{\nsim} \pi'$.*

**Definition 3.** *We call a set of policies $\Pi$ be type-independent under oracle $\psi$ if for all $\pi \in \Pi$ and $\pi' \in \Pi$ such that $\pi \neq \pi'$, we have $\pi \overset{\psi}{\nsim} \pi'$.*

The $\psi$-type characterization gives rise to a measurement of quantity for the set of policies $\mathcal{H}$, denoted by $n^\psi(\mathcal{H})$.

**Definition 4.** *Given a hypothesis set $\mathcal{H}$, the type number $n^\psi(\mathcal{H})$ under oracle $\psi$ is defined as the size of a largest type-independent subset of $\mathcal{H}$, i.e.,*

$$n^\psi(\mathcal{H}) = \max |\Pi|,$$

*where $\Pi \subset \mathcal{H}$ and $\Pi$ is type-independent under oracle $\psi$.*

## 1.2 Regret Analysis

In this subsection, we restrict our discussion to cases where the cardinality of the hypothesis set is finite, i.e., $|\mathcal{H}| < \infty$. This condition is emphasized through the notation $\mathcal{H}_{\text{fin}}$. We also assume that the realization assumption holds, i.e., $\pi^* \in \mathcal{H}_{\text{fin}}$.

**Theorem 5.** *Given an MDP with generalized eluder coefficient $d_{\text{GEC}}(\cdot)$ and a finite hypothesis class $\mathcal{H}_{\text{fin}}$ with $\pi^* \in \mathcal{H}_{\text{fin}}$, by setting*

$$\eta = \sqrt{\frac{d_{\text{GEC}}(1/\sqrt{HK})}{\log(Hn^\psi(\mathcal{H}_{\text{fin}})/\delta) \cdot HK}},$$

*the regret of the MEX algorithm applying on $\mathcal{H}_{\text{fin}}$ with oracle $\psi$ after $K$ episodes is upper bounded by, with probability at least $1 - \delta$,*

$$\text{Regret}(K) \lesssim \sqrt{d_{\text{GEC}}(1/\sqrt{HK}) \cdot \log(Hn^\psi(\mathcal{H}_{\text{fin}})/\delta) \cdot HK}.$$

**Proof.** See Appendix 3.2 □

The sole term pertaining to the size of the hypothesis set is $n^\psi(\mathcal{H}_{\text{fin}})$. Consequently, the magnitude of regret is solely influenced by the type number associated with a hypothesis set, as opposed to the cardinality of the hypothesis set. This phenomenon occurs because policies that are categorized under the same type by policy $\psi$ yield identical rewards when implemented in the MEX algorithm.

The type number $n^\psi(\mathcal{H}_{\text{fin}})$ depends on the choice of the oracle $\psi$, which makes it hard to verify when the explicit form of $\psi$ is not given. However, we can introduce a stronger notion of type and verify the upper bound of $n^\psi(\mathcal{H}_{\text{fin}})$.

**Definition 6.** *(Strong) We call two policies $\pi$ and $\pi'$ to be of the same s-type if*

$$V(\mu, \pi) = V(\mu, \pi') = V(\mu', \pi) = V(\mu', \pi')$$

*for all $\mu \in \text{BR}(\pi)$ and $\mu' \in \text{BR}(\pi')$. The relationship are denoted as $\pi \overset{s}{\sim} \pi'$.*

Similar to the definition of type number under oracle $\psi$, we can define strong type number $n_{\text{stype}}(\mathcal{H})$.

**Lemma 7.** $n^\psi(\mathcal{H}) \leq n_{\text{stype}}(\mathcal{H})$ *for all $\psi$ be a best response oracle.*

## 2 Regret Analysis for Infinite Hypothesis Set

In this subsection, we discuss the cases where the cardinality of the hypothesis set is infinite, i.e., $|\mathcal{H}| = \infty$. This condition is emphasized through the notation $\mathcal{H}_{\text{inf}}$. We keep assume that the realization assumption holds, i.e., $\pi^* \in \mathcal{H}_{\text{inf}}$.

### 2.1 Approximate an Infinite Hypothesis Set by a Finite Hypothesis Set

A direct approach to handling an infinite hypothesis set is to approximate it as a finite hypothesis set. First, we outline what makes a good approximation.

**Definition 8.** *($\varepsilon_\psi$-optimal approximation) A finite hypothesis set $\mathcal{H}_{\text{fin}}$ is called an $\varepsilon_\psi$-optimal approximation of $\mathcal{H}_{\text{inf}}$ if*

$$\min_{\pi \in \mathcal{H}} |V(\psi(\pi), \pi) - V(\psi(\pi^*), \pi^*)| \leq \varepsilon_\psi$$

We denote $\pi_{\text{det}}^*$ as $\varepsilon_\psi$-optimal policy defined as

$$\pi_{\text{det}}^* \triangleq \underset{\pi \in \mathcal{H}}{\operatorname{argmin}} |V(\psi(\pi), \pi) - V(\psi(\pi^*), \pi^*)|.$$

In the following, we list some examples of the $\varepsilon_\psi$-optimal approximation. These examples are established based on the following lemma.

**Lemma 9.** *For any best response oracle $\psi$, if $\|\pi - \pi^*\| \le \varepsilon$, then*

$$|V(\psi(\pi), \pi) - V(\psi(\pi^*), \pi^*)| \le L_\psi \varepsilon,$$

*where $L_\psi > 0$ is a constant.*

**Example 10.** Given a finite hypothesis set $\mathcal{H}_{\text{fin}}$. For $\varepsilon > 0$, Define an infinite hypothesis $\mathcal{H}_{\text{inf}}$ as

$$\mathcal{H}_{\text{inf}} \triangleq \{\pi | \|\pi - \pi'\| \le \varepsilon, \pi' \in \mathcal{H}_{\text{fin}}\}.$$

For the constructed $\mathcal{H}_{\text{inf}}$, $\mathcal{H}_{\text{fin}}$ is an $L_\psi \varepsilon$-optimal approximation of $\mathcal{H}_{\text{inf}}$.

**Example 11.** Given a specific neural network structure $\mathcal{N}$, we define $\mathcal{H}_{\text{inf}}$ as the set comprising all neural networks characterized by the set all possible parameters $\Theta$ that is in accordance with the specified structure, formally represented as,

$$\mathcal{H}_{\text{inf}} = \{\pi | \pi \in \mathcal{N}(\theta), \theta \in \Theta\}.$$

We proceed to create a discretization of $\Theta$, denoted as $\hat{\Theta}$. The finite approximation set $\mathcal{H}_{\text{fin}}$ is defined as

$$\mathcal{H}_{\text{fin}} = \{\pi | \pi \in \mathcal{N}(\theta), \theta \in \hat{\Theta}\}.$$

By the choice of discretization interval, we can ensure that $\|\pi - \pi^*\| \le \varepsilon$. Consequently, the discretized set $\mathcal{H}_{\text{fin}}$ serves as $L_\psi \varepsilon$-optimal approximation of $\mathcal{H}_{\text{inf}}$.

## 2.2 Regret Analysis

Now, given an infinite hypothesis set $\mathcal{H}_{\text{inf}}$ with an $\varepsilon_\psi$-optimal approximation set $\mathcal{H}_{\text{fin}}$, we are prepared to execute the MEX algorithm within the confines of $\mathcal{H}_{\text{fin}}$. The regret analysis is given in the following Theorem.

**Theorem 12.** *Given an MDP with generalized eluder coefficient $d_{\text{GEC}}(\cdot)$ and an infinite hypothesis class $\mathcal{H}_{\text{inf}}$ with $\pi^* \in \mathcal{H}_{\text{inf}}$. For any $\varepsilon_\psi$-optimal approximation $\mathcal{H}_{\text{fin}}$ of $\mathcal{H}_{\text{inf}}$, by setting*

$$\eta = \sqrt{\frac{d_{\text{GEC}}(1/\sqrt{HK})}{\log(Hn^\psi(\mathcal{H}_{\text{fin}})/\delta) \cdot HK}},$$

*the regret of the MEX algorithm applying on $\mathcal{H}_{\text{fin}}$ with oracle $\psi$ after $K$ episodes is upper bounded by, with probability at least $1 - \delta$,*

$$\operatorname{Regret}(K) \lesssim \sqrt{d_{\text{GEC}}(1/\sqrt{HK}) \cdot \log(Hn^\psi(\mathcal{H}_{\text{fin}})/\delta) \cdot HK} + K\varepsilon_\psi.$$

3

**Proof.** By the choice of $\pi^k$, we have

$$V(\psi(\pi_{\det}^*), \pi_{\det}^*) - \eta \sum_{h=1}^{H} L_h^{k-1}(\pi_{\det}^*) \leq V(\mu^k, \pi^k) - \eta \sum_{h=1}^{H} L_h^{k-1}(\pi^k)$$

for all $k \in [K]$. By Definition 8,

$$V(\psi(\pi_{\det}^*), \pi_{\det}^*) \geq V(\psi(\pi^*), \pi^*) - \varepsilon_\psi.$$

Thus,

$$V(\psi(\pi^*), \pi^*) - V(\mu^k, \pi^k) \leq \eta \sum_{h=1}^{H} L_h^{k-1}(\pi_{\det}^*) - \eta \sum_{h=1}^{H} L_h^{k-1}(\pi^k) + \varepsilon_\psi.$$

Follow the same procedure in the proof of Theorem 5 leads to the proof. $\qquad\square$

**Remark 13.** The linear term $K\varepsilon_\psi$ cannot be eliminated. Consider the best case where $\pi^k = \pi_{\det}^*$ for all $k \in [K]$. The regret is

$$\begin{aligned}
\mathrm{Regret}(K) &= \sum_{k=1}^{K} V(\psi(\pi_{\det}^*), \pi_{\det}^*) - V(\psi(\pi^*), \pi^*) \\
&= K(V(\psi(\pi_{\det}^*), \pi_{\det}^*) - V(\psi(\pi^*), \pi^*)) \\
&\leq K\varepsilon_\psi.
\end{aligned}$$

# 3  Appendix

## 3.1  Lemmas

**Lemma 14.** *With probability at least $1 - \delta$, for any $(h, k) \in [H] \times [K]$, $\mu^s \in \mathrm{BR}(\pi^s)$, and $\pi \in \Pi$*

$$L_h^{k-1}(\pi^*) - L_h^{k-1}(\pi) \leq -2 \sum_{s=1}^{k-1} \mathbb{E}_{\xi_h \sim \mu^s}[\ell_{\pi^s}(\pi; \xi_h)] + 2\log(H |\Pi| / \delta).$$

**Proof.** Given $\pi \in \mathcal{H}$, we denote the random variable $X_{h,\pi}^k$ as

$$X_{h,\pi}^k = \log\left( \frac{\mathbb{P}_{h,\pi^*}(s_{h+1}^k | s_h^k, a_h^k)}{\mathbb{P}_{h,\pi}(s_{h+1}^k | s_h^k, a_h^k)} \right).$$

Now we define a filtration $\{\mathcal{F}_{h,k}\}_{k=1}^K$ as (B.25) in [1]. Thus we have $X_{h,\pi}^k \in \mathcal{F}_{h,k}$. Therefore, by applying Lemma D.1 in [1], we have that with probability at least $1 - \delta$, for any $(h, k) \in [H] \times [K]$, and $\pi \in \Pi$, we have

$$-\frac{1}{2} \sum_{s=1}^{k-1} X_{h,\pi}^s \leq \sum_{s=1}^{k-1} \log\mathbb{E}\left[ \exp\left\{ -\frac{1}{2} X_{h,\pi}^s \right\} | \mathcal{F}_{h,s-1} \right] + \log(H |\Pi| / \delta). \tag{1}$$

Meanwhile, by (B.27) in [1], for any $\mu^s \in \mathrm{BR}(\pi^s)$, the conditional expectation equals to

$$\mathbb{E}\left[ \exp\left\{ -\frac{1}{2} X_{h,\pi}^s \right\} | \mathcal{F}_{h,s-1} \right] = 1 - \mathbb{E}_{(s_h^s, a_h^s) \sim \mu^s}[D_H(\mathbb{P}_{h,\pi^*}(\cdot | s_h^s, a_h^s) || \mathbb{P}_{h,\pi}(\cdot | s_h^s, a_h^s))]. \tag{2}$$

Denote $\mathbb{E}_{(s_h^s, a_h^s) \sim \mu^s}[D_H(\mathbb{P}_{h, \pi^*}(\cdot \mid s_h^s, a_h^s) \| \mathbb{P}_{h, \pi^s}(\cdot \mid s_h^s, a_h^s))]$ as $\mathbb{E}_{\xi_h \sim \mu^s}[\ell_{\pi^s}(\pi; \xi_h)]$. Using the fact $\log(x) \leq x - 1$ and substituting (2) into (1) finishes the proof. $\qquad\square$

## 3.2 Proof of Theorem 5

**Proof.** We decompose the regret into two terms,

$$
\mathrm{Regret}(K) \triangleq \sum_{k=1}^{K} V(\psi(\pi^*), \pi^*) - V(\psi(\pi^k), \pi^*)
$$

$$
= \underbrace{\sum_{k=1}^{K} V(\psi(\pi^*), \pi^*) - V(\psi(\pi^k), \pi^k)}_{\text{Term (i)}} + \underbrace{\sum_{k=1}^{K} V(\psi(\pi^k), \pi^k) - V(\psi(\pi^k), \pi^*)}_{\text{Term (ii)}}.
$$

**Term (i).** By the choice of $\pi^k$, we have

$$
V(\psi(\pi^*), \pi^*) - \eta \sum_{h=1}^{H} L_h^{k-1}(\pi^*) \leq V(\mu^k, \pi^k) - \eta \sum_{h=1}^{H} L_h^{k-1}(\pi^k)
$$

for all $k \in [K]$. Thus,

$$
V(\psi(\pi^*), \pi^*) - V(\mu^k, \pi^k) \leq \eta \sum_{h=1}^{H} L_h^{k-1}(\pi^*) - \eta \sum_{h=1}^{H} L_h^{k-1}(\pi^k). \tag{3}
$$

for any $\pi^k \overset{\psi}{\sim} \pi^{k'}$, we have

$$
V(\psi(\pi^k), \pi^k) = V(\psi(\pi^{k'}), \pi^{k'}).
$$

Thus, an upper bound for $V(\psi(\pi^*), \pi^*) - V(\mu^k, \pi^k)$ is also an upper bound for $V(\psi(\pi^*), \pi^*) - V(\mu^{k'}, \pi^{k'})$. Applying Lemma 14, we have that with probability at least $1 - \delta$, for any $(h, k) \in [H] \times [K]$, $\mu^s = \psi(\pi^s)$ and $\pi^k \in \mathcal{H}_{\mathrm{fin}}$,

$$
L_h^{k-1}(\pi^*) - L_h^{k-1}(\pi^k) \leq -2 \sum_{s=1}^{k-1} \mathbb{E}_{\xi_h \sim \mu^s}[\ell_{\pi^s}(\pi; \xi_h)] + 2\log(Hn^\psi(\mathcal{H}_{\mathrm{fin}})/\delta).
$$

Substituting the above equation into (3) gives us that with probability at least $1 - \delta$, for any $k \in [K]$, $\mu^s = \psi(\pi^s)$ and $\pi^k \in \mathcal{H}_{\mathrm{fin}}$

$$
V(\psi(\pi^*), \pi^*) - V(\mu^k, \pi^k) \leq -2\eta \sum_{h=1}^{H} \sum_{s=1}^{k-1} \mathbb{E}_{\xi_h \sim \mu^s}[\ell_{\pi^s}(\pi; \xi_h)] + 2H\eta\log(Hn^\psi(\mathcal{H}_{\mathrm{fin}})/\delta).
$$

Summing over $[K]$ gives us

$$
\text{Term (i)} \leq -2\eta \sum_{k=1}^{K} \sum_{h=1}^{H} \sum_{s=1}^{k-1} \mathbb{E}_{\xi_h \sim \mu^s}[\ell_{\pi^s}(\pi; \xi_h)] + 2\eta KH\log(Hn^\psi(\mathcal{H}_{\mathrm{fin}})/\delta).
$$

**Term (ii).** Follow the proof of Theorem 4.4 in [1], we have that for all $\mu^s = \psi(\pi^s)$

$$
\text{Term (ii)} \leq 2\eta \sum_{k=1}^{K} \sum_{h=1}^{H} \sum_{s=1}^{k-1} \mathbb{E}_{\xi_h \sim \mu^s}[\ell_{\pi^s}(\pi; \xi_h)] + \frac{d_{\mathrm{GEC}}(\varepsilon_{\mathrm{conf}})}{8\eta} + \sqrt{d_{\mathrm{GEC}}(\varepsilon_{\mathrm{conf}})HK} + \varepsilon_{\mathrm{conf}}HK.
$$

**Combining Term (i) and Term (ii).**

$$\text{Regret}(K) = \text{Term (i)} + \text{Term (ii)}$$
$$\leq 2\eta K H \log(H n^{\psi}(\mathcal{H}_{\text{fin}})/\delta) + \frac{d_{\text{GEC}}(\varepsilon_{\text{conf}})}{8\eta} + \sqrt{d_{\text{GEC}}(\varepsilon_{\text{conf}}) H K} + \varepsilon_{\text{conf}} H K.$$

Set $\varepsilon_{\text{conf}} = 1/\sqrt{HK}$ and

$$\eta = \sqrt{\frac{d_{\text{GEC}}(1/\sqrt{HK})}{\log(H n^{\psi}(\mathcal{H}_{\text{fin}})/\delta) \cdot HK}}$$

leads to the proof. $\qquad\square$

# Bibliography

[1] Zhihan Liu, Miao Lu, Wei Xiong, Han Zhong, Hao Hu, Shenao Zhang, Sirui Zheng, Zhuoran Yang, and Zhaoran Wang. One Objective to Rule Them All: A Maximization Objective Fusing Estimation and Planning for Exploration. may 2023.