## A Pure Exploration perspective on Game Tree Search

Wouter M. Koolen



Delft, Monday 19<sup>th</sup> November, 2018

#### Menu

- What is Pure Exploration?
- Relation to Reinforcement Learning?
- Why care?
  - Pure Exploration problems occur as sub-problems in RL!
  - Novel/interesting/powerful PE learning algorithms! "Fresh"
- My focus: Pure Exploration Renaissance (2016)
  - Track-and-Stop algorithm for Best Arm Identification
  - 2 BAI-MCTS approach for Game Tree Search
  - Murphy Sampling for Games Trees of "depth 1.5"

- Introduction
- 2 Relation of RL and PE

#### Oure Exploration Intro: Best Arm Identification

- Model
- Sample Complexity Lower Bound
- Algorithms

#### Game Tree Search

- Game Trees of Arbitrary Depth
- Confidence Intervals on Min/Max
- Game Trees of Depth 1.5 (Maximum/Minimum)
   Results



## Relation of RL and PE



## Example 1: Phased Q-Learning

[Even-Dar, Mannor, and Mansour, 2002]

Early example(s) of Pure Exploration as sub-module in RL

```
Initialise V_0(s) := 0 for each state s.

for phase i = 1, 2, ... do

for each state s do

Run Best Arm Identification algorithm on

a \mapsto r + \gamma V_i(s') where (r, s') \sim \mathbb{P}(r, s'|s, a)

Store estimate in V_{i+1}(s).

end for

end for
```

## Example 2: AlphaZero

#### MCTS as Policy/Value Improvement Operator





K distributions parameterised by their means  $\mu = (\mu_1, \dots, \mu_K)$ . The **best arm** is

$$i^* = \underset{i \in [K]}{\operatorname{argmax}} \mu_i$$



K distributions parameterised by their means  $\mu = (\mu_1, \dots, \mu_K)$ . The **best arm** is

$$i^* = \underset{i \in [K]}{\operatorname{argmax}} \mu_i$$

#### Strategy

- Stopping rule  $\tau \in \mathbb{N}$
- In round  $t \leq \tau$  sampling rule picks  $I_t \in [K]$ . See  $X_t \sim \mu_{I_t}$ .
- Recommendation rule  $\hat{l} \in [K]$ .



K distributions parameterised by their means  $\mu = (\mu_1, \dots, \mu_K)$ . The **best arm** is

$$i^* = \underset{i \in [K]}{\operatorname{argmax}} \mu_i$$

#### Strategy

- Stopping rule  $\tau \in \mathbb{N}$
- In round  $t \leq \tau$  sampling rule picks  $I_t \in [K]$ . See  $X_t \sim \mu_{I_t}$ .
- Recommendation rule  $\hat{I} \in [K]$ .

Realisation of interaction:  $(I_1, X_1), \ldots, (I_{\tau}, X_{\tau}), \hat{I}$ .



K distributions parameterised by their means  $\mu = (\mu_1, \dots, \mu_K)$ . The **best arm** is

$$i^* = \underset{i \in [K]}{\operatorname{argmax}} \mu_i$$

#### Strategy

- Stopping rule  $\tau \in \mathbb{N}$
- In round  $t \leq \tau$  sampling rule picks  $I_t \in [K]$ . See  $X_t \sim \mu_{I_t}$ .
- Recommendation rule  $\hat{l} \in [K]$ .

Realisation of interaction:  $(I_1, X_1), \ldots, (I_{\tau}, X_{\tau}), \hat{I}$ .

**Two** objectives: sample efficiency  $\tau$  and correctness  $\hat{l} = i^*$ .

Koolen (CW

#### Objective

On bandit  $\mu$ , strategy  $(\tau, (I_t)_t, \hat{I})$  has • error probability  $\mathbb{P}_{\mu}(\hat{I} \neq i^*(\mu))$ , and

• sample complexity  $\mathbb{E}_{\mu}[\tau]$ .

Idea: constrain one, optimise the other.

## Objective



On bandit  $\mu$ , strategy  $(\tau, (I_t)_t, \hat{I})$  has

- error probability  $\mathbb{P}_{\boldsymbol{\mu}} (\hat{l} \neq i^*(\boldsymbol{\mu}))$ , and
- sample complexity  $\mathbb{E}_{\mu}[\tau]$ .

Idea: constrain one, optimise the other.

#### Definition

Fix small confidence  $\delta \in (0,1)$ . A strategy is  $\delta$ -correct if

 $\mathbb{P}_{\boldsymbol{\mu}}ig(\hat{l} 
eq i^*(\boldsymbol{\mu})ig) \leq \delta$  for every bandit model  $\boldsymbol{\mu}$ .

(Generalisation: output  $\epsilon$ -best arm)

## Objective



On bandit  $\mu$ , strategy  $(\tau, (I_t)_t, \hat{I})$  has

- error probability  $\mathbb{P}_{\mu}(\hat{l} \neq i^{*}(\mu))$ , and
- sample complexity  $\mathbb{E}_{\mu}[\tau]$ .

Idea: constrain one, optimise the other.

#### Definition

Fix small confidence  $\delta \in (0,1)$ . A strategy is  $\delta$ -correct if

 $\mathbb{P}_{\boldsymbol{\mu}}ig(\hat{l} 
eq i^*(\boldsymbol{\mu})ig) \leq \delta$  for every bandit model  $\boldsymbol{\mu}$ .

(Generalisation: output  $\epsilon$ -best arm)

#### Goal: minimise $\mathbb{E}_{\mu}[\tau]$ over all $\delta$ -correct strategies.

#### Model

## Algorithms



- Sampling rule  $I_t$ ?
- Stopping rule  $\tau$ ?
- Recommendation rule  $\hat{I}$ ?

$$\hat{l} = \mathop{\mathrm{argmax}}_{i \in [K]} \hat{\mu}_i( au)$$

where  $\hat{\mu}(t)$  is **empirical mean**.

#### Instance-Dependent Sample Complexity Lower bound

Define the alternatives to  $\mu$  by  $Alt(\mu) = \{\lambda | i^*(\lambda) \neq i^*(\mu)\}.$ 

Instance-Dependent Sample Complexity Lower bound

Define the alternatives to  $\mu$  by  $Alt(\mu) = \{\lambda | i^*(\lambda) \neq i^*(\mu)\}.$ 

Theorem (Castro 2014, Garivier and Kaufmann 2016) Fix a  $\delta$ -correct strategy. Then for every bandit model  $\mu$ 

$$\mathbb{E}_{oldsymbol{\mu}}[ au] \; \geq \; \mathcal{T}^*(oldsymbol{\mu}) \ln rac{1}{\delta}$$

where the characteristic time  $T^*(\mu)$  is given by

$$\frac{1}{\mathcal{T}^*(\boldsymbol{\mu})} = \max_{\boldsymbol{w} \in \Delta_K} \min_{\boldsymbol{\lambda} \in \mathsf{Alt}(\boldsymbol{\mu})} \sum_{i=1}^K w_i \mathsf{KL}(\mu_i \| \lambda_i).$$

Instance-Dependent Sample Complexity Lower bound

Define the **alternatives** to  $\mu$  by  $Alt(\mu) = \{\lambda | i^*(\lambda) \neq i^*(\mu)\}.$ 

Theorem (Castro 2014, Garivier and Kaufmann 2016) Fix a  $\delta$ -correct strategy. Then for every bandit model  $\mu$ 

$$\mathbb{E}_{oldsymbol{\mu}}[ au] \ \geq \ \mathcal{T}^*(oldsymbol{\mu}) \ln rac{1}{\delta}$$

where the characteristic time  $T^*(\mu)$  is given by

$$\frac{1}{T^*(\boldsymbol{\mu})} = \max_{\boldsymbol{w} \in \Delta_K} \min_{\boldsymbol{\lambda} \in \mathsf{Alt}(\boldsymbol{\mu})} \sum_{i=1}^K w_i \mathsf{KL}(\mu_i \| \lambda_i).$$

Intuition (going back to Lai and Robbins [1985]): if observations are likely under both  $\mu$  and  $\lambda$ , yet  $i^*(\mu) \neq i^*(\lambda)$ , then learner cannot stop and be correct in both.

Koolen (CWI

#### Example

K = 5 arms, Bernoulli  $\mu = (0, 0.1, 0.2, 0.3, 0.4)$ .

$$T^*(\mu) = 200.4$$
  $w^*(\mu) = (0.45, 0.46, 0.06, 0.02, 0.01)$ 

At  $\delta = 0.05$ , the time gets multiplied by  $\ln \frac{1}{\delta} = 3.0$ .

## Sampling Rule

Look at the lower bound again. Any good algorithm **must** sample with optimal (**oracle**) proportions

$$m{w}^*(m{\mu}) \;=\; rgmax_{m{w}\in riangle_K} \min_{m{\lambda}\in \mathsf{Alt}(m{\mu})} \; \sum_{i=1}^K w_i \, \mathsf{KL}(\mu_i \| \lambda_i)$$

## Sampling Rule

Look at the lower bound again. Any good algorithm **must** sample with optimal (**oracle**) proportions

$$w^*(\mu) = rgmax_{oldsymbol{w}\in riangle_K} \min_{oldsymbol{\lambda}\in \mathsf{Alt}(\mu)} \sum_{i=1}^K w_i \, \mathsf{KL}(\mu_i \| \lambda_i)$$

#### Track-and-Stop

Idea: draw  $I_t \sim w^*(\hat{\mu}(t))$ .

- Ensure  $\hat{\mu}(t) 
  ightarrow \mu$  hence  $N_i(t)/t 
  ightarrow w_i^*$  by "forced exploration"
- Draw arm with  $N_i(t)/t$  below  $w_i^*$  (tracking)
- Computation of  $w^*$  (reduction to 1d line search)

### All in all

Final result: lower and upper bound meet on every problem instance.

Theorem (Garivier and Kaufmann 2016)

For Track-and-Stop algorithm, for any bandit  $\mu$ 

$$\limsup_{\delta \to 0} \frac{\mathbb{E}_{\mu}[\tau]}{\ln \frac{1}{\delta}} = T^{*}(\mu)$$

## All in all

Final result: lower and upper bound meet on every problem instance.

Theorem (Garivier and Kaufmann 2016)

For Track-and-Stop algorithm, for any bandit  $\mu$ 

$$\limsup_{\delta \to 0} \frac{\mathbb{E}_{\mu}[\tau]}{\ln \frac{1}{\delta}} = T^{*}(\mu)$$

Very similar optimality result for **Top Two Thompson Sampling** by Russo [2016]. Here  $N_i(t)/t \rightarrow w_i^*$  result of posterior sampling.



2 Relation of RL and PE

#### Oure Exploration Intro: Best Arm Identification

- Model
- Sample Complexity Lower Bound
- Algorithms

#### Game Tree Search

- Game Trees of Arbitrary Depth
- Confidence Intervals on Min/Max
- Game Trees of Depth 1.5 (Maximum/Minimum)
   Results



### Challenge Environment: Stochastic Game Tree Search



## Challenge Environment: Stochastic Game Tree Search



## Model [Teraoka et al., 2014]



Maximin Action Identification Problem

Find best move at root from samples of leaves.

## Model [Teraoka et al., 2014]



Maximin Action Identification Problem Find **best move at root** from samples of **leaves**.



## Methods for Best Arm Identification





## Methods for Best Arm Identification



. . .



## Methods for Best Arm Identification





LUCB, UGapE,

. . .















Correctness  $(\epsilon, \delta)$ -PAC algorithms

#### Efficiency

Sample complexity function of leaf gaps  $\Delta_\ell$ 

$$O\left(\sum_{\ell \in \mathcal{L}} \frac{1}{\Delta_\ell^2 \vee \epsilon^2} \log\left(\frac{1}{\delta}\right)\right)$$

#### How to build confidence intervals on min/max nodes



### How to build confidence intervals on min/max nodes



More principled approach in [Kaufmann, Koolen, and Garivier, 2018]. Many equal intervals  $\Rightarrow$  higher lower bound.



2 Relation of RL and PE

#### Oure Exploration Intro: Best Arm Identification

- Model
- Sample Complexity Lower Bound
- Algorithms

#### Game Tree Search

- Game Trees of Arbitrary Depth
- Confidence Intervals on Min/Max
- Game Trees of Depth 1.5 (Maximum/Minimum)
   Results



## Simplify





Best Arm Identification [Garivier and Kaufmann, 2016] Solved Depth 2 [Garivier, Kaufmann, and Koolen, 2016] Open

## Simple Instance: Minimum Threshold Identification



Fix threshold  $\gamma$ .



For 
$$t = 1, \ldots, \tau$$

- Pick leaf A<sub>t</sub>
- See  $X_t \sim \mu_{A_t}$

Recommend  $\hat{m} \in \{<,>\}$ 

Goal: fixed confidence  $\mathbb{P}_{\mu}$  {error}  $< \delta$ and small sample complexity  $\mathbb{E}_{\mu}[\tau]$ 

#### Lower Bound

Generic lower bound [Castro, 2014, Garivier and Kaufmann, 2016] shows sample complexity for any  $\delta$ -correct algorithm is at least

 $\mathbb{E}_{\mu}[\tau] \geq T^{*}(\mu) \ln \frac{1}{\delta}.$ 

#### Lower Bound

Generic lower bound [Castro, 2014, Garivier and Kaufmann, 2016] shows sample complexity for any  $\delta$ -correct algorithm is at least

$$\mathbb{E}_{oldsymbol{\mu}}[ au] \ \geq \ T^*(oldsymbol{\mu}) \ln rac{1}{\delta}.$$

For our problem the characteristic time and oracle weights are

$$T^{*}(\boldsymbol{\mu}) = \begin{cases} \frac{1}{\mathsf{KL}(\boldsymbol{\mu}^{*},\boldsymbol{\gamma})} & \boldsymbol{\mu}^{*} < \boldsymbol{\gamma}, \\ \sum_{a} \frac{1}{\mathsf{KL}(\boldsymbol{\mu}_{a},\boldsymbol{\gamma})} & \boldsymbol{\mu}^{*} > \boldsymbol{\gamma}, \end{cases} \quad \mathbf{w}^{*}_{a}(\boldsymbol{\mu}) = \begin{cases} \mathbf{1}_{a=a^{*}} & \boldsymbol{\mu}^{*} < \boldsymbol{\gamma}, \\ \frac{1}{\mathsf{KL}(\boldsymbol{\mu}_{a},\boldsymbol{\gamma})} & \frac{1}{\mathsf{KL}(\boldsymbol{\mu}_{j},\boldsymbol{\gamma})} & \boldsymbol{\mu}^{*} > \boldsymbol{\gamma}. \end{cases}$$

#### Dichotomous Oracle Behaviour! Sampling Rule?



## Sampling Rules

- Lower Confidence Bounds Play  $A_t = \arg \min_a \operatorname{LCB}_a(t)$
- **Thompson Sampling** ( $\Pi_{t-1}$  is posterior after t-1 rounds) Sample  $\theta \sim \Pi_{t-1}$ , then play  $A_t = \arg \min_a \theta_a$ .
- Murphy Sampling condition on low minimum mean Sample  $\theta \sim \prod_{t=1} (\cdot |\min_a \theta_a < \gamma)$ , then play  $A_t = \arg \min_a \theta_a$ .



## Intuition for Murphy Sampling

- When  $\mu^* < \gamma$  conditioning is immaterial:  $\theta \approx \mu$  and MS  $\equiv$  TS.
- When μ\* > γ conditioning results in θ ≈ (μ<sub>1</sub>,..., γ,..., μ<sub>K</sub>). Index a lowered to γ with probability ∝ <sup>1</sup>/<sub>KL(μ<sub>2</sub>,γ)</sub> [Russo, 2016].

## Murphy Sampling Rule [KKG, NIPS'18]

#### Theorem

Asymptotic optimality:  $N_{\mathsf{a}}(t)/t o w^*_{\mathsf{a}}(\mu)$  for all  $\mu$ 

| Sampling rule           | $\leq$       | $\geq$       |
|-------------------------|--------------|--------------|
| Thompson Sampling       | $\checkmark$ | ×            |
| Lower Confidence Bounds | ×            | $\checkmark$ |
| Murphy Sampling         | $\checkmark$ | $\checkmark$ |

#### Lemma

Any anytime sampling strategy  $(A_t)_t$  ensuring  $\frac{N_t}{t} \to \boldsymbol{w}^*(\boldsymbol{\mu})$  and good stopping rule  $\tau_{\delta}$  guarantee  $\limsup_{\delta \to 0} \frac{\tau_{\delta}}{\ln \frac{1}{\lambda}} \leq T^*(\boldsymbol{\mu})$ .

#### Conclusion

- Pure Exploration currently going through a renaissance
- Instance-optimal identification algorithms
  - Best Arm
  - Game Tree Search
  - ▶ ...
- Moving toward more complex queries. RL on the horizon ....
- Useful submodules

### Conclusion

- Pure Exploration currently going through a renaissance
- Instance-optimal identification algorithms
  - Best Arm
  - Game Tree Search
  - ▶ ...
- Moving toward more complex queries. RL on the horizon ....
- Useful submodules

# Thank you! And let's talk!