

# Monte-Carlo Tree Search by Best Arm Identification

Emilie Kaufmann and Wouter M. Koolen

Inria Lille  
SequeL team

CWI  
Machine Learning Group



Centrum Wiskunde & Informatica

Inria-CWI workshop  
Amsterdam, September 20th, 2017

... a new Associate Team proposal

## 6 PAC

involving

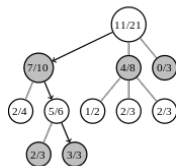
- Peter Grünwald (CWI, Machine Learning Group)
- Wouter M. Koolen (CWI, Machine Learning Group)
- Benjamin Guedj (Inria Lille, MODAL project-team)
- Emilie Kaufmann (Inria Lille, Sequel project-team)

Broader goal: Probably Approximately Correct - Learning

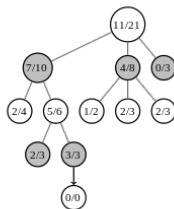
<sup>6</sup> Safe, Efficient, Sequential, Active, Structured, Ideal

# Monte-Carlo Tree Search for games

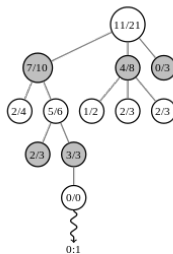
Selection



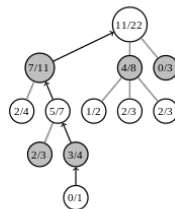
Expansion



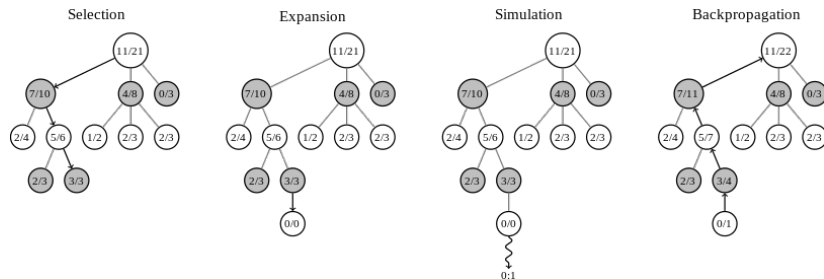
Simulation



Backpropagation



# Monte-Carlo Tree Search for games



We introduce an idealized model:

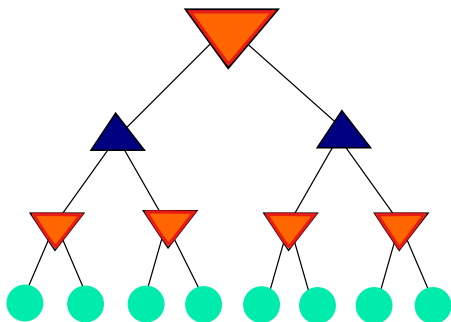
- *fixed* maximin tree
- *i.i.d.* payouts starting from each leaf

and propose **new algorithms** with **sample complexity guarantees**

- 1 Problem formulation
- 2 The BAI-MCTS architecture
- 3 UGapE-MCTS and LUCB-MCTS
- 4 Towards optimal algorithms

- 1 Problem formulation
- 2 The BAI-MCTS architecture
- 3 UGapE-MCTS and LUCB-MCTS
- 4 Towards optimal algorithms

# A simple model for MCTS



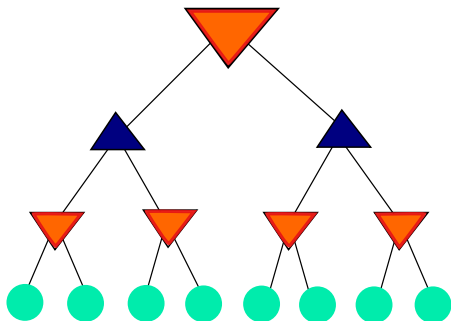
A fixed MAXMIN game tree  $\mathcal{T}$ , with leaves  $\mathcal{L}$ .

Orange inverted triangle MAX node (= your move)

Blue triangle MIN node (= adversary move)

Green circle Leaf  $l$ : stochastic oracle  $\mathcal{O}_l$  that evaluates the position

# A simple model for MCTS



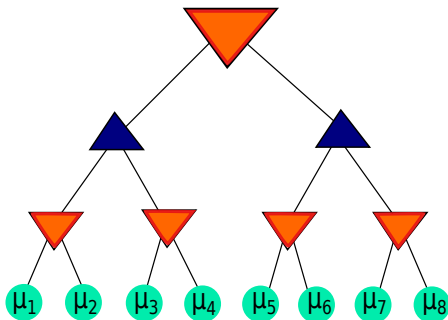
At round  $t$  a **MCTS algorithm**:

- picks a path down to a leaf  $L_t$
- get an evaluation of this leaf  $X_t \sim \mathcal{O}_{L_t}$

Assumption: i.i.d. successive evaluations,  $\mathbb{E}_{X \sim \mathcal{O}_\ell}[X] = \mu_\ell$



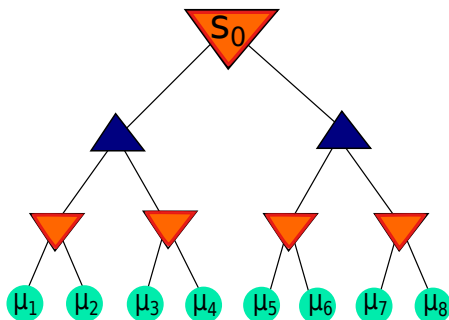
# A simple model for MCTS



At round  $t$  a **MCTS algorithm**:

- picks a path down to a leaf  $L_t$
- get an evaluation of this leaf  $X_t \sim \mathcal{O}_{L_t}$

Assumption: i.i.d. successive evaluations,  $\mathbb{E}_{X \sim \mathcal{O}_\ell}[X] = \mu_\ell$

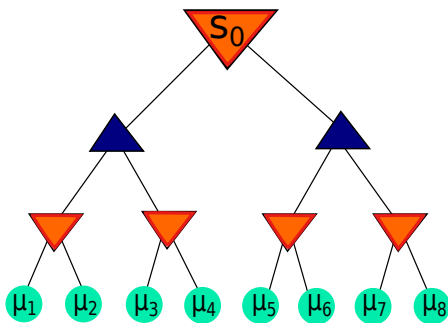


A MCTS algorithm should find the **best move at the root**:

$$V_s = \begin{cases} \mu_s & \text{if } s \in \mathcal{L}, \\ \max_{c \in \mathcal{C}(s)} V_c & \text{if } s \text{ is a MAX node,} \\ \min_{c \in \mathcal{C}(s)} V_c & \text{if } s \text{ is a MIN node.} \end{cases}$$

$$s^* = \operatorname{argmax}_{s \in \mathcal{C}(s_0)} V_s$$

# A PAC learning framework



MCTS algorithm:  $(L_t, \tau, \hat{s}_T)$ , where

- $L_t$  is the **sampling rule**
- $\tau$  is the **stopping rule**
- $\hat{s}_T \in \mathcal{C}(s_0)$  is the **recommendation rule**

is  $(\epsilon, \delta)$ -PAC if  $\mathbb{P}(V_{\hat{s}_T} \geq V_{s^*} - \epsilon) \geq 1 - \delta$ .

Goal:  $(\epsilon, \delta)$ -PAC algorithm with a small **sample complexity**  $\tau$ .

# A simpler problem: best arm identification

Reminiscent of a **bandit model**:



A Best Arm Identification algorithm:  $(A_t, \tau, \hat{s}_\tau)$ , where

- $A_t$  is the **sampling rule**
- $\tau$  is the **stopping rule**
- $\hat{s}_\tau \in \mathcal{C}(s_0)$  is the **recommendation rule**

is  $(\epsilon, \delta)$ -PAC if

$$\mathbb{P}(\mu_{\hat{s}_\tau} \geq \mu^* - \epsilon) \geq 1 - \delta.$$

# A simpler problem: best arm identification

Reminiscent of a **bandit model**:



A Best Arm Identification algorithm:  $(A_t, \tau, \hat{s}_T)$ , where

- $A_t$  is the **sampling rule**
- $\tau$  is the **stopping rule**
- $\hat{s}_T \in \mathcal{C}(s_0)$  is the **recommendation rule**

is  $(\epsilon, \delta)$ -PAC if

$$\mathbb{P}(\mu_{\hat{s}_T} \geq \mu^* - \epsilon) \geq 1 - \delta.$$

## The BAI problem:

How to adaptively sample the arms so as to identify as quickly as possible the arm with highest mean ?

# MCTS: a structured BAI problem

Reminiscent of a bandit model:



A Best Arm Identification algorithm:  $(L_t, \tau, \hat{s}_\tau)$ , where

- $L_t$  is the **sampling rule**
- $\tau$  is the **stopping rule**
- $\hat{s}_\tau \in \mathcal{C}(s_0)$  is the **recommendation rule**

is  $(\epsilon, \delta)$ -PAC if

$$\mathbb{P}(V_{\hat{s}_\tau} \geq V_{s^*} - \epsilon) \geq 1 - \delta.$$

## The MCTS problem:

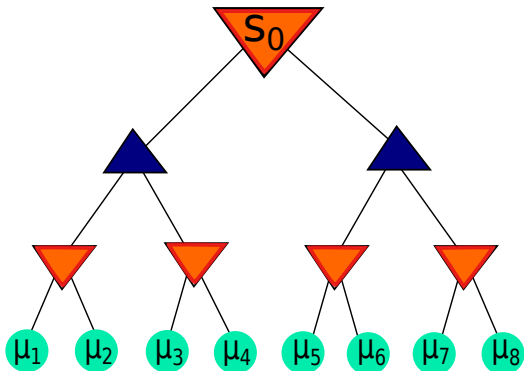
How to adaptively sample the **leaves of a maxmin tree** so as to identify as quickly as possible the **best action at the root** ?

- 1 Problem formulation
- 2 The BAI-MCTS architecture
- 3 UGapE-MCTS and LUCB-MCTS
- 4 Towards optimal algorithms

# A key building block: confidence intervals

Using the samples collected for the leaves, one can build, for  $\ell \in \mathcal{L}$ ,

$[\text{LCB}_\ell(t), \text{UCB}_\ell(t)]$  a confidence interval on  $\mu_\ell$

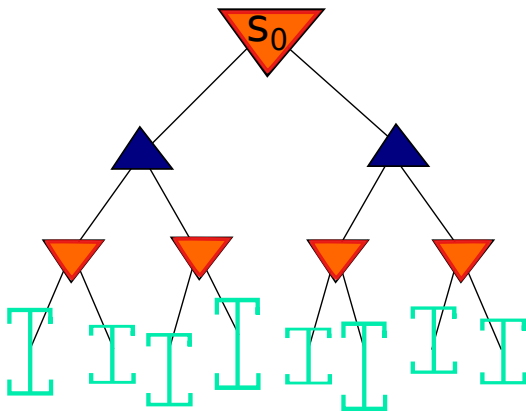




# A key building block: confidence intervals

Using the samples collected for the leaves, one can build, for  $\ell \in \mathcal{L}$ ,

$[\text{LCB}_\ell(t), \text{UCB}_\ell(t)]$  a confidence interval on  $\mu_\ell$

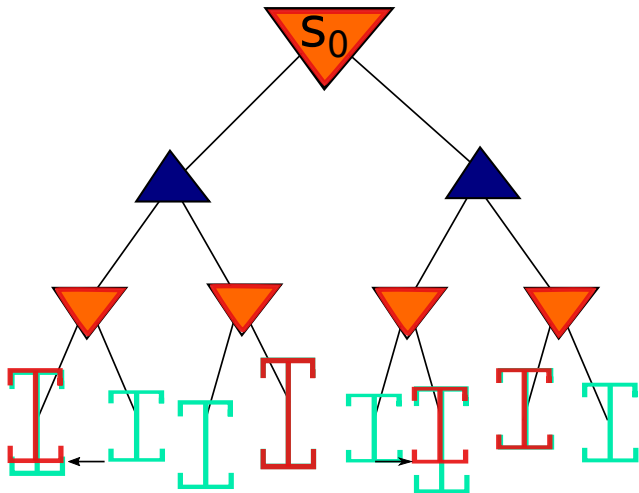


**Idea:** Propagate these confidence intervals up in the tree

# A key building block: confidence intervals

MAX node:

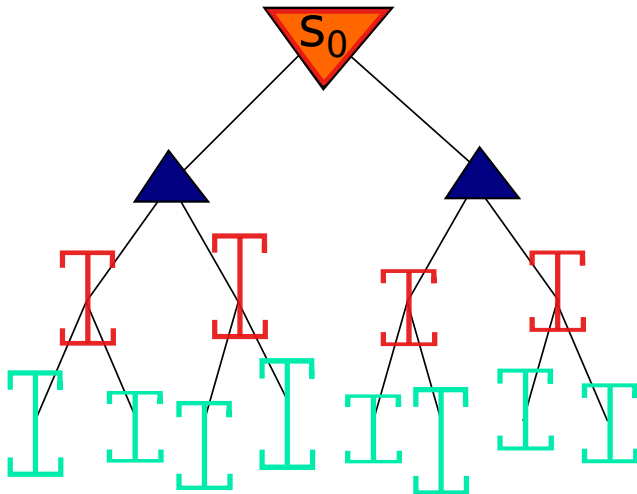
$$UCB_s(t) = \max_{c \in \mathcal{C}(s)} UCB_c(t) \quad LCB_s(t) = \max_{c \in \mathcal{C}(s)} LCB_c(t)$$



# A key building block: confidence intervals

MAX node:

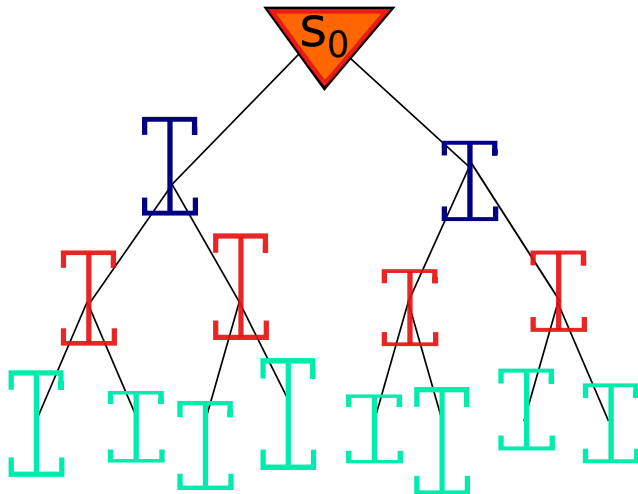
$$UCB_s(t) = \max_{c \in \mathcal{C}(s)} UCB_c(t) \quad LCB_s(t) = \max_{c \in \mathcal{C}(s)} LCB_c(t)$$



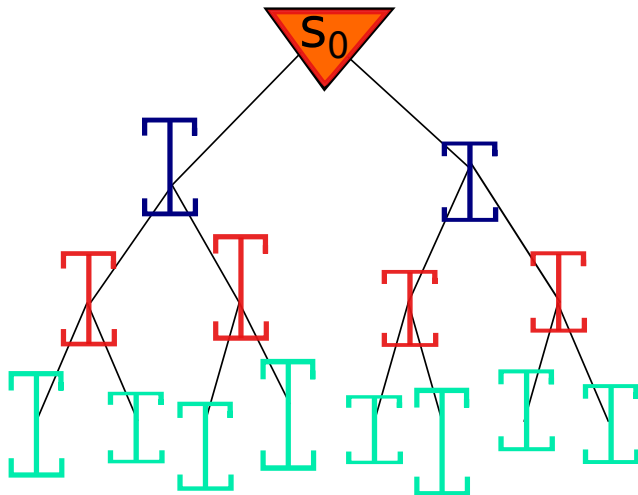
# A key building block: confidence intervals

MIN node:

$$UCB_s(t) = \min_{c \in \mathcal{C}(s)} UCB_c(t) \quad LCB_s(t) = \min_{c \in \mathcal{C}(s)} LCB_c(t)$$



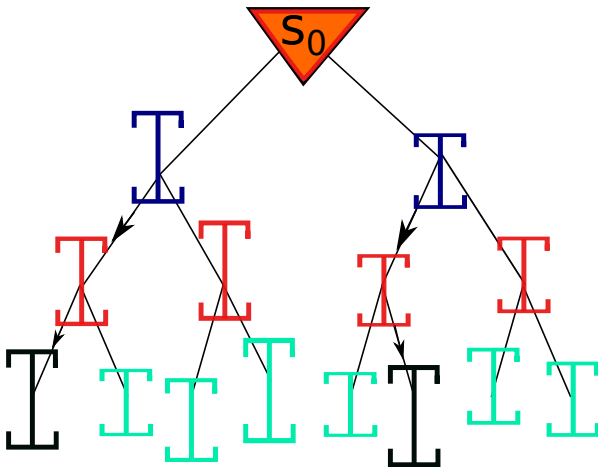
# Property of this construction



$$\bigcap_{\ell \in L} (\mu_\ell \in \mathcal{I}_\ell(t)) \Rightarrow \bigcap_{s \in T} (V_s \in \mathcal{I}_s(t))$$

# Representative leaves

$l_s(t)$ : **representative leaf** of internal node  $s \in \mathcal{T}$ .



**Idea:** alternate optimistic/pessimistic moves starting from  $s$

# Generic BAI-MCTS algorithm

**Input:** a BAI algorithm

**Initialization:**  $t = 0$ .

**while not** BAIStop ( $\{s \in \mathcal{C}(s_0)\}$ ) **do**

$R_{t+1} = \text{BAIStep}(\{s \in \mathcal{C}(s_0)\})$

    Sample the **representative leaf**  $L_{t+1} = \ell_{R_{t+1}}(t)$

    Update the information about the arms.  $t = t + 1$ .

**end**

**Output:** BAIReco ( $\{s \in \mathcal{C}(s_0)\}$ )

# Generic BAI-MCTS algorithm

**Input:** a BAI algorithm

**Initialization:**  $t = 0$ .

**while not** BAIStop ( $\{s \in \mathcal{C}(s_0)\}$ ) **do**

$R_{t+1} = \text{BAIStep}(\{s \in \mathcal{C}(s_0)\})$

    Sample the representative leaf  $L_{t+1} = \ell_{R_{t+1}}(t)$

    Update the information about the arms.  $t = t + 1$ .

**end**

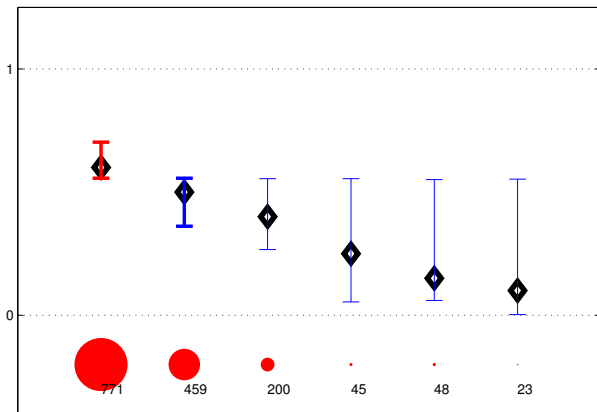
**Output:** BAIReco ( $\{s \in \mathcal{C}(s_0)\}$ )

... sometimes reduces to updating confidence intervals!



- 1 Problem formulation
- 2 The BAI-MCTS architecture
- 3 UGapE-MCTS and LUCB-MCTS**
- 4 Towards optimal algorithms

# An example of BAI algorithm: LUCB



The (KL)-LUCB algorithm

[Kalyanakrishnan et al. 12, Kaufmann and Kalyanakrishnan 13]

based on the UGapE algorithm [Gabillon et al. 12]

- **Sampling rule:**  $R_{t+1}$  is the least sampled among two promising depth-one nodes:

$$\underline{a}_t = \operatorname{argmin}_{a \in \mathcal{C}(s_0)} B_a(t) \quad \text{and} \quad \underline{b}_t = \operatorname{argmax}_{b \in \mathcal{C}(s_0) \setminus \{\underline{a}_t\}} \operatorname{UCB}_b(t),$$

where

$$B_s(t) = \max_{s' \in \mathcal{C}(s_0) \setminus \{s\}} \operatorname{UCB}_{s'}(t) - \operatorname{LCB}_s(t).$$

- **Stopping rule:**

$$\tau = \inf \{ t \in \mathbb{N} : \operatorname{UCB}_{\underline{b}_t}(t) - \operatorname{LCB}_{\underline{a}_t}(t) < \epsilon \}$$

- **Recommendation rule:**  $\hat{s}_\tau = \underline{a}_\tau$

# Theoretical guarantees

We choose confidence intervals of the form

$$\text{LCB}_\ell(t) = \hat{\mu}_\ell(t) - \sqrt{\frac{\beta(N_\ell(t), \delta)}{2N_\ell(t)}}$$

$$\text{UCB}_\ell(t) = \hat{\mu}_\ell(t) + \sqrt{\frac{\beta(N_\ell(t), \delta)}{2N_\ell(t)}}$$

where  $\beta(s, \delta)$  is some **exploration function**.

## Correctness

If  $\delta \leq \max(0.1|\mathcal{L}|, 1)$ , for the choice

$$\beta(s, \delta) = \log(|\mathcal{L}|/\delta) + 3 \log \log(|\mathcal{L}|/\delta) + (3/2) \log(\log s + 1)$$

**UGapE-MCTS is  $(\epsilon, \delta)$ -PAC.**

$$H_\epsilon^*(\mu) := \sum_{\ell \in \mathcal{L}} \frac{1}{\Delta_\ell^2 \vee \Delta_*^2 \vee \epsilon^2}$$

where

$$\Delta_* := V(s^*) - V(s_2^*)$$

$$\Delta_\ell := \max_{s \in \text{Ancestors}(\ell) \setminus \{s_0\}} |V_{\text{Parent}(s)} - V_s|$$

## Sample complexity

With probability larger than  $1 - \delta$ , the total number of leaves explorations performed by UGapE-MCTS is upper bounded as

$$\tau = O \left( H_\epsilon^*(\mu) \log \left( \frac{1}{\delta} \right) \right).$$

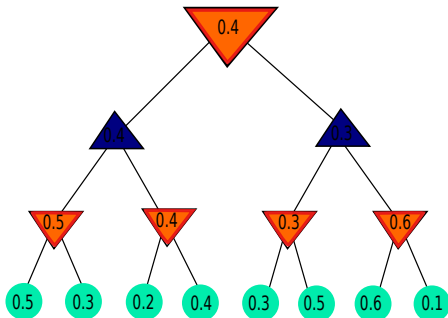
# Theoretical guarantees

$$H_\epsilon^*(\mu) := \sum_{\ell \in \mathcal{L}} \frac{1}{\Delta_\ell^2 \vee \Delta_*^2 \vee \epsilon^2}$$

where

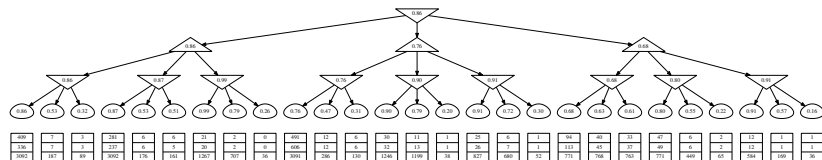
$$\Delta_* := V(s^*) - V(s_2^*)$$

$$\Delta_\ell := \max_{s \in \text{Ancestors}(\ell) \setminus \{s_0\}} |V_{\text{Parent}(s)} - V_s|$$



# Numerical results

$\epsilon = 0, \delta = 0.1 \cdot 27$  ( $N = 10^6$  simulations)



LUCB-MCTS (0.72% errors, 1551 samples)

UGapE-MCTS (0.75% errors, 1584 samples)

FindTopWinner (0% errors, 20730 samples) [Teraoka et al. 14]

- 1 Problem formulation
- 2 The BAI-MCTS architecture
- 3 UGapE-MCTS and LUCB-MCTS
- 4 Towards optimal algorithms



# A sample complexity lower bound

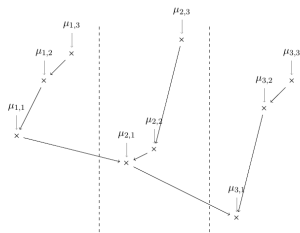
## Theorem

Let  $\epsilon = 0$ . Any  $\delta$ -correct algorithm satisfies

where 
$$\mathbb{E}_{\mu}[\tau] \geq T^*(\mu) \log(1/(3\delta))$$

$$T^*(\mu)^{-1} := \sup_{\mathbf{w} \in \Sigma_{|\mathcal{L}|}} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{\ell \in \mathcal{L}} w_{\ell} \text{KL}(\mathcal{B}(\mu_{\ell}), \mathcal{B}(\lambda_{\ell})).$$

Depth-two tree:



The optimal proportions satisfy

$$w_{i,j}^*(\mu) = 0$$

if  $i \geq 2$  and  $j \geq 2$ .

# A sample complexity lower bound

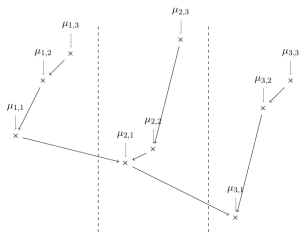
## Theorem

Let  $\epsilon = 0$ . Any  $\delta$ -correct algorithm satisfies

where 
$$\mathbb{E}_{\mu}[\tau] \geq T^*(\mu) \log(1/(3\delta))$$

$$T^*(\mu)^{-1} := \sup_{\mathbf{w} \in \Sigma_{|\mathcal{L}|}} \inf_{\lambda \in \text{Alt}(\mu)} \sum_{\ell \in \mathcal{L}} w_{\ell} \text{KL}(\mathcal{B}(\mu_{\ell}), \mathcal{B}(\lambda_{\ell})).$$

Depth-two tree:



The optimal proportions satisfy

$$w_{i,j}^*(\mu) = 0$$

if  $i \geq 2$  and  $j \geq 2$ .

**A more general sparsity pattern?**

## Our contributions:

- a generic way to use a BAI algorithm for MCTS
- PAC and sample complexity guarantees for UGapE-MCTS and LUCB-MCTS...
- ... that also displays good empirical performance

## Future work:

- identify the *optimal* sample complexity of the MCTS problem... (i.e. matching upper and lower bounds)
- ... and that of other structured Best Arm Identification problems [Ajallooeian et al., ALT 17]

## Our contributions:

- a generic way to use a BAI algorithm for MCTS
- PAC and sample complexity guarantees for UGapE-MCTS and LUCB-MCTS...
- ... that also displays good empirical performance

## Future work:

- identify the *optimal* sample complexity of the MCTS problem... (i.e. matching upper and lower bounds)
- ... and that of other structured Best Arm Identification problems [Ajallooeian et al., ALT 17]

## Reference:

E. Kaufmann & W.M. Koolen,  
*Monte-Carlo Tree Search by Best Arm Identification*  
to appear in NIPS 2017