


 Viktor Bengs^{a,b}, Björn Haddendorst^c, Eyke Hüllermeier^{a,b}
^aInstitute of Informatics, University of Munich, Germany

^bMunich Center for Machine Learning, Germany

^cDepartment of Computer Science, Paderborn University

viktor.bengs@lmu.de, eyke@lmu.de, willem.waegeman@UGent.be



TL;DR

Extension of Copeland winner identification in dueling bandits for indifference feedback with novel lower bounds and a worst-case nearly optimal learning algorithm

DUELING BANDITS WITH INDIFFERENCES

Setting

- **Given:** Different arms (options) $a_1, \dots, a_n \iff 1, \dots, n \iff \mathcal{A}$

- **Action at time t :** Choose a pair of arms $i_t \in \mathcal{A}$ and $j_t \in \mathcal{A} \setminus \{i_t\}$

- **Observation at time t :**

 either $i_t \succ j_t$, i.e., arm i_t is strictly preferred over arm j_t

 or $i_t \prec j_t$, i.e., arm j_t is strictly preferred over arm i_t

 or $i_t \cong j_t$, i.e., neither i_t is strictly preferred over j_t nor the opposite (*indifference between i_t and j_t*)

- **Stochastic feedback assumption:** Each possible explicit observations is determined by one of the following matrices $P^{\succ}, P^{\prec}, P^{\cong} \in [0, 1]^{n \times n}$:

$$P_{i_t, j_t}^{\succ} = \mathbb{P}(i_t \succ j_t) \quad P_{i_t, j_t}^{\prec} = \mathbb{P}(i_t \prec j_t) \quad P_{i_t, j_t}^{\cong} = \mathbb{P}(i_t \cong j_t)$$

 \rightsquigarrow A problem instance is characterized by $\mathbf{P} = ((P_{i,j}^{\succ}, P_{i,j}^{\cong}, P_{i,j}^{\prec}))_{i < j}$

Goal

(i) Finding a Copeland winner (COWI), i.e., an element of

$$\mathcal{C}(\mathbf{P}) = \{i \in \mathcal{A} \mid \text{CP}(\mathbf{P}, i) = \max_j \text{CP}(\mathbf{P}, j)\},$$

where

$$\text{CP}(\mathbf{P}, i) = \sum_{j \neq i} 1_{[P_{i,j}^{\succ} > \max\{P_{i,j}^{\prec}, P_{i,j}^{\cong}\}]} + \frac{1}{2} \sum_{j \neq i} 1_{[P_{i,j}^{\cong} > \max\{P_{i,j}^{\prec}, P_{i,j}^{\succ}\}]},$$

 is the Copeland score of arm $i \in \mathcal{A}$

(ii) Conducting as few as possible duels (low sample complexity)

Formal Goal: For a given error bound $\delta \in (0, 1)$ design algorithm A which

- uses $\tau^A(\mathbf{P})$ duels in total such that $\mathbb{E}[\tau^A(\mathbf{P})]$ is small
- returns $\hat{i} \in \mathcal{A}$ such that $\mathbb{P}(\hat{i} \notin \mathcal{C}(\mathbf{P})) \leq \delta$

 for any problem instance \mathbf{P} .

REFERENCES

[1] Róbert Busa-Fekete, Balázs Szörényi, Paul Weng, Weiwei Cheng, and Eyke Hüllermeier. Top- k selection based on adaptive sampling of noisy preferences. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 1094–1102, 2013.

[2] Shubham Anand Jain, Rohan Shah, Sanit Gupta, Denil Mehta, Inderjeet J Nair, Jian Vora, Sushil Khyalia, Sourav Das, Vinay J Ribeiro, and Shivaram Kalyanakrishnan. PAC mode estimation using PPR martingale confidence sequences. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 5815–5852. PMLR, 2022.

[3] Tanguy Urvoy, Fabrice Clerot, Raphael Féraud, and Sami Naamane. Generic exploration and k -armed voting bandits. In *Proceedings of International Conference on Machine Learning (ICML)*, pages 91–99, 2013.



LEARNING ALGORITHM

POCOWISTA

Idea of POTential COpeland Winner STays Algorithm (POCOWISTA):

1. Duel arm i_t having highest potentially Copeland score with arm j_t having highest current Copeland score
2. Conduct duel via efficient PPR-1V1 routine [2] to find mode of $(P_{i_t, j_t}^{\succ}, P_{i_t, j_t}^{\cong}, P_{i_t, j_t}^{\prec})$

Algorithm POCOWISTA

```

1: Input: Set of arms  $\mathcal{A}$ , error prob.  $\delta \in (0, 1)$ 
2: Initialization:  $e \leftarrow 1$  and for each  $i \in \mathcal{A}$  set
    $D(i) \leftarrow \{i\}$  (set of already compared arms)
    $\widehat{\text{CP}}(i) \leftarrow 0$  (current Copeland score)
    $\overline{\text{CP}}(i) \leftarrow n - 1$  (potential Copeland score)
3: while  $\nexists i$  s.t.  $\widehat{\text{CP}}(i) \geq \overline{\text{CP}}(j) \forall j \in \mathcal{A} \setminus \{i\}$  do
4:    $i_e = \text{argmax}_{i \in \mathcal{A}} \widehat{\text{CP}}(i)$ 
5:    $j_e = \text{argmax}_{j \in \mathcal{A} \setminus D(i_e)} \widehat{\text{CP}}(j)$ 
6:    $k \leftarrow \text{PPR-1V1}(i_e, j_e, \delta / \binom{n}{2})$ 
7:   SCORES-UPDATE( $i_e, j_e, k$ )
8:    $e \leftarrow e + 1$ 
9: end while
10: return  $\text{argmax}_{i \in \mathcal{A}} \widehat{\text{CP}}(i)$ 
    
```

Algorithm SCORES-UPDATE

```

1: Input: Arms  $i, j$ , ternary decision  $k \in \{1, 2, 3\}$ 
2: if  $k = 1$  then
3:    $\widehat{\text{CP}}(i) \leftarrow \widehat{\text{CP}}(i) + 1$ 
4: else if  $k = 2$  then
5:    $\widehat{\text{CP}}(i) \leftarrow \widehat{\text{CP}}(i) + 1/2$ ,  $\widehat{\text{CP}}(j) \leftarrow \widehat{\text{CP}}(j) + 1/2$ 
6: else
7:    $\widehat{\text{CP}}(j) \leftarrow \widehat{\text{CP}}(j) + 1$ 
8: end if
9:  $D(i) \leftarrow D(i) \cup \{j\}$ ,  $D(j) \leftarrow D(j) \cup \{i\}$ 
10:  $\overline{\text{CP}}(i) \leftarrow n - |D(i)| + \widehat{\text{CP}}(i)$ 
11:  $\overline{\text{CP}}(j) \leftarrow n - |D(j)| + \widehat{\text{CP}}(j)$ 
    
```

TRA-POCOWISTA

 What if the problem instance \mathbf{P} is transitive?

Definition. \mathbf{P} is *transitive* if for each distinct $i, j, k \in \mathcal{A}$ holds:

1. Transitivity of strict preference.
 - If $P_{i,j}^{\succ} > \max(P_{i,j}^{\prec}, P_{i,j}^{\cong})$ and $P_{j,k}^{\succ} > \max(P_{j,k}^{\prec}, P_{j,k}^{\cong})$, then $P_{i,k}^{\succ} > \max(P_{i,k}^{\prec}, P_{i,k}^{\cong})$.
2. IP-transitivity.
 - If $P_{i,j}^{\cong} > \max(P_{i,j}^{\prec}, P_{i,j}^{\succ})$ and $P_{j,k}^{\succ} > \max(P_{j,k}^{\prec}, P_{j,k}^{\cong})$, then $P_{i,k}^{\succ} > \max(P_{i,k}^{\prec}, P_{i,k}^{\cong})$.
3. PI-transitivity.
 - If $P_{i,j}^{\succ} > \max(P_{i,j}^{\prec}, P_{i,j}^{\cong})$ and $P_{j,k}^{\cong} > \max(P_{j,k}^{\prec}, P_{j,k}^{\succ})$, then $P_{i,k}^{\succ} > \max(P_{i,k}^{\prec}, P_{i,k}^{\cong})$.
4. Transitivity of indifference.
 - If $P_{i,j}^{\cong} > \max(P_{i,j}^{\prec}, P_{i,j}^{\succ})$ and $P_{j,k}^{\succ} > \max(P_{j,k}^{\prec}, P_{j,k}^{\cong})$, then $P_{i,k}^{\cong} > \max(P_{i,k}^{\prec}, P_{i,k}^{\succ})$.

 \Rightarrow Updates can be made more efficient

Algorithm TRA-POCOWISTA

```

1: Input: Set of arms  $\mathcal{A}$ , error prob.  $\delta \in (0, 1)$ 
2: Initialization:  $e \leftarrow 1$  and for each  $i \in \mathcal{A}$ 
    $D(i) \leftarrow \{i\}$ ,  $\widehat{\text{CP}}(i) \leftarrow 0$ ,  $\overline{\text{CP}}(i) \leftarrow n - 1$ 
    $W(i) \leftarrow \emptyset$  (set of defeated arms)
    $I(i) \leftarrow \emptyset$  (set of indifferent arms)
    $L(i) \leftarrow \emptyset$  (set of superior arms)
3: while  $\nexists i$  s.t.  $\widehat{\text{CP}}(i) \geq \overline{\text{CP}}(j) \forall j \in \mathcal{A} \setminus \{i\}$  do
4:    $i_e = \text{argmax}_{i \in \mathcal{A}} \widehat{\text{CP}}(i)$ 
5:    $j_e = \text{argmax}_{j \in \mathcal{A} \setminus D(i_e)} \widehat{\text{CP}}(j)$ 
6:    $k \leftarrow \text{PPR-1V1}(i_e, j_e, \delta / n)$ 
7:   TRANSITIVE-SCORE-UPDATE( $i_e, j_e, k$ )
8:    $e \leftarrow e + 1$ 
9: end while
10: return  $\text{argmax}_{i \in \mathcal{A}} \widehat{\text{CP}}(i)$ 
    
```

Algorithm TRANSITIVE-SCORE-UPDATE

```

1: Input: Arms  $i, j, k \in \{1, 2, 3\}$ 
2: if  $k = 1$  then
3:    $\widehat{\text{CP}}(i) \leftarrow \widehat{\text{CP}}(i) + |W(j) \cup I(j)| + 1$ 
4:    $W(i) \leftarrow W(i) \cup W(j) \cup I(j) \cup \{j\}$ 
5:    $D(i) \leftarrow D(i) \cup W(j) \cup I(j) \cup \{j\}$ 
6:    $L(j) \leftarrow L(j) \cup L(i) \cup I(i) \cup \{i\}$ 
7:    $D(j) \leftarrow D(j) \cup L(i) \cup I(i) \cup \{i\}$ 
8: else if  $k = 2$  then
9:    $\widehat{\text{CP}}(i) \leftarrow \widehat{\text{CP}}(i) + |W(j)| + 1/2(1 + |I(j)|)$ 
10:   $\widehat{\text{CP}}(j) \leftarrow \widehat{\text{CP}}(j) + |W(i)| + 1/2(1 + |I(i)|)$ 
11:   $W(i) \leftarrow W(i) \cup W(j)$ ,  $W(j) \leftarrow W(i)$ 
12:   $L(i) \leftarrow L(i) \cup L(j)$ ,  $L(j) \leftarrow L(i)$ 
13:   $I(i) \leftarrow I(i) \cup I(j) \cup \{j\}$ ,  $I(j) \leftarrow I(i) \cup I(j) \cup \{i\}$ 
14:   $D(i) \leftarrow D(i) \cup D(j)$ ,  $D(j) \leftarrow D(i)$ 
15: else
16:  Same as for  $k = 1$  with  $i$  and  $j$  reversed
17: end if
18: Same steps as line 10 and 11 in SCORE-UPDATE
    
```

THEORETICAL RESULTS

Lower bounds

Informal Version: For \mathbf{P} with $\min_{i < j} |P_{i,j}^{(1)} - P_{i,j}^{(2)}| > \Delta$ the lower bounds are $\Omega(n^2 / \Delta^2 \ln 1/\delta)$,

 where $P_{i,j}^{(1)}, P_{i,j}^{(2)}, P_{i,j}^{(3)}$ are the order statistics of $P_{i,j}^{\succ}, P_{i,j}^{\cong}$ and $P_{i,j}^{\prec}$.

Formal Version: If A correctly identifies the COWI with confidence $1 - \delta$, then

$$\mathbb{E}[\tau^A(\mathbf{P})] \geq \ln \frac{1}{2.4\delta} \sum_{j \in \mathcal{A} \setminus \{i^*\}} C_j \min_{k \in L(j) \cup I(j)} \frac{1}{D_{j,k}(\mathbf{P})},$$

 where $\mathcal{C}(\mathbf{P}) = \{i^*\}$ and in the case with indifferences

$$D_{j,k}(\mathbf{P}) := \max\{\text{KL}_{j,k}^{(1)}, \text{KL}_{j,k}^{(2)}\}$$

$$\text{KL}_{j,k}^{(1)} = \text{KL}((P_{j,k}^{\succ}, P_{j,k}^{\cong}, P_{j,k}^{\prec}), (P_{j,k}^{\cong}, P_{j,k}^{\succ}, P_{j,k}^{\prec})),$$

$$\text{KL}_{j,k}^{(2)} = \text{KL}((P_{j,k}^{\succ}, P_{j,k}^{\cong}, P_{j,k}^{\prec}), (P_{j,k}^{\prec}, P_{j,k}^{\cong}, P_{j,k}^{\succ})),$$

$$C_j = \max_{(i,l) \in \Psi(j)} \frac{\binom{|I(i)|}{i} \binom{|L(j)|}{l} 1_{|i| \geq 1} + \binom{|I(j)|}{i} \binom{|L(i)|}{l} 1_{|j| \geq 1}}{\binom{|I(i)|}{i} \binom{|L(j)|}{l} 1_{|i| \geq 1} + \binom{|I(j)|}{i} \binom{|L(i)|}{l} 1_{|j| \geq 1}},$$

$$\Psi(j) := \{(i, l) \in \{0, \dots, |I(j)|\} \times \{0, \dots, |L(j)|\} \mid i + 2l \geq 2d_j + 1\}$$

 for any \mathbf{P} with $\min_{j,k} \min\{P_{j,k}^{\succ}, P_{j,k}^{\cong}, P_{j,k}^{\prec}\} > 0$.

Upper bounds

Informal Version: Worst-case sample complexities have the order

$$\frac{\text{POCOWISTA}}{\frac{n^2}{\Delta_{i,j}^2} \ln \left(\frac{n}{\sqrt{\delta}} \cdot \frac{1}{\Delta_{i,j}}\right)} \quad \frac{\text{TRA-POCOWISTA}^*}{\frac{n}{\Delta_{i,j}^2} \ln \left(\frac{n}{\sqrt{\delta}} \cdot \frac{1}{\Delta_{i,j}}\right)} \quad \frac{\text{SAVAGE}^{**} [3]}{\frac{n^2}{\Delta_{i,j}^2} \ln \left(\frac{n}{\delta} \cdot \frac{1}{\Delta_{i,j}}\right)} \quad \frac{\text{PBR-CCSO}^{**} [1]}{\frac{n^2}{\Delta_{i,j}^2} \ln \left(\frac{n}{\delta} \cdot \frac{1}{\Delta_{i,j}}\right)}$$

 *if \mathbf{P} is transitive

** if there are no indifferences

Formal Version: For any $\mathbf{P} = ((P_{i,j}^{\succ}, P_{i,j}^{\cong}, P_{i,j}^{\prec}))_{i < j}$, such that there exists no pair $i, j \in \mathcal{A}$ with $i \neq j$ and $P_{i,j}^{\succ} = P_{j,i}^{\prec} = 1/3$, it holds

(i) for A := POCOWISTA that

$$\mathbb{P}(\hat{i}_A \in \mathcal{C}(\mathbf{P}) \text{ and } \tau^A(\mathbf{P}) \leq t(\mathbf{P}, \delta)) \geq 1 - \delta,$$

 where $t(\mathbf{P}, \delta) \leq \sum_{i < j} t_0((P_{i,j}^{\succ}, P_{i,j}^{\cong}, P_{i,j}^{\prec}), \delta / \binom{n}{2})$,

$$t_0((p_1, p_2, p_3), \delta) = \frac{c_1 p_{(1)}}{(p_{(1)} - p_{(2)})^2} \ln \left(\frac{\sqrt{2} c_2 p_{(1)}}{\sqrt{\delta} (p_{(1)} - p_{(2)})} \right), \quad (1)$$

 $p_{(1)} \geq p_{(2)} \geq p_{(3)}$ is the order statistic of p_1, p_2, p_3 , $c_1 = 194.07$, and $c_2 = 79.86$.

 (ii) for A := TRA-POCOWISTA if \mathbf{P} transitive that

$$\mathbb{P}(\hat{i}_A \in \mathcal{C}(\mathbf{P}) \text{ and } \tau^A(\mathbf{P}) \leq \tilde{t}(\mathbf{P}, \delta)) \geq 1 - \delta,$$

 where $\tilde{t}(\mathbf{P}, \delta) = \sum_{e=1}^E t_0((P_{i_e, j_e}^{\succ}, P_{i_e, j_e}^{\cong}, P_{i_e, j_e}^{\prec}), \delta/n)$, t_0 is as in (1) and $E \leq n$.