

# Ranking Interruptus: When Truncated Rankings Are Better and How to Measure That

Enrique Amigó, Stefano Mizzaro, Damiano Spina









### Which Ranking is Better?



### Why Truncated Rankings?





#### **Conversational Search**

**Small Screens** 

## How to Pick the Right Metric?



"My system is the best one!"





#### **Axiomatic Analysis**

System Effectiveness vs. User Satisfaction

#### **Axiomatic Analysis of Effectiveness Metrics**



### Contributions

- 1. Formal properties for truncated rankings
- 2. Extension of traditional metrics (ERR, NDCG, and RBP) by adding a user effort factor
- 3. Theoretical analysis of effectiveness metrics:
  - De-facto standard metrics do not satisfy desirable properties to evaluate truncated rankings
  - Observational Information Effectiveness (OIE) satisfies them all
- 4. Empirical evidence using 9 TREC test collections
- 5. Guidelines for metric selection in different ranking scenarios

	>~8
BASIC	PREMIUM
Ø	
Ø	<b>Ø</b>
Ø	
×	
×	•
	<b>•</b>
8	0



#### **Desirable Properties**

		Non-Truncated					Truncated			
		Metric	Priority	Top-W.	DEEP. TH.	SHALL. TH.	CONFIDENCE	Recall	Redundancy	
	FULL RANKING METRICS	Spearman Kendall AUC-ROC	$\checkmark$ $\checkmark$ $\checkmark$			$\checkmark$ $\checkmark$	√ √ √			] /
Metrics GR REI PRO USI BA TEI DO BA	BINARY RELEVANCE	P@N R@N R-p RR F-measure AP	✓	✓	$\checkmark$	$\checkmark$	$\checkmark$	~ ~ ~ ~ ~	$\checkmark$	
	GRADED RELEVANCE	NDCG Q-measure BPref	$\checkmark$ $\checkmark$ $\checkmark$	$\frac{1}{\sqrt{2}}$		$\checkmark$ $\checkmark$ $\checkmark$		$\checkmark$ $\checkmark$ $\checkmark$		
	PROBABILISTIC USER MODEL BASED	ERR RBP iRBU	$\checkmark$ $\checkmark$ $\checkmark$	$\checkmark$ $\checkmark$ $\checkmark$	$\checkmark$ $\checkmark$ $\checkmark$	$\checkmark$			✓ ✓	
	TERMINAL DOCUMENT BASED	NDCGT ERRT RBPT	$\checkmark$ $\checkmark$ $\checkmark$	$\checkmark$ $\checkmark$ $\checkmark$	√ √	√ √	√ √ √	~	~	
	UTILITY BASED	Flat Utility RBU DCGU ERRU RBPU	$\begin{array}{c} \checkmark \checkmark$	$\langle \cdot \rangle$	✓ ✓ ✓	$\checkmark$ $\checkmark$ $\checkmark$			√ √	
	INFORMATION BASED	OIE	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	

Designed for non-truncated rankings

Add one 'relevant' terminal document at the bottom of the ranking

Add a component to take "effort" into account

Based on Information Theory: quantity of information obtained from ranking and ground truth

### **Empirical Validation**

- Nine TREC test collections
- Optimal Truncation Points
- Truncated vs. Full Ranking Effectiveness
- Validation of the aspects captured by truncated rankings

# (Some) Results

- Confirmation of theoretical analysis
  - Observational Information Effectiveness (OIE)
  - RBPU: extension of RBP that incorporates effort penalty
- For most metrics, optimal truncation is rather substantial



Figure 2: The distributions of optimal ranking lengths according to different truncation metrics. Results are averaged across the nine datasets.

#### **Guidelines for Metric selection**

To help practitioners identifying which effectiveness metric should be used for different ranking tasks / conditions





### Summary



Ranking Interruptus: When Truncated Rankings Are Better and How to Measure That

Stefano Mizzaro

Enrique Amigó

Damiano Spina







Evaluation of Truncated Rankings

Axiomatic Analysis + Empirical Validation Guidelines to Choose the Appropriate Metric



#### Future work (long road ahead!)

- Experiments with systems which incorporate a stopping criterion
- User studies to calibrate penalty effort

- ...