# Recommending ordered sections with similar section filtering to help structuring Wikipedia articles

*Supervisor*
Prof. Dr. Philippe Cudré-Mauroux
eXascale Infolab, Department of Informatics, University of Fribourg

*Co-Supervisor*
Natalia Ostapuk
eXascale Infolab, Department of Informatics, University of Fribourg

Sergiy Goloviatinski, University of Neuchâtel
Master thesis
03.02.2022

# Outline

- Context and motivation
- Related work
- Improvements of existing method
- Redundant sections filtering
- Section ordering
- Prototype demo
- Conclusion

# Context and motivation

- Wikipedia articles are structured with sections
- No official guidelines on structure, only few community made
  - → Need to look at examples
- Good quality articles have more sections (7.32 on average compared to 4.63 for all articles)
  - →More sections would improve quality
- 37% of articles are marked as stubs
- 87% of sections are used only in one article (e.g. "Early life of Dionysius the Elder")
  - → Need for standardization
- →An algorithm for section recommendation would respond to an existing need

# Structuring Wikipedia Articles with Section Recommendations



- Use Wikipedia category network
- Recommended sections → ranked by P(Section | Category)

Tiziano Piccardi, et al. Structuring wikipedia articles with section recommendations. *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval,* 2018
https://dl.acm.org/doi/pdf/10.1145/3209978.3209984

# Structuring Wikipedia Articles with Section Recommendations: Idea

People

Early life
Personal life
Death

is a

is a

Scientists

Actors

Early life
Personal life
Death
Awards

Early life
Personal life
Death
Filmography

- Use Wikipedia category network as ontology
- A scientist is a person
- An actor is also a person
- Therefore, sections from articles about scientists and actors can be recommended for articles about persons

# Structuring Wikipedia Articles with Section Recommendations: Problem

## Category:French artists

From Wikipedia, the free encyclopedia

### Subcategories

- ▶ French cartoonists (4 C, 110 P) ← is a(french cartoonist, french artist) ✅
- ▶ French ceramists (2 C, 31 P)
- ▶ French cinematographers (3 C, 117 P)
- ▶ French comics artists (3 C, 168 P)
- ▶ French conceptual artists (28 P)
- ▶ French contemporary artists (1 C, 164 P)

● Wikipedia category network is not guaranteed to be an ontology

is a (french artists, french people in arts occupations) ✅

is a(french artists, french art) ❌

Categories: French people in arts occupations │ Artists by nationality │ French art

# Structuring Wikipedia Articles with Section Recommendations: Method (1)

55 DBpedia types

About: Pablo Picasso

An Entity of Type: person, from Named Graph: http://dbpedia.org, within Data Space: dbpedia.org

Pablo Ruiz Picasso (25 October 1881 – 8 April 1973) was a Spanish painter, sculptor, printmaker, ceramicist and theatre designer who spent most of his adult life in France. Regarded as one of the most influential artists of the 20th century, he is known for co-founding the Cubist movement, the invention of constructed sculpture, the co-invention of collage, and for the wide variety of styles that he helped develop and explore. Among his most famous works are the proto-Cubist Les Demoiselles d'Avignon (1907), and Guernica (1937), a dramatic portrayal of the bombing of Guernica by German and Italian air forces during the Spanish Civil War.

- Need a way to measure if a category relationship is ontological
- Using an external source: DBpedia
- Entities on DBpedia about Wikipedia articles
- Use type attribute

# Structuring Wikipedia Articles with Section Recommendations: Method (2)

Person + Location + Organization



Person                Location            Organization

# Structuring Wikipedia Articles with Section Recommendations: Method (3a)



Article type distribution

- Mark which categories are pure, based on article type distribution
- Bottom-up approach
- If pure, add articles to parent
- Continue while categories pure in the hierarchy

# Structuring Wikipedia Articles with Section Recommendations: Method (3b)



Article type distribution

- Because child categories are pure, add their articles to parent category
- Article type distribution of parent contains articles from children
- But now, there are too many different types in the parent category

# Structuring Wikipedia Articles with Section Recommendations: Method (3c)

- The category is marked as not pure, we prune the graph

11

# Structuring Wikipedia Articles with Section Recommendations: Method (4)

Gini coefficient to quantify purity of article type distribution: minimum 0.966



gini ≥ 0.966

gini < 0.966

# Backup slide: gini coefficient



| 1. Perfect Equality | 2. Unequal | 3. More Unequal | 4. Total Inequality |
|---|---|---|---|
| $G_1 = 0$ | $G_2 > 0$ | $G_3 > G_2$ | $G_4 = 1$ |

Cumulative %Participation (y-axis), Cumulative %Population (x-axis), 0% to 100%

# Structuring Wikipedia Articles with Section Recommendations: Pipeline

14

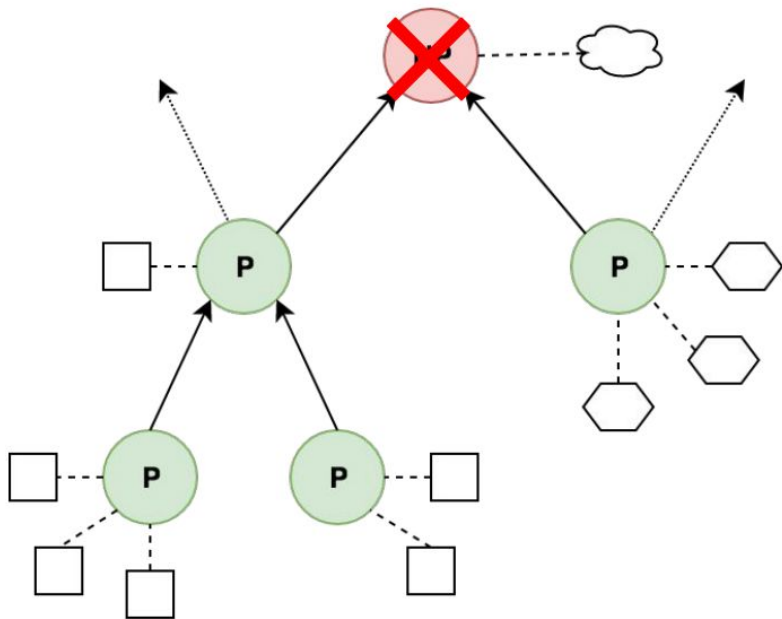# Structuring Wikipedia Articles with Section Recommendations: Limitations

### Semantically redundant sections

1. Gameplay
2. Story
3. Plot
4. Development and release
5. Development
6. Reception
7. Sequel

*Article about video game*

### No logical section ordering

1. First round (0.8) ← P(Section | Category)
2. Second round (0.8)
3. Third round (0.8)
4. Final (0.8)
5. Fourth round (0.6)
6. Quarter-finals (0.6)
7. Semi-finals (0.6)
8. Fifth round (0.4)

*Article about Scottish football cup*

# Improvement of existing method

# Improvement of existing method: unknown article types (1)

About: Albedo
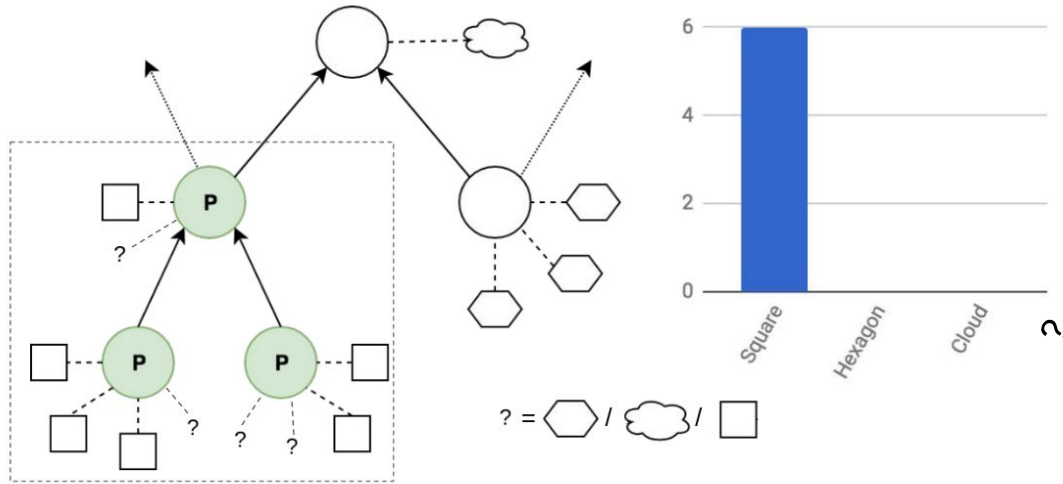
?

An Entity of Type: Thing, from Named Graph: http://dbpedia.org, within Data Space: dbpedia.org

Albedo (/ælˈbiːdoʊ/; from Latin albedo 'whiteness') is the measure of the diffuse reflection of solar radiation out of the total solar radiation and measured on a scale from 0, corresponding to a black body that absorbs all incident radiation, to 1, corresponding to a body that reflects all incident radiation. The term albedo was introduced into optics by Johann Heinrich Lambert in his 1760 work Photometria.

- "Thing" is too generic for a type, not used in method
- Therefore we don't know the type of the article

# Improvement of existing method: unknown article types (2)



- Unknown article types not included in the article type distribution
- Those articles could be of any type
- Categories containing unknown article types are marked as pure, but maybe it would not be the case if those article types would be known

# Improvement of existing method: example of noise in recommendations

1. Early life
2. Biography
3. Career
4. History
5. Death

*Article about "Monument types"*

- Sections in red not relevant for monuments
- Those sections in red come from the category "Names inscribed under the Arc de Triomphe"

# Solution

- Filter out unknown types articles from dataset
- Reduced number of articles from ~3.6 million to ~2 million
- 56% of articles remained

1. Early life
2. Biography
3. Career
4. History
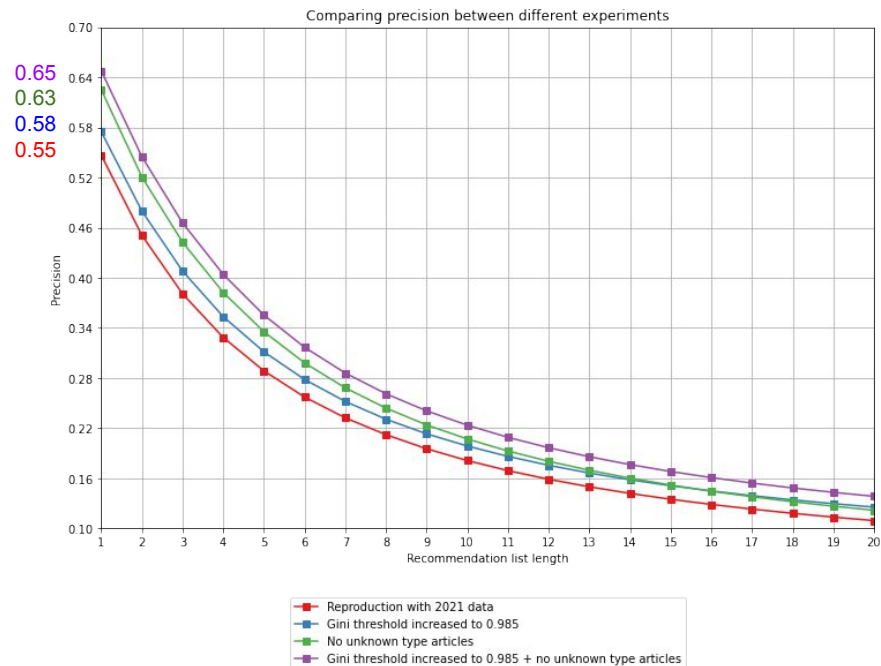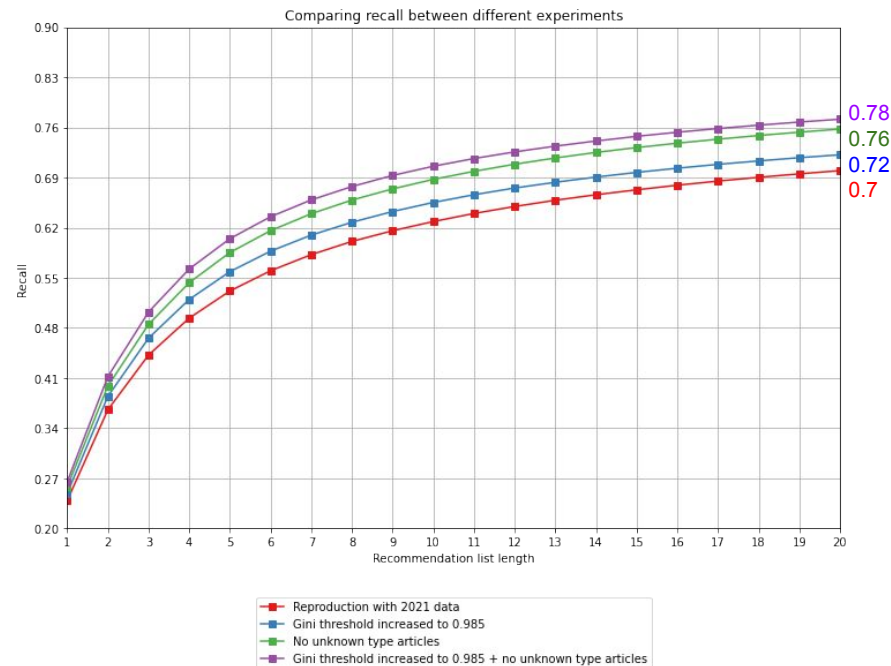5. Death

remove articles with unknown types

1. History
2. Gallery
3. Architecture
4. Location
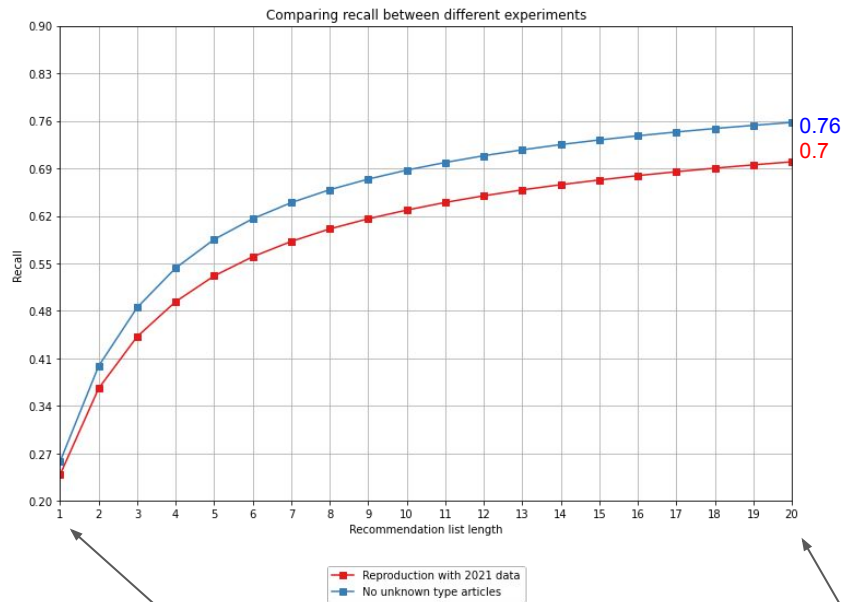5. Description

*Article about "Monument types"*

# (backup slide) Improvement to existing method: increasing gini threshold

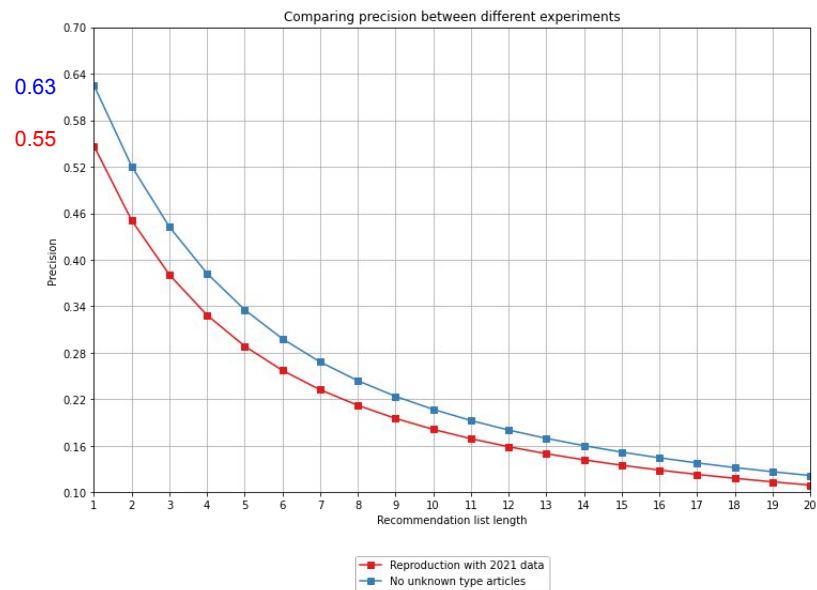- Gini threshold increased from 0.966 to 0.985 gave us a better performance

# (backup slide) Existing method improvement: results

# Existing method improvement: results



Comparing recall between different experiments

0.76
0.7

most frequent section in category



Comparing precision between different experiments

0.63
0.55

20 most frequent sections in category

23

# Redundant sections filtering

# Redundant sections filtering: idea

Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks

Nils Reimers and Iryna Gurevych

Redundant sections = semantically similar section contents

Plot [edit]

Section content A → embedding → **A**

Synopsis [edit]

Section content B → embedding → **B**

Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*
https://arxiv.org/abs/1908.10084

semantic similarity(Plot, Synopsis) = cosine similarity(**A**, **B**)

# Redundant sections filtering: problem

- Which sections from which article to compare ?
- e.g. "History" section has different content depending on context
- Group categories which are likely to have articles in common
- →group sections that are likely to appear together in recommendation lists
- Find semantically similar sections from articles in each group

(backup slide) Redundant sections filtering: group categories by context (1)



$$\frac{|Articles(C1) \cap Articles(C3)|}{|Articles(C1) \cup Articles(C3)|} = \frac{|\{A2\}|}{|\{A1, A2\}|} = \frac{1}{2}$$

Articles(x): articles belonging to category x

# (backup slide) Redundant sections filtering: group categories by context (2)



Community detection

Louvain method

# Redundant sections filtering: method (1)



**Category group A**

Article 1
Section X
Content X$_1$

Article 2
Section X
Content X$_2$

Article 3
Section X
Content X$_3$

Article 4
Section Y
Content Y$_1$

Article 5
Section Y
Content Y$_2$

embedding →

Embedding X$_1$

Embedding X$_2$

Embedding X$_3$

Embedding Y$_1$

Embedding Y$_2$

find similar pairs →

**Category group A**

| **Section X** | **cosine sim>0.5** | **Section Y** |
|---|---|---|
| Embedding X$_1$ | 0.8 | Embedding Y$_1$ |
| Embedding X$_2$ | 0.6 | Embedding Y$_2$ |
| Embedding X$_3$ | 0.7 | |

mean cosine sim(section X, section Y, A) = 0.7
# similar section contents(section X, section Y, A) = 3

# Redundant sections filtering: method (2)



- Distribution of cosine similarities different in each category group
- Defined semantic filtering level to decide if redundant sections will be filtered
- Given level corresponds to a cosine similarity threshold different in each category group
- Section pairs with cosine similarity above threshold considered as redundant

- Level 0: no filtering
- Level 1: filter top 1/3 most similar
- Level 2: filter top 2/3 most similar
- Level 3: filter all sections detected as similar

# Redundant sections filtering: method (3)

Similar section pairs from same category group with cosine sim > threshold given by semantic filtering level

1. Early life
2. Biography
3. Life
4. Career
5. Personal
6. Later career
7. Appreciation

*Article about classical musician*

Pairs of semantically similar sections

(Biography, Life)
(Life, Early Life)
(Later career, Life)
(Career, Later career)
(Career, Appreciation)

$$\frac{\#\ similar\ section\ contents(Biography, Life)}{\#\ section\ contents(Biography) \cdot \#\ section\ contents(Life)}$$

Biography — 0.0012 — Life

Later_career

0.0001

0.0005

Early_life

0.0003

Career — 0.0003 — Appreciation

Community detection with Louvain method

1. Biography
2. Career
3. Personal

Keep section x in each similar section group with highest
$$\#\ section\ contents(x) \cdot mean(edge\ weights(x))$$

# Redundant sections filtering: example

1. Gameplay
2. Story
3. Plot
4. Development and release
5. Development
6. Reception
7. Sequel

*Article about video game*

redundant sections filtering

1. Gameplay
2. Plot
3. Development
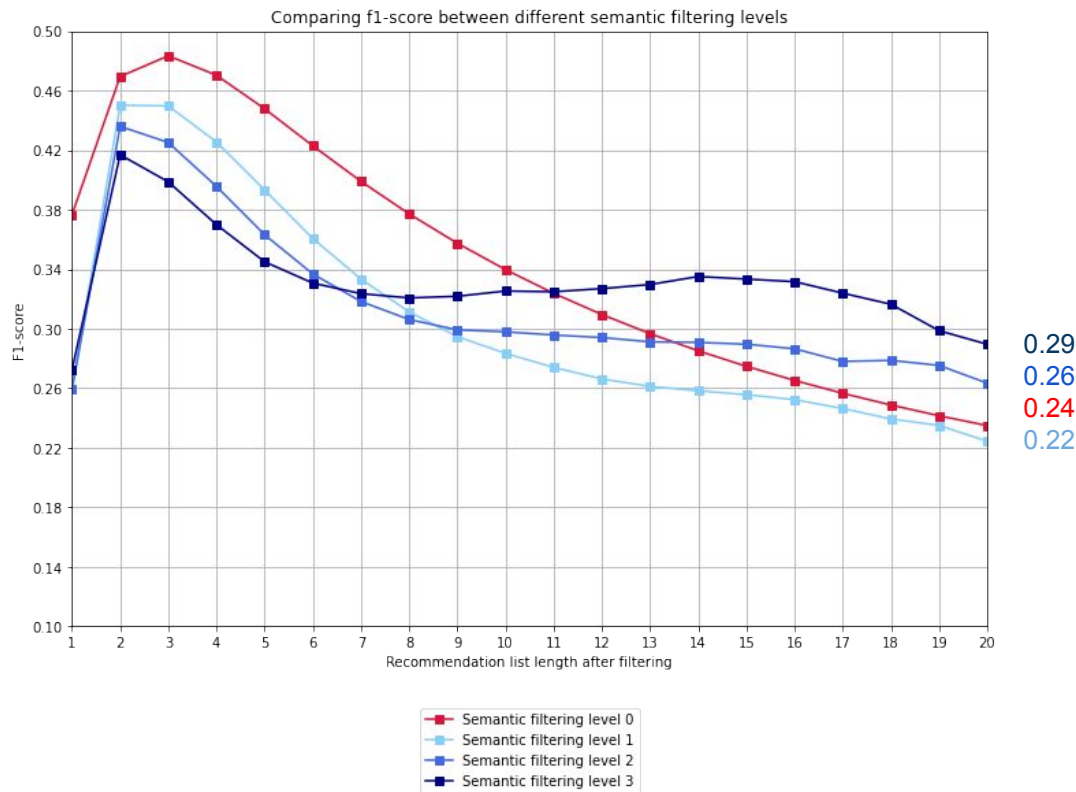4. Reception
5. Sequel

# Redundant sections filtering: results



Comparing f1-score between different semantic filtering levels

0.29
0.26
0.24
0.22

# Section ordering

# Section ordering: method (1)

- Order value: relative position of section in article
- Beginning to end order value: span of the section inside article

**Category X**

**Article 1**

1. Section A — 0
2. Section B — 0.25
3. Section C — 0.5
4. Section D — 0.75
   — 1

**Article 2**

1. Section A — 0
2. Section B — 0.5
   — 1

Beginning order value for A

End order value for A, beginning order value for B

Average beginning and end order values by section in each category:
Beginning order value(A) = (0 + 0)/2 = 0
End order value(A) = (0.25 + 0.5)/2 = 0.375

35

# Section ordering: method (2)

begin(x) : average beginning order value of section x
end(x) : average end order value of section x

**Category 1**

begin(A) = 0
end(A) = 0.2

begin(B) = 0.1
end(B) = 0.3

**Category 2**

begin(B) = 0.2
end(B) = 0.6

begin(C) = 0.5
end(C) = 1

$$midpoint(x) = \frac{begin(x) + end(x)}{2}$$

**Recommended sections**

- Section B
- Section C
- Section A

| X | begin(X) | end(X) | midpoint(X) |
|---|----------|--------|-------------|
| B | (0.1+0.2)/2 = 0.15 | (0.3+0.6)/2 = 0.45 | (0.15+0.45)/2 = 0.3 |
| C | 0.5 | 1 | (0.5+1)/2 = 0.75 |
| A | 0 | 0.2 | (0+0.2)/2 = 0.1 |

**Ordered recommended sections**

1. Section A (0.1)
2. Section B (0.3)
3. Section C (0.75)

# Section ordering: example

ordered by P(Section | Category)

ordered by midpoint(Section)

1. First round (0.8)
2. Second round (0.8)
3. Third round (0.8)
4. Final (0.8)
5. Fourth round (0.6)
6. Quarter-finals (0.6)
7. Semi-finals (0.6)
8. Fifth round (0.4)

ordering

1. First round (0.22)
2. Second round (0.33)
3. Third round (0.44)
4. Fourth round (0.53)
5. Fifth round (0.59)
6. Quarter-finals (0.7)
7. Semi-finals (0.81)
8. Final (0.9)

*Article about Scottish football cup*

# (backup slide) Section ordering: evaluation, kendall's tau

Mirella Lapata. Automatic evaluation of information ordering: Kendall's tau. *Computational Linguistics*, *Volume 32, Number 4*, 2006
https://aclanthology.org/J06-4002/

$$\tau = 1 - \frac{2 \cdot I}{0.5 \cdot n(n-1)}$$

A = [1,2,3,4]

B = [1,3,2,4]

with:
- I : nb of intersections
- n : nb of elements

C = [4,3,2,1]

tau(A,A) = 1

$$tau(A,B) = 1 - \frac{2 \cdot 1}{0.5 \cdot 4(4-1)} = 1 - \frac{2}{6} = \frac{2}{3}$$

tau(A,C) = -1

# (backup slide) Section ordering: experiments

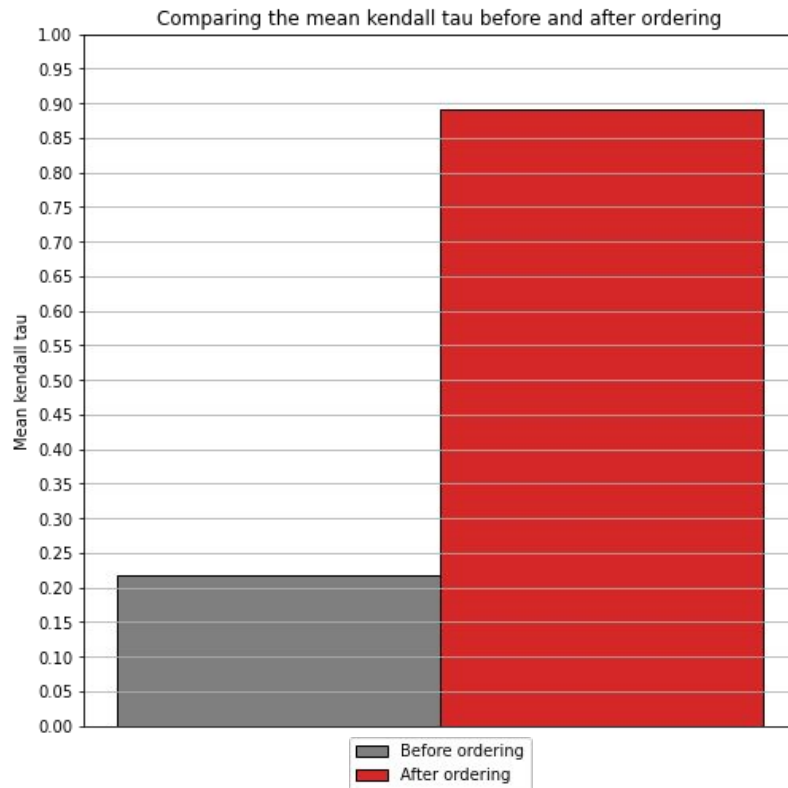- Order by ascending mean(beginning order value,end order value), different for each category
- Order by ascending beginning order value, different for each category
- Order by ascending mean(beginning order value,end order value), among all categories (as if all articles in same category)

# (backup slide) Section ordering: results for different experiments



Comparing the mean kendall tau before and after ordering

# Section ordering: results

Averaged for recommendation
lists of size 20 which had at
least 2 sections in common with
ground truth



Comparing the mean kendall tau before and after ordering

Legend:
- Before ordering
- After ordering

# Prototype demo

http://127.0.0.1:5000

# Conclusion

- Improved performance of existing method (precision@1 0.55→0.65 recall@20 0.7 →0.77)
- Added two features
  - Redundant section filtering (f1 score@20 0.24→0.29)
  - Order sections logically (average Kendall's tau 0.22→0.89)
- Implemented prototype to demonstrate use case