

Linear Score Normalization Revisited

Ilya Markov

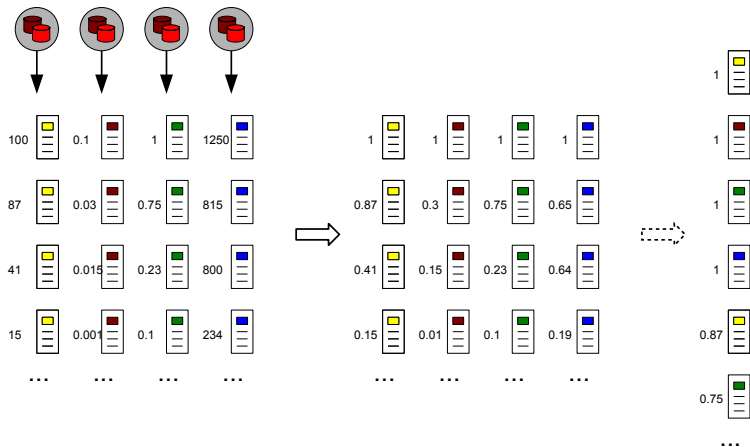
ilya.markov@usi.ch

University of Lugano

Outline

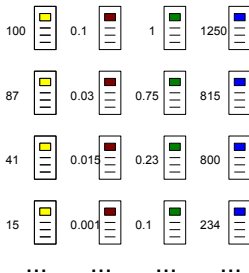
- 1 Score Normalization
- 2 Linear Score Normalization
- 3 Improving Linear Score Normalization

Score Normalization

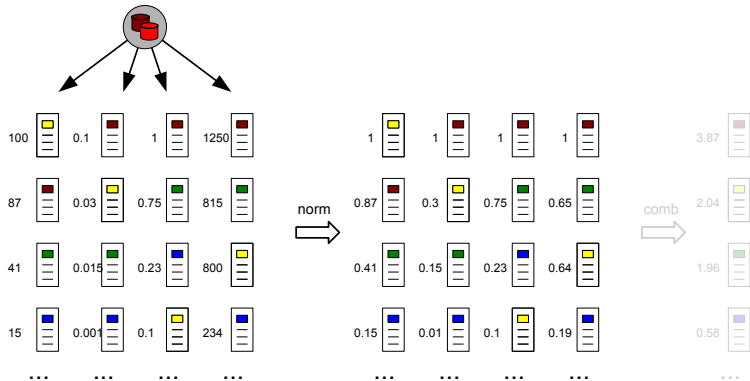


Assumptions

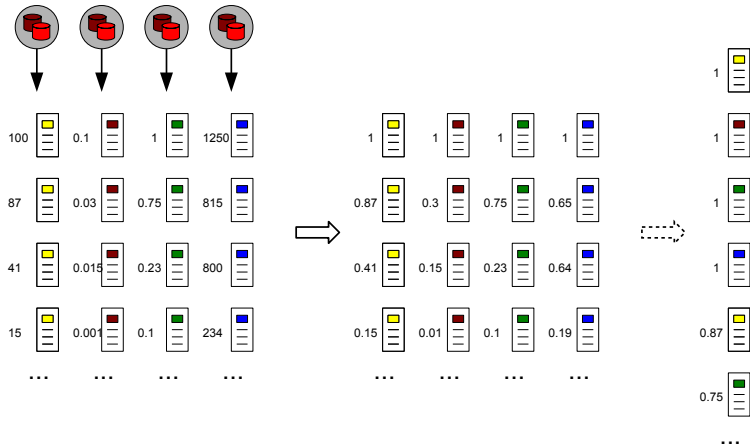
- 1 Document scores are provided
- 2 Document lists are disjoint



Data Fusion



Score Normalization



Outline

- 1 Score Normalization
- 2 Linear Score Normalization**
- 3 Improving Linear Score Normalization

Linear Score Normalization

▶ MinMax

$$s_{MinMax} = \frac{s - s_{min}}{s_{max} - s_{min}}$$

▶ Z-Score

$$s_{ZScore} = \frac{s - \mu}{\sigma}$$

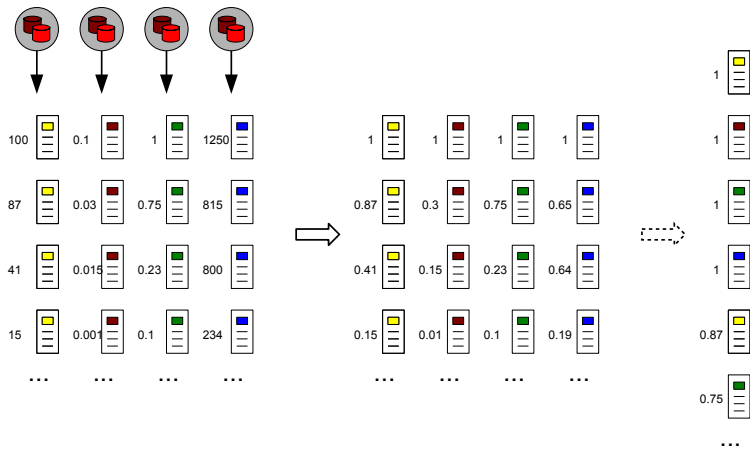
▶ Sum

$$s' = s - s_{min}, \quad s_{Sum} = \frac{s'}{\sum_i s'_i}$$

MinMax

$$S_{MinMax} = \frac{S - S_{min}}{S_{max} - S_{min}}$$

MinMax



MinMax Assumptions

- A1 Each collection contains at least 1 relevant document
- A2 This document is most likely to be ranked 1st

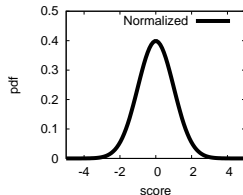
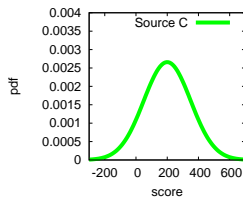
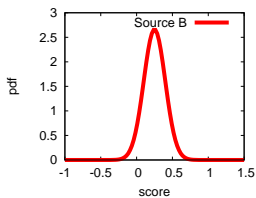
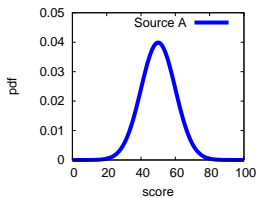
MinMax Results

# sources	top-10		top-100		top-1000	
	MAP	P@10	MAP	P@10	MAP	P@10
30	0.063	0.213	0.110	0.213	0.111	0.213
1000	0.015	0.037	0.017	0.037	0.016	0.037

Z-Score

$$S_{ZScore} = \frac{s - \mu}{\sigma}$$

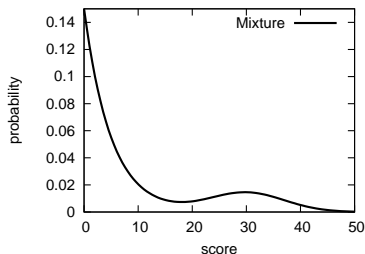
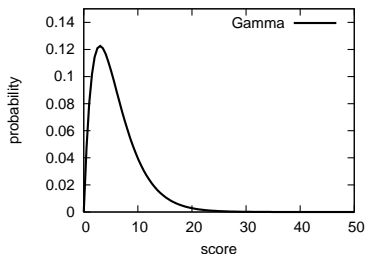
Z-Score



Z-Score Assumptions

- A1 Each collection contains relevant documents
- A2 These documents are likely to be ranked high
- A3 Document score distributions have the bell shape

Document Score Distributions



Z-Score Results

# sources	top-10		top-100		top-1000	
	MAP	P@10	MAP	P@10	MAP	P@10
30	0.057	0.168	0.059	0.072	0.066	0.133
1000	0.016	0.047	0.010	0.034	0.011	0.035

Sum

$$s' = s - s_{min}, \quad s_{Sum} = \frac{s'}{\sum_i s'_i}$$

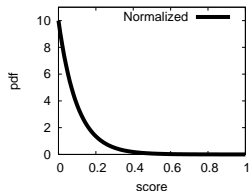
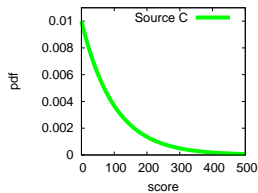
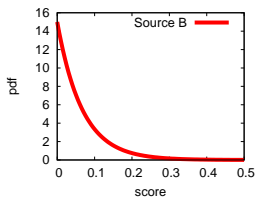
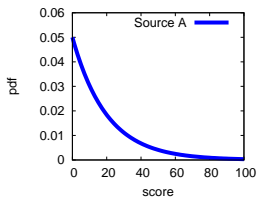
Sum

$$s \sim \mathcal{E}(s; \lambda), \quad \lambda = \frac{n}{\sum_i s_i}$$

$$s_{Sum} \sim \mathcal{E}(s_{Sum}; \lambda_{Sum})$$

$$\lambda_{Sum} = \frac{n}{\sum_i s_{Sum}} = n$$

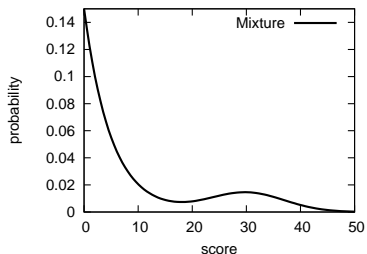
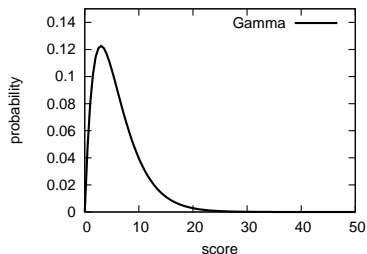
Sum



Sum Assumptions

- A1 Each collection contains relevant documents
- A2 These documents are likely to be ranked high
- A3 Document scores are distributed exponentially

Document Score Distributions



Sum Results

# sources	top-10		top-100		top-1000	
	MAP	P@10	MAP	P@10	MAP	P@10
30	0.056	0.158	0.052	0.054	0.035	0.041
1000	0.015	0.039	0.006	0.024	0.004	0.019

Linear Score Normalization Assumptions

- A1 Each collection contains relevant documents
- A2 These documents are likely to be ranked high
- A3 [Z-Score and Sum] Document scores are distributed in a certain way

Linear Score Normalization Results

	top-10		top-100		top-1000	
	MAP	P@10	MAP	P@10	MAP	P@10
MinMax	0.063 ^Δ	0.213 ^Δ	0.110 ^Δ	0.213 ^Δ	0.111 ^Δ	0.213 ^Δ
Z-Score	0.057	0.168	0.059	0.072	0.066	0.133
Sum	0.056 [∇]	0.158 [∇]	0.052 [∇]	0.054 [∇]	0.035 [∇]	0.041 [∇]

Table: 30 sources

	top-10		top-100		top-1000	
	MAP	P@10	MAP	P@10	MAP	P@10
MinMax	0.015	0.037	0.017 ^Δ	0.037	0.016 ^Δ	0.037
Z-Score	0.016	0.047	0.010	0.034	0.011	0.035
Sum	0.015	0.039	0.006 [∇]	0.024 [▼]	0.004 [∇]	0.019 [∇]

Table: 1000 sources

$$S_{MinMax} = \frac{S - S_{min}}{S_{max} - S_{min}}$$

Outline

- 1 Score Normalization
- 2 Linear Score Normalization
- 3 Improving Linear Score Normalization**

Linear Score Normalization Assumptions

- A1 Each collection contains relevant documents
- A2 These documents are likely to be ranked high
- A3 [Z-Score and Sum] Document scores are distributed in a certain way

Dealing with A1 Assumption

- 1 Force data to satisfy A1
- 2 *Change methods to relax A1*

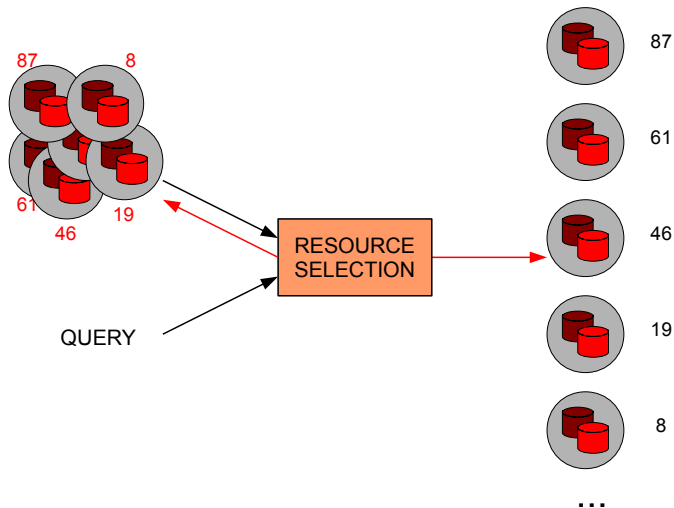
Weighted Linear Score Normalization

$$s_{weighted} = w(C|q) \cdot s_{linear}$$

Calculating Weights

- 1 Based on a ranked list itself
- 2 *Based on external evidence*

Resource Selection



CORI Results Merging

$$s_{CORI} = \frac{1 + 0.4 \cdot s_{MinMax}(C|q)}{1.4} \cdot s_{MinMax}$$

Parameterized CORI Results Merging

$$s_{CORI} = \frac{1 + \lambda \cdot s_{MinMax}(C|q)}{1 + \lambda} \cdot s_{MinMax}$$

CORI Results

	top-10		top-100		top-1000	
	MAP	P@10	MAP	P@10	MAP	P@10
CORI(cv)	0.092 ^Δ	0.376 [▲]	0.154 ^Δ	0.353	0.157	0.361
CORI	0.073	0.349	0.129	0.369 [▲]	0.162	0.371
SAFE	0.068	0.273	0.087	0.241	0.085	0.234
MinMax	0.063	0.213	0.110	0.213	0.111	0.213

Table: 30 sources

	top-10		top-100		top-1000	
	MAP	P@10	MAP	P@10	MAP	P@10
CORI(cv)	0.049 ^Δ	0.281	0.065 ^Δ	0.289	0.065	0.282
CORI	0.028	0.269	0.042	0.287	0.056	0.292 [▲]
SAFE	0.031	0.116	0.030	0.114	0.026	0.105
MinMax	0.015	0.039	0.006	0.024	0.004	0.019

Table: 1000 sources

$$S_{MinMax} = \frac{S - S_{min}}{S_{max} - S_{min}}$$

$$s_{MinMax} = w(C|q) \cdot s_{MinMax}$$

$$s_{CORI} = \frac{1 + \lambda \cdot s_{MinMax}(C|q)}{1 + \lambda} \cdot s_{MinMax}$$

Answers