

Want to Normalize Scores? Ask Me How!

Unsupervised Linear Score Normalization:
Intuition, Assumptions and Performance

*Ilya Markov*¹, *Avi Arampatzis*², *Fabio Crestani*¹

¹University of Lugano,
Lugano, Switzerland

²Democritus University of Thrace,
Xanthi, Greece

July 2, 2012

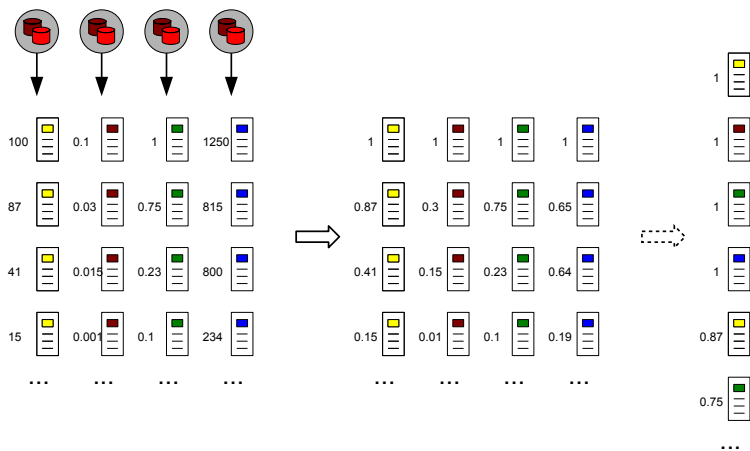
$$S_{norm} = \frac{s}{max}$$

$$S_{norm} = W \frac{s}{max}$$

Outline

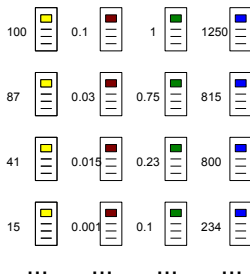
- 1 Introduction
- 2 Linear Score Normalization
- 3 Improving Linear Score Normalization

Score Normalization

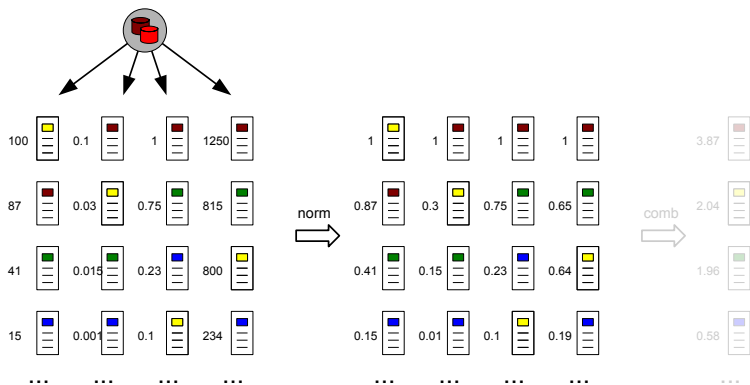


Assumptions

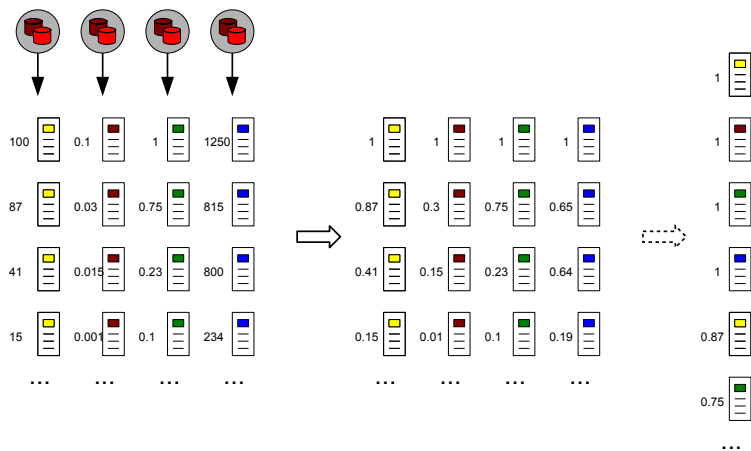
- 1 Document scores are provided
- 2 Document lists are disjoint



Data Fusion



Score Normalization



Outline

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

Linear Score Normalization

- MinMax

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

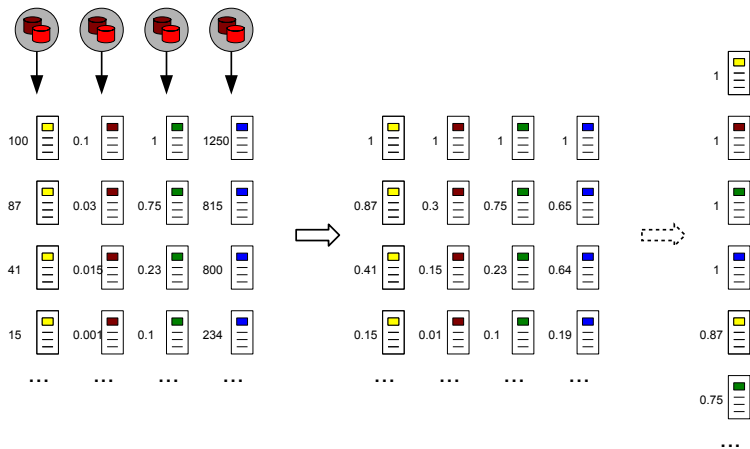
- Z-Score

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

- Sum

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

MinMax



$$S_{norm} = \frac{s - \min}{\max - \min}$$

MinMax

- Formula

$$S_{norm} = \frac{s - \min}{\max - \min}$$

- Assumptions

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

- Discussion

- 1st documents are ranked before any other in the merged list.
- High early precision is achieved when assumptions are hold.

- Results

	gov2.30		gov2.1000	
	p@10	p@20	p@10	p@20
top-10	0.1966	0.1936	0.0161	0.0148
top-1000	0.1953	0.1919	0.0161	0.0148

Max

- Observations
 - The lowest theoretical score of many scoring functions is 0.
- Formula

$$s_{norm} = \frac{s}{max}$$

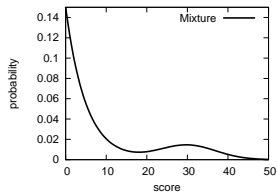
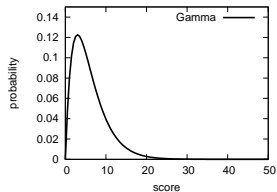
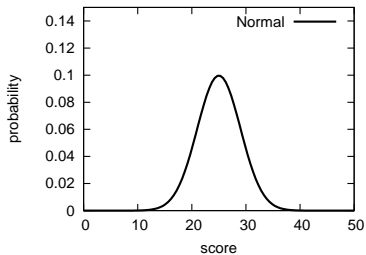
- Results

		gov2.30		gov2.1000	
		p@10	p@20	p@10	p@20
top-10	MinMax	0.1966	0.1936	0.0161	0.0148
	Max	0.1953	0.1926	0.0161	0.0148
top-1000	MinMax	0.1953	0.1919	0.0161	0.0148
	Max	0.1953	0.1919	0.0161	0.0148

- Discussion
 - Performance similar to MinMax.
 - Not affected by minimum score.

Z-Score

$$S_{norm} = \frac{S - \mu}{\sigma}$$



Z-Score

- Formula

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

- Assumptions

- Score distribution is a bell-shape curve.

- Discussion

- May be true for top documents.
- Low ranked documents affect normalized scores of top ranked documents.

- Results

	gov2.30		gov2.1000	
	p@10	p@20	p@10	p@20
top-10	0.1839	0.1852	0.0007	0.0037
top-1000	0.0819	0.0953	0.0013	0.0023

Unit Variance (UV)

- Formula

$$S_{norm} = \frac{S}{\sigma}$$

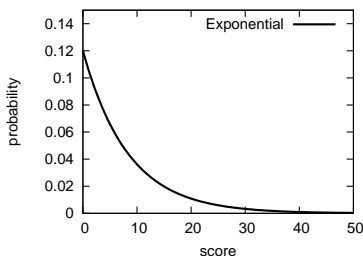
- Results

		gov2.30		gov2.1000	
		p@10	p@20	p@10	p@20
top-10	Z-Score	0.1839	0.1852	0.0007	0.0037
	UV	<i>0.1980</i>	<i>0.1936</i>	<i>0.0228</i>	<i>0.0188</i>
top-1000	Z-Score	0.0819	0.0953	0.0013	0.0023
	UV	<i>0.1248</i>	<i>0.1144</i>	0.0013	0.0007

- Discussion
 - More robust and less affected by low scores.

Sum

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



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

$$p(s_{norm}) = p\left(\frac{s'}{\sum_i s'_i}\right) \sim \mathcal{E}(s; n)$$

Sum

- Formula

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

- Assumptions

- Document scores are distributed exponentially.
- The mean of different exponentials is made the same.

- Discussion

- Should not hold especially for top documents.

- Results

	gov2.30		gov2.1000	
	p@10	p@20	p@10	p@20
top-10	0.1785	0.1758	0.0013	0.0023
top-1000	0.0134	0.0114	0.0013	0.0017

Overall Results

	gov2.30		gov2.1000	
	p@10	p@20	p@10	p@20
MinMax	0.1966	0.1936	0.0161	0.0148
Max	0.1953	0.1926	0.0161	0.0148
Z-Score	0.1839	0.1852	0.0007	0.0037
UV	0.1980	0.1936	<i>0.0228</i>	<i>0.0188</i>
Sum	0.1785	0.1758	0.0013	0.0023

Table: Top 10 documents are retrieved.

	gov2.30		gov2.1000	
	p@10	p@20	p@10	p@20
MinMax	<i>0.1953</i>	<i>0.1919</i>	<i>0.0161</i>	<i>0.0148</i>
Max	<i>0.1953</i>	<i>0.1919</i>	<i>0.0161</i>	<i>0.0148</i>
Z-Score	0.0819	0.0953	0.0013	0.0023
UV	0.1248	0.1144	0.0013	0.0007
Sum	0.0134	0.0114	0.0013	0.0017

Table: Top 1000 documents are retrieved.

Outline

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

$$S_{norm} = \frac{s}{max}$$

$$S_{norm} = W \frac{s}{max}$$

Possible Weights for MinMax

... calculated directly from a list of document scores.

- Min and max
- Mean and standard deviation
- Median and median absolute deviation
- Sum and mean

MinMax * Stdv

- Formula

$$s_{norm} = \sigma \frac{s - min}{max - min}$$

- Results

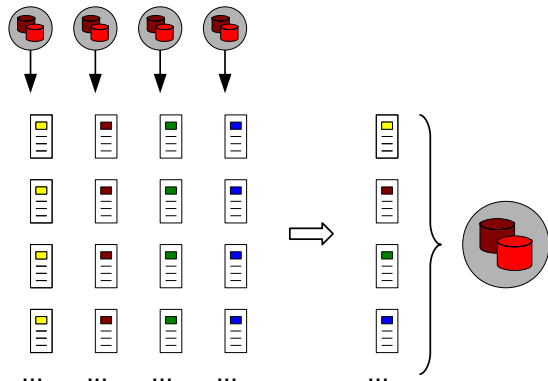
		gov2.30		gov2.1000	
		p@10	p@20	p@10	p@20
top-10	MinMax	<i>0.1966</i>	<i>0.1936</i>	<i>0.0161</i>	<i>0.0148</i>
	MM-Stdv	0.1604	0.1728	0.0054	0.0081
top-1000	MinMax	<i>0.1953</i>	<i>0.1919</i>	<i>0.0161</i>	<i>0.0148</i>
	MM-Stdv	0.1738	0.1617	0.0114	0.0148

- Discussion
 - The σ weight does not help (neither do other weights).

Resource Selection for MinMax

- + Removes non-relevant collections.
- + Provides collection relevant scores.
- Complex and requires centralized sample index.

Centralized Sample Index



Centralized sample index can be used for

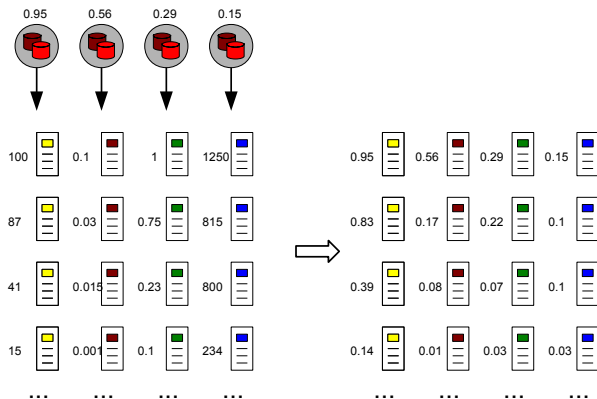
- resource selection
- score normalization

Sometimes cannot be built.

MinMax and Resource Selection

	gov2.30		gov2.1000	
	p@10	p@20	p@10	p@20
no selection	0.1966	0.1936	0.0161	0.0148
10 collections	0.2148	0.2060	0.1369	0.1275
3 collections	<i>0.3342</i>	<i>0.2752</i>	<i>0.2389</i>	<i>0.1973</i>

Weighted MinMax



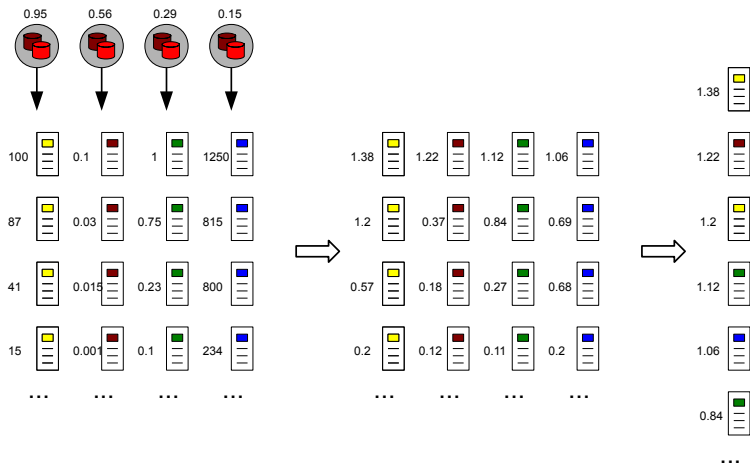
$$S_{norm} = W_{coll} \frac{s - \min}{\max - \min}$$

Weighted MinMax

		gov2.30		gov2.1000	
		p@10	p@20	p@10	p@20
no selection	MinMax	0.1966	0.1936	0.0161	0.0148
	WeightedMM	<i>0.3658</i>	<i>0.2802</i>	<i>0.2705</i>	<i>0.2205</i>
10 collections	MinMax	0.2148	0.2060	0.1369	0.1275
	WeightedMM	<i>0.3664</i>	<i>0.2782</i>	<i>0.2705</i>	<i>0.2228</i>
3 collections	MinMax	0.3342	0.2752	0.2389	0.1973
	WeightedMM	<i>0.3638</i>	<i>0.2812</i>	<i>0.2678</i>	<i>0.2027</i>

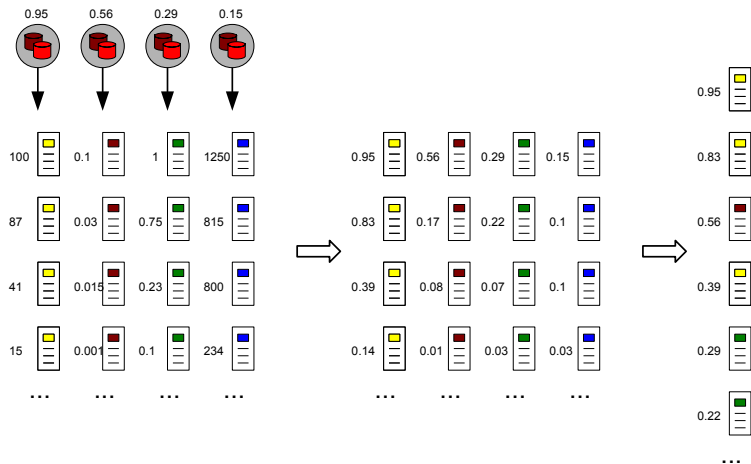
- No need in resource selection.
- Removing some collections hurts weighted MinMax.

Collection Retrieval Inference Network (CORI)



$$s_{norm} = (1 + 0.4 \cdot w_{coll}) \frac{s - \min}{\max - \min}$$

Weighted MinMax



$$s_{norm} = w_{coll} \frac{s - \min}{\max - \min}$$

Collection Retrieval Inference Network (CORI)

		gov2.30		gov2.1000	
		p@10	p@20	p@10	p@20
no selection	WeightedMM	<i>0.3658</i>	<i>0.2802</i>	<i>0.2705</i>	<i>0.2205</i>
	CORI	0.2906	0.2295	0.1879	0.1131
10 collections	WeightedMM	<i>0.3664</i>	<i>0.2782</i>	<i>0.2705</i>	<i>0.2228</i>
	CORI	0.2718	0.2181	0.1966	0.1544
3 collections	WeightedMM	<i>0.3638</i>	<i>0.2812</i>	<i>0.2678</i>	<i>0.2027</i>
	CORI	0.3396	0.2768	0.2409	0.1987

- CORI depends on how many collections are selected.
- No clear intuition behind the formula.

Other Score Normalization

	gov2.30		gov2.1000	
	p@10	p@20	p@10	p@20
MinMax	0.1966	0.1936	0.0161	0.0148
WeightedMM	<i>0.3658</i>	<i>0.2802</i>	<i>0.2705</i>	<i>0.2205</i>
CORI	0.2906	0.2295	0.1879	0.1131
SAFE	0.1879	0.1698	0.0698	0.0718
HIS	0.1728	0.1845	0.0819	0.0822

Table: No resource selection.

	gov2.30		gov2.1000	
	p@10	p@20	p@10	p@20
MinMax	0.3342	0.2752	0.2389	0.1973
WeightedMM	<i>0.3638</i>	0.2812	<i>0.2678</i>	0.2027
CORI	0.3396	0.2768	0.2409	0.1987
SAFE	0.2691	0.2490	0.2477	0.1983
HIS	0.3054	0.2816	0.2070	0.1883

Table: 3 collections are selected.

Open Question

Is it possible to calculate collection weights simpler,
with minimum additional information
and without centralized sample index?

$$S_{norm} = \frac{s}{max}$$

$$S_{norm} = W \frac{s}{max}$$

Alternative ways to calculate collection weights?