# Axiometrics: An Axiomatic Approach to Evaluation Metrics

**Stefano Mizzaro**

Udine University, Italy

mizzaro@uniud.it

essir 2015
thessaloniki.greece

---

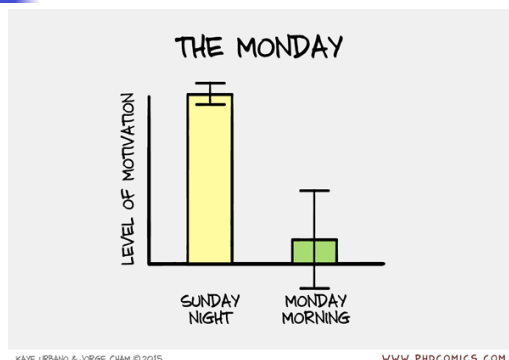## Let's be clear from the start: I. Won't. Go. Overtime.



---

## Aims – 1

- To present both:
  - basic material (what you find in books)
  - advanced material (recently published, even not yet published!)
- Links with Julio, Enrique, Evangelos
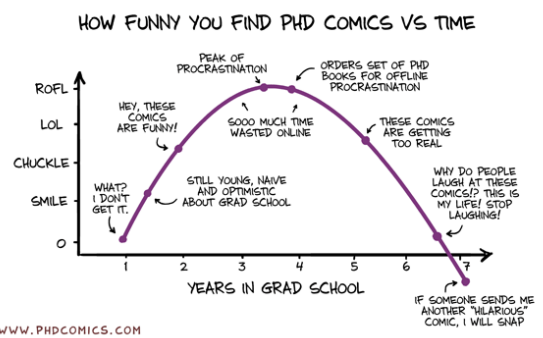  - not (yet!) fully integrated
  - a bit disorganized…

---

## Aims – 2

- Intro to IR Evaluation
- Intro to IR Evaluation Measures / Metrics
- The Link between Measurement Theory and Metrics
  - Intro to Measurement Theory (Scales)
  - Metrics Analysis
  - The Axiometrics Framework (*)

---

## The Monday effect…



---

## You don't find it funny?!

## Outline

- Evaluation [5']
- Measures / metrics [15']
- Measurement theory [15']
- Metrics analysis [5']
- Axiometrics framework [15']

## Outline

- Evaluation [5']
- Measures / metrics [15']
- Measurement theory [15']
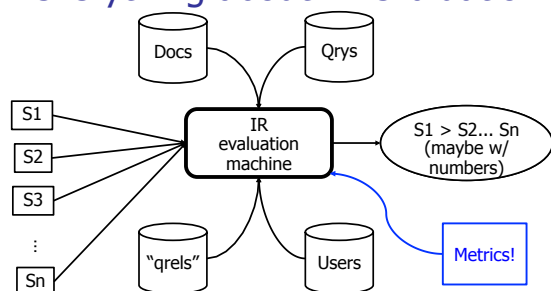- Metrics analysis [5']
- Axiometrics framework [15']

## What is evaluation (in IR)?

- Eh…
- Ideally: a machine telling you how good an IR system is
- "Good": effective, capable to retrieve relevant (useful?!) documents
  - (efficiency is also studied, but focus is on effectiveness)

## The importance of evaluation in IR

- Everybody agrees that evaluation is of paramount importance in IR
- One of the most evaluation-oriented disciplines in computer/information sciences
- We're busy doing a lot of evaluation since the 60s
  - So this talk is relevant. I do not know if it is useful :-)

## A short history of nearly everything about IR evaluation



## The importance of evaluation in IR

- Everybody agrees that evaluation is of paramount importance in IR
- One of the most evaluation-oriented disciplines in computer/information sciences
- We're busy doing a lot of evaluation since the 60s
- And we don't know (agree on) how to evaluate

## (Other) Issues in IR evaluation

- Relevance
  - "Topicality"?
  - "Utility"?
  - ...
- Methodology
  - Test collection, benchmark ,TREC-like
  - User study (--> Diane)
  - Large log analysis
  - ...

## Outline

- Evaluation [5′]
- **Measures / metrics [15′]**
- Measurement theory [15′]
- Metrics analysis [5′]
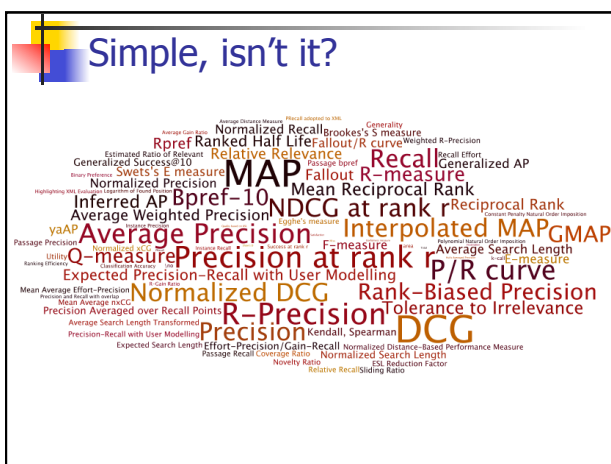- Axiometrics framework [15′]

## We go down the dark evaluation metrics rabbit hole



## IR effectiveness metric

- (or, measure)
- ""A number telling us how effective an IR system is""
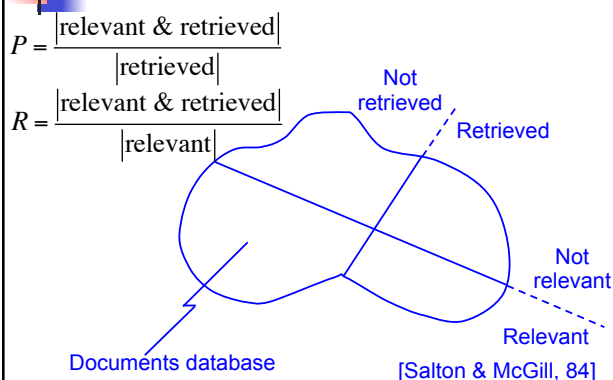- Simple, isn't it?

## Simple, isn't it?



| 1. | Precision | | 45. | Discounted Cumulative Gain |
|---|---|---|---|---|
| 2. | Recall | | 46. | Normalized DCG |
| 3. | Swets's E measure | | 47. | NDCG at rank r |
| 4. | Fallout | | 48. | Average Weighted Precision |
| 5. | Normalized Recall | | 49. | Weighted R-Precision |
| 6. | Normalized Precision | | 50. | Average Distance Measure |
| 7. | P/R curve | | 51. | PRecall adopted to XML |
| 8. | Fallout/R curve | | 52. | Precision and Recall with overlap |
| 9. | Brookes's S measure | | 53. | \no |
| 10. | Expected Search Length | | 54. | area |
| 11. | ESL Reduction Factor | | 55. | Success at rank r |
| 12. | Precision Averaged over Recall Points | | 56. | Average Gain Ratio |
| 13. | Sliding Ratio | | 57. | R-Gain Ratio |
| 14. | Coverage Ratio | | 58. | Binary Preference |
| 15. | Novelty Ratio | | 59. | Bpref-10 |
| 16. | Relative Recall | | 60. | Q-measure |
| 17. | E-measure | | 61. | R-measure |
| 18. | F-measure | | 62. | Tolerance to Irrelevance |
| 19. | Utility | | 63. | Estimated Ratio of Relevant |
| 20. | Average Precision | | 64. | Egghe's measure |
| 21. | Mean Average Precision | | 65. | Passage Recall |
| 22. | Interpolated MAP | | 66. | Passage Precision |
| 23. | Precision at rank r | | 67. | Passage bpref |
| 24. | R-Precision | | 68. | Kendall, Spearman |
| 25. | Generality | | 69. | Mean Average nxCG |
| 26. | Miss | | 70. | Normalized xCG |
| 27. | Shaw's D | | 71. | Effort-Precision/Gain-Recall |
| 28. | PRecall | | 72. | Mean Average Effort-Precision |
| 29. | Satisfaction | | 73. | Precision-Recall with User Modelling |
| 30. | Frustration | | 74. | Geometric MAP |
| 31. | Total | | 75. | Generalized AP |
| 32. | Usefulness measure | | 76. | Inferred AP |
| 33. | Hull's Averaged Precision | | 77. | Expected Precision-Recall with User Modelling |
| 34. | Average Search Length | | 78. | Rpref |
| 35. | Quality based on ASL | | 79. | Generalized Success@10 |
| 36. | Normalized Distance-Based Performance Measure | | 80. | k-call |
| 37. | Recall Effort | | 81. | Highlighting XML Evaluation |
| 38. | Relative Relevance | | 82. | Ranking Efficiency |
| 39. | Ranked Half Life | | 83. | Rank-Biased Precision |
| 40. | Reciprocal Rank | | 84. | Average Search Length Transformed |
| 41. | Mean Reciprocal Rank | | 85. | Logarithm of Found Position |
| 42. | Instance Precision | | 86. | Normalized Search Length |
| 43. | Instance Recall | | 87. | Polynomial Natural Order Imposition |
| 44. | Classification Accuracy | | 88. | Constant Penalty Natural Order Imposition |

## A shorter list

- Precision, Recall
- Precision-Recall curve
- MAP (Mean Average Precision)
- P@n (Precision at n)
- NDCG (Normalized Discounted Cumulative Gain)
- MRR (Mean Reciprocal Rank)
- RBP (Rank Biased Precision)
- TBG (Time Based Gain)

## Precision & Recall

$$P = \frac{|relevant\ \&\ retrieved|}{|retrieved|}$$

$$R = \frac{|relevant\ \&\ retrieved|}{|relevant|}$$

Not retrieved

Retrieved

Not relevant

Relevant

Documents database

[Salton & McGill, 84]

## P & R: probabilistic definition

- $P$ = p(relevant | retrieved)
- $R$ = p(retrieved | relevant)
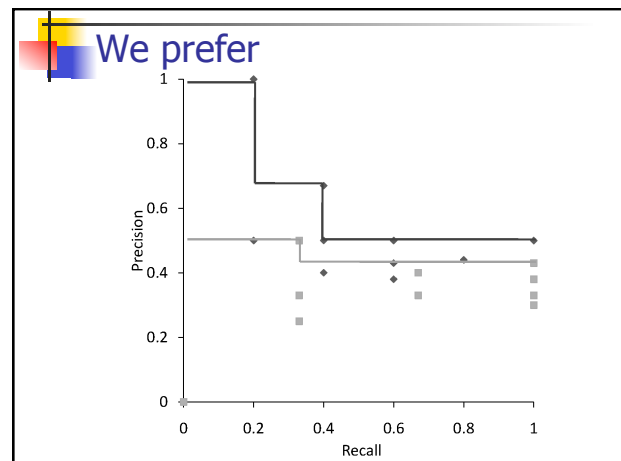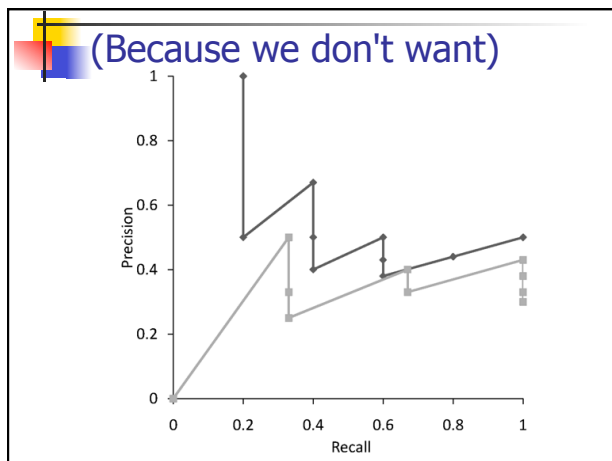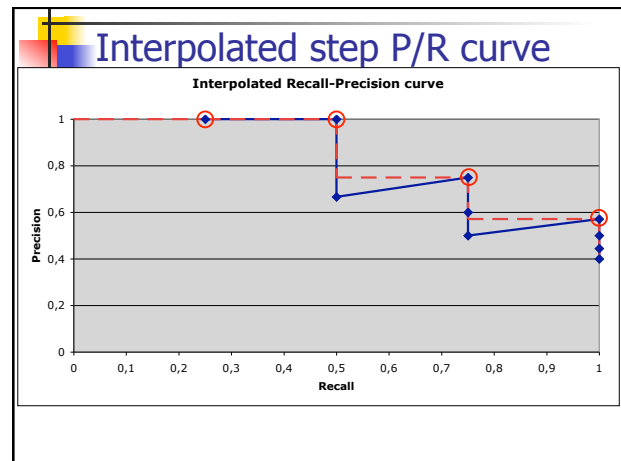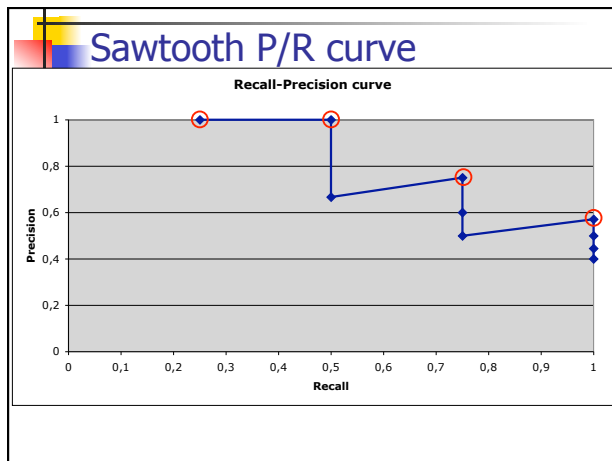
## Ranking!

- But today all IR systems **rank** the documents
- Limitations of P&R
  - 2 numbers, not just one
  - Not affected by the rank of retrieved docs.
- Solutions: (too?) many.
  - Precision/Recall curve
  - MAP (Mean Average Precision)
  - …
- Let us see some examples

## Rank

| Rank | Rel? | R | P |
|------|------|------|------|
| 1 | 1 | 0,25 | 1 |
| 2 | 1 | 0,5 | 1 |
| 3 | 0 | 0,5 | 0,67 |
| 4 | 1 | 0,75 | 0,75 |
| 5 | 0 | 0,75 | 0,6 |
| 6 | 0 | 0,75 | 0,5 |
| 7 | 1 | 1 | 0,57 |
| 8 | 0 | 1 | 0,5 |
| 9 | 0 | 1 | 0,44 |
| 10 | 0 | 1 | 0,4 |

## Rank

| Rank | Rel? | R | P |
|------|------|------|------|
| 1 | 1 | **0,25** | **1** |
| 2 | 1 | **0,5** | **1** |
| 3 | 0 | 0,5 | 0,67 |
| 4 | 1 | **0,75** | **0,75** |
| 5 | 0 | 0,75 | 0,6 |
| 6 | 0 | 0,75 | 0,5 |
| 7 | 1 | **1** | **0,57** |
| 8 | 0 | 1 | 0,5 |
| 9 | 0 | 1 | 0,44 |
| 10 | 0 | 1 | 0,4 |

## Sawtooth P/R curve



## Interpolated step P/R curve



## (Because we don't want)



## We prefer



## (We could do but we don't do)



## So we average the step curves

## Over N queries and we get



## And of course on 11 levels of recall



## We happily compare systems ?



## Although often…



## P/R curve --> MAP

- P/R curve
  - It is not a number
  - It can be transformed into a number by measuring **the area below the curve**
- --> **AP (Average Precision)**
- --> **MAP (Mean Average Precision)**
- Good property: top-heavyness

| Rank | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | … |
|------|---|---|---|---|---|---|---|---|---|
| Rel  | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | … |

- User model?

## P@n

- Simply count how many relevant documents are retrieved in the first n positions of the rank
- P@10 useful for classical Web search engines
- P@1 for "Feeling lucky"

# Non binary relevance

- Some documents are "more relevant" than others
- Discounted Cumulative Gain (DCG, NDCG)
  - Different relevance --> different gain for the user
    - E.g., H --> 3, R --> 2, P --> 1, N --> 0
  - Sum of the gains while walking down the rank
  - Discounting more and more: late rank positions give less gain even if of equal relevance ("top-heaviness")

# Example

| Rank | Rel |
|------|-----|
| 1 | H |
| 2 | P |
| 3 | N |
| 4 | R |
| 5 | H |
| 6 | R |
| 7 | N |
| 8 | P |
| 9 | N |
| 10 | P |

# Example

| Rank | Rel | Gain |
|------|-----|------|
| 1 | H | 3 |
| 2 | P | 1 |
| 3 | N | 0 |
| 4 | R | 2 |
| 5 | H | 3 |
| 6 | R | 2 |
| 7 | N | 0 |
| 8 | P | 1 |
| 9 | N | 0 |
| 10 | P | 1 |

# Example

| Rank | Rel | Gain | CG |
|------|-----|------|-----|
| 1 | H | 3 | 3 |
| 2 | P | 1 | 4 |
| 3 | N | 0 | 4 |
| 4 | R | 2 | 6 |
| 5 | H | 3 | 9 |
| 6 | R | 2 | 11 |
| 7 | N | 0 | 11 |
| 8 | P | 1 | 12 |
| 9 | N | 0 | 12 |
| 10 | P | 1 | 13 |

# Example

| Rank | Rel | Gain | CG | Discount |
|------|-----|------|-----|----------|
| 1 | H | 3 | 3 | *log(1)* 1 |
| 2 | P | 1 | 4 | log(2) |
| 3 | N | 0 | 4 | log(3) |
| 4 | R | 2 | 6 | log(4) |
| 5 | H | 3 | 9 | log(5) |
| 6 | R | 2 | 11 | log(6) |
| 7 | N | 0 | 11 | log(7) |
| 8 | P | 1 | 12 | log(8) |
| 9 | N | 0 | 12 | log(9) |
| 10 | P | 1 | 13 | log(10) |

# Example

| Rank | Rel | Gain | CG | Discount | DG |
|------|-----|------|-----|----------|-----|
| 1 | H | 3 | 3 | *log(1)* 1 | 3/1=3 |
| 2 | P | 1 | 4 | log(2) | 1/log(2)=1 |
| 3 | N | 0 | 4 | log(3) | 0/log(3)=0 |
| 4 | R | 2 | 6 | log(4) | 2/log(4)=1 |
| 5 | H | 3 | 9 | log(5) | 3/log(5)=1.3 |
| 6 | R | 2 | 11 | log(6) | 2/log(6)=.8 |
| 7 | N | 0 | 11 | log(7) | 0/log(7)=0 |
| 8 | P | 1 | 12 | log(8) | 1/log(8)=.3 |
| 9 | N | 0 | 12 | log(9) | 0/log(9)=0 |
| 10 | P | 1 | 13 | log(10) | 1/log(10)=.3 |

## Example

| Rank | Rel | Gain | CG | Discount | DG | DCG |
|---|---|---|---|---|---|---|
| 1 | H | 3 | 3 | *log(1)* 1 | 3/1=3 | 3.0 |
| 2 | P | 1 | 4 | log(2) | 1/log(2)=1 | 4.0 |
| 3 | N | 0 | 4 | log(3) | 0/log(3)=0 | 4.0 |
| 4 | R | 2 | 6 | log(4) | 2/log(4)=1 | 5.0 |
| 5 | H | 3 | 9 | log(5) | 3/log(5)=1.3 | 6.3 |
| 6 | R | 2 | 11 | log(6) | 2/log(6)=.8 | 7.1 |
| 7 | N | 0 | 11 | log(7) | 0/log(7)=0 | 7.1 |
| 8 | P | 1 | 12 | log(8) | 1/log(8)=.3 | 7.4 |
| 9 | N | 0 | 12 | log(9) | 0/log(9)=0 | 7.4 |
| 10 | P | 1 | 13 | log(10) | 1/log(10)=.3 | 7.7 |

## Example

| Rank | Rel | Gain | CG | Discount | DG | DCG | DCG Ideal |
|---|---|---|---|---|---|---|---|
| 1 | H | 3 | 3 | *log(1)* 1 | 3/1=3 | 3.0 | (H) 3.0 |
| 2 | P | 1 | 4 | log(2) | 1/log(2)=1 | 4.0 | (H) 6.0 |
| 3 | N | 0 | 4 | log(3) | 0/log(3)=0 | 4.0 | (R) 7.3 |
| 4 | R | 2 | 6 | log(4) | 2/log(4)=1 | 5.0 | (R) 8.3 |
| 5 | H | 3 | 9 | log(5) | 3/log(5)=1.3 | 6.3 | (P) 8.7 |
| 6 | R | 2 | 11 | log(6) | 2/log(6)=.8 | 7.1 | (P) 9.1 |
| 7 | N | 0 | 11 | log(7) | 0/log(7)=0 | 7.1 | (P) 9.4 |
| 8 | P | 1 | 12 | log(8) | 1/log(8)=.3 | 7.4 | (N) 9.4 |
| 9 | N | 0 | 12 | log(9) | 0/log(9)=0 | 7.4 | (N) 9.4 |
| 10 | P | 1 | 13 | log(10) | 1/log(10)=.3 | 7.7 | (N) 9.4 |

## Example

| Rank | Rel | Gain | CG | Discount | DG | DCG | DCG Ideal | NDCG |
|---|---|---|---|---|---|---|---|---|
| 1 | H | 3 | 3 | *log(1)* 1 | 3/1=3 | 3.0 | (H) 3.0 | 1.00 |
| 2 | P | 1 | 4 | log(2) | 1/log(2)=1 | 4.0 | (H) 6.0 | 0.67 |
| 3 | N | 0 | 4 | log(3) | 0/log(3)=0 | 4.0 | (R) 7.3 | 0.55 |
| 4 | R | 2 | 6 | log(4) | 2/log(4)=1 | 5.0 | (R) 8.3 | 0.61 |
| 5 | H | 3 | 9 | log(5) | 3/log(5)=1.3 | 6.3 | (P) 8.7 | 0.72 |
| 6 | R | 2 | 11 | log(6) | 2/log(6)=.8 | 7.1 | (P) 9.1 | 0.78 |
| 7 | N | 0 | 11 | log(7) | 0/log(7)=0 | 7.1 | (P) 9.4 | 0.75 |
| 8 | P | 1 | 12 | log(8) | 1/log(8)=.3 | 7.4 | (N) 9.4 | 0.78 |
| 9 | N | 0 | 12 | log(9) | 0/log(9)=0 | 7.4 | (N) 9.4 | 0.78 |
| 10 | P | 1 | 13 | log(10) | 1/log(10)=.3 | 7.7 | (N) 9.4 | 0.82 |

## Example

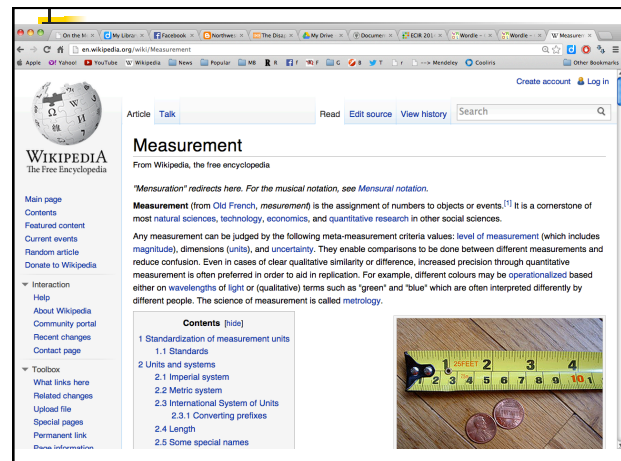| Rank | Rel | Gain | CG | Discount | DG | DCG | DCG Ideal | NDCG |
|---|---|---|---|---|---|---|---|---|
| 1 | H | 3 | 3 | *log(1)* 1 | 3/1=3 | 3.0 | (H) 3.0 | 1.00 |
| 2 | P | 1 | 4 | log(2) | 1/log(2)=1 | 4.0 | (H) 6.0 | 0.67 |
| 3 | N | 0 | 4 | log(3) | 0/log(3)=0 | 4.0 | (R) 7.3 | 0.55 |
| 4 | R | 2 | 6 | log(4) | 2/log(4)=1 | 5.0 | (R) 8.3 | 0.61 |
| 5 | H | 3 | 9 | log(5) | 3/log(5)=1.3 | 6.3 | (P) 8.7 | 0.72 |
| 6 | R | 2 | 11 | log(6) | 2/log(6)=.8 | 7.1 | (P) 9.1 | 0.78 |
| 7 | N | 0 | 11 | log(7) | 0/log(7)=0 | 7.1 | (P) 9.4 | 0.75 |
| 8 | P | 1 | 12 | log(8) | 1/log(8)=.3 | 7.4 | (N) 9.4 | 0.78 |
| 9 | N | 0 | 12 | log(9) | 0/log(9)=0 | 7.4 | (N) 9.4 | 0.78 |
| 10 | P | 1 | 13 | log(10) | 1/log(10)=.3 | 7.7 | (N) 9.4 | 0.82 |

## A shorter list

- Precision, Recall
- Precision-Recall curve
- MAP
- P@n
- NDCG
- MRR (Mean Reciprocal Rank)
- RBP (Rank Biased Precision)
- TBG (Time Based Gain)

## … even continuous relevance

- Dynamometer, Hand force grip, even physiological data
- IR System showing a "relevance bar" close to each document
  - "estimation of the amount of relevance"
- Magnitude estimation (paper @ last SIGIR)
- …

## Metrics classification

- Metrics could be classified on the basis of underlying notions of:
  - relevance (binary, ranking, continuous)
  - retrieval (binary, ranking, continuous)
- 3 x 3 (or N x N) grid, …

## Classification (incomplete!)

**Relevance**

| | Binary | Rank | Continuous |
|---|---|---|---|
| Continuous | | Sliding ratio | ADM |
| Rank | | NDCG, Kendall/ Spearman? | |
| Binary | P&R, E, F, … | RP curve, MAP, R-prec, … | |

**Retrieval**

## So, if you understand it, I'm not being clear enough :-)

Average Distance Measure
Average Gain Ratio      Normalized Recall Brookes's S measure      Generality
Rpref Ranked Half Life Fallout/R curve Weighted R-Precision
Estimated Ratio of Relevant      Relative Relevance            Recall Recall Effort
Generalized Success@10                          Passage bpref      Generalized AP
Binary Preference Swets's E measure Fallout R-measure
Normalized Precision                      Mean Reciprocal Rank
Highlighting XML Evaluation Logarithm of found Positions
Inferred AP Bpref-10      NDCG at rank r Reciprocal Rank
Average Weighted Precision                                Constant Penalty Natural Order Imposition
yaAP      Average Precision…      Interpolated MAP GMAP
Passage Precision                            F-measure                Polynomial Natural Order Imposition
Utility Eoghe's measure e-measure
Ranking Efficiency Q-measure Precision at rank r P/R curve Average Search Length
                        Classification Accuracy      VAP
Expected Precision–Recall with User Modelling
Mean Average Effort–Precision R-Gain Ratio
Precision and Recall with overlap      Rank–Biased Precision
Mean Average nxCG Normalized DCG
Precision Averaged over Recall Points                  Tolerance to Irrelevance
Average Search Length Transformed      R–Precision
Precision–Recall with User Modelling Precision Kendall, Spearman DCG
Expected Search Length Effort–Precision/Gain–Recall Normalized Distance–Based Performance Measure
Passage Recall Coverage Ratio      Normalized Search Length
      Novelty Ratio            ESL Reduction Factor
Relative Recall Sliding Ratio

--> Julio, Enrique, Evan…

## Take home messages so far

- In IR, Evaluation is important.
- There are **many (100+)** metrics
  - System-oriented metrics only
    - (R, P, MAP, NDCG, …)
  - Let alone human/user-oriented metrics
    - (user satisfaction, fatigue, …)
  - And still IR only, no clustering, filtering…
- Just a few are indeed used
  - (which ones? Why those? --> Julio)
- Definitions of some of them

## And now for something completely different…



*And now for something completely different*

## Outline

- Evaluation [5']
- Measures / metrics [15']
- Measurement theory [15']
- Metrics analysis [5']
- Axiometrics framework [15']

## Measurement

- Definition: A process aimed at determining a relationship between a physical quantity and a unit of measurement
- Typically, one assigns numbers to objects/ events
- Studied in Measurement Theory
  - Reasonably settled



---



## Measurement Theory

- One important concept: which numbers can I select when measuring? What properties do they have?
- Which **measurement scale**?
  - (or "**level**")
  - E.g., to measure length:
    - Meters
    - Inches
    - "Longer than" (?)



---

## Measurement scales

- Standard set of scales:
  1. Nominal
  2. Ordinal
  3. Interval
  4. Ratio
- (other proposals exist)

---

## Perhaps not so well settled…



---

## And indeed

- Nicholas Chrisman
- 1998
- Proposes 10 scales, not just 4

*[Embedded paper image: "Rethinking Levels of Measurement for Cartography" by Nicholas R. Chrisman]*

## Anyway "Good Old Fashioned" Measurement scales

1. Nominal
2. Ordinal
3. Interval
4. Ratio

## 4. Ratio scale

- I'm twice taller than him
- He is twice richer than me
  - (Both "how much" & "how many")
- I'm twice older than you
  - Years, months (*12), days, …
- Zero
  - Age starts from zero!
- But not
  - Today is twice as hot as yesterday?

## 3. Interval scale

- Today is twice as hot as yesterday?
- In the last two days we had the same increase of 5ºC in temperature
  - The difference between today and yesterday temperature is the same as …
- Dates are another example
  - Difference: OK (2014 – 2012 = 2006 – 2004)
  - Ratio: KO (2000 is not 2 * 1000)
  - (ratios of differences: OK)

## 2. Ordinal scale

- A measure is not an amount but a **rank**
- It is a form of measurement!
- Ex: Lines ranked according to their length
- It does **not** mean that:
  - the first is twice as long than the second
  - the length difference between the 1st and the 2nd is the same as the 2nd and the 3rd

1st
2nd
3rd
4th

## 1. Nominal scale

- Qualitative
- Categories
- Names, gender, nationality, …
- Can be Dichotomous or Non-dichotomous
- No assumptions on ratios, distances, ranks.

## Legit operations

- Given a scale, only some operations make sense
  - Arithmetic: +, -, *, /
  - Statistic: Mean, median, mode, …
- Ex:
  - Average height, weight, …: OK
  - Average gender: KO

## Legit relational/math operations

|          | =, ≠ | >, < | +, - | ×, ÷ |
|----------|:----:|:----:|:----:|:----:|
| Nominal  | ✓    | ✗    | ✗    | ✗    |
| Ordinal  | ✓    | ✓    | ✗    | ✗    |
| Interval | ✓    | ✓    | ✓    | ✗    |
| Ratio    | ✓    | ✓    | ✓    | ✓    |

## Permissible transformations

|          | a · x | a · x + b | Monotonic | 1-to-1 |
|----------|:-----:|:---------:|:---------:|:------:|
| Nominal  | ✓     | ✓         | ✓         | ✓      |
| Ordinal  | ✓     | ✓         | ✓         | ✗      |
| Interval | ✓     | ✓         | ✗         | ✗      |
| Ratio    | ✓     | ✗         | ✗         | ✗      |

## Meaning

- If you transform the measure, are you still measuring the same thing?
  - Nationality
  - Rank
  - Temperature
  - Money

## Examples

- Nominal scale for nationality
  - Greek= 1
  - Italian = 2
  - Spanish = 3
  - Japanese = 4
  - …
  - 3 − 1 = 4 − 2 Uh?!?
- Whereas interval scale for temperature ºC
  - 30º − 10º  = 40º − 20º: ok

## Legit statistics

|          | Mode | Median | Mean |  |
|----------|:----:|:------:|:--------:|:---------------------:|
|          |      |        | Arithmetic | Geometric, Harmonic |
| Nominal  | ✓    | ✗      | ✗        | ✗                     |
| Ordinal  | ✓    | ✓      | ✗        | ✗                     |
| Interval | ✓    | ✓      | ✓        | ✗                     |
| Ratio    | ✓    | ✓      | ✓        | ✓                     |

## Examples

- Nationality, mode, ok
- Mean rank? Uh?!
  - (think of ranking ten students...)
- Mean temperature, ok

---

METHODOLOGIST'S CORNER

**Likert scales, levels of measurement and the "laws" of statistics**

**Geoff Norman**

**Abstract**  Reviewers of research reports frequently criticize the choice of statistical methods. While some of these criticisms are well-founded, frequently the use of various parametric methods such as analysis of variance, regression, correlation are faulted because: (a) the sample size is too small, (b) the data may not be normally distributed, or (c) The data are from Likert scales, which are ordinal, so parametric statistics cannot be used. In this paper, I dissect these arguments, and show that many studies, dating back to the 1930s consistently show that parametric statistics are robust with respect to violations of these assumptions. Hence, challenges like those above are unfounded, and parametric methods can be utilized without concern for "getting the wrong answer".

---

## Why are you telling me this?



---



---

## Why are you telling me this?

- Two reasons
  - 1. Measurement theory and scales can be used to directly analyze IR metrics
  - 2. Because IR can be seen as a measurement. Of relevance --> Language to define axioms on metrics
- Let's see 1. first

---

## Outline

- Evaluation [5′]
- Measures / metrics [15′]
- Measurement theory [15′]
- **Metrics analysis [5′]**
- Axiometrics framework [15′]

---

## MRR (Mean Reciprocal Rank)

- Take the rank of the 1st relevant doc.
- Take the reciprocal (…)

$$RR = \frac{1}{rank(i)}$$

- Then take the mean (…) over some topics

$$MRR = \frac{1}{|Q|} \sum_{q \in Q} \frac{1}{rank(i_q)}$$
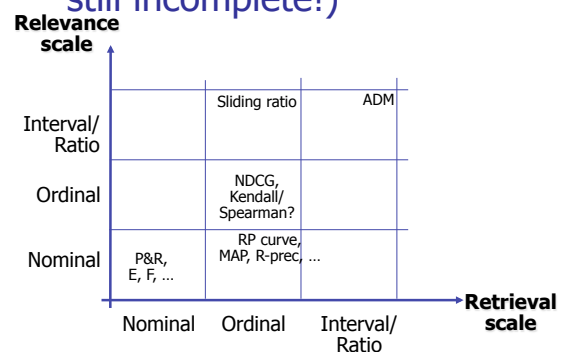
## MRR?!

- It is used
- In many papers
- Even in some TREC tracks

- When analyzed with "measurement theory glasses", is is not "measurement theory proof"
- A **reciprocal** is taken on an **ordinal** scale…
- … then it is **averaged**…

## Even worse than that… NDCG

- H,R,P,N --> 3, 2, 1, 0 ?
  - Linear, most common
- H,R,P,N --> 100, 10, 1, 0 ?! (or 4, 2, 1, 0)
  - Exponential, sometimes used, actually
- H,R,P,N --> 100, 99, 90, 0 ?!?
  - "Crazy", never heard of… why?
- Arbitrary choice!
- By transforming relevance levels into gains, we transform an 2. Ordinal scale --> 4. Ratio scale!
- And we also discount
  - Dividing by log(rank)…
  - All metrics that can be modeled as gain/discount…

## Classification (more principled, still incomplete!)

| Relevance scale | Nominal | Ordinal | Interval/Ratio |
|---|---|---|---|
| Interval/Ratio | | Sliding ratio | ADM |
| Ordinal | | NDCG, Kendall/Spearman? | |
| Nominal | P&R, E, F, … | RP curve, MAP, R-prec, … | |
| | Nominal | Ordinal | Interval/Ratio |

Retrieval **scale**

## To summarize, but not to conclude

- 100+ metrics
- Measurement theory seems a useful tool
  - Metric classification
  - Some arbitrary choices
  - Some metrics are not "Measurement theory proof"
- Metric **engineering** seems more an **art** (artisan) than a **science**…
- A more principled approach?

## Outline

- Evaluation [5′]
- Measures / metrics [15′]
- Measurement theory [15′]
- Metrics analysis [5′]
- Axiometrics framework [15′]

## So, back to metrics…



## Hm. One hundred.

- What do they have in common?
- Research question:
  **Are there some axioms that any effectiveness metric should satisfy?**
- "Axioms"? ("Axiometrics"! [Arjen])
  - Properties, constraints, Laws, …
- Maybe a not so crazy idea:
  - [van Rijsbergen 1979], [Bollman 1984], [Amigó, Gonzalo 2009] (& more), [Sebastiani 2015], SWIRL 2012, …

## 2. IR as measurement of relevance

- Both an IR system and a human assessor **measure** the **relevance** of a document
- Maybe on different **scales**:
  - Ranked SE output: ordinal scale
  - Human TREC-like relevance assessor: nominal (ordinal) scale (binary, relevant/nonrelevant)
  - … IR related tasks like categorization, filtering, …

## So, relevance measurement(s)

- This is general
  - Change the scale and you get the corresponding relevance measurement
  - Allows to take into account binary, graded, ranked,… relevance and retrieval

## Notation: $\sigma$ and $\alpha$

- 2 relevance measurements
  - by a system ($\sigma$)
    - $\sigma(q,d)$, $\sigma(q,D)$, $\sigma(Q,D)$
    - e.g. SE ranked output
  - by a human (user/assessor) ($\alpha$)
    - $\alpha(q,d)$, $\alpha(q,D)$, $\alpha(Q,D)$
    - e.g., TREC qrels

## Similarity

- Given two relevance measurements $\alpha$ and $\sigma$, we can define a notion of similarity between them (given a query/topic $q$ and a document $d$)

$$\underset{q,d}{\text{sim}}\,(\alpha, \sigma)$$

- And an IR system should provide a measurement $\sigma$ that is similar to the human measurement $\alpha$

## Now, similarity…

# Psychological Review

Copyright © 1977 by the American Psychological Association, Inc.

**VOLUME 84   NUMBER 4   JULY 1977**

Features of Similarity

Amos Tversky
Hebrew University
Jerusalem, Israel

The metric and dimensional assumptions that underlie the geometric represen-
tation of similarity are questioned on both theoretical and empirical grounds.
A new set-theoretical approach to similarity is developed in which objects are
represented as collections of features, and similarity is described as a feature-
matching process. Specifically, a set of qualitative assumptions is shown to

---

## Similarity comparison

- And we can compare similarities
- For example,

$$\sin_{q,d} (\alpha, \sigma) < \sin_{q,d} (\alpha, \sigma')$$

- means that on the query $q$ and the document $d$, and given the human relevance judgment $\alpha$, system $\sigma$ is worse than $\sigma'$ (i.e., less similar to $\alpha$)

---

## Similarity --> Metric

- On the basis of the notion of similarity, we can define an IR effectiveness metric
- The more $\sigma$ is similar to $\alpha$, the higher the metric value

---

## So, to summarize (but not to conclude!)

- **Measurement** theory, Measurement scales, IR as relevance measurement
  - $\sigma(q,d)$, $\sigma(q,D)$, $\sigma(Q,D)$, $\alpha(q,d)$, $\alpha(q,D)$, $\alpha(Q,D)$
- **Similarity**

$$\sin_{q,d} (\alpha, \sigma) < \sin_{q,d} (\alpha, \sigma')$$

- **Metric**

$$\text{metric}_{Q,D}(\alpha, \sigma)$$
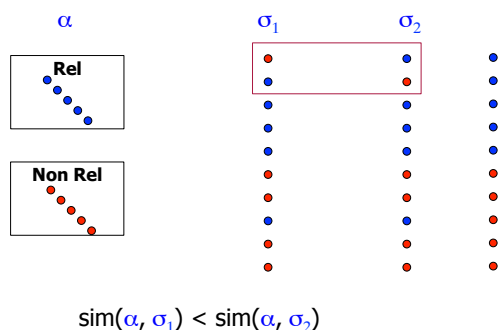
---

## Generality: Across different measurement scales

- We can compute the similarity of two relevance judgments when they are on the same scale
- …but more than that…
- … also when they are on different scales, in some cases
  - E.g., the classical ad-hoc retrieval
    - scale($\sigma$) = [[ordinal]]
    - scale($\alpha$) = [[nominal]] (binary relevance. Ordinal)

---

## Same scales: binary IR

$\alpha$ $\sigma_1$ $\sigma_2$

| Rel | Ret | Ret |
| Non Rel | Non Ret | Non Ret |

$\text{sim}(\alpha, \sigma_1) > \text{sim}(\alpha, \sigma_2)$

## Different scales: ad hoc IR

$$\text{sim}(\alpha, \sigma_1) < \text{sim}(\alpha, \sigma_2)$$

## Ok, but we want Axioms!

## Some details…

- I do not trust our axioms too much yet…
- … preliminary work…
- Actually:
  - I think axioms are correct and consistent
  - I don't know if they are complete
- Stating axioms is also useful to "test the framework"
  - Measurement theory is an effective language to state them!
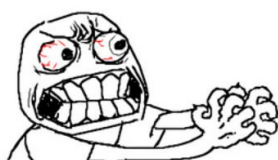
## A First Axiom

**Axiom 3** (Similarity of two systems). *Let q be a query, d a document, α a human relevance measurement and σ and σ′ two system relevance measurements such that*

$$\sigma(q, d) = \sigma'(q, d). \tag{1}$$

*Then*

$$\sim_{q,d}(\alpha, \sigma) = \sim_{q,d}(\alpha, \sigma'). \tag{2}$$

## Come on you're kidding!

- All this and then such a stupid axiom????

## Ok, ok. Second Axiom

**Axiom 6** (Overestimated documents). *Let q be a query, d and d′ two document, α a human relevance measurement and σ a system relevance measurements such that*
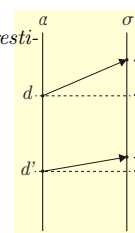
$$\alpha(d) > \alpha(d'),$$
$$\sim_{q,d}(\alpha, \sigma) < \sim_{q,d'}(\alpha, \sigma)$$

*and (6) and (8) hold (i.e., both d and d′ are overestimated), then*

$$\text{metric}_{q,d}(\alpha, \sigma) < \text{metric}_{q,d'}(\alpha, \sigma).$$

- (d is more relevant, sim on d is lower, then metric value on d has to be lower)
- d is both "more wrong" and "more visible" to the user

## Ok, ok. Third Axiom

**Axiom 8** (System relevance). *Let $q$ be a query, $d$ and $d'$ two documents, $\alpha$ a human relevance measurement and $\sigma$ a system relevance measurement such that $\mathsf{sim}_{q,d}\,(\alpha,\sigma) = \mathsf{sim}_{q,d'}\,(\alpha,\sigma)$, $\sigma(d) > \sigma(d')$, and*

$$\alpha(d) \geq \alpha(d'). \tag{10}$$

*Then*

$$d \sqsupset_{metric(\alpha,\sigma)} d'.$$

(document $d$ affects the metric value more than $d'$)

## Meaning?

- Corollary:
- By taking scale($\sigma$) = [[Rank]] we derive that:
  - Early rank positions affect a metric value more than later rank positions
  - IR metrics should be "**top-heavy**"
- Previous Axiom 8 states a more abstract/ general principle, independent of the scales

## Now, a last Axiom

**Axiom 9** (User relevance). *Let $q$ be a query, $d$ and $d'$ two documents, $\alpha$ a human relevance measurement and $\sigma$ a system relevance measurement such that: $\mathsf{sim}_{q,d}\,(\alpha,\sigma) = \mathsf{sim}_{q,d'}\,(\alpha,\sigma)$, $\alpha(d) > \alpha(d')$, and*

$$\sigma(d) \geq \sigma(d'). \tag{11}$$

*Then*

$$d \sqsupset_{metric(\alpha,\sigma)} d'.$$

(document $d$ affects the metric value more than $d'$)

## Meaning?

- A metric should weigh more, and be more affected, by more relevant documents
- "$\alpha$ top heavyness", "human top heavyness"
  - Perhaps less intuitive than previous axiom,
  - but it does indeed seem natural in the framework
  - by symmetry (treat $\alpha$ as $\sigma$)
  - To evaluate a nonrelevant document as nonrelevant is an easy job (the vast majority of documents in a collection are nonrelevant)

## Meaning?

- Consequence: linear gain values of 3, 2, 1, 0 in NDCG (for H, R, P, N) can be questioned
  - Exponential 100, 10, 1, 0 (or 4, 2, 1, 0) might be better
  - (already proposed in the original paper, but not much used…)
  - And "crazy" 100, 99, 90, 0 is wrong!

## A theorem

**Theorem 2** (Unbalanced query). *Let $Q$ be a query set, $q \notin Q$ a query, $D$ a document set, $\alpha$ a human relevance measurement and $\sigma$ and $\sigma'$ two system relevance measurements such that*

$$\underset{Q,D}{\mathsf{metric}}\,(\alpha,\sigma) > \underset{Q,D}{\mathsf{metric}}\,(\alpha,\sigma')$$

*and*

$$\underset{Q\cup\{q\},D}{\mathsf{metric}}\,(\alpha,\sigma) \leq \underset{Q\cup\{q\},D}{\mathsf{metric}}\,(\alpha,\sigma').$$

*Then*

$$\underset{q,D}{\mathsf{metric}}\,(\alpha,\sigma) < \underset{q,D}{\mathsf{metric}}\,(\alpha,\sigma').$$

## A theorem

**Theorem 2** (*Unbalanced query*). *Let $Q$ be a query set, $q \notin Q$ a query, $D$ a document set, $\alpha$ a human relevance measurement and $\sigma$ and $\sigma'$ two system relevance measurements such that*

$$\underset{Q,D}{\mathrm{metric}}\,(\alpha, \sigma) > \underset{Q,D}{\mathrm{metric}}\,(\alpha, \sigma')$$

*and*

$$\underset{Q \cup \{q\},D}{\mathrm{metric}}\,(\alpha, \sigma) \leq \underset{Q \cup \{q\},D}{\mathrm{metric}}\,(\alpha, \sigma').$$

*Then*

$$\underset{q,D}{\mathrm{metric}}\,(\alpha, \sigma) < \underset{q,D}{\mathrm{metric}}\,(\alpha, \sigma').$$

on the query set $Q$ and the document collection $D$, and given the human relevance judgment $\alpha$, $\sigma$ is more effective than $\sigma'$

we add a new query $q$ and then $\sigma$ becomes less effective than $\sigma'$
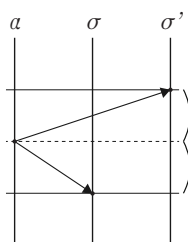
$\sigma$ is less effective than $\sigma'$ also on the new query $q$

## When we can't say anything (i.e., we can't state an axiom)
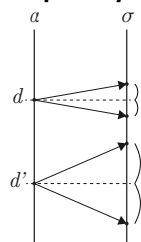


## Why we can't say anything

"P-oriented"
"R-oriented"          "Top-heavy"



## Hopefully you got the point

- By relying on **measurement theory**,
- one can define **relevance measurement,**
- that in turn allows to define **similarity** between (human and system) measurements,
- that in turn allows to define **Axioms** and **Theorems** on **metrics**
- that seem somehow interesting
  - not (always) trivial, more general, even inspiring …

## Outline + wrap-up

- Evaluation [5′]
  - Yes, it's complex!
- Measures / metrics [15′]
  - Oh dear, so many metrics?
- Measurement theory [15′]
  - Maybe a useful tool…
- Metrics analysis [5′]
  - Are we doing it wrong?!
- Axiometrics framework [15′]
  - Attempt to shift focus from metric A vs. B to metric properties

## Biblio

- [Amigó et al. 2009] A comparison of extrinsic clustering evaluation metrics based on formal constraints. IRJ 2009.
- [Bollman 1984] Two axioms for evaluation measures in information retrieval, SIGIR 1984.
- [Busin & Mizzaro 2013] Axiometrics: An Axiomatic Approach to Information Retrieval Effectiveness Metrics, ICTIR 2013.
- [Maddalena & Mizzaro 2014] Axiometrics: Axioms of Information Retrieval Effectiveness Metrics, EVIA@NTCIR 2014.
- [Sebastiani 2015] An Axiomatically Derived Measure for the Evaluation of Classification Algorithms, ICTIR 2015, to appear.
- [van Rijsbergen 1979] Information Retrieval, 1979.