



Avaliação de Sistemas de RI

Prof. Luciano Barbosa

(Parte do material retirado dos slides dos livros adotados)



Contexto

- Objetivo:
 - Medir a qualidade dos resultados dos engenhos de busca -> quanto o sistema satisfaz a necessidade do usuário
- Desafio:
 - Resultados podem ser intepretados de forma diferente pelos usuários -> avaliação por grupos diferentes
- Associa métrica quantitativa com o desempenho do sistema
- Compara diferentes sistemas de RI



2 Componentes Principais

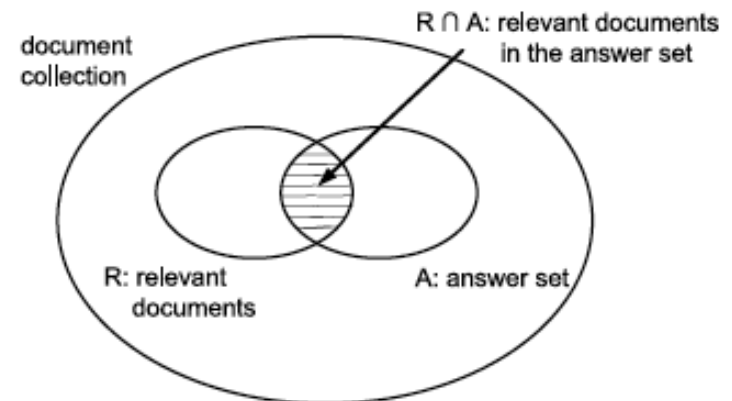
- Métrica: avaliação associada com a relevância dos resultados do usuário
- Dados rotulados: comparar resultados produzidos pelo sistema com resultados sugeridos por humanos para as mesmas consultas



Métricas

- Precisão: fração de docs relevantes do total recuperados

$$Precision = \frac{|R \cap A|}{|A|}$$



- Revocação: fração dos docs relevantes recuperados

$$Recall = \frac{|R \cap A|}{|R|}$$



Precisão e Revocação

- Todos os docs da coleção avaliados
- Em engenhos de busca:
 - Docs não apresentados de uma única vez
 - Lista ranqueada dos docs
- Varia com a posição do ranqueamento
- Mais apropriado: curva de precisão e revocação



Precisão e Revocação

- Considere o conjunto de 10 docs relevantes para a consulta q_1

$$R_{q_1} = \{d_3, d_5, d_9, d_{25}, d_{39}, d_{44}, d_{56}, d_{71}, d_{89}, d_{123}\}$$

- Considere um algoritmo que gere o ranqueamento:

01. d_{123} •	06. d_9 •	11. d_{38}
02. d_{84}	07. d_{511}	12. d_{48}
03. d_{56} •	08. d_{129}	13. d_{250}
04. d_6	09. d_{187}	14. d_{113}
05. d_8	10. d_{25} •	15. d_3 •



Precisão e Revocação

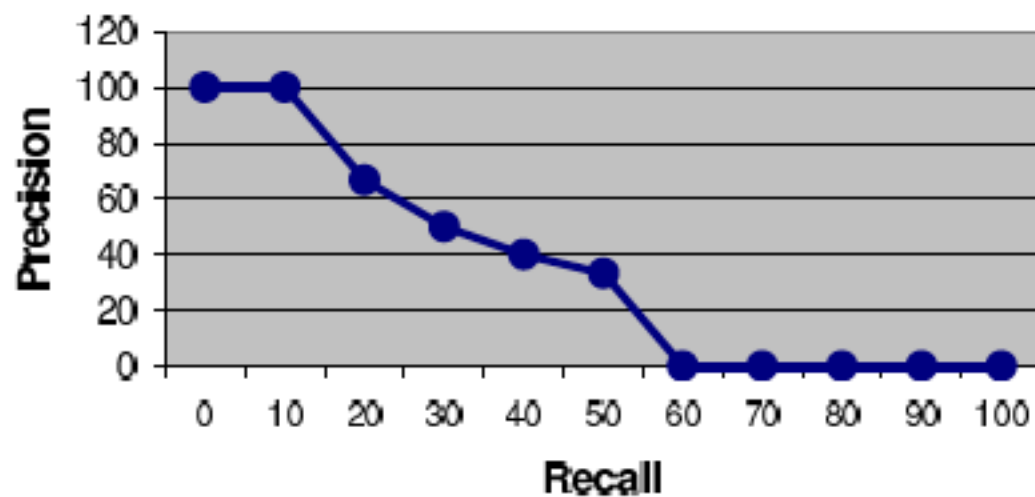
- Documento d_{123} na posição 1 é relevante
 - 10% dos relevantes (revocação) e precisão igual a 100%
- Documento d_{56} na posição 3 é relevante
 - Revocação = 20% e precisão = 66.6%

01. d_{123} •	06. d_9 •	11. d_{38}
02. d_{84}	07. d_{511}	12. d_{48}
03. d_{56} •	08. d_{129}	13. d_{250}
04. d_6	09. d_{187}	14. d_{113}
05. d_8	10. d_{25} •	15. d_3 •



Exemplo de Curva

- | | | |
|-----------------|----------------|---------------|
| 01. d_{123} • | 06. d_9 • | 11. d_{38} |
| 02. d_{84} | 07. d_{511} | 12. d_{48} |
| 03. d_{56} • | 08. d_{129} | 13. d_{250} |
| 04. d_6 | 09. d_{187} | 14. d_{113} |
| 05. d_8 | 10. d_{25} • | 15. d_3 • |

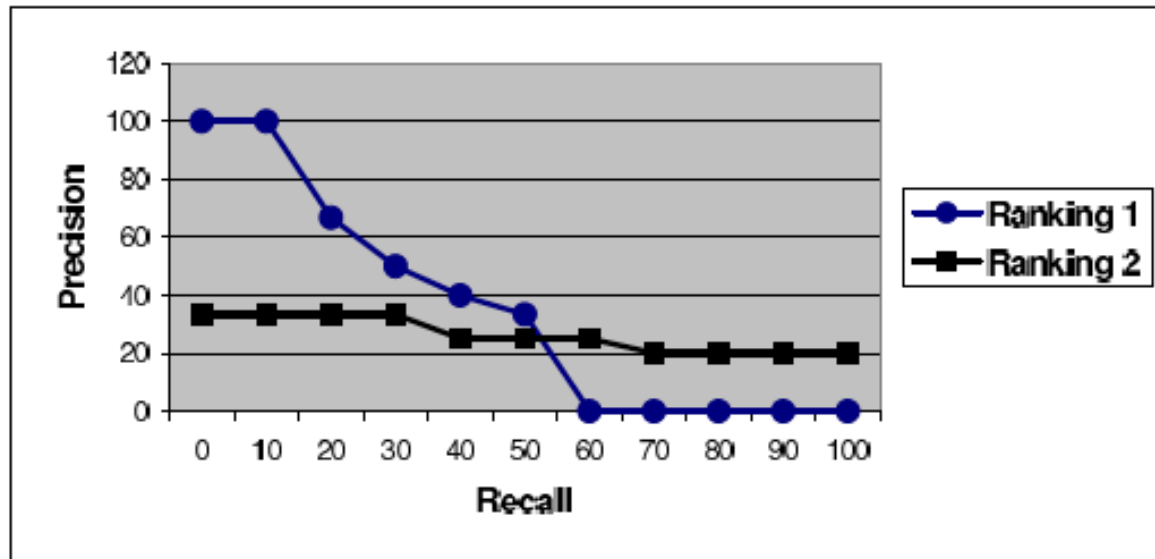


Recall	Precision
0	100
10	100
20	66.6
30	50
40	40
50	33.3
60	0
70	0
80	0
90	0
100	0



Precisão e Revocação

- Comparando algoritmos



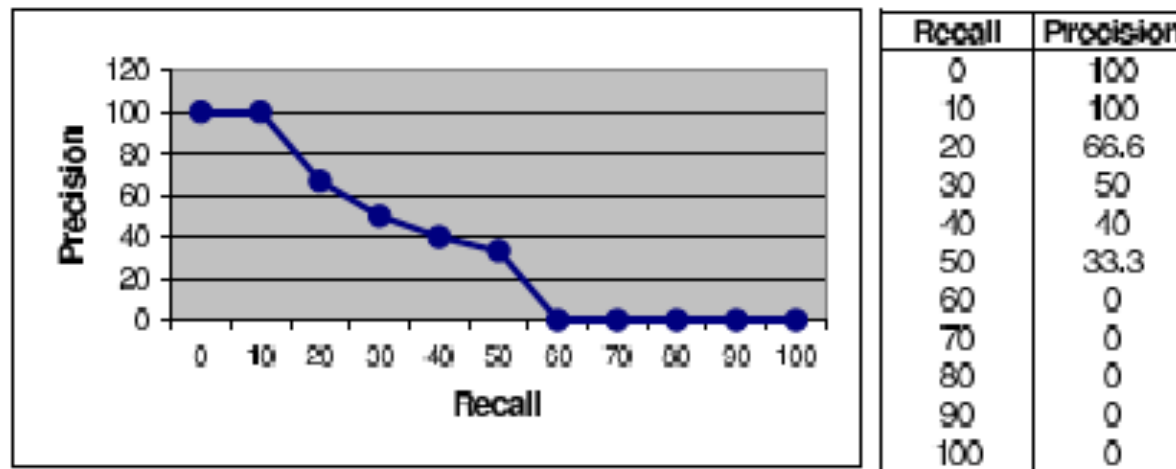
- Precisão média em um nível de revocação r

$$\overline{P}(r_j) = \sum_{i=1}^{N_q} \frac{P_i(r_j)}{N_q}$$



MAP: Mean Average Precision

- Média da precisão à medida que novos documentos são apresentados
- Exemplo:



$$MAP_1 = \frac{1 + 0.66 + 0.5 + 0.4 + 0.33 + 0 + 0 + 0 + 0 + 0}{10} = 0.28$$



Valor de Precisão

- Independente de recall
- Número de docs relevantes no topo do ranqueamento
- Precisão em 5 ($P@5$): mede a precisão qdo 5 docs são apresentados



Exemplo

- Resultado de consulta:

01. d_{123} •	06. d_9 •	11. d_{38}
02. d_{84}	07. d_{511}	12. d_{48}
03. d_{56} •	08. d_{129}	13. d_{250}
04. d_6	09. d_{187}	14. d_{113}
05. d_8	10. d_{25} •	15. d_3 •

- $P@5 = 40\%$ e $P@10 = 40\%$



Mean Reciprocal Rank

- Medir a primeira resposta correta
 - Q&A
 - Busca por sites ou URLs

Google bndes

Web News Images Videos Maps More Search tools

About 8,160,000 results (0.32 seconds)

BNDES - O banco nacional do desenvolvimento
www.bndes.gov.br/ [Translate this page](#)
Um órgão vinculado ao Ministério do Desenvolvimento, Indústria e Comércio Exterior e tem como objetivo apoiar empreendimentos que contribuam para o ...
[Apoio Financeiro](#) - [BNDES Transparente](#) - [A Empresa](#) - [Acesso à Informação](#)

BNDES - Brazilian Development Bank
www.bndes.gov.br/SiteBNDES/bndes/bndes_en/ [Translate this page](#)
The Brazilian Development Bank (BNDES) and the New Development Bank of the BRICS signed a memorandum of understanding aimed at strengthening ...

Cartão BNDES
<https://www.cartaoibndes.gov.br/> [Translate this page](#)
Clique aqui para conhecer o Cartão BNDES. Parceiros. Banco do Brasil Banco do Nordeste Banco Santander Banrisul Bradesco BRDE Caixa Econômica ...
[Simulador](#) - [Solicite seu Cartão BNDES](#) - [Busca de Produtos](#) - [Dúvidas](#)

Quora what is the highest

Home Write Notifications Luciano Ask Question

16 FOLLOWERS

Last asked: 6 Jan

QUESTION TOPICS

- Mount Everest
- Mountains
- Geography

Edit Topics

QUESTION STATS

Views	2,898
Followers	16
Merged Questions	5
Edits	

What is the highest mountain in the world?

[Re-Ask](#) Follow 16 Comment Share Downvote

Have this question too? Re-Ask to get an answer.

Luciano Barbosa
Add Bio • Make Anonymous

Write your answer, or answer later

6 Answers

Ben Cherry, Entrepreneur, Engineer
1.1k Views

The simplest measurement is the point on the Earth's surface that is farthest from the Earth's center. By this measure, Chimborazo in Ecuador is the highest mountain in the world.

Chimborazo's summit is 3,967.1 miles from the center of the earth, vs. just 3,965.8 miles for Mount Everest (a difference of 1.3 miles).

<http://en.wikipedia.org/wiki/Chimborazo>

Written 13 Dec 2011 • View Upvotes

[Upvote](#) 10 [Downvote](#) [Comments](#) 14 [Share](#)

Christian Van Den Berge, Photographer, analytics, SEO, web, travel
577 Views

actually it depends where you start measuring; at the base of the mountain or at sea level?

From sea level it is indeed Mount Everest, but the highest mountain really is

There's more on Quora...

Pick new people and topics to follow and see the best answers on Quora.

[Update Your Interests](#)

RELATED QUESTIONS

- What does a mountain contain?
- What is the highest mountain in Ireland?
- What is the view like from the summit of the highest mountains of the world?
- What is the highest mountain in Australia?
- What is the world's tallest mountain (base to peak)?
- What is the world's largest (volume) mountain?
- What rock type makes up the third highest mountain in the world, Kanchenjunga Mountain?
- Where is the world's largest mountain located?
- What area has the highest density of Starbucks in the world?
- What is the highest mountain peak of Western Ghats of India?

[More Related Questions](#)



MRR: Mean Reciprocal Rank

- Reciprocal ranking
 - R_i : ranqueamento relativo à consulta q_i
 - $S_{correct}(R_i)$: posição da primeira resposta correta em R_i
 - S_h : limiar para posição no ranqueamento

$$RR(\mathcal{R}_i) = \begin{cases} \frac{1}{S_{correct}(\mathcal{R}_i)} & \text{if } S_{correct}(\mathcal{R}_i) \leq S_h \\ 0 & \text{otherwise} \end{cases}$$

- Para um conjunto de consultas

$$MRR(Q) = \sum_i^{N_q} RR(\mathcal{R}_i)$$



E-Measure

- Combina precisão e revocação

$$E(j) = 1 - \frac{1 + b^2}{\frac{b^2}{r(j)} + \frac{1}{P(j)}}$$

- b : parâmetro para dar mais peso para precisão ou revocação (definido pelo usuário)
- $r(j)$: revocação na posição j do ranqueamento
- $P(j)$: precisão na posição j do ranqueamento



E-Measure

$$E(j) = 1 - \frac{1 + b^2}{\frac{b^2}{r(j)} + \frac{1}{P(j)}}$$

- $b=0$
 - $E(j) = 1 - P(j)$
 - $E(j)$ se torna uma função de precisão
- $b \rightarrow \infty$
 - $\lim_{b \rightarrow \infty} E(j) = 1 - r(j)$
 - $E(j)$ se torna uma função de revocação
- $b=1$: F-measure (média harmônica)



F-Measure

- Valor único que combina precisão e revocação

$$F(j) = \frac{2}{\frac{1}{r(j)} + \frac{1}{P(j)}}$$

- Assume valores entre 0 e 1
- $F(j)=0$: nenhum documento relevante é recuperado
- $F(j)=1$: todos documentos relevantes são recuperados
- Assume valor alto apenas quando ambos precisão e revocação são altos
- Para maximizar F-measure, é preciso encontrar o melhor compromisso entre precisão e revocação



Chamada



Cumulated Gain (CG)

- Precisão e revocação medem apenas avaliações binárias
- Sem distinção de graus de relevância
- CG combina diferentes graus de relevância



Cumulated Gain (CG)

- Considere resultados com notas de 0 a 3
- 0 não relevante e 3 muito relevante
- Por exemplo

$$\begin{aligned} R_{q_1} &= \{ [d_3, 3], [d_5, 3], [d_9, 3], [d_{25}, 2], [d_{39}, 2], \\ &\quad [d_{44}, 2], [d_{56}, 1], [d_{71}, 1], [d_{89}, 1], [d_{123}, 1] \} \\ R_{q_2} &= \{ [d_3, 3], [d_{56}, 2], [d_{129}, 1] \} \end{aligned}$$



Cumulated Gain (CG)

- Do topo para o último, soma os valores cumulativos em cada ponto do ranqueamento

$$G_1 = (1, 0, 1, 0, 0, 3, 0, 0, 0, 2, 0, 0, 0, 0, 3)$$

$$G_2 = (0, 0, 2, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 3)$$

$$CG_1 = (1, 1, 2, 2, 2, 5, 5, 5, 5, 7, 7, 7, 7, 7, 10)$$

$$CG_2 = (0, 0, 2, 2, 2, 2, 2, 3, 3, 3, 3, 3, 3, 3, 6)$$



Cumulated Gain (CG)

- De forma mais formal:

$$CG_j[i] = \begin{cases} G_j[1] & \text{if } i = 1; \\ G_j[i] + CG_j[i - 1] & \text{otherwise} \end{cases}$$



Discounted Cumulated Gain (DCG)

- Modificação: menor peso para documentos mais baixo no ranking
- Logaritmo da posição no ranqueamento
- Mais formalmente:

$$DCG_j[i] = \begin{cases} G_j[1] & \text{if } i = 1; \\ \frac{G_j[i]}{\log_2 i} + DCG_j[i - 1] & \text{otherwise} \end{cases}$$



Discounted Cumulated Gain (DCG)

$$G_1 = (1, \boxed{0, 1}, 0, 0, 3, 0, 0, 0, 2, 0, 0, 0, 0, 3)$$

$$G_2 = (0, 0, 2, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 3)$$

$$CG_1 = (1, \boxed{1, 2}, 2, 2, 5, 5, 5, 5, 7, 7, 7, 7, 7, 10)$$

$$CG_2 = (0, 0, 2, 2, 2, 2, 2, 3, 3, 3, 3, 3, 3, 3, 6)$$

$$DCG_1 = (1.0, \boxed{1.0, 1.6}, 1.6, 1.6, 2.8, 2.8, 2.8, 2.8, 3.4, 3.4, 3.4, 3.4, 3.4, 4.2)$$

$$DCG_2 = (0.0, 0.0, 1.3, 1.3, 1.3, 1.3, 1.3, 1.6, 1.6, 1.6, 1.6, 1.6, 1.6, 1.6, 2.4)$$

- Menos afetados por docs relevantes no fim do ranqueamento
- Calcula-se média sobre conjunto de consultas para comparação

$$\overline{CG}[i] = \sum_{j=1}^{N_q} \frac{CG_j[i]}{N_q}; \quad \overline{DCG}[i] = \sum_{j=1}^{N_q} \frac{DCG_j[i]}{N_q}$$



Discounted Cumulated Gain (DCG)

$$G_1 = (1, 0, 1, 0, 0, 3, 0, 0, 0, 2, 0, 0, 0, 0, 3)$$

$$G_2 = (0, 0, 2, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 3)$$

$$CG_1 = (1, 1, 2, 2, 2, 5, 5, 5, 5, 7, 7, 7, 7, 7, 10)$$

$$CG_2 = (0, 0, 2, 2, 2, 2, 2, 3, 3, 3, 3, 3, 3, 3, 6)$$

$$DCG_1 = (1.0, 1.0, 1.6, 1.6, 1.6, 2.8, 2.8, 2.8, 2.8, 3.4, 3.4, 3.4, 3.4, 3.4, 4.2)$$

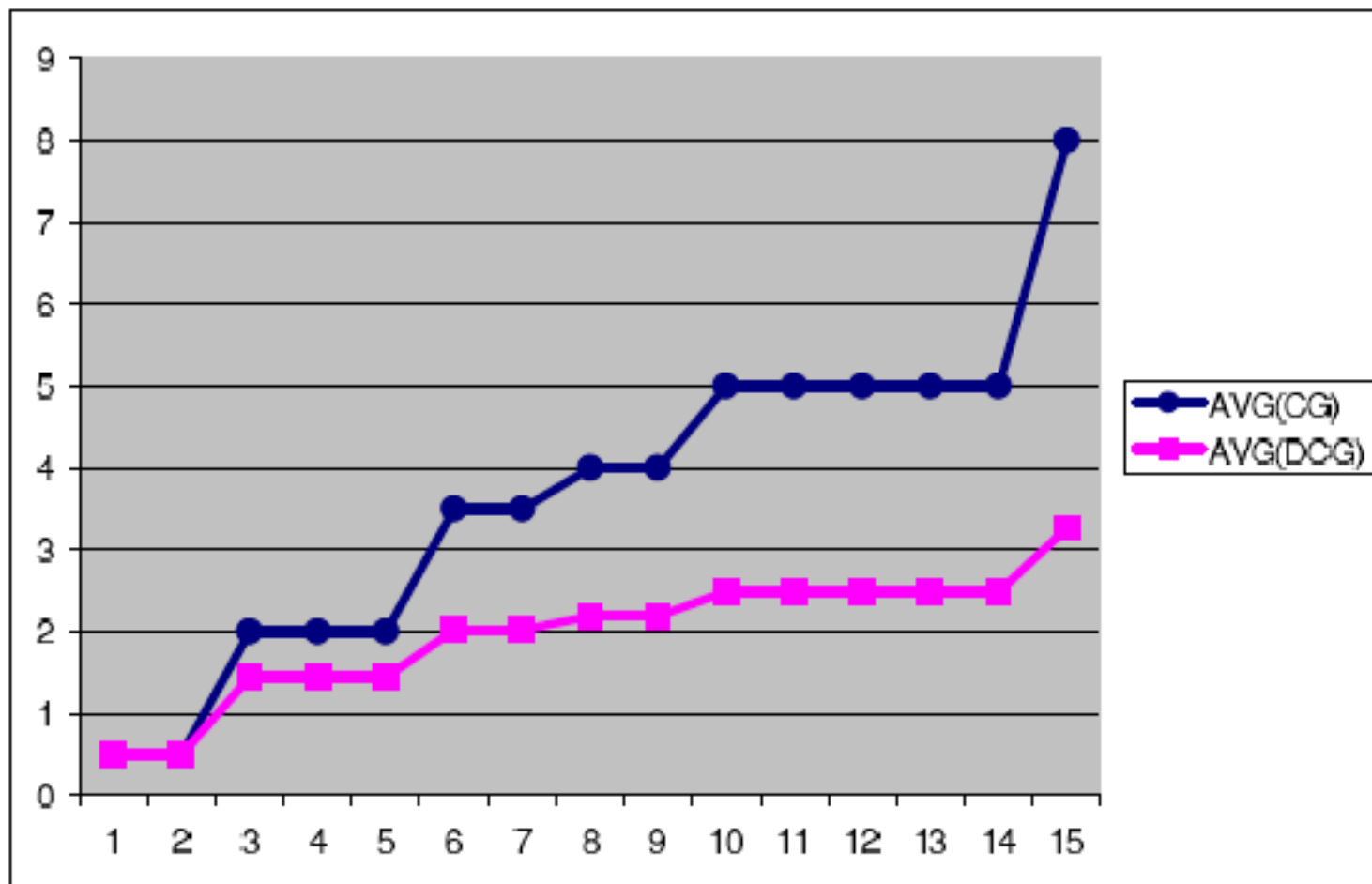
$$DCG_2 = (0.0, 0.0, 1.3, 1.3, 1.3, 1.3, 1.3, 1.6, 1.6, 1.6, 1.6, 1.6, 1.6, 1.6, 2.4)$$

- Menos afetados por docs relevantes no fim do ranqueamento
- Calcula-se média sobre conjunto de consultas para comparação

$$\overline{CG}[i] = \sum_{j=1}^{N_q} \frac{CG_j[i]}{N_q}; \quad \overline{DCG}[i] = \sum_{j=1}^{N_q} \frac{DCG_j[i]}{N_q}$$



Ganho Acumulado Descontado (DCG)





DCG e CG Ideais

- Comparar DCG e CG com um baseline
- Exemplo para notas de 0 a 3
- Ordena as notas: ideal gain vector
- Ranqueamento ideal

$$R_{q_1} = \{ [d_3, 3], [d_5, 3], [d_9, 3], [d_{25}, 2], [d_{39}, 2], \\ [d_{44}, 2], [d_{56}, 1], [d_{71}, 1], [d_{89}, 1], [d_{123}, 1] \}$$

$$R_{q_2} = \{ [d_3, 3], [d_{56}, 2], [d_{129}, 1] \}$$

$$IG_1 = (3, 3, 3, 2, 2, 2, 1, 1, 1, 1, 0, 0, 0, 0, 0)$$

$$IG_2 = (3, 2, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0)$$

- Quanto o ranqueamento está distante do melhor caso



DCG e CG Ideais

$$IG_1 = (3, 3, 3, 2, 2, 2, 1, 1, 1, 1, 0, 0, 0, 0, 0)$$

$$IG_2 = (3, 2, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0)$$

$$ICG_1 = (3, 6, 9, 11, 13, 15, 16, 17, 18, 19, 19, 19, 19, 19, 19)$$

$$ICG_2 = (3, 5, 6, 6, 6, 6, 6, 6, 6, 6, 6, 6, 6, 6, 6)$$

$$IDCG_1 = (3.0, 6.0, 7.9, 8.9, 9.8, 10.5, 10.9, 11.2, 11.5, 11.8, 11.8, 11.8, 11.8, 11.8, 11.8)$$

$$IDCG_2 = (3.0, 5.0, 5.6, 5.6, 5.6, 5.6, 5.6, 5.6, 5.6, 5.6, 5.6, 5.6, 5.6, 5.6, 5.6)$$



DCG Normalizado

- Comparar diferentes curvas de DCG para algoritmos de ranqueamento diferentes
- Algoritmo de ranqueamento usado em N_q consultas
- Fórmulas:

$$NCG[i] = \frac{\overline{CG}[i]}{\overline{ICG}[i]}; \quad NDCG[i] = \frac{\overline{DCG}[i]}{\overline{IDCG}[i]}$$

$$\overline{ICG}[i] = \sum_{j=1}^{N_q} \frac{ICG_j[i]}{N_q}; \quad \overline{IDCG}[i] = \sum_{j=1}^{N_q} \frac{IDCG_j[i]}{N_q}$$



DCG Normalizado

$$NCG[i] = \frac{\overline{CG}[i]}{\overline{ICG}[i]}; \quad NDCG[i] = \frac{\overline{DCG}[i]}{\overline{IDCG}[i]}$$

$$\overline{CG} = (0.5, 0.5, 2.0, 2.0, 2.0, 3.5, 3.5, 4.0, 4.0, 5.0, 5.0, 5.0, 5.0, 5.0, 8.0)$$

$$\overline{DCG} = (0.5, 0.5, 1.5, 1.5, 1.5, 2.1, 2.1, 2.2, 2.2, 2.5, 2.5, 2.5, 2.5, 2.5, 3.3)$$

$$\overline{ICG} = (3.0, 5.5, 7.5, 8.5, 9.5, 10.5, 11.0, 11.5, 12.0, 12.5, 12.5, 12.5, 12.5, 12.5, 12.5)$$

$$\overline{IDCG} = (3.0, 5.5, 6.8, 7.3, 7.7, 8.1, 8.3, 8.4, 8.6, 8.7, 8.7, 8.7, 8.7, 8.7, 8.7)$$

$$NCG = (0.17, 0.09, 0.27, 0.24, 0.21, 0.33, 0.32, \\ 0.35, 0.33, 0.40, 0.40, 0.40, 0.40, 0.40, 0.64)$$

$$NDCG = (0.17, 0.09, 0.21, 0.20, 0.19, 0.25, 0.25, \\ 0.26, 0.26, 0.29, 0.29, 0.29, 0.29, 0.29, 0.38)$$



Discussão sobre Métricas DCG

- Levam em consideração vários níveis de relevância
- DCG permite diminuir o peso de docs relevantes com baixo valor de ranqueamento
- Limitação: difícil definir vários níveis de relevância



Correlação de Ranqueamento

- Em alguns casos
 - Não conseguimos avaliar relevância diretamente
 - Interessados em determinar quanto um ranqueamento diferencia de outro já conhecido
- Comparar ordem relativa de dois ranqueamentos
- Técnicas estatísticas de correlação de ranqueamento



Correlação de Spearman

$$S(\mathcal{R}_1, \mathcal{R}_2) = 1 - \frac{6 \times \sum_{j=1}^K (s_{1,j} - s_{2,j})^2}{K \times (K^2 - 1)}$$

- \mathcal{R}_1 e \mathcal{R}_2 : ranqueamentos a serem comparados
- K : número de elementos no ranqueamento
- s_{ij} : posição do documento no ranqueamento i na posição j
- Métricas
 - Correlação=1: ranqueamentos iguais
 - Correlação=0: completamente independentes
 - Correlação=-1: ranqueamento inverso



Correlação de Spearman

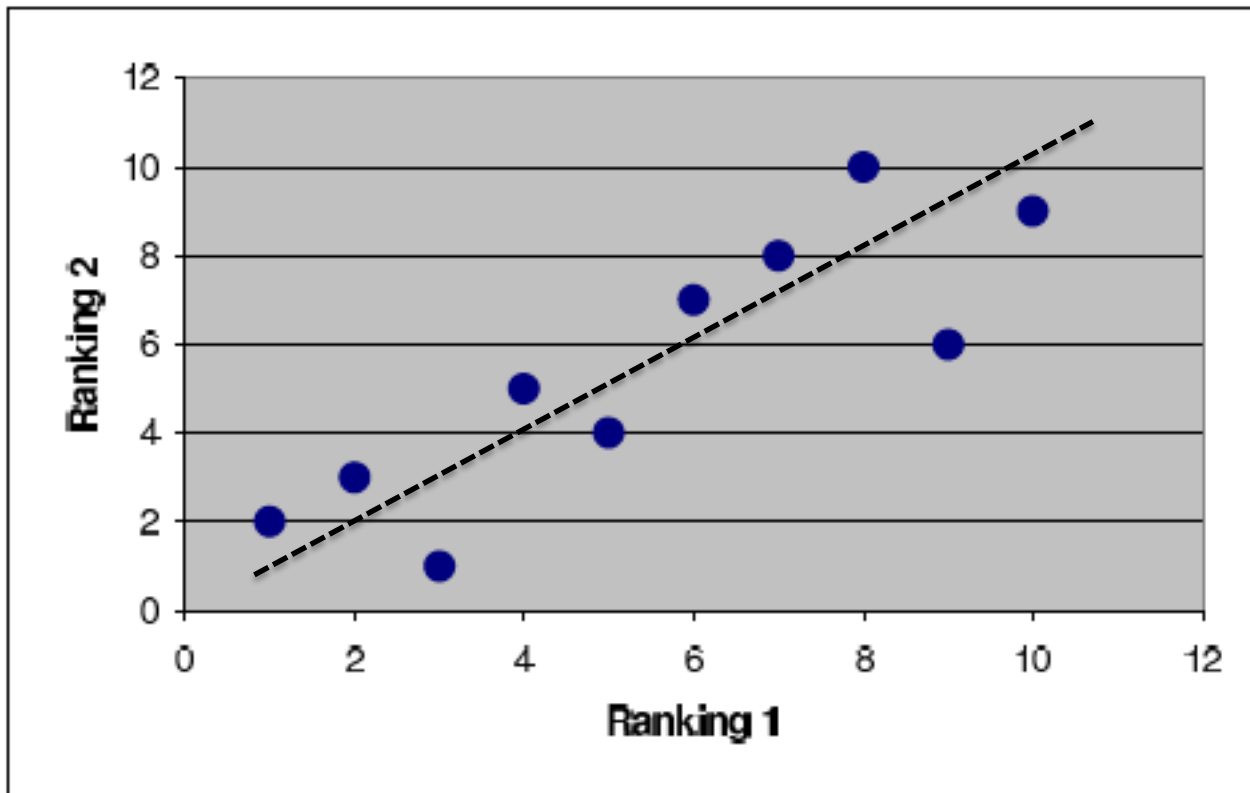
$$S(\mathcal{R}_1, \mathcal{R}_2) = 1 - \frac{6 \times \sum_{j=1}^K (s_{1,j} - s_{2,j})^2}{K \times (K^2 - 1)}$$

$$S(\mathcal{R}_1, \mathcal{R}_2) = 1 - \frac{6 \times 24}{10 \times (10^2 - 1)} = 1 - \frac{144}{990} = 0.854$$

documents	$s_{1,j}$	$s_{2,j}$	$s_{1,j} - s_{2,j}$	$(s_{1,j} - s_{2,j})^2$
d_{123}	1	2	-1	1
d_{84}	2	3	-1	1
d_{56}	3	1	+2	4
d_6	4	5	-1	1
d_8	5	4	+1	1
d_9	6	7	-1	1
d_{511}	7	8	-1	1
d_{129}	8	10	-2	4
d_{187}	9	6	+3	9
d_{25}	10	9	+1	1
Sum of Square Distances				24

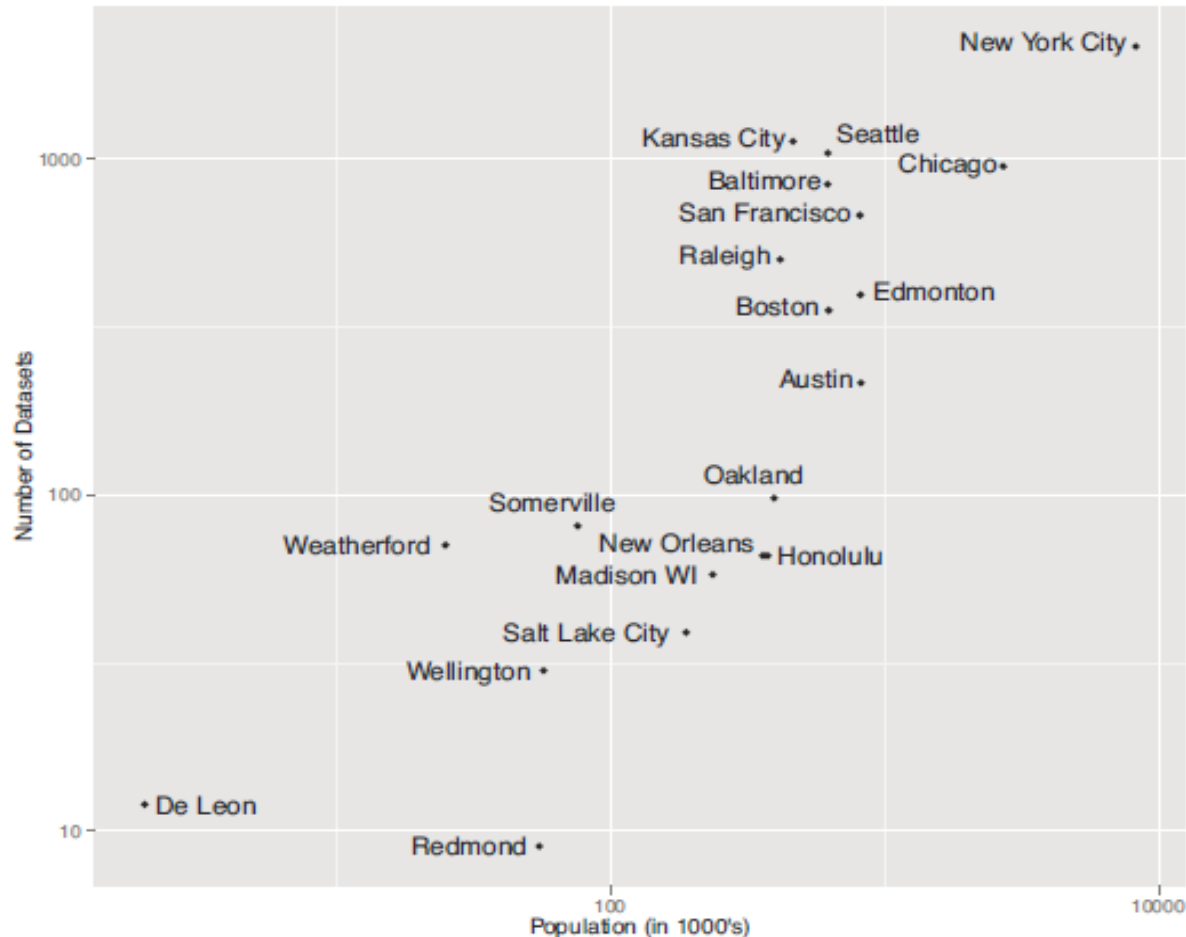


Correlação de Spearman





Uso em Outros Domínios: Dados Urbanos





Kendal Tau

- Correlação de Spearman não tem uma interpretação clara
- Kendal Tau: baseado na ideia de pares de docs concordantes e discordantes em dois ranqueamentos



Kendal Tau

- Considere os ranqueamentos dados por R_1 e R_2 :

documents	$s_{1,j}$	$s_{2,j}$	$s_{1,j} - s_{2,j}$
d_{123}	1	2	-1
d_{84}	2	3	-1
d_{56}	3	1	+2
d_6	4	5	-1
d_8	5	4	+1



Ranqueamentos de docs por R_1 e R_2 :

documents	$s_{1,j}$	$s_{2,j}$	$s_{1,j} - s_{2,j}$
d_{123}	1	2	-1
d_{84}	2	3	-1
d_{56}	3	1	+2
d_6	4	5	-1
d_8	5	4	+1

Pares de docs ordenados de R_1

$[d_{123}, d_{84}]$, $[d_{123}, d_{56}]$, $[d_{123}, d_6]$, $[d_{123}, d_8]$,
 $[d_{84}, d_{56}]$, $[d_{84}, d_6]$, $[d_{84}, d_8]$,
 $[d_{56}, d_6]$, $[d_{56}, d_8]$,
 $[d_6, d_8]$

Pares de docs ordenados de R_2

$[d_{56}, d_{123}]$, $[d_{56}, d_{84}]$, $[d_{56}, d_8]$, $[d_{56}, d_6]$,
 $[d_{123}, d_{84}]$, $[d_{123}, d_8]$, $[d_{123}, d_6]$,
 $[d_{84}, d_8]$, $[d_{84}, d_6]$,
 $[d_8, d_6]$



Kendal Tau

- Contar os pares concordantes e discordantes
- No exemplo anterior:
 - Para um total de 20 pares ordenados: $K(K-1)$
 - 14 concordantes e 6 discordantes
- Coeficiente de Kendal Tau

$$\tau(\mathcal{R}_1, \mathcal{R}_2) = P(\mathcal{R}_1 = \mathcal{R}_2) - P(\mathcal{R}_1 \neq \mathcal{R}_2)$$

- No exemplo:

$$\begin{aligned}\tau(\mathcal{R}_1, \mathcal{R}_2) &= \frac{14}{20} - \frac{6}{20} \\ &= 0.4\end{aligned}$$



Kendal Tau

- Seja:
 - Número de pares discordantes: $\Delta(\mathcal{R}_1, \mathcal{R}_2)$
 - Total de pares: $K(K - 1)$
- Coeficiente kendal tau

$$\tau(\mathcal{R}_1, \mathcal{R}_2) = 1 - \frac{2 \times \Delta(\mathcal{R}_1, \mathcal{R}_2)}{K(K-1)}$$

- No exemplo anterior:

$$\tau(\mathcal{R}_1, \mathcal{R}_2) = 1 - \frac{2 \times 6}{5(5 - 1)} = 0.4$$



Coleções de Referência: TREC

- TREC: conferência dedicada a experimentos em grandes coleções de teste
- Um conjunto de experimentos são apresentados
- Grupos de pesquisa comparam resultados
- 3 partes:
 - Documentos
 - Tópicos: consultas descritas em linguagem natural
 - Conjunto de documentos relevantes para cada tópico



Coleções de Referência: TREC

- Vem crescendo com o tempo
 - TREC-3: 2 Gb
 - TREC-6: 5.8 Gb
 - Terabyte (GOV2): 25 milhões de páginas no .gov



Coleções de Referência: TREC

- Fontes

WSJ	→ <i>Wall Street Journal</i>
AP	→ Associated Press (news wire)
ZIFF	→ Computer Selects (articles), Ziff-Davis
FR	→ Federal Register
DOE	→ US DOE Publications (abstracts)
SJMN	→ <i>San Jose Mercury News</i>
PAT	→ US Patents
FT	→ <i>Financial Times</i>
CR	→ Congressional Record
FBIS	→ Foreign Broadcast Information Service
LAT	→ <i>LA Times</i>



Documentos na Coleção: TREC-6 (Discos 1 e 2)

Disk	Contents	Size Mb	Number Docs	Words/Doc. (median)	Words/Doc. (mean)
1	WSJ, 1987-1989	267	98,732	245	434.0
	AP, 1989	254	84,678	446	473.9
	ZIFF	242	75,180	200	473.0
	FR, 1989	260	25,960	391	1315.9
	DOE	184	226,087	111	120.4
2	WSJ, 1990-1992	242	74,520	301	508.4
	AP, 1988	237	79,919	438	468.7
	ZIFF	175	56,920	182	451.9
	FR, 1988	209	19,860	396	1378.1



Documentos na Coleção: TREC-6 (Discos 3-6)

Disk	Contents	Size Mb	Number Docs	Words/Doc. (median)	Words/Doc. (mean)
3	SJMN, 1991	287	90,257	379	453.0
	AP, 1990	237	78,321	451	478.4
	ZIFF	345	161,021	122	295.4
	PAT, 1993	243	6,711	4,445	5391.0
4	FT, 1991-1994	564	210,158	316	412.7
	FR, 1994	395	55,630	588	644.7
	CR, 1993	235	27,922	288	1373.5
5	FBIS	470	130,471	322	543.6
	LAT	475	131,896	351	526.5
6	FBIS	490	120,653	348	581.3



Coleções

- Documentos no formato SGML
- Estruturas comuns
 - Número do documento <DOCNO>
 - Campo com o texto do doc <TEXT>
- Estrutura depende da coleção



Example de um Documento da TREC (WSJ)

```
<doc>
```

```
<docno> WSJ880406-0090 </docno>
```

```
<hl> AT&T Unveils Services to Upgrade Phone Networks  
Under Global Plan </hl>
```

```
<author> Janet Guyon (WSJ Staff) </author>
```

```
<dateline> New York </dateline>
```

```
<text>
```

```
American Telephone & Telegraph Co introduced the first  
of a new generation of phone services with broad ...
```

```
</text>
```

```
</doc>
```




Coleções da Web

- Tarefa de RI na Web (introduzido na TREC-9)
 - VLC2: Internet Archive de 1997
 - WT2g e WT10g: subconjuntos do VLC2
 - .GOV e .GOV2: páginas do .gov (2002 e 2004)

Collection	# Docs	Avg Doc Size	Collection Size
VLC2 (WT100g)	18,571,671	5.7 KBytes	100 GBytes
WT2g	247,491	8.9 KBytes	2.1 GBytes
WT10g	1,692,096	6.2 KBytes	10 GBytes
.GOV	1,247,753	15.2 KBytes	18 GBytes
.GOV2	27 million	15 KBytes	400 GBytes



TREC: Consulta

- Representa a necessidade de informação do usuário
- Consulta = tópico
- Coleção da TREC -> conjunto de tópicos
- Em linguagem natural
- Conversão de um tópico em uma consulta



Exemplo de Tópico

<top>

<num> Number: 168

<title> Topic: Financing AMTRAK

<desc> Description:

A document will address the role of the Federal Government in financing the operation of the National Railroad Transportation Corporation (AMTRAK)

<narr> Narrative: A relevant document must provide information on the government's responsibility to make AMTRAK an economically viable entity. It could also discuss the privatization of AMTRAK as an alternative to continuing government subsidies. Documents comparing government subsidies given to air and bus transportation with those provided to AMTRAK would also be relevant

</top>



Documentos Relevantes

- Obtido baseado no método de pooling
 - Coleção com os melhores resultados produzidos por vários sistemas (ex. 100)
 - Documentos na coleção são avaliados por pessoas
- Suposições
 - A grande maioria dos documentos relevantes está na coleção
 - Os documentos não presentes na coleção são irrelevantes



Tarefas

- Ad hoc: consultas feitas a um base fixa de documentos
- Roteamento: consultas sobre a base em constante mudança
- Outras tarefas: chinês, Q&A, entre-línguas, legal, IR com fala, blog, spam etc