

# CIÊNCIA DE DADOS (BIG DATA PROCESSING AND ANALYTICS)

RECUPERAÇÃO DA INFORMAÇÃO NA WEB  
E EM REDES SOCIAIS



Professor curador  
Luciano Moreira Camilo e Silva





# **TRILHA 1**

## **INTRODUÇÃO À RECUPERAÇÃO DA INFORMAÇÃO**

---

# Recuperação de dados *versus* Recuperação da informação

## Recuperação de dados

determinístico • estruturado • exato



Fonte: GettyImages.

## Recuperação da informação

probabilístico • não-estruturado • relevante



Fonte: GettyImages.

# Recuperação de dados *versus* Recuperação da informação



**Fonte:** Elaborada pelo autor.

# Sistemas de Recuperação da informação

- Disciplina da grande área das ciências sociais (biblioteconomia).
- Evolução em parceria com a disciplina de sistemas da informação.
- Tornou-se dominante após a década de 1990, como a ferramenta *de facto* de recuperação de documentos na web.



Sistema ZATOR 800 de recuperação da informação de 1951.  
Fonte: MOOERS, C. N. *Aslib Proceedings*, v. 8, n. 1, p. 3-22, [s.d.].

# Modelos de sistema de Recuperação da informação

## Modelo Booleano

São Paulo

Santo

Paulo

Cidade

Fonte: Elaborada pelo autor.

## Modelo Vetorial

Consulta:  
"São Paulo" E "Cidade"

7

$\theta_7$

9

$\theta_7$

$\theta_8$

8

Cidade





# Índices Invertidos

*No meio do caminho tinha uma pedra  
Tinha uma pedra no meio do caminho  
Tinha uma pedra  
No meio do caminho tinha uma pedra  
Nunca me esquecerei desse acontecimento  
Na vida de minhas retinas tão fatigadas  
Nunca me esquecerei que no meio do caminho  
Tinha uma pedra  
Tinha uma pedra no meio do caminho  
No meio do caminho tinha uma pedra.*

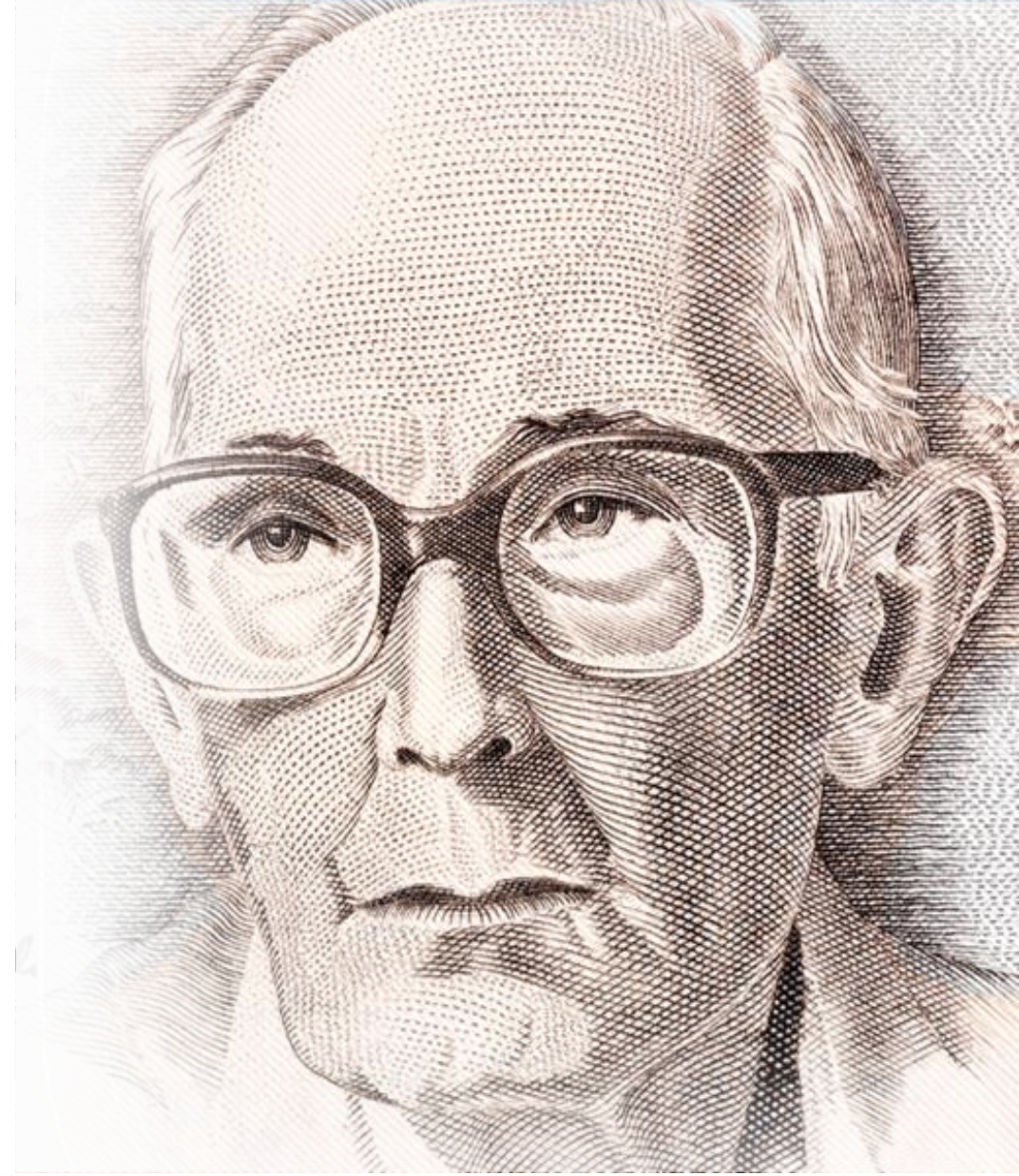
Carlos Drummond de Andrade





# Índices Invertidos

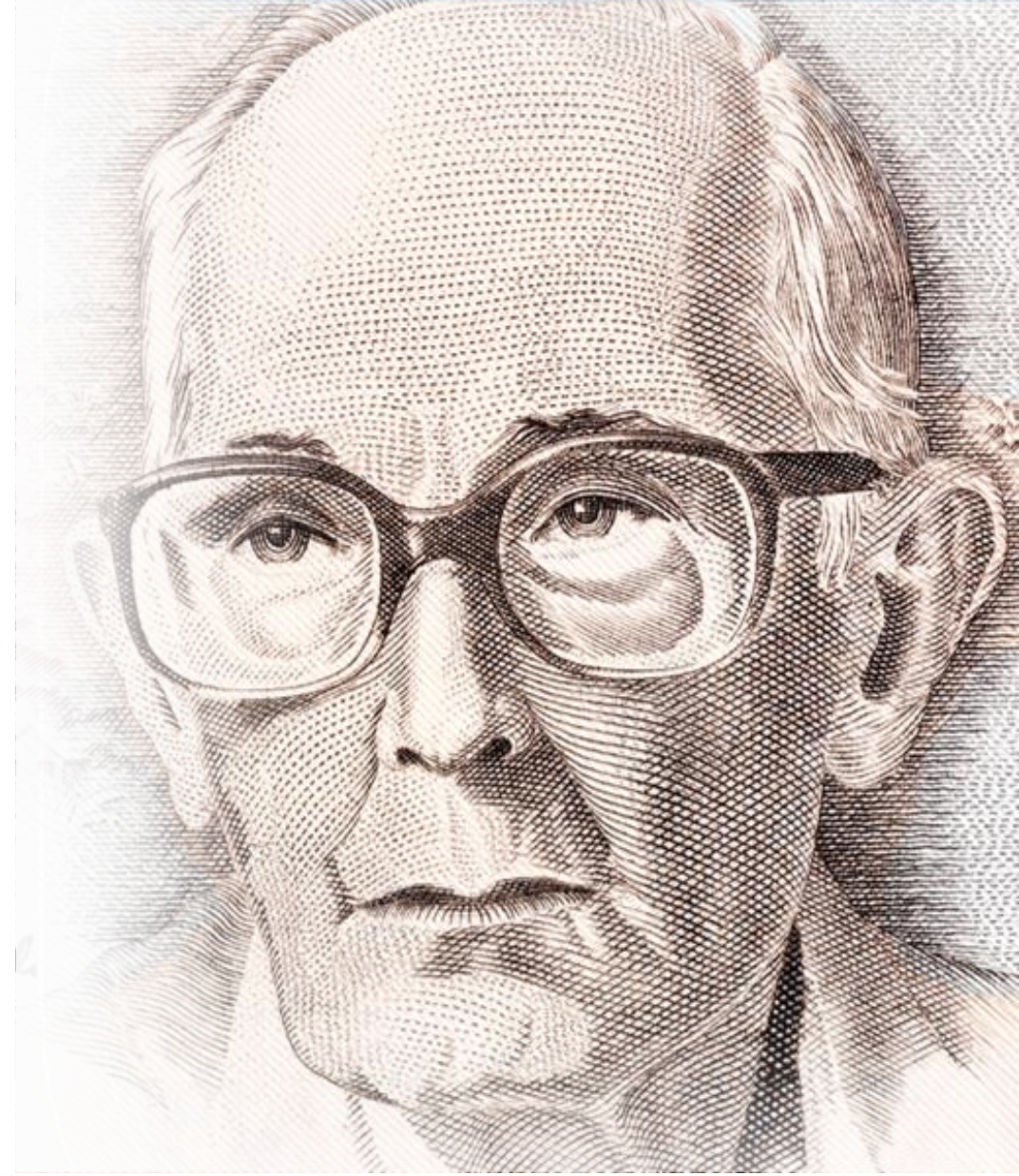
1. *No meio do caminho tinha uma pedra*
2. *Tinha uma pedra no meio do caminho*
3. *Tinha uma pedra*
4. *No meio do caminho tinha uma pedra*
5. *Nunca me esquecerei desse acontecimento*
6. *Na vida de minhas retinas tão fatigadas*
7. *Nunca me esquecerei que no meio do caminho*
8. *Tinha uma pedra*
9. *Tinha uma pedra no meio do caminho*
10. *No meio do caminho tinha uma pedra.*





# Índices Invertidos

1. *no meio do caminho tinha uma pedra*
2. *tinha uma pedra no meio do caminho*
3. *tinha uma pedra*
4. *no meio do caminho tinha uma pedra*
5. *nunca me esquecerei desse acontecimento*
6. *na vida de minhas retinas tão fatigadas*
7. *nunca me esquecerei que no meio do caminho*
8. *tinha uma pedra*
9. *tinha uma pedra no meio do caminho*
10. *no meio do caminho tinha uma pedra.*



## Índices Invertidos

1. *no meio do caminho tinha uma pedra*
2. *tinha uma pedra no meio do caminho*
3. *tinha uma pedra*
4. *no meio do caminho tinha uma pedra*
5. *nunca me esquecerei desse acontecimento*
6. *na vida de minhas retinas tão fatigadas*
7. *nunca me esquecerei que no meio do caminho*
8. *tinha uma pedra*
9. *tinha uma pedra no meio do caminho*
10. *no meio do caminho tinha uma pedra.*

Word	Doc_Iids
	5
	1 2 4 7 9 10

## Índices Invertidos

1. *no meio do caminho tinha uma pedra*
2. *tinha uma pedra no meio do caminho*
3. *tinha uma pedra*
4. *no meio do caminho tinha uma pedra*
5. *nunca me esquecerei desse acontecimento*
6. *na vida de minhas retinas tão fatigadas*
7. *nunca me esquecerei que no meio do caminho*
8. *tinha uma pedra*
9. *tinha uma pedra no meio do caminho*
10. *no meio do caminho tinha uma pedra.*

Word	Doc_ids
<i><b>acontecimento</b></i>	5
<i><b>caminho</b></i>	1,2,4,7,9,10

# Índices Invertidos

1. *no meio do caminho tinha uma pedra*
2. *tinha uma pedra no meio do caminho*
3. *tinha uma pedra*
4. *no meio do caminho tinha uma pedra*
5. *nunca me esquecerei desse acontecimento*
6. *na vida de minhas retinas tão fatigadas*
7. *nunca me esquecerei que no meio do caminho*
8. *tinha uma pedra*
9. *tinha uma pedra no meio do caminho*
10. *no meio do caminho tinha uma pedra.*

Word	Doc_Ids
<i>acontecimento</i>	5
<i>caminho</i>	1,2,4,7,9,10
<i>esquecerei</i>	5,7



# Índices Invertidos

1. *no meio do caminho tinha uma pedra*
2. *tinha uma pedra no meio do caminho*
3. *tinha uma pedra*
4. *no meio do caminho tinha uma pedra*
5. *nunca me esquecerei desse acontecimento*
6. *na vida de minhas retinas tão fatigadas*
7. *nunca me esquecerei que no meio do caminho*
8. *tinha uma pedra*
9. *tinha uma pedra no meio do caminho*
10. *no meio do caminho tinha uma pedra.*

Word	Doc_Ids
<i>acontecimento</i>	5
<i>caminho</i>	1,2,4,7,9,10
<i>esquecerei</i>	5,7
<i>fatigados</i>	6

## Índices Invertidos

1. *no meio do caminho tinha uma pedra*
2. *tinha uma pedra no meio do caminho*
3. *tinha uma pedra*
4. *no meio do caminho tinha uma pedra*
5. *nunca me esquecerei desse acontecimento*
6. *na vida de minhas retinas tão fatigadas*
7. *nunca me esquecerei que no meio do caminho*
8. *tinha uma pedra*
9. *tinha uma pedra no meio do caminho*
10. *no meio do caminho tinha uma pedra.*

Word	Doc_Iids
<i>acontecimento</i>	5
<i>caminho</i>	1,2,4,7,9,10
<i>esquecerei</i>	5,7
<i>fatigados</i>	6
<i>meio</i>	1,2,4,7,9,10

# Índices Invertidos

1. *no meio do caminho tinha uma pedra*
2. *tinha uma pedra no meio do caminho*
3. *tinha uma pedra*
4. *no meio do caminho tinha uma pedra*
5. *nunca me esquecerei desse acontecimento*
6. *na vida de minhas retinas tão fatigadas*
7. *nunca me esquecerei que no meio do caminho*
8. *tinha uma pedra*
9. *tinha uma pedra no meio do caminho*
10. *no meio do caminho tinha uma pedra.*

Word	Doc_Ids
<i>acontecimento</i>	5
<i>caminho</i>	1,2,4,7,9,10
<i>esquecerei</i>	5,7
<i>fatigados</i>	6
<i>meio</i>	1,2,4,7,9,10
<i>nunca</i>	5,7

# Índices Invertidos

1. *no meio do caminho tinha uma pedra*
2. *tinha uma pedra no meio do caminho*
3. *tinha uma pedra*
4. *no meio do caminho tinha uma pedra*
5. *nunca me esquecerei desse acontecimento*
6. *na vida de minhas retinas tão fatigadas*
7. *nunca me esquecerei que no meio do caminho*
8. *tinha uma pedra*
9. *tinha uma pedra no meio do caminho*
10. *no meio do caminho tinha uma pedra.*

Word	Doc_Ids
<i>acontecimento</i>	5
<i>caminho</i>	1,2,4,7,9,10
<i>esquecerei</i>	5,7
<i>fatigados</i>	6
<i>meio</i>	1,2,4,7,9,10
<i>nunca</i>	5,7
<i>pedra</i>	1,2,3,4,8,9,10



# Índices Invertidos

1. *no meio do caminho tinha uma pedra*
2. *tinha uma pedra no meio do caminho*
3. *tinha uma pedra*
4. *no meio do caminho tinha uma pedra*
5. *nunca me esquecerei desse acontecimento*
6. *na vida de minhas retinas tão fatigadas*
7. *nunca me esquecerei que no meio do caminho*
8. *tinha uma pedra*
9. *tinha uma pedra no meio do caminho*
10. *no meio do caminho tinha uma pedra.*

Word	Doc_Ids
<i>acontecimento</i>	5
<i>caminho</i>	1,2,4,7,9,10
<i>esquecerei</i>	5,7
<i>fatigados</i>	6
<i>meio</i>	1,2,4,7,9,10
<i>nunca</i>	5,7
<i>pedra</i>	1,2,3,4,8,9,10
<i>retinas</i>	6

# Índices Invertidos

1. *no meio do caminho tinha uma pedra*
2. *tinha uma pedra no meio do caminho*
3. *tinha uma pedra*
4. *no meio do caminho tinha uma pedra*
5. *nunca me esquecerei desse acontecimento*
6. *na vida de minhas retinas tão fatigadas*
7. *nunca me esquecerei que no meio do caminho*
8. *tinha uma pedra*
9. *tinha uma pedra no meio do caminho*
10. *no meio do caminho tinha uma pedra.*

Word	Doc_Ids
<i>acontecimento</i>	5
<i>caminho</i>	1,2,4,7,9,10
<i>esquecerei</i>	5,7
<i>fatigados</i>	6
<i>meio</i>	1,2,4,7,9,10
<i>nunca</i>	5,7
<i>pedra</i>	1,2,3,4,8,9,10
<i>retinas</i>	6
<i>tinha</i>	1,2,3,4,8,9,10

---

# Índices Invertidos

1. *no meio do caminho tinha uma pedra*
2. *tinha uma pedra no meio do caminho*
3. *tinha uma pedra*
4. *no meio do caminho tinha uma pedra*
5. *nunca me esquecerei desse acontecimento*
6. *na vida de minhas retinas tão fatigadas*
7. *nunca me esquecerei que no meio do caminho*
8. *tinha uma pedra*
9. *tinha uma pedra no meio do caminho*
10. *no meio do caminho tinha uma pedra.*

Word	Doc_Ids
<i>acontecimento</i>	5
<i>caminho</i>	1,2,4,7,9,10
<i>esquecerei</i>	5,7
<i>fatigados</i>	6
<i>meio</i>	1,2,4,7,9,10
<i>nunca</i>	5,7
<i>pedra</i>	1,2,3,4,8,9,10
<i>retinas</i>	6
<i>tinha</i>	1,2,3,4,8,9,10
<i>vida</i>	6

# Algoritmo PageRank

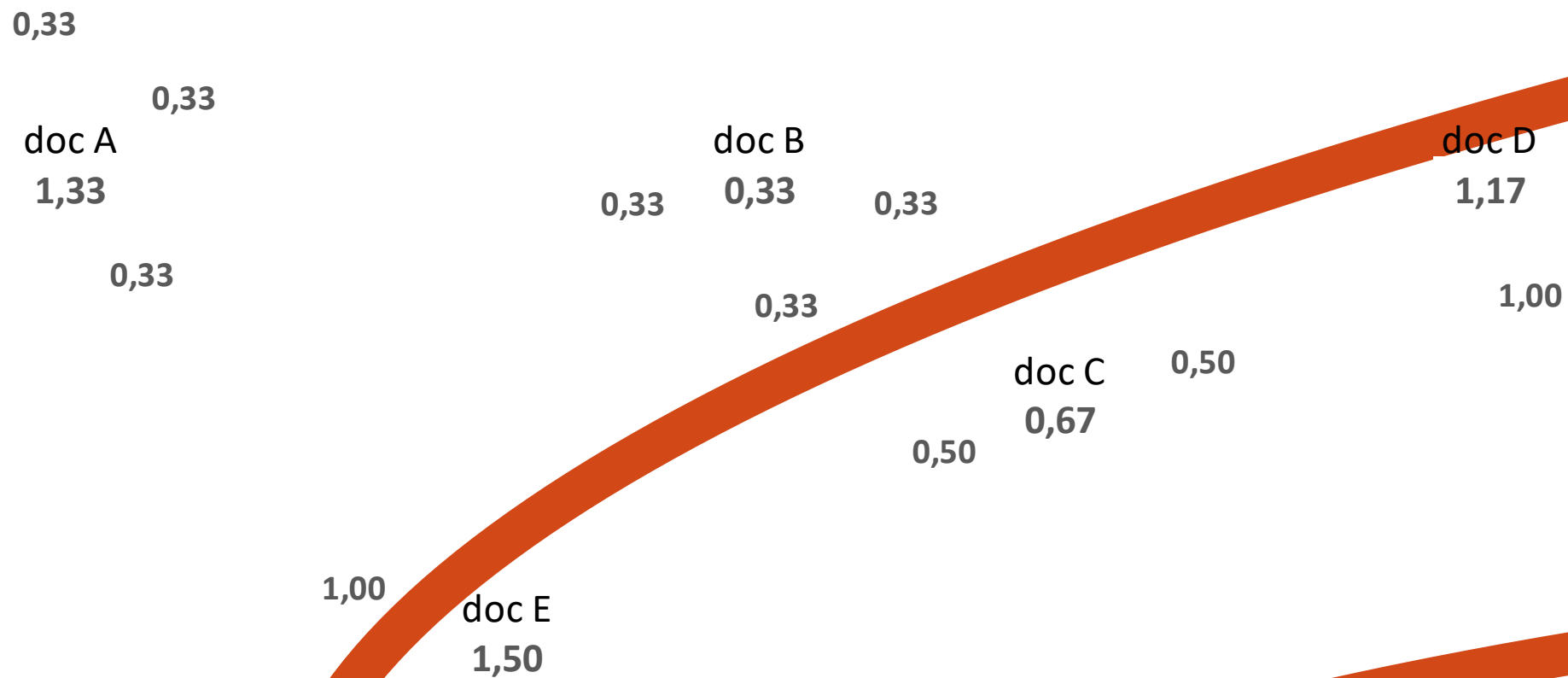
- Classificação das páginas da web baseada em reputação.
- Cada documento recebe e doa pontos baseados em hiperlinks e em sua reputação.
- Criado pelo cofundador do Google, Larry Page.



Fonte: Wikimedia.



# Algoritmo PageRank



Fonte: Elaborada pelo autor.

# Algoritmo PageRank

doc A  
1,33

doc B  
0,33

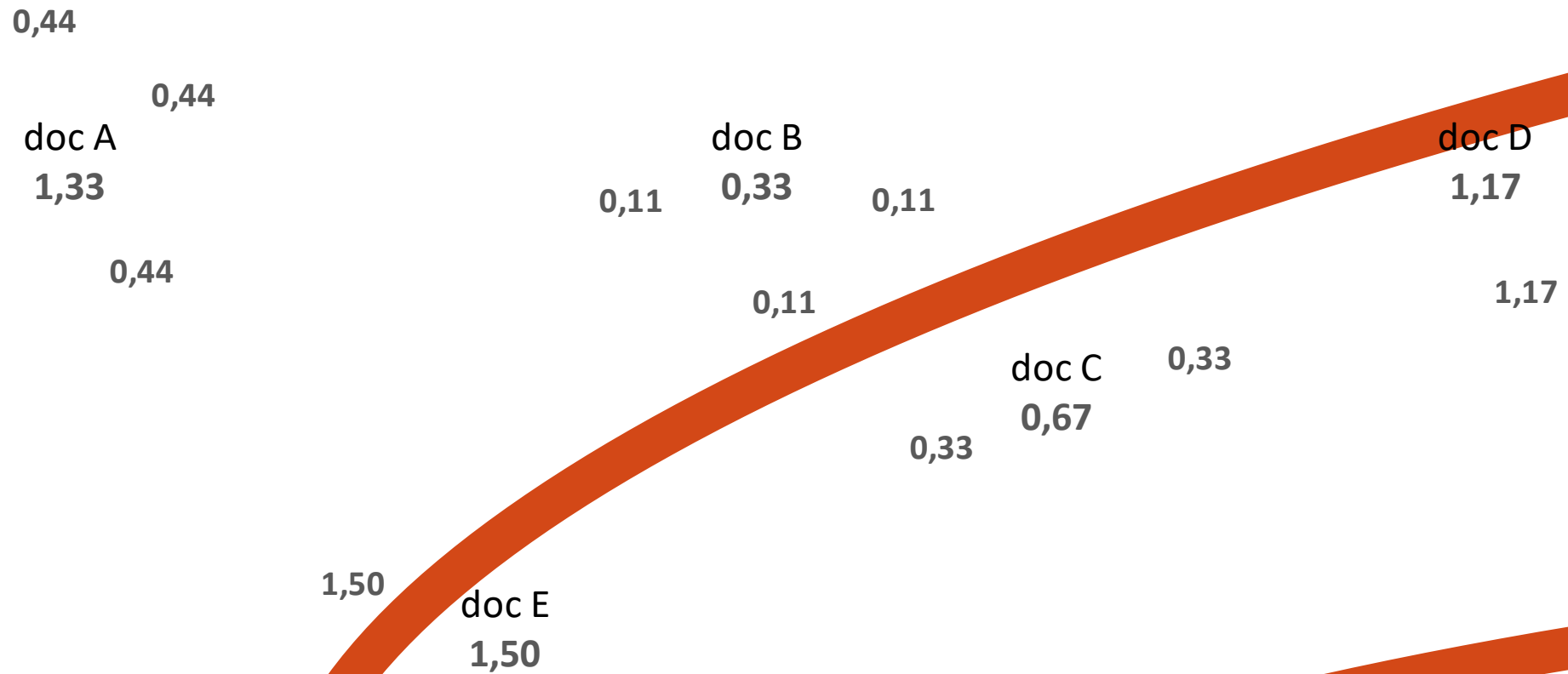
doc D  
1,17

doc C  
0,67

doc E  
1,50

Fonte: Elaborada pelo autor.

# Algoritmo PageRank



Fonte: Elaborada pelo autor.

# Algoritmo PageRank

doc A  
1,61

doc B  
0,44

doc D  
0,89

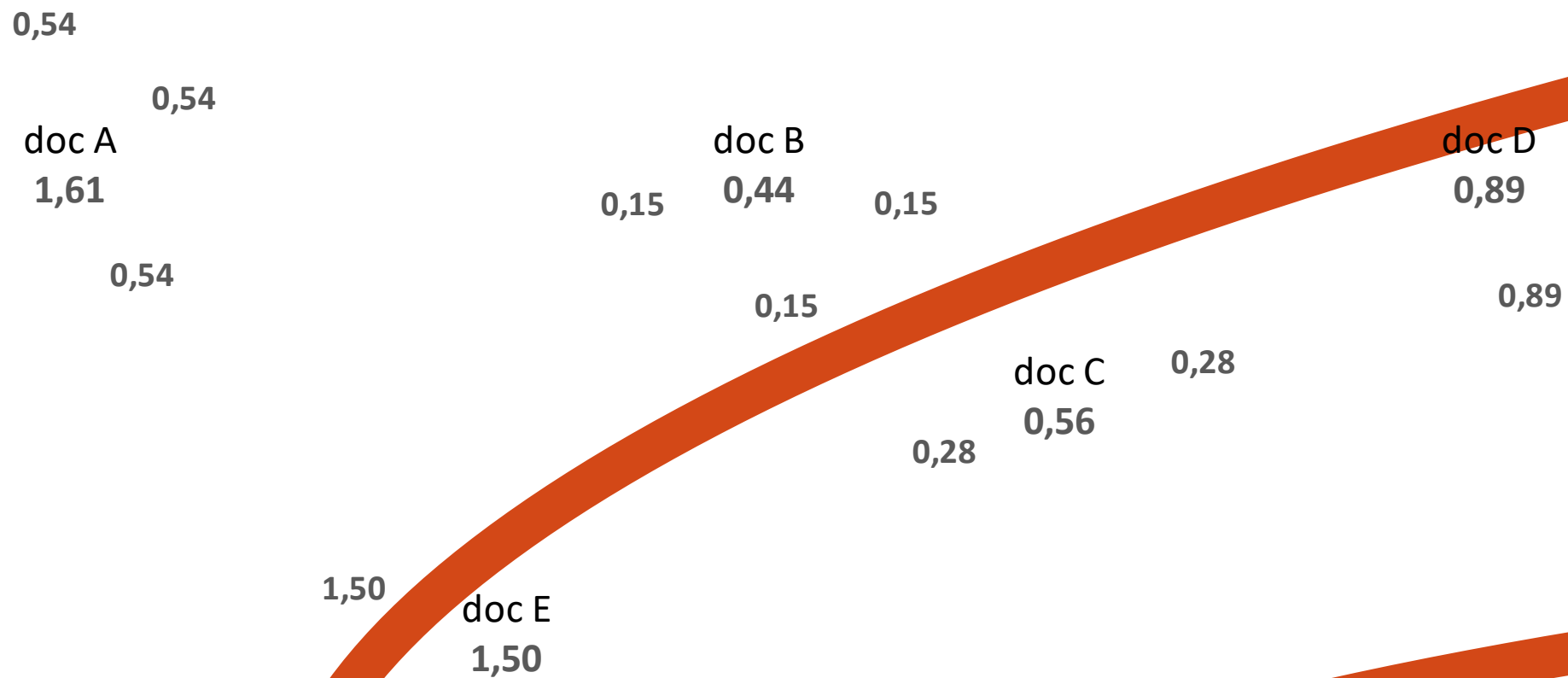
doc C  
0,56

doc E  
1,50

Fonte: Elaborada pelo autor.



# Algoritmo PageRank



Fonte: Elaborada pelo autor.

# Algoritmo PageRank

doc A  
1,65

doc B  
0,54

doc D  
0,96

doc C  
0,69

doc E  
1,17

Fonte: Elaborada pelo autor.

# Algoritmo PageRank

doc A  
1,35

doc B  
0,55

doc D  
1,07

doc C  
0,73

doc E  
1,31

Fonte: Elaborada pelo autor.

# Algoritmo PageRank

doc A  
1,49

doc B  
0,45

doc D  
1,00

doc C  
0,63

doc E  
1,44

Fonte: Elaborada pelo autor.

# Algoritmo PageRank – Convergência

doc A  
1,50

doc B  
0,50

doc D  
1,00

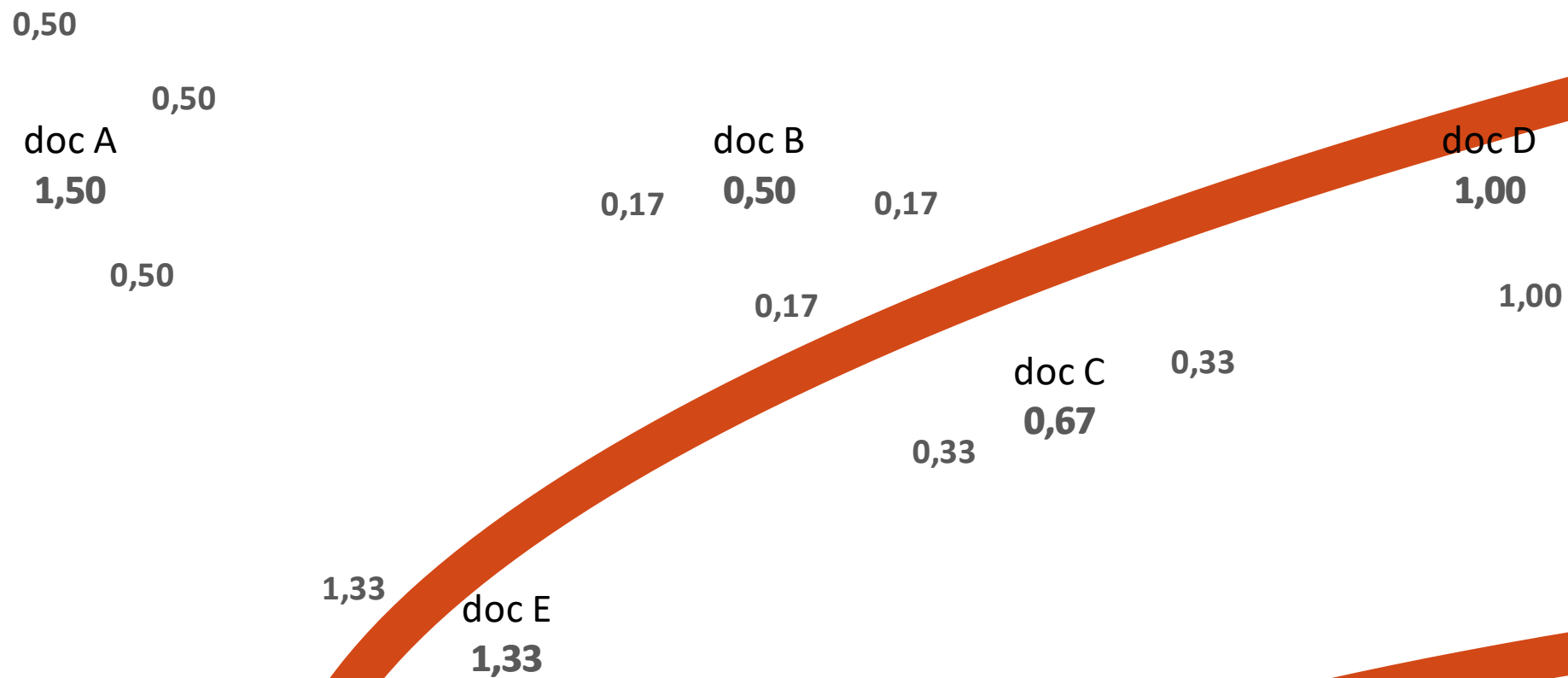
doc C  
0,67

doc E  
1,33

Fonte: Elaborada pelo autor.



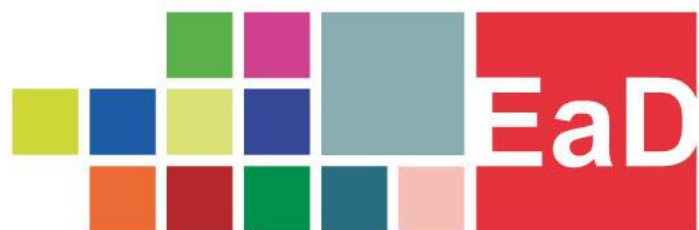
# Algoritmo PageRank – Convergência



Fonte: Elaborada pelo autor.

# Operadores avançados em sistemas de Recuperação da informação

Operador	Para que serve	Exemplo
<b><i>-palavra</i></b>	Exclui os resultados que possuem esta palavra.	“ontologia”-filosofia
<b><i>“palavra”</i></b>	Busca pela frase exata e não considera sinônimos ou conjunções de gênero, número e grau.	“worked at night”
<b><i>“dentro de uma frase”</i></b>	Busca por uma frase exata, e não pela combinação de palavras.	“melhor carro da formula 1”
<b><i>..</i></b>	Procura considerando valores numéricos dentro do intervalo de números.	presidente 2002..2020
<b><i>intitle:</i></b>	Pesquisa apenas no título da página do documento.	Election <b><i>intitle:</i></b> biden
<b><i>inurl:</i></b>	Pesquisa apenas na URL de identificação do documento.	eleição <b><i>inurl:</i></b> 2018
<b><i>filetype:</i></b>	Especifica o formato de documentos procurado (por exemplo: PDF, PPT, DOC, XLS).	<b><i>filetype:</i></b> pdf
<b><i>site:</i></b>	Faz a pesquisa exclusivamente dentro do site especificado.	<b><i>site:</i></b> senado.gov.br



Universidade Presbiteriana  
**Mackenzie**