

# Inverse Document Frequency (idf) weighting

# Document Frequency (df)



# Document Frequency (df)

$df_t$



# Document Frequency (df)

$df_t$  = Number of documents containing the term  $t$

# Inverse Document Frequency (idf)

$$\text{idf}_t = \log_{10} (N/\text{df}_t)$$



# Inverse Document Frequency (idf)

$$\text{idf}_t = \log_{10} (N/\text{df}_t)$$

N = Total number of documents

# Inverse Document Frequency (idf)

$$\text{idf}_t = \log_{10} (N/\text{df}_t)$$

N = Total number of documents

N = 1000

Term	$\text{df}_t$	$\text{idf}_t$
Analytics		
book		
is		
of		
the		

# Inverse Document Frequency (idf)

$$\text{idf}_t = \log_{10} (N/\text{df}_t)$$

N = Total number of documents

N = 1000

Term	$\text{df}_t$	$\text{idf}_t$
Analytics	100	
book		
is		
of		
the		



# Inverse Document Frequency (idf)

$$\text{idf}_t = \log_{10} (N/\text{df}_t)$$

N = Total number of documents

N = 1000

Term	$\text{df}_t$	$\text{idf}_t$
Analytics	100	
book	10	
is		
of		
the		

# Inverse Document Frequency (idf)

$$\text{idf}_t = \log_{10} (N/\text{df}_t)$$

N = Total number of documents

N = 1000

Term	$\text{df}_t$	$\text{idf}_t$
Analytics	100	
book	10	
is	800	
of	500	
the	1000	

# Inverse Document Frequency (idf)

$$\text{idf}_t = \log_{10} (N/\text{df}_t)$$

N = Total number of documents

N = 1000

Term	$\text{df}_t$	$\text{idf}_t$
Analytics	100	$\log_{10} (N/\text{df}_t)$
book	10	
is	800	
of	500	
the	1000	

# Inverse Document Frequency (idf)

$$\text{idf}_t = \log_{10} (N/\text{df}_t)$$

N = Total number of documents

N = 1000

Term	$\text{df}_t$	$\text{idf}_t$
Analytics	100	$\log_{10} (10)$
book	10	
is	800	
of	500	
the	1000	

# Inverse Document Frequency (idf)

$$\text{idf}_t = \log_{10} (N/\text{df}_t)$$

N = Total number of documents

N = 1000

Term	$\text{df}_t$	$\text{idf}_t$
Analytics	100	1
book	10	
is	800	
of	500	
the	1000	

# Inverse Document Frequency (idf)

$$\text{idf}_t = \log_{10} (N/\text{df}_t)$$

N = Total number of documents

N = 1000

Term	$\text{df}_t$	$\text{idf}_t$
Analytics	100	1
book	10	$\log_{10} (100)$
is	800	
of	500	
the	1000	

# Inverse Document Frequency (idf)

$$\text{idf}_t = \log_{10} (N/\text{df}_t)$$

N = Total number of documents

N = 1000

Term	$\text{df}_t$	$\text{idf}_t$
Analytics	100	1
book	10	2
is	800	
of	500	
the	1000	

# Inverse Document Frequency (idf)

$$\text{idf}_t = \log_{10} (N/\text{df}_t)$$

N = Total number of documents

N = 1000

Term	$\text{df}_t$	$\text{idf}_t$
Analytics	100	1
book	10	2
is	800	0.09
of	500	0.3
the	1000	0



# Inverse Document Frequency for Ranked Retrieval

Query (q): book of Analytics

Doc 1 (d1): This book is of Analytics

Doc 2 (d2): Big-Data Analytics is the process of examining large volume of data



# Inverse Document Frequency for Ranked Retrieval

Query (q): book of Analytics

Doc 1 (d1): This book is of Analytics

Doc 2 (d2): Big-Data Analytics is the process of examining large volume of data

Term	df(t)
Analytics	2
Big-Data	1
book	1
data	1
examining	1
is	2
large	1
of	2
process	1
the	1
this	1
volume	1

# Inverse Document Frequency for Ranked Retrieval

Query (q): book of Analytics

Doc 1 (d1): This book is of Analytics

Doc 2 (d2): Big-Data Analytics is the process of examining large volume of data

Term	df(t)	idf(t)
Analytics	2	0
Big-Data	1	0.3
book	1	0.3
data	1	0.3
examining	1	0.3
is	2	0
large	1	0.3
of	2	0
process	1	0.3
the	1	0.3
this	1	0.3
volume	1	0.3

# Inverse Document Frequency for Ranked Retrieval

Query (q): book of Analytics

Doc 1 (d1): This book is of Analytics

Doc 2 (d2): Big-Data Analytics is the process of examining large volume of data

Term	df(t)	idf(t)
Analytics	2	0
Big-Data	1	0.3
book	1	0.3
data	1	0.3
examining	1	0.3
is	2	0
large	1	0.3
of	2	0
process	1	0.3
the	1	0.3
this	1	0.3
volume	1	0.3

Score (d1) = Sum over terms in both query and doc 1

# Inverse Document Frequency for Ranked Retrieval

Query (q): book of Analytics

Doc 1 (d1): This book is of Analytics

Doc 2 (d2): Big-Data Analytics is the process of examining large volume of data

Term	df(t)	idf(t)
Analytics	2	0
Big-Data	1	0.3
book	1	0.3
data	1	0.3
examining	1	0.3
is	2	0
large	1	0.3
of	2	0
process	1	0.3
the	1	0.3
this	1	0.3
volume	1	0.3

$$\text{Score (d1)} = 0 + 0.3 + 0$$

# Inverse Document Frequency for Ranked Retrieval

Query (q): book of Analytics

Doc 1 (d1): This book is of Analytics

Doc 2 (d2): Big-Data Analytics is the process of examining large volume of data

Term	df(t)	idf(t)
Analytics	2	0
Big-Data	1	0.3
book	1	0.3
data	1	0.3
examining	1	0.3
is	2	0
large	1	0.3
of	2	0
process	1	0.3
the	1	0.3
this	1	0.3
volume	1	0.3

Score (d1) = 0.3

Score (d2) = Sum over terms in both query and doc 2

# Inverse Document Frequency for Ranked Retrieval

Query (q): book of Analytics

Doc 1 (d1): This book is of Analytics

Doc 2 (d2): Big-Data Analytics is the process of examining large volume of data

Term	df(t)	idf(t)
Analytics	2	0
Big-Data	1	0.3
book	1	0.3
data	1	0.3
examining	1	0.3
is	2	0
large	1	0.3
of	2	0
process	1	0.3
the	1	0.3
this	1	0.3
volume	1	0.3

Score (d1) = 0.3

Score (d2) = 0 + 0

# Inverse Document Frequency for Ranked Retrieval

Query (q): book of Analytics

Doc 1 (d1): This book is of Analytics

Doc 2 (d2): Big-Data Analytics is the process of examining large volume of data

Term	df(t)	idf(t)
Analytics	2	0
Big-Data	1	0.3
book	1	0.3
data	1	0.3
examining	1	0.3
is	2	0
large	1	0.3
of	2	0
process	1	0.3
the	1	0.3
this	1	0.3
volume	1	0.3

Score (d1) = 0.3

Score (d2) = 0



# Term Frequency - Inverse Document Frequency (tf-idf)



# Term Frequency - Inverse Document Frequency (tf-idf)

- Product of Term Frequency (tf) and Inverse Document Frequency (idf)



# Term Frequency - Inverse Document Frequency (tf-idf)

- Product of Term Frequency (tf) and Inverse Document Frequency (idf)


$$W_{t,d} = [1 + \log_{10}(tf_{t,d})] \times [\log_{10}(N/df_t)]$$

# Term Frequency - Inverse Document Frequency (tf-idf)

Query (q): book of Analytics

Doc 1 (d1): This book is of Analytics

Doc 2 (d2): Big-Data Analytics is the process of examining large volume of data

Analytics  
Vidhya

# Term Frequency - Inverse Document Frequency (tf-idf)

Query (q): book of Analytics

Doc 1 (d1): This book is of Analytics

Doc 2 (d2): Big-Data Analytics is the process of examining large volume of data

Term	d1(tf-idf)	d2(tf-idf)
Analytics	0	0
Big-Data	0	0.3
book	0.3	0
data	0	0.3
examining	0	0.3
is	0	0
large	0	0.3
of	0	0
process	0	0.3
the	0	0.3
this	0.3	0
volume	0	0.3

# Term Frequency - Inverse Document Frequency (tf-idf)

Query (q): book of Analytics

Doc 1 (d1): This book is of Analytics

Doc 2 (d2): Big-Data Analytics is the process of examining large volume of data

Term	d1(tf-idf)	d2(tf-idf)
Analytics	0	0
Big-Data	0	0.3
book	0.3	0
data	0	0.3
examining	0	0.3
is	0	0
large	0	0.3
of	0	0
process	0	0.3
the	0	0.3
this	0.3	0
volume	0	0.3

Score (d1) = 0.3

# Term Frequency - Inverse Document Frequency (tf-idf)

Query (q): book of Analytics

Doc 1 (d1): This book is of Analytics

Doc 2 (d2): Big-Data Analytics is the process of examining large volume of data

Term	d1(tf-idf)	d2(tf-idf)
Analytics	0	0
Big-Data	0	0.3
book	0.3	0
data	0	0.3
examining	0	0.3
is	0	0
large	0	0.3
of	0	0
process	0	0.3
the	0	0.3
this	0.3	0
volume	0	0.3

Score (d1) = 0.3

Score (d2) = 0

# Term Frequency - Inverse Document Frequency (tf-idf)

Query (q): book of Analytics

Doc 1 (d1): This book is of Analytics

Doc 2 (d2): Big-Data Analytics is the process of examining large volume of data

Term	d1(tf-idf)	d2(tf-idf)
Analytics	0	0
Big-Data	0	0.3
book	0.3	0
data	0	0.3
examining	0	0.3
is	0	0
large	0	0.3
of	0	0
process	0	0.3
the	0	0.3
this	0.3	0
volume	0	0.3

Score (d1) = 0.3

Score (d2) = 0



# Term Frequency - Inverse Document Frequency (tf-idf)

- Product of Term Frequency (tf) and Inverse Document Frequency (idf)



# Term Frequency - Inverse Document Frequency (tf-idf)

- Product of Term Frequency (tf) and Inverse Document Frequency (idf)
- Increases with the occurrences of a term within the document due to tf part


$$W_{t,d} = [1 + \log_{10}(tf_{t,d})] \times [\log_{10}(N/df_t)]$$

# Term Frequency - Inverse Document Frequency (tf-idf)

- Product of Term Frequency (tf) and Inverse Document Frequency (idf)
- Increases with the occurrences of a term within the document due to tf part
- Increases with the rarity of a term in the collection due to idf part

$$W_{t,d} = [1 + \log_{10}(tf_{t,d})] \times [\log_{10}(N/df_t)]$$



Thank You