Introduction to Information Retrieval INF 141
Donald J. Patterson

Content adapted from Hinrich Schütze http://www.informationretrieval.org

Overview

- Boolean Retrieval
- Weighted Boolean Retrieval
- Zone Indices
- Term Frequency Metrics
- The full vector space model





From the bottom

• "Grep"



- "Grep"
 - Querying without an index or a crawl



- "Grep"
 - Querying without an index or a crawl
 - Whenever you want to find something you look through the entire document for it.



- "Grep"
 - Querying without an index or a crawl
 - Whenever you want to find something you look through the entire document for it.
 - Example:



- "Grep"
 - Querying without an index or a crawl
 - Whenever you want to find something you look through the entire document for it.
 - Example:
 - You have the collected works of Shakespeare on disk



- "Grep"
 - Querying without an index or a crawl
 - Whenever you want to find something you look through the entire document for it.
 - Example:
 - You have the collected works of Shakespeare on disk
 - You want to know which play contains the words



- "Grep"
 - Querying without an index or a crawl
 - Whenever you want to find something you look through the entire document for it.
 - Example:
 - You have the collected works of Shakespeare on disk
 - You want to know which play contains the words
 - "Brutus AND Caesar"

• "Grep"



- "Grep"
 - "Brutus AND Caesar" is the query.



- "Grep"
 - "Brutus AND Caesar" is the query.
 - This is a boolean query. Why?



- "Grep"
 - "Brutus AND Caesar" is the query.
 - This is a boolean query. Why?
 - What other operators could be used?



- "Grep"
 - "Brutus AND Caesar" is the query.
 - This is a boolean query. Why?
 - What other operators could be used?
 - The grep solution:



- "Grep"
 - "Brutus AND Caesar" is the query.
 - This is a boolean query. Why?
 - What other operators could be used?
 - The grep solution:
 - Read all the files and all the text and output the intersection of the files



• "Grep"



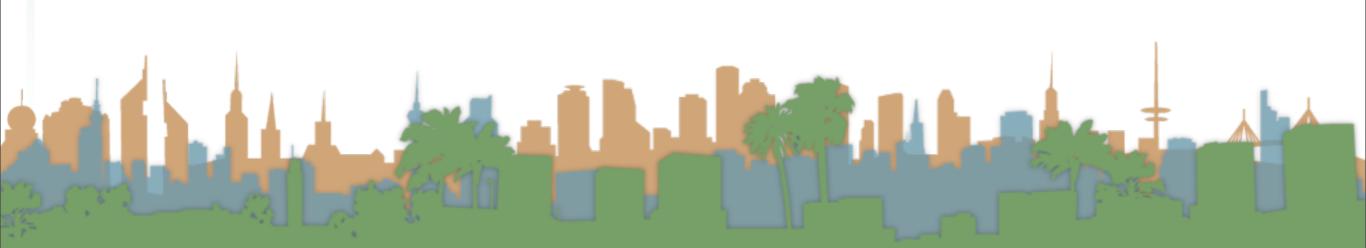
- "Grep"
 - Slow for large corpora



- "Grep"
 - Slow for large corpora
 - Calculating "NOT" requires exhaustive scanning



- "Grep"
 - Slow for large corpora
 - Calculating "NOT" requires exhaustive scanning
 - Some operations not feasible



- "Grep"
 - Slow for large corpora
 - Calculating "NOT" requires exhaustive scanning
 - Some operations not feasible
 - Query: "Romans NEAR Countrymen"



- "Grep"
 - Slow for large corpora
 - Calculating "NOT" requires exhaustive scanning
 - Some operations not feasible
 - Query: "Romans NEAR Countrymen"
 - Doesn't support ranked retrieval



- "Grep"
 - Slow for large corpora
 - Calculating "NOT" requires exhaustive scanning
 - Some operations not feasible
 - Query: "Romans NEAR Countrymen"
 - Doesn't support ranked retrieval
- Moving beyond grep is the motivation for the inverted index.