

Information Retrieval

CS5154/6054

Yizong Cheng

8/23/2022

CS5154/6054 Information Retrieval
Course Syllabus
Fall Semester 2022

Instructor Information:

Instructor: Yizong Cheng
Office Location: 812A Rhodes
Email Address: chengy@ucmail.uc.edu
Phone Number: 513-556-1809
Office Hours: 3:00-4:00 MW

General Course Information:

Course Number:	CS5154/6054
Course Title:	Information Retrieval
Credit Hours:	3
Contact Hours:	3
Meeting Days/Times:	9:30-10:50 TH, Swift 500
Prerequisites:	Python programming
Course Description:	Information storage and retrieval with unstructured data. Inverted index, tf-idf and cosine similarity, and relevance-based evaluation. Probabilistic information retrieval with text classification, feature selection, and clustering.
Course Delivery Mode:	This course is an in-person course. There may be frequent in-class quizzes.
Course Location and/or Access:	This course can be found on Canvas (accessible via canopy.uc.edu). Lecture notes, reference links, and homework assignments will be posted on the course site in Canvas.

Course Resources:

Required Textbooks and Materials:	Introduction to Information Retrieval by Manning, Raghavan, and Schutze, Cambridge 2008 (IIR)
Supplemental Textbooks and Materials:	Blueprints for Text Analytics Using Python by Albrecht, Ramachandran, and Winkler, O'Reilly, 2020 (Blueprints)
Required Software or Hardware:	Python

Course Assignments:

Quizzes:	in-class
Assignments:	mostly Python programming
Exams:	in-class, midterm and final

Grading:

Final numerical grades for the course are determined based on the following weights for assignments.

Assignment:	Weighting:
Quizzes and assignments:	50%
Exams:	50%
Total:	100%

Final letter grades for the course will be determined using your final numerical grade. The range will be different for undergraduate (CS5154) and graduate (CS6054) sections.

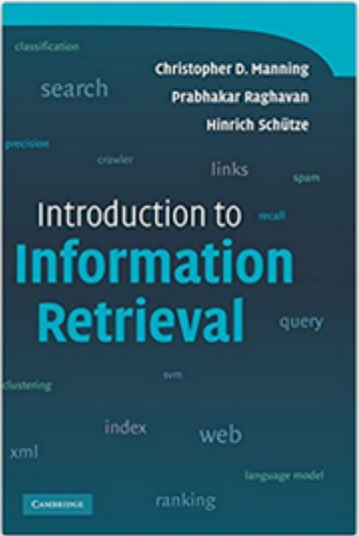
FileEditViewHistoryBookmarksToolsHelp

Introduction to Information Retrieval

https://www.amazon.com/Introduction-Information-Retrieval-Christopher-Manning/

Back to results

Look inside



See all 3 images

Introduction to Information Retrieval Illustrated Edition

by Christopher D. Manning (Author), & 2 more

★★★★★ 166 ratings

See all formats and editions

eTextbook \$22.22 - \$38.00	Hardcover \$19.99 - \$67.55
--------------------------------	--------------------------------

Read with Our Free App

Class-tested and coherent, this groundbreaking new textbook teaches web-era information retrieval, including web search and the related areas of text classification and text clustering from basic concepts. Written from a computer science perspective by three leading experts in the field, it gives an up-to-date treatment of all aspects of the design and

Read more

Rent
\$19.99

List Price: \$70.99
Save: \$51.00 (72%)

Due Date: Dec 14, 2022 Rental Details

- FREE return shipping at the end of the semester.
- Access codes and supplements are not guaranteed with rentals.

In Stock.

Rented from Apex_media

Fulfilled by Amazon

FREE delivery Wednesday, August 24

Or fastest delivery Saturday, August 20. Order within 9 hrs

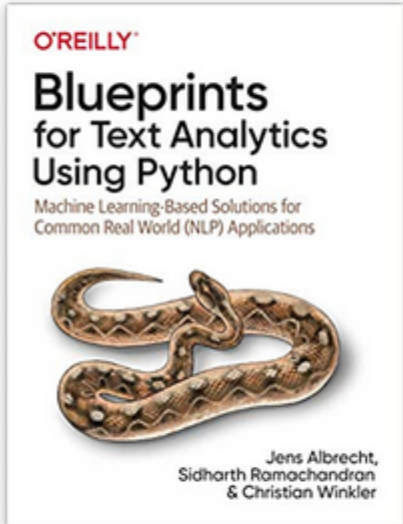
FileEditViewHistoryBookmarksToolsHelp

Blueprints for Text Analytics Usi X

←→↻🏠🔒https://www.amazon.com/Blueprints-Text-Analytics-Using-Python/dp/149207408X/★📧📦🔔☰


Back to results

Look inside



O'REILLY
**Blueprints
for Text Analytics
Using Python**
Machine Learning-Based Solutions for
Common Real World (NLP) Applications

Jens Albrecht,
Sidharth Ramachandran
& Christian Winkler



See all 2 images

Blueprints for Text Analytics Using Python: Machine Learning-Based Solutions for Common Real World (NLP) Applications 1st Edition

by [Jens Albrecht](#) (Author), & 2 more

★★★★★ 24 ratings

See all formats and editions

Kindle \$44.99	Paperback \$51.99
Read with Our Free App	6 Used from \$32.00 22 New from \$39.38

Turning text into valuable information is essential for businesses looking to gain a competitive advantage. With recent improvements in natural language processing (NLP)

Buy new: **\$51.99**

List Price: ~~\$69.99~~ ⓘ
Save: \$18.00 (26%)

& **FREE Returns**

FREE delivery Wednesday, August 24

Or fastest delivery **Saturday, August 20**. Order within **7 hrs 52 mins**

📍 [Select delivery location](#)

Only 13 left in stock (more on the way).

As an alternative, the [Kindle eBook](#) is available now and can be read on any device with the free Kindle app.

Qty: 1

Course Schedule:

Week	Class Period	Topic	Reading
1	Aug 23	term-document matrix	1 IIR
	Aug 25	tokenization with regular expressions	2 IIR and 4 Blueprints
2	Aug 30	set similarity and ranking	3 IIR
	Sep 1	spelling correction	3 IIR
3	Sep 6	tf-idf and cosine similarity	6 IIR
	Sep 8	sklearn vectorizer and ngram	5 Blueprints
4	Sep 13	precision and recall	8 IIR
	Sep 15	evaluating IR	8 IIR

What to Learn Today

- Information retrieval
 - Documents and terms
 - retrieve all documents containing a term
 - Boolean query : AND, OR, NOT between terms
- The ad hoc retrieval: fixed documents and changing query
 - Filtering or classification: fixed query and changing documents
- Documents as sets of terms
 - document-term matrix
 - Inverted index
 - Python implementation with list, set, and dictionary (dict)

Information Retrieval

- Information retrieval (IR) is finding material (usually documents) of an unstructured nature (usually text) that satisfies an information need from within large collections (usually stored on computers).
- The term **information retrieval** was coined by Calvin Mooers (1919-1994) in 1948/1950.



ICM 1950

UNIVERSITY OF CINCINNATI,
CINCINNATI, OHIO, U. S. A.

INFORMATION RETRIEVAL VIEWED AS TEMPORAL SIGNALLING

CALVIN N. MOOERS

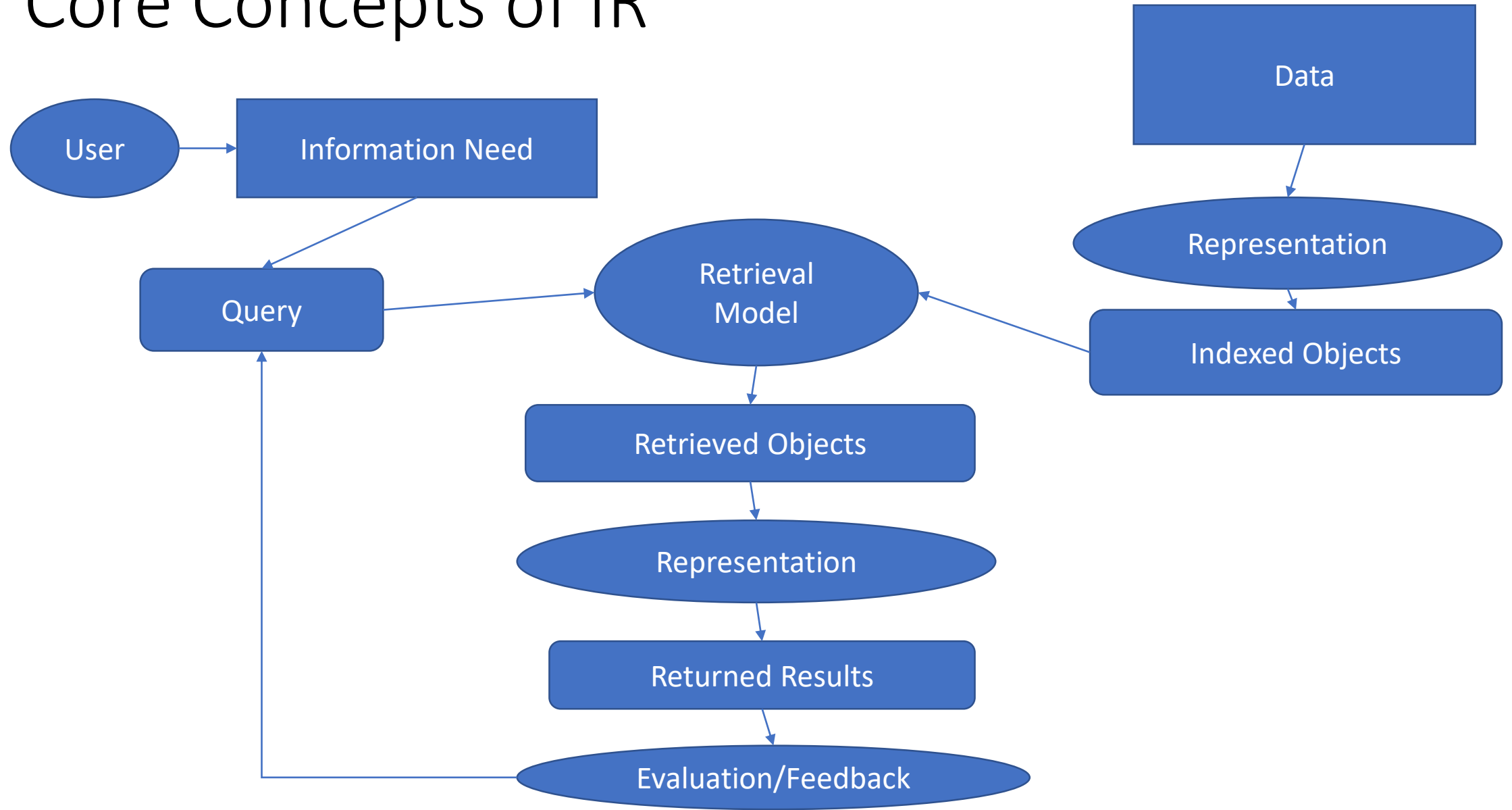
The problem of directing a user to stored information, some of which may be unknown to him, is the problem of "information retrieval". Signalling theories can be applied, though this is a form of temporal signalling, which distinguishes it from the point-to-point signalling currently under study by others.

In information retrieval, the addressee or receiver rather than the sender is the active party. Other differences are that communication is temporal from one epoch to a later epoch in time, though possibly at the same point in space; communication is in all cases unidirectional; the sender cannot know the particular message that will be of later use to the receiver and must send all possible

The *ad hoc* retrieval problem

- Given a user information need and a collection of documents, the IR system determines how well the documents satisfy the query and returns a subset of relevant documents to the user.
 - Collection is relatively stable.
 - Queries are created and used dynamically; change fast
 - “Ad-hoc”: formed or used for specific or immediate problems or needs – Merriam-Webster’s collegiate Dictionary
- Filtering: Queries are stable, while the collection changes.

Core Concepts of IR



Quoting ACM Survey 2019 Article 15

- In their broader sense, information retrieval systems (IRSs) aim to fulfill users' information needs expressed in a keyword-based query.
- More precisely, information retrieval (IR) refers to the process of selection, from a document repository, of the documents *likely* to be relevant to a particular information need, formulated by a query.
- Based on this definition, an IRS has to deal with a collection of documents, with users' information needs, and with the notion of relevance.

Documents and the Query

- In IR, documents are carriers of information;
 - in their original forms, they are human-understandable objects (e.g., Web pages, articles, books, and images),
 - which an IRS must transform into machine-understandable objects.
- This process is called indexing,
 - and its outcome is the association of a set of features (terms in textual documents) with documents.
 - These features constitute the basic elements employed to formally represent a document.
- A user's information needs are motivated by a user's information gap;
 - a query is a representation of these needs.
 - Once formal representations have been provided for both documents and queries, the system compares them to assess the relevance of each document to the considered query.

The IR Model

- Relevance is a complex notion composed of several dimensions, such as topicality, popularity, and novelty.
- An IRS can only estimate relevance, and generally topicality is the core relevance dimension.
- The assessment of topical relevance relies on the definition of a model, the IR model, which provides a formal means to represent and compare both documents and queries.
- Different mathematical theories have been employed to define IR models, which include **set theory**, linear algebra, probability theory, and formal logics.

Boolean Retrieval

- The Boolean model is arguably the simplest model to base an information retrieval system on.
- Queries are Boolean expressions, e.g., Caesar and Brutus
- The search engine returns all documents that satisfy the Boolean expression.
- Windows OS has files indexed, too.
- Document representation for Boolean retrieval: sets of words
 - Each document is a set of words. (words in general are substrings of the text)
 - The collection is a collection of sets.
 - A collection of sets is called a hypergraph in mathematics. Each word is a vertex and each document (a set of words) is a hyperedge.

Brutus AND Caesar AND NOT Calpurnia

Antony and Cleopatra, Act III, Scene ii

Agrippa [Aside to Domitius Enobarbus]: Why, Enobarbus,
 When Antony found Julius Caesar dead,
 He cried almost to roaring; and he wept
 When at Philippi he found Brutus slain.

Hamlet, Act III, Scene ii

Lord Polonius: I did enact Julius Caesar: I was killed i' the
 Capitol; Brutus killed me.

► **Figure 1.2** Results from Shakespeare for the query **Brutus AND Caesar AND NOT Calpurnia**.

FileEditViewHistoryBookmarksToolsHelp

The Python Tutorial — Python 3 ×

←→↻🏠🔒https://docs.python.org/3/tutorial/📄☆📧🐘☰

☰🐍3.10.6🔍Go

The Python Tutorial

Python is an easy to learn, powerful programming language. It has efficient high-level data structures and a simple but effective approach to object-oriented programming. Python's elegant syntax and dynamic typing, together with its interpreted nature, make it an ideal language for scripting and rapid application development in many areas on most platforms.

The Python interpreter and the extensive standard library are freely available in source or binary form for all major platforms from the Python web site, <https://www.python.org/>, and may be freely distributed. The same site also contains distributions of and pointers to many free third party Python modules, programs and tools, and additional documentation.

The Python interpreter is easily extended with new functions and data types implemented in C or C++ (or other languages callable from C). Python is also suitable as an extension language for customizable applications.

This tutorial introduces the reader informally to the basic concepts and features of the Python language and system. It helps to have a Python interpreter handy for hands-on experience, but all examples are self-contained, so the tutorial can be read off-line as well.

For a description of standard objects and modules, see [The Python Standard Library](#). [The Python Language Reference](#) gives a more formal definition of the language. To write extensions in C or C++, read [Extending and Embedding the Python Interpreter](#) and [Python/C API Reference Manual](#). There are also several books covering Python in depth.

This tutorial does not attempt to be comprehensive and cover every single feature, or even every commonly used feature. Instead, it introduces many of Python's most noteworthy features, and will give you a good idea of the language's flavor and style. After reading it, you will be able to read and write Python modules and programs, and you will be ready to learn more about the various Python library modules described in [The Python Standard Library](#).

File Edit View History Bookmarks Tools Help

5. Data Structures — Python 3.10

← → ↻ 🏠 🔒 https://docs.python.org/3/tutorial/datastructures.html 📄 ☆ 📧 🐘 ☰

☰ 🐍 3.10.6 🔍 Go

5. Data Structures

This chapter describes some things you've learned about already in more detail, and adds some new things as well.

5.1. More on Lists

The list data type has some more methods. Here are all of the methods of list objects:

list.append(x)
Add an item to the end of the list. Equivalent to `a[len(a):] = [x]`.

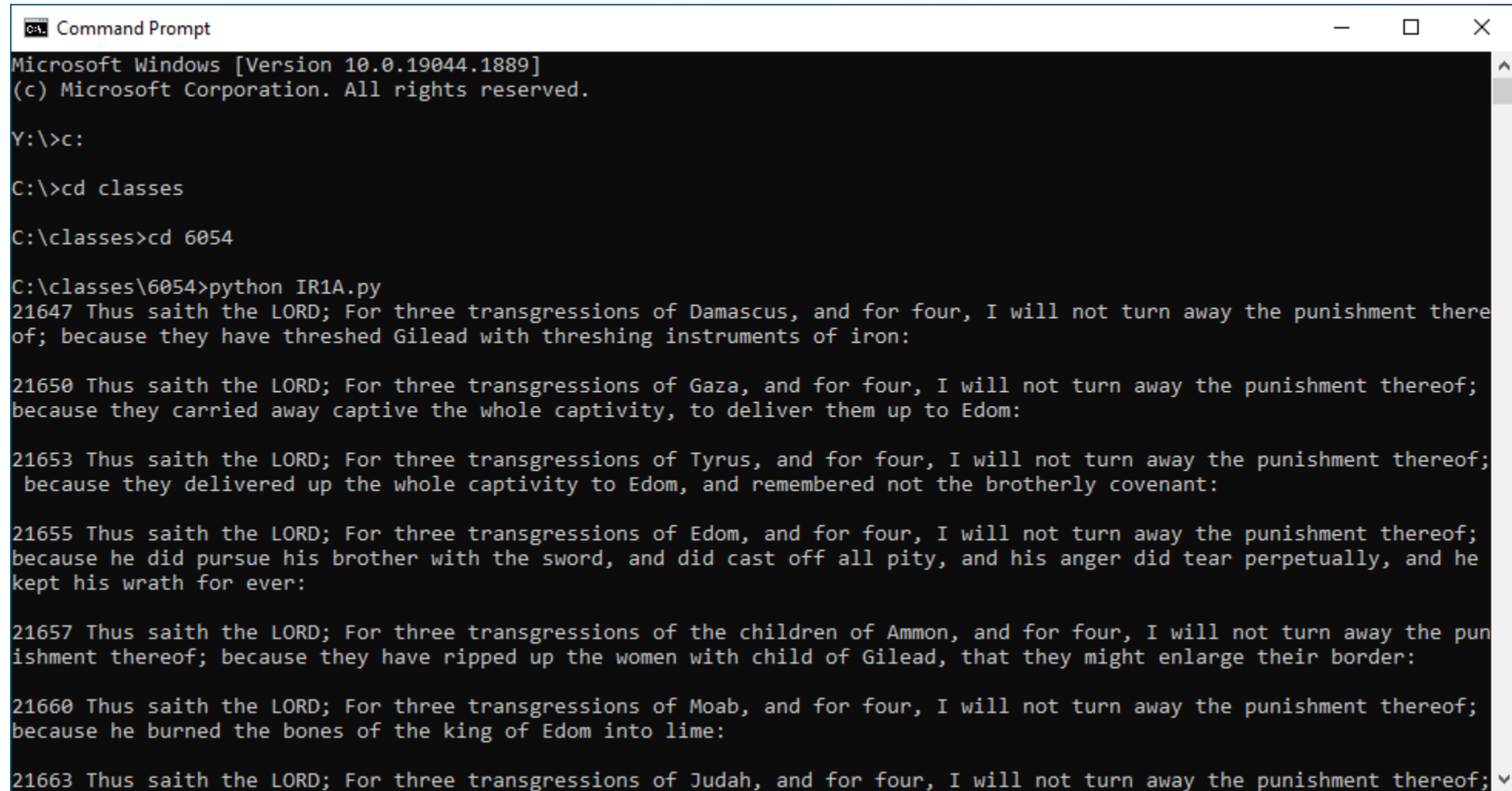
list.extend(iterable)
Extend the list by appending all the items from the iterable. Equivalent to `a[len(a):] = iterable`.

list.insert(i, x)
Insert an item at a given position. The first argument is the index of the element before which to insert, so `a.insert(0, x)` inserts at the front of the list, and `a.insert(len(a), x)` is equivalent to `a.append(x)`.

list.remove(x)
Remove the first item from the list whose value is equal to x. It raises a `ValueError` if there is no such item.

list.pop([i])

Preferred Way to Run Python Programs



```
Ca\ Command Prompt
Microsoft Windows [Version 10.0.19044.1889]
(c) Microsoft Corporation. All rights reserved.

Y:\>c:

C:\>cd classes

C:\classes>cd 6054

C:\classes\6054>python IR1A.py
21647 Thus saith the LORD; For three transgressions of Damascus, and for four, I will not turn away the punishment thereof; because they have threshed Gilead with threshing instruments of iron:

21650 Thus saith the LORD; For three transgressions of Gaza, and for four, I will not turn away the punishment thereof; because they carried away captive the whole captivity, to deliver them up to Edom:

21653 Thus saith the LORD; For three transgressions of Tyrus, and for four, I will not turn away the punishment thereof; because they delivered up the whole captivity to Edom, and remembered not the brotherly covenant:

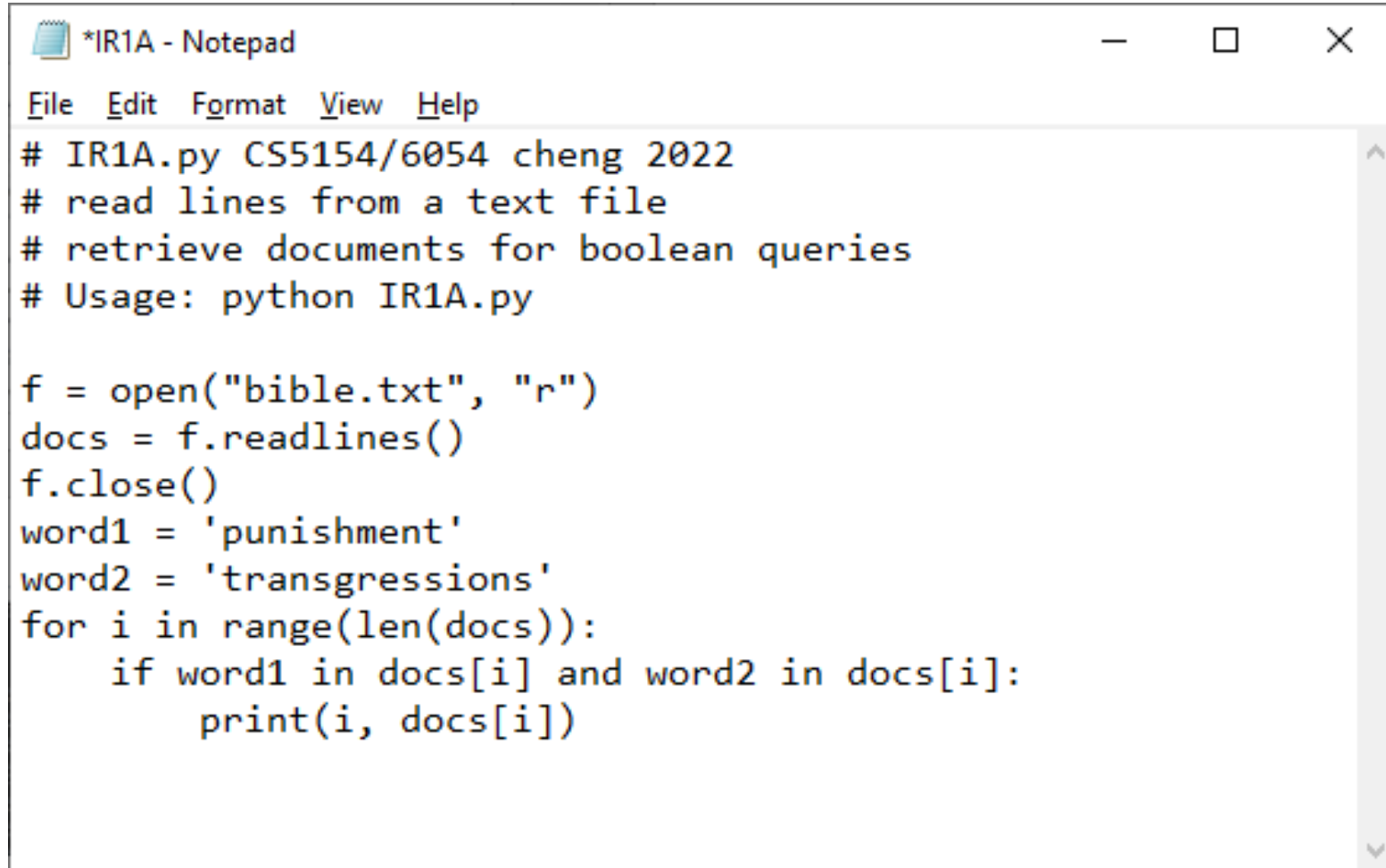
21655 Thus saith the LORD; For three transgressions of Edom, and for four, I will not turn away the punishment thereof; because he did pursue his brother with the sword, and did cast off all pity, and his anger did tear perpetually, and he kept his wrath for ever:

21657 Thus saith the LORD; For three transgressions of the children of Ammon, and for four, I will not turn away the punishment thereof; because they have ripped up the women with child of Gilead, that they might enlarge their border:

21660 Thus saith the LORD; For three transgressions of Moab, and for four, I will not turn away the punishment thereof; because he burned the bones of the king of Edom into lime:

21663 Thus saith the LORD; For three transgressions of Judah, and for four, I will not turn away the punishment thereof;
```

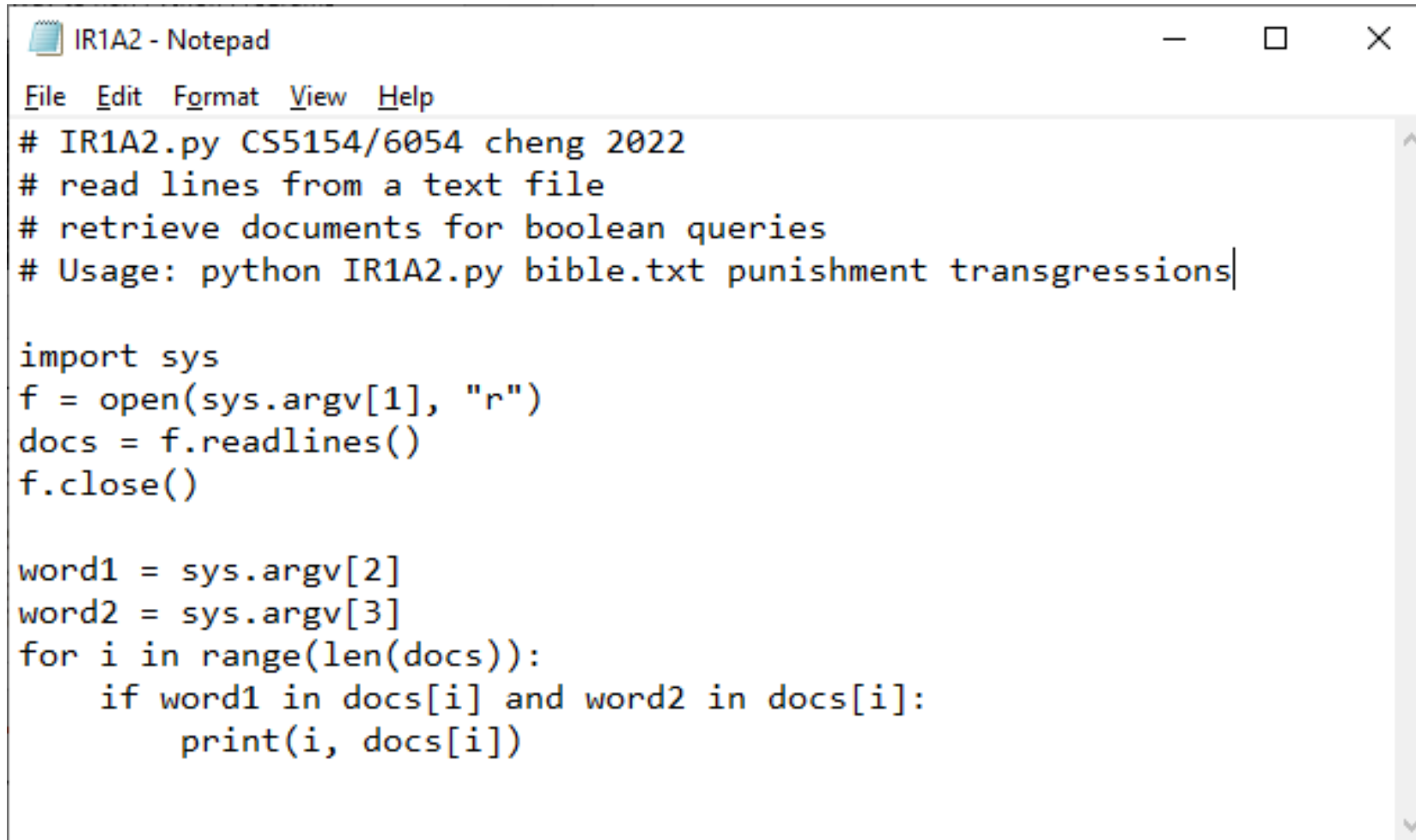
Editing Python Program Using Notepad



```
*IR1A - Notepad
File Edit Format View Help
# IR1A.py CS5154/6054 cheng 2022
# read lines from a text file
# retrieve documents for boolean queries
# Usage: python IR1A.py

f = open("bible.txt", "r")
docs = f.readlines()
f.close()
word1 = 'punishment'
word2 = 'transgressions'
for i in range(len(docs)):
    if word1 in docs[i] and word2 in docs[i]:
        print(i, docs[i])
```

Using Command Line Arguments



```
IR1A2 - Notepad
File Edit Format View Help
# IR1A2.py CS5154/6054 cheng 2022
# read lines from a text file
# retrieve documents for boolean queries
# Usage: python IR1A2.py bible.txt punishment transgressions

import sys
f = open(sys.argv[1], "r")
docs = f.readlines()
f.close()

word1 = sys.argv[2]
word2 = sys.argv[3]
for i in range(len(docs)):
    if word1 in docs[i] and word2 in docs[i]:
        print(i, docs[i])
```

```
ca. Command Prompt

C:\classes\6054>python IR1A2.py bible.txt punishment transgressions
21647 Thus saith the LORD; For three transgressions of Damascus, and for four, I will not turn away the punishment thereof; because they have threshed Gilead with threshing instruments of iron:

21650 Thus saith the LORD; For three transgressions of Gaza, and for four, I will not turn away the punishment thereof; because they carried away captive the whole captivity, to deliver them up to Edom:

21653 Thus saith the LORD; For three transgressions of Tyrus, and for four, I will not turn away the punishment thereof; because they delivered up the whole captivity to Edom, and remembered not the brotherly covenant:

21655 Thus saith the LORD; For three transgressions of Edom, and for four, I will not turn away the punishment thereof; because he did pursue his brother with the sword, and did cast off all pity, and his anger did tear perpetually, and he kept his wrath for ever:

21657 Thus saith the LORD; For three transgressions of the children of Ammon, and for four, I will not turn away the punishment thereof; because they have ripped up the women with child of Gilead, that they might enlarge their border:

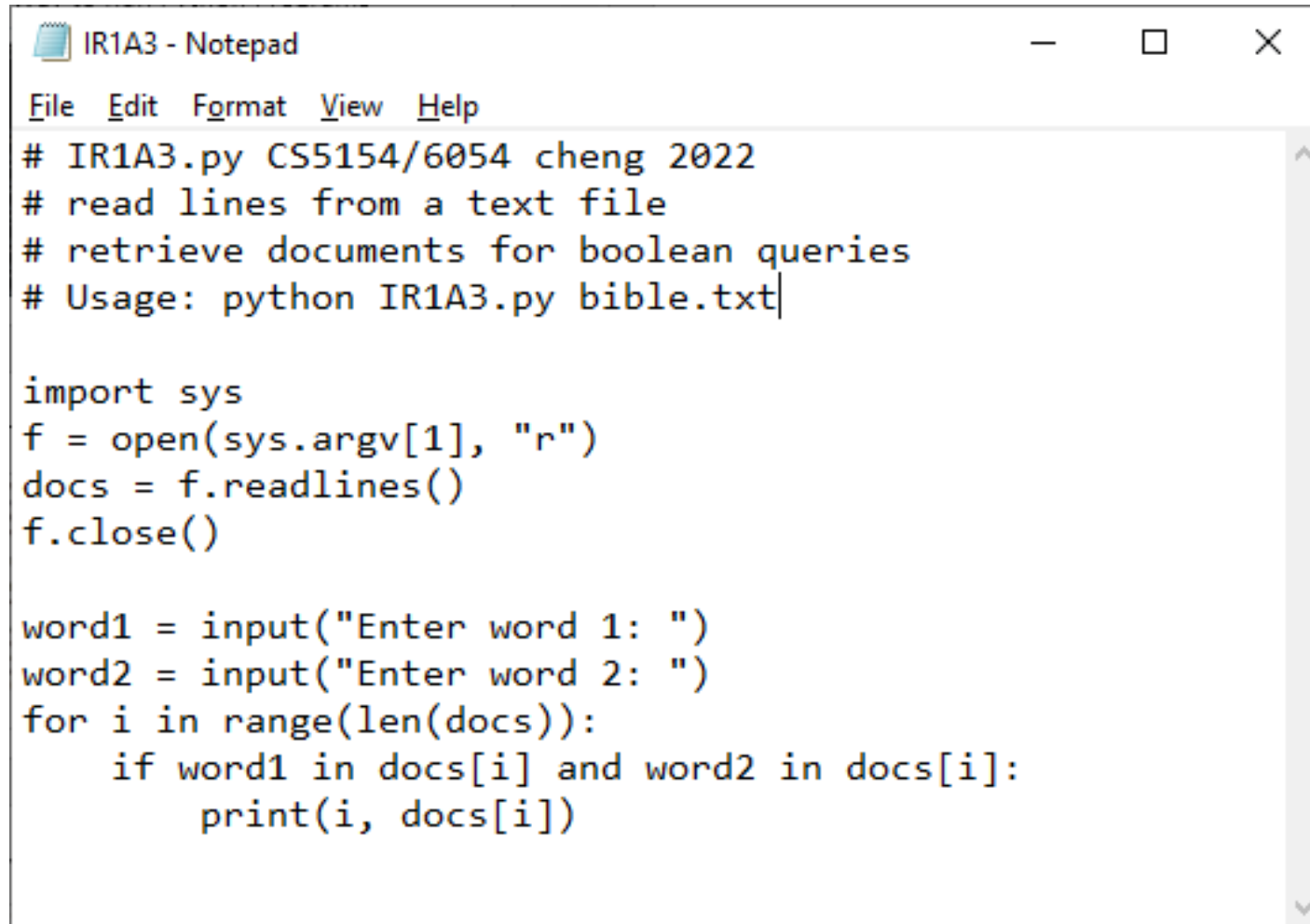
21660 Thus saith the LORD; For three transgressions of Moab, and for four, I will not turn away the punishment thereof; because he burned the bones of the king of Edom into lime:

21663 Thus saith the LORD; For three transgressions of Judah, and for four, I will not turn away the punishment thereof; because they have despised the law of the LORD, and have not kept his commandments, and their lies caused them to err, after the which their fathers have walked:

21665 Thus saith the LORD; For three transgressions of Israel, and for four, I will not turn away the punishment thereof; because they sold the righteous for silver, and the poor for a pair of shoes;

C:\classes\6054>
```

Using User Input



```
IR1A3 - Notepad
File Edit Format View Help
# IR1A3.py CS5154/6054 cheng 2022
# read lines from a text file
# retrieve documents for boolean queries
# Usage: python IR1A3.py bible.txt

import sys
f = open(sys.argv[1], "r")
docs = f.readlines()
f.close()

word1 = input("Enter word 1: ")
word2 = input("Enter word 2: ")
for i in range(len(docs)):
    if word1 in docs[i] and word2 in docs[i]:
        print(i, docs[i])
```



```
ca. Command Prompt

C:\classes\6054>python IR1A3.py bible.txt
Enter word 1: punishment
Enter word 2: transgressions
21647 Thus saith the LORD; For three transgressions of Damascus, and for four, I will not turn away the punishment thereof; because they have threshed Gilead with threshing instruments of iron:

21650 Thus saith the LORD; For three transgressions of Gaza, and for four, I will not turn away the punishment thereof; because they carried away captive the whole captivity, to deliver them up to Edom:

21653 Thus saith the LORD; For three transgressions of Tyrus, and for four, I will not turn away the punishment thereof; because they delivered up the whole captivity to Edom, and remembered not the brotherly covenant:

21655 Thus saith the LORD; For three transgressions of Edom, and for four, I will not turn away the punishment thereof; because he did pursue his brother with the sword, and did cast off all pity, and his anger did tear perpetually, and he kept his wrath for ever:

21657 Thus saith the LORD; For three transgressions of the children of Ammon, and for four, I will not turn away the punishment thereof; because they have ripped up the women with child of Gilead, that they might enlarge their border:

21660 Thus saith the LORD; For three transgressions of Moab, and for four, I will not turn away the punishment thereof; because he burned the bones of the king of Edom into lime:

21663 Thus saith the LORD; For three transgressions of Judah, and for four, I will not turn away the punishment thereof; because they have despised the law of the LORD, and have not kept his commandments, and their lies caused them to err, after the which their fathers have walked:

21665 Thus saith the LORD; For three transgressions of Israel, and for four, I will not turn away the punishment thereof; because they sold the righteous for silver, and the poor for a pair of shoes;
```

Printed Anno Dom. 1635.

119. 145. thy word is ve yf *Pr. 15*

THE EXHAUSTIVE
CONCORDANCE
TO THE
UNITED STATES
CONSTITUTION

WITH TOPICAL INDEX AND
RAPID REFERENCE CONSTITUTION

DENNIS EIZZOCO
EDITOR

Firm Foundation Press

Term-Document Incidence Matrix Fig 1.1 IIR

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
ANTHONY	1	1	0	0	0	1
BRUTUS	1	1	0	1	0	0
CAESAR	1	1	0	1	1	1
CALPURNIA	0	1	0	0	0	0
CLEOPATRA	1	0	0	0	0	0
MERCY	1	0	1	1	1	1
WORSER	1	0	1	1	1	0

...

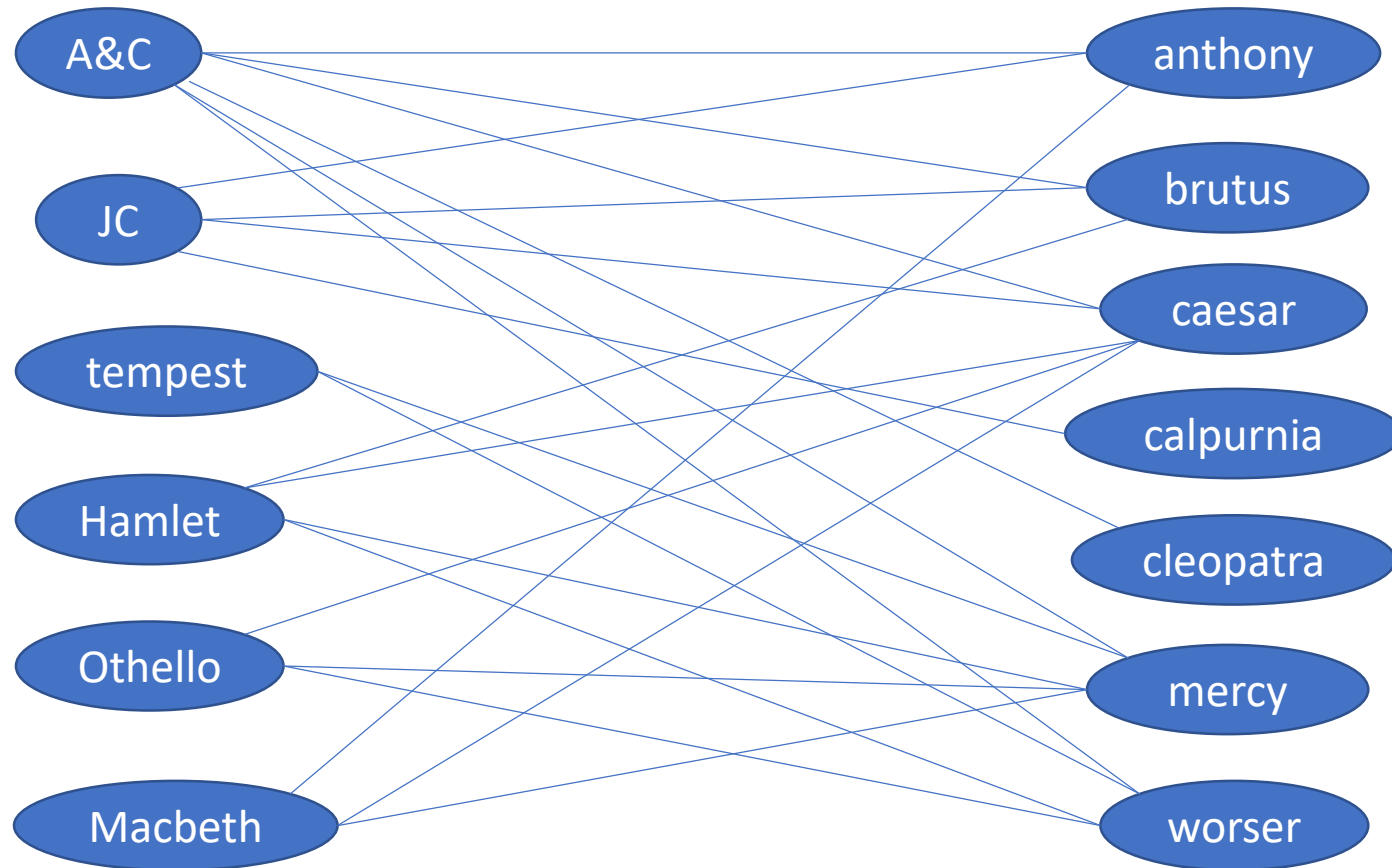
Entry is 1 if term occurs. Example: CALPURNIA occurs in *Julius Caesar*.

Entry is 0 if term doesn't occur. Example: CALPURNIA doesn't occur in *The tempest*.

A Large but Sparse Matrix

- Size of incidence matrix: number of documents times number terms ! too large for large collections
- But the matrix is very sparse – mostly 0s, few 1s.
 - Unless, a feature is shared by the most (the word “a”).
 - In that case, we may flip the feature (“occur” becomes “does not occur”)
- Mathematically, it is a bipartite graph.
- The 1s in the incidence matrix are edges of the graph.
- All we need is an edge list file for storage.

The Bipartite Graph



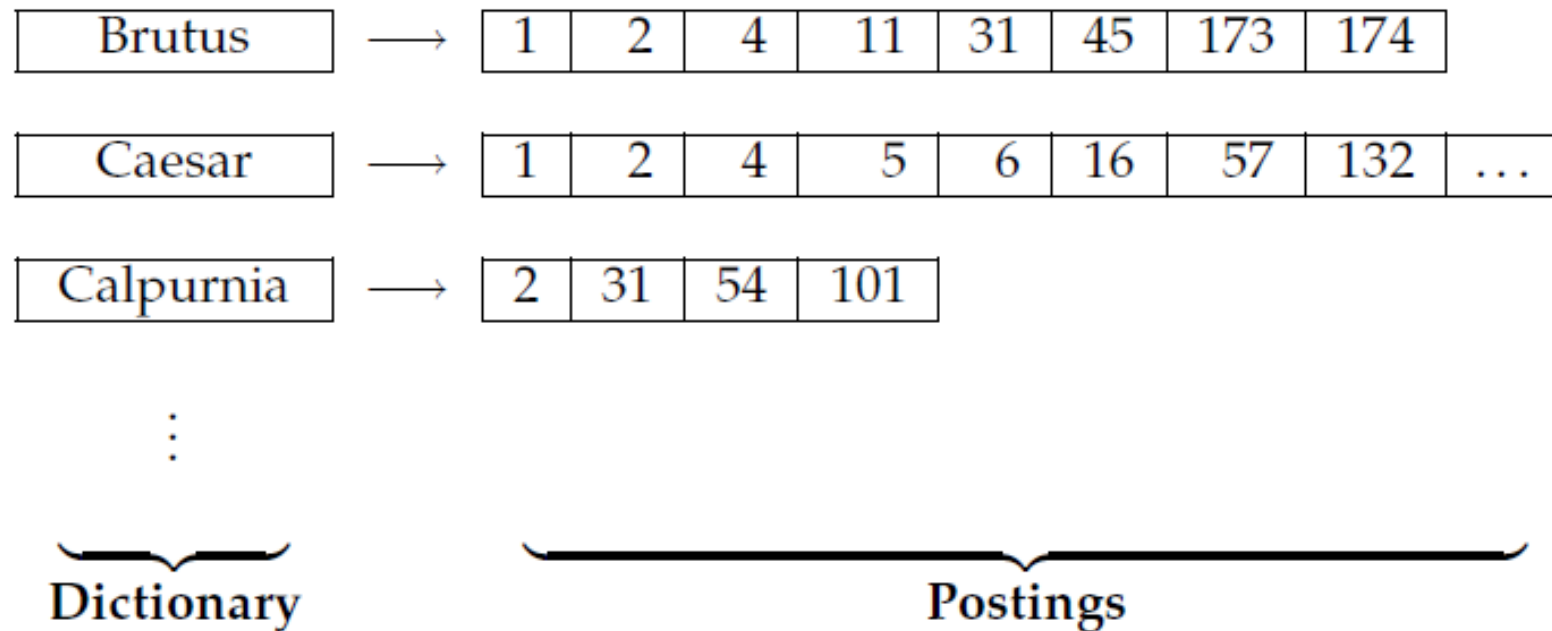
Representation of a Bipartite Graph

- The adjacency matrix (the term-document incidence matrix)
- The edges list
 - (A&C, anthony), (A&C, brutus), (JC, brutus),...
- A map from documents to sets of terms.
 - Hamlet : {brutus, caesar, mercy, worser}
 - Othello: {caesar, mercy, wroser}
- A map from terms to sets of documents. (Inverted index)
 - anthony: {A&C, JC, Macbeth}
 - brutus: {A&C, JC, Hamlet}

Inverted Index

- For each term t , we store a list of all documents that contain t .
 - A list of sets of terms
 - = For each term t , we store the 1s in its row in the incidence matrix
- It allows us to access all documents containing a term, instead of all terms in a document, which is the raw data, in an efficient way.
- In information retrieval, we may want to find or rank all documents containing a subset of terms.
- For this purpose, the inverted index is a necessary preprocessing step of the raw data, or a list of sets of terms.

Dictionary and Postings



► **Figure 1.3** The two parts of an inverted index. The dictionary is commonly kept in memory, with pointers to each postings list, which is stored on disk.

Example Documents

Doc 1

I did enact Julius Caesar: I was killed
i' the Capitol; Brutus killed me.

Doc 2

So let it be with Caesar. The noble Brutus
hath told you Caesar was ambitious:

Making the Inverted Index

1. Collect the documents to be indexed:

Friends, Romans, countrymen. So let it be with Caesar ...

2. Tokenize the text, turning each document into a list of tokens:

Friends Romans countrymen So ...

3. Do linguistic preprocessing, producing a list of normalized tokens, which are the indexing terms:

friend roman countryman so ...

4. Index the documents that each term occurs in by creating an inverted index, consisting of a dictionary and postings.

Tokenization (Edges of the Bipartite Graph)

term	docID		
I	1	so	2
did	1	let	2
enact	1	it	2
julius	1	be	2
caesar	1	with	2
I	1	caesar	2
was	1	the	2
killed	1	noble	2
i'	1	brutus	2
the	1	hath	2
capitol	1	told	2
brutus	1	you	2
killed	1	caesar	2
me	1	was	2
		ambitious	2

Sorting

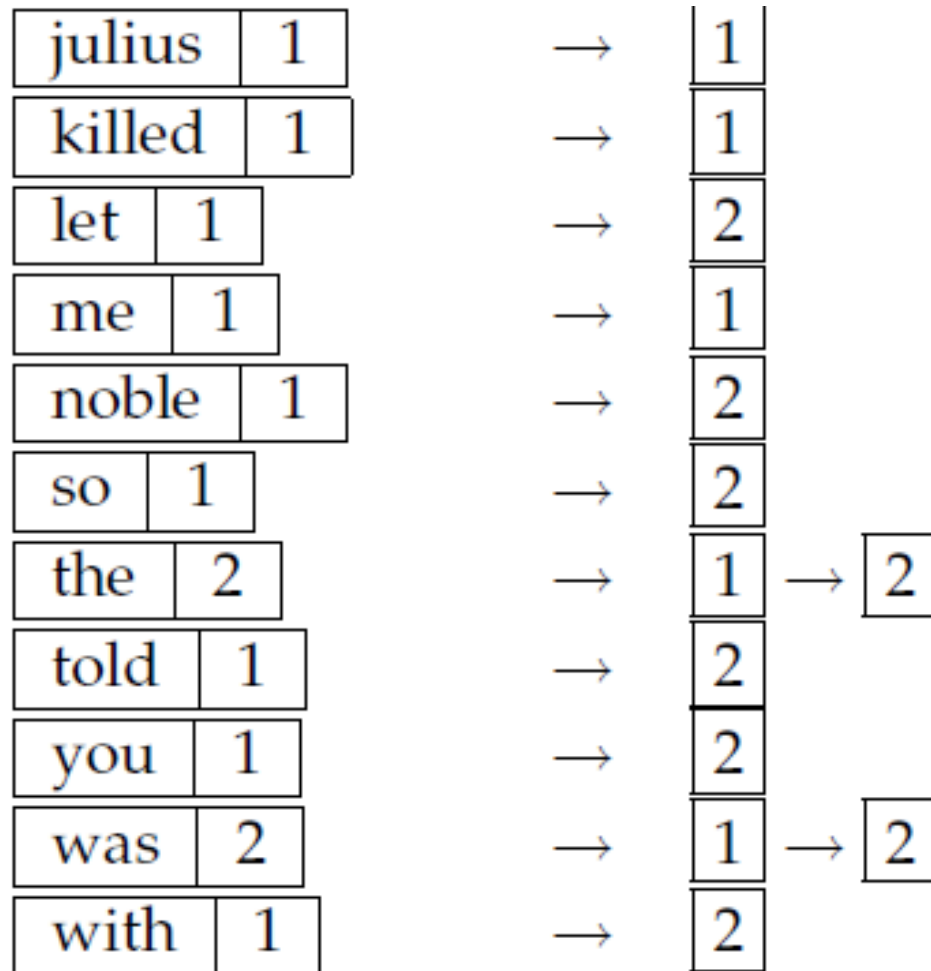
term	docID
ambitious	2
be	2
brutus	1
brutus	2
capitol	1
caesar	1
caesar	2
caesar	2
did	1
enact	1
hath	1
I	1
I	1

i'	1
it	2
julius	1
killed	1
killed	1
let	2
me	1
noble	2
so	2
the	1
the	2
told	2
you	2
was	1
was	2
with	2

Dictionary and Postings Lists

term	doc. freq.	→	postings lists
ambitious	1	→	2
be	1	→	2
brutus	2	→	1 → 2
capitol	1	→	1
caesar	2	→	1 → 2
did	1	→	1
enact	1	→	1
hath	1	→	2
I	1	→	1
i'	1	→	1
it	1	→	2

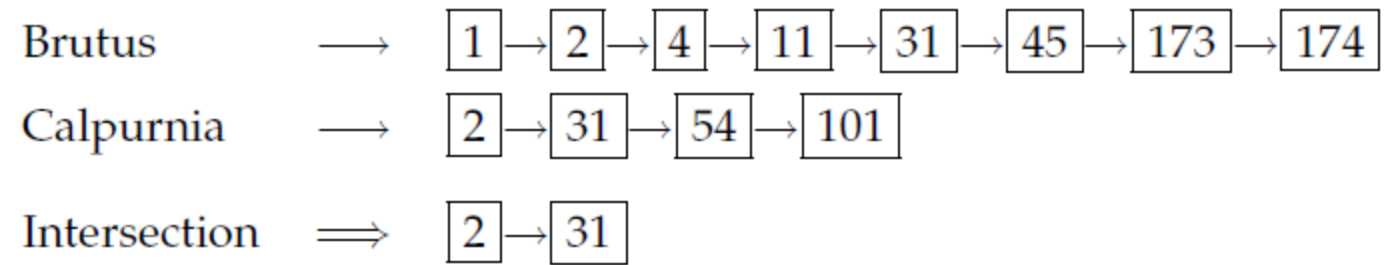
A Map from Terms to Sets of DocIDs



Processing Boolean Expressions

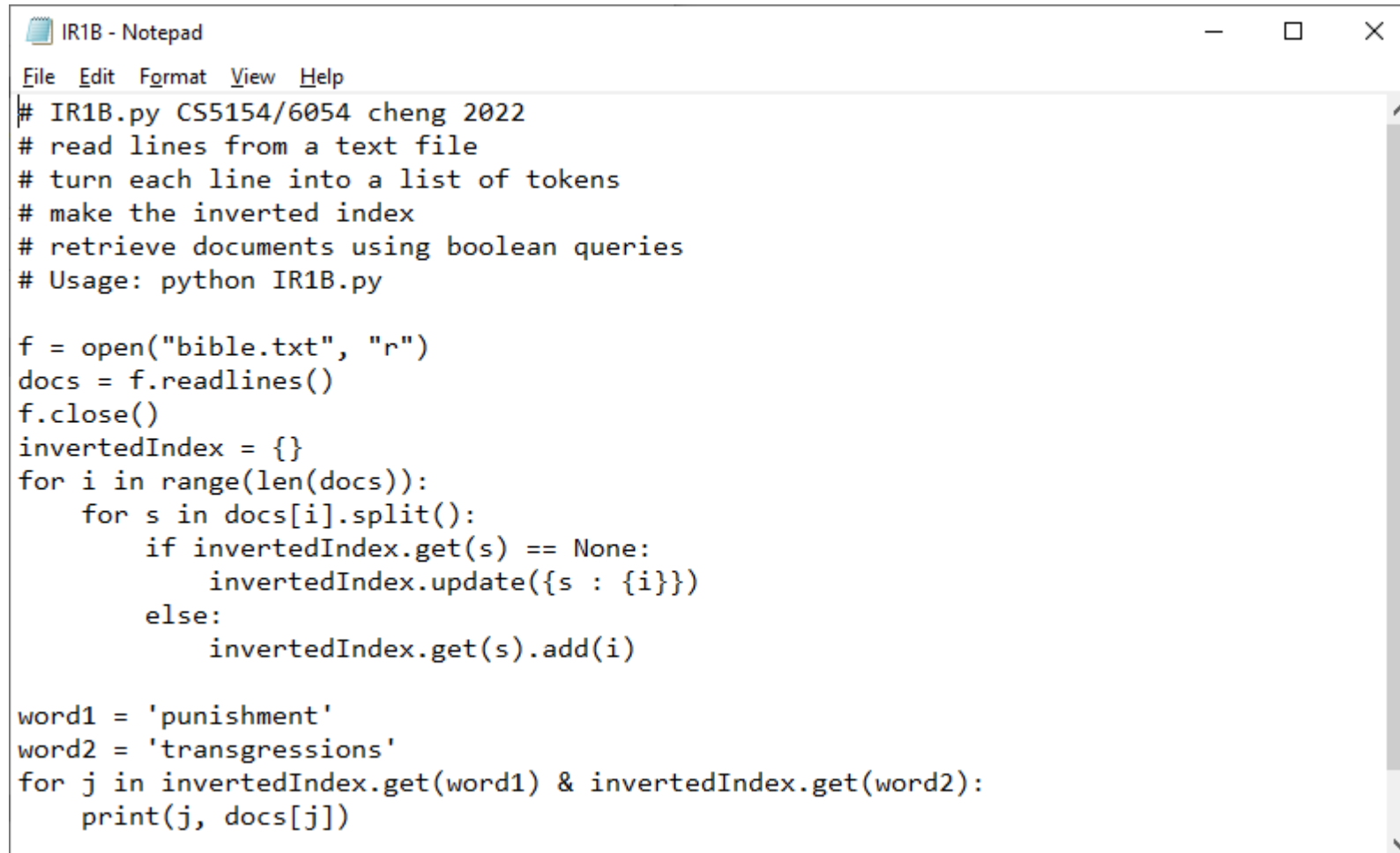
- Consider the query: Brutus AND Calpurnia
- To find all matching documents using inverted index:
 - 1 Locate Brutus in the dictionary
 - 2 Retrieve its postings list from the postings file
 - 3 Locate Calpurnia in the dictionary
 - 4 Retrieve its postings list from the postings file
 - 5 Intersect the two postings lists
 - 6 Return intersection to user

Intersection of Postings Lists



► Figure 1.5 Intersecting the postings lists for Brutus and Calpurnia from Figure 1.3.

Inverted Index as a dict (map)



```
IR1B - Notepad
File Edit Format View Help
# IR1B.py CS5154/6054 cheng 2022
# read lines from a text file
# turn each line into a list of tokens
# make the inverted index
# retrieve documents using boolean queries
# Usage: python IR1B.py

f = open("bible.txt", "r")
docs = f.readlines()
f.close()
invertedIndex = {}
for i in range(len(docs)):
    for s in docs[i].split():
        if invertedIndex.get(s) == None:
            invertedIndex.update({s : {i}})
        else:
            invertedIndex.get(s).add(i)

word1 = 'punishment'
word2 = 'transgressions'
for j in invertedIndex.get(word1) & invertedIndex.get(word2):
    print(j, docs[j])
```

Boolean Queries

- The example was a simple conjunctive query . . .
- . . . the Boolean retrieval model can answer any query that is a Boolean expression.
 - Boolean queries are queries that use and, or and not to join query terms.
 - Views each document as a set of terms.
 - Is precise: Document matches condition or not.
- Primary commercial retrieval tool for 3 decades
- Many professional searchers (e.g., lawyers) still like Boolean queries.
 - You know exactly what you are getting.
- Many search systems you use are also Boolean: search system on your laptop, in your email reader, on the intranet etc

Pros and Cons of the Boolean Model

- Key property: Documents either match or don't.
- Good for expert users with precise understanding of their needs and of the collection.
- Also good for applications: Applications can easily consume 1000s of results.
- Not good for the majority of users
- Most users are not capable of writing Boolean queries . . .
 - . . . or they are, but they think it's too much work.
- Most users don't want to wade through 1000s of results.
- This is particularly true of web search.

Problems with Boolean Search

- Boolean queries often result in either too few (=0) or too many (1000s) results.
- Query 1 (boolean conjunction): [standard user dlink 650]
 - → 200,000 hits – feast
- Query 2 (boolean conjunction): [standard user dlink 650 no card found]
 - → 0 hits – famine
- In Boolean retrieval, it takes a lot of skill to come up with a query that produces a manageable number of hits.

No Problems with Ranking

- With ranking, large result sets are not an issue.
- Just show the top 10 results and the user won't be overwhelmed
- Premise: the ranking algorithm works: More relevant results are ranked higher than less relevant results.

Assignment 1: due 8/26/2022

- Add “import time” at the beginning of IR1A.py and IR1B.py.
- Run `t1 = time.process_time_ns()` (time in nanoseconds) before and after a section of the code and output the difference between two times.
 - Feel free to use variables `t2`, `t3`, `t4`,... so you can `print(t2-t1)` etc. at the end of the program.
- Find the process times for the following tasks.
 - Reading the collection of documents (both IR1A and IR1B).
 - Making the inverted index (IR1B).
 - Boolean retrieval with an expression with two terms (IR1A without and IR1B with the inverted index).
 - You may want to have a loop of 100 iterations of this block of code to see the difference.

Hints

- How to do the same thing 100 times?
 - for rep in range(100):
- Using the timer (example: number of nanoseconds reading data)

```
import time

t1 = time.process_time_ns()
f = open("bible.txt", "r")
docs = f.readlines()
f.close()
t2 = time.process_time_ns()
print(t2 - t1)
```