

# Regular Expressions

CS5154/6054

Yizong Cheng

8/25/2022

```
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)
```

import re

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

import re

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

O'REILLY®

# Blueprints for Text Analytics Using Python

Machine Learning-Based Solutions for  
Common Real World (NLP) Applications



Jens Albrecht,  
Sidharth Ramachandran  
& Christian Winkler

## **Chapter 4. Preparing Textual Data for Statistics and Machine Learning**

---

## REGULAR EXPRESSIONS

Regular expressions are an essential tool for text data preparation. They can be used not only for tokenization and data cleaning but also for the identification and treatment of email addresses, salutations, program code, and more.

Python has the standard library `re` for regular expressions and the newer, backward-compatible library `regex` that offers support for POSIX character classes and some more flexibility.

A good overview about the available meta-characters like `^` as well as character classes like `\w` is available at [W3Schools](https://www.w3schools.com/python/python_regex.asp). There is also a number of interactive websites to develop and test regular expressions, e.g., <https://regex101.com> (make sure to set the flavor to Python).

In many packages, you will find the precompiled regular expressions like this:

---

```
RE_BRACKET = re.compile('\[[^\[\]]*\]')
text = RE_BRACKET.sub(' ', text)
```

---



Tutorials ▼

References ▼

Exercises ▼

Videos

Website

Paid Courses

Log in



HTML

CSS

JAVASCRIPT

SQL

PYTHON

PHP

BOOTSTRAP

HOW TO



## Python Tutorial

Python HOME

Python Intro

Python Get Started

Python Syntax

Python Comments

Python Variables

Python Data Types

Python Numbers

Python Casting

Python Strings

Python Booleans

Python Operators

Python Lists

Python Tuples

Python Sets

Python Dictionaries

# Python RegEx

&lt; Previous

Next &gt;

A RegEx, or Regular Expression, is a sequence of characters that forms a search pattern.

RegEx can be used to check if a string contains the specified search pattern.



HTML

CSS

JAVASCRIPT

SQL

PYTHON

PHP

BOOTSTRAP

HOW TO



## Python Tutorial

Python HOME

Python Intro

Python Get Started

Python Syntax

Python Comments

Python Variables

Python Data Types

Python Numbers

Python Casting

Python Strings

Python Booleans

Python Operators

Python Lists

Python Tuples

Python Sets

Python Dictionaries

Python If...Else

Python While Loops

Python For Loops

## RegEx Functions

The `re` module offers a set of functions that allows us to search a string for a match:

Function	Description
<u><a href="#">findall</a></u>	Returns a list containing all matches
<u><a href="#">search</a></u>	Returns a <u><a href="#">Match object</a></u> if there is a match anywhere in the string
<u><a href="#">split</a></u>	Returns a list where the string has been split at each match
<u><a href="#">sub</a></u>	Replaces one or many matches with a string

LIKE US



Get certified  
by completing  
a course  
today!



Get started

CODE GAME

FileEditViewHistoryBookmarksToolsHelp

Python RegEx

←→↻🏠🔒https://www.w3schools.com/python/python\_regex.asp📄☆📧📧📧☰


🏠HTMLCSSJAVASCRIPTSQLPYTHONPHPBOOTSTRAPHOW TO🌙🌐🔍^

Python TutorialPython HOMEPython IntroPython Get StartedPython SyntaxPython CommentsPython VariablesPython Data TypesPython NumbersPython CastingPython StringsPython BooleansPython OperatorsPython ListsPython TuplesPython SetsPython DictionariesPython If...ElsePython While LoopsPython For LoopsPython FunctionsPython LambdaPython ArraysPython Classes/ObjectsPython InheritancePython Iterators

# Metacharacters

Metacharacters are characters with a special meaning:

Character	Description	Example	Try it
[]	A set of characters	"[a-m]"	<a href="#">Try it »</a>
\	Signals a special sequence (can also be used to escape special characters)	"\d"	<a href="#">Try it »</a>
.	Any character (except newline character)	"he..o"	<a href="#">Try it »</a>
^	Starts with	"^hello"	<a href="#">Try it »</a>
\$	Ends with	"world\$"	<a href="#">Try it »</a>
*	Zero or more occurrences	"aix*"	<a href="#">Try it »</a>
+	One or more occurrences	"aix+"	<a href="#">Try it »</a>
{}	Exactly the specified number of occurrences	"al{2}"	<a href="#">Try it »</a>
	Either or	"falls stays"	<a href="#">Try it »</a>
()	Capture and group		



Play Game





HTML

CSS

JAVASCRIPT

SQL

PYTHON

PHP

BOOTSTRAP

HOW TO



## Python Tutorial

Python HOME

Python Intro

Python Get Started

Python Syntax

Python Comments

Python Variables

Python Data Types

Python Numbers

Python Casting

Python Strings

Python Booleans

Python Operators

Python Lists

Python Tuples

Python Sets

Python Dictionaries

Python If...Else

Python While Loops

Python For Loops

Python Functions

Python Lambda

Python Arrays

Python Classes/Objects

## Special Sequences

A special sequence is a `\` followed by one of the characters in the list below, and has a special meaning:

Character	Description	Example	Try it
<code>\A</code>	Returns a match if the specified characters are at the beginning of the string	<code>"\AThe"</code>	<a href="#">Try it »</a>
<code>\b</code>	Returns a match where the specified characters are at the beginning or at the end of a word (the "r" in the beginning is making sure that the string is being treated as a "raw string")	<code>r"\bain"</code> <code>r"ain\b"</code>	<a href="#">Try it »</a> <a href="#">Try it »</a>
<code>\B</code>	Returns a match where the specified characters are present, but NOT at the beginning (or at the end) of a word (the "r" in the beginning is making sure that the string is being treated as a "raw string")	<code>r"\Bain"</code> <code>r"ain\B"</code>	<a href="#">Try it »</a> <a href="#">Try it »</a>

[HTML](#)[CSS](#)[JAVASCRIPT](#)[SQL](#)[PYTHON](#)[PHP](#)[BOOTSTRAP](#)[HOW TO](#)

## Python Tutorial

[Python HOME](#)[Python Intro](#)[Python Get Started](#)[Python Syntax](#)[Python Comments](#)[Python Variables](#)[Python Data Types](#)[Python Numbers](#)[Python Casting](#)[Python Strings](#)[Python Booleans](#)[Python Operators](#)[Python Lists](#)[Python Tuples](#)[Python Sets](#)[Python Dictionaries](#)[Python If...Else](#)[Python While Loops](#)[Python For Loops](#)[Python Functions](#)[Python Lambda](#)[Python Arrays](#)[Python Classes/Objects](#)

\d

Returns a match where the string contains digits (numbers from 0-9)

\d

[Try it »](#)

\D

Returns a match where the string DOES NOT contain digits

\D

[Try it »](#)

\s

Returns a match where the string contains a white space character

\s

[Try it »](#)

\S

Returns a match where the string DOES NOT contain a white space character

\S

[Try it »](#)

\w

Returns a match where the string contains any word characters (characters from a to Z, digits from 0-9, and the underscore \_ character)

\w

[Try it »](#)

\W

Returns a match where the string DOES NOT contain any word characters

\W

[Try it »](#)

\Z

Returns a match if the specified characters are at the end of the string

"Spain\Z"

[Try it »](#)

FileEditViewHistoryBookmarksToolsHelp

Python RegEx

https://www.w3schools.com/python/python\_regex.asp

HTMLCSSJAVASCRIPTSQLPYTHONPHPBOOTSTRAPHOW TO

Python Tutorial

Python HOME

Python Intro

Python Get Started

Python Syntax

Python Comments

Python Variables

Python Data Types

Python Numbers

Python Casting

Python Strings

Python Booleans

Python Operators

Python Lists

Python Tuples

Python Sets

Python Dictionaries

Python If...Else

Python While Loops

Python For Loops

Python Functions

Python Lambda

Python Arrays

Python Classes/Objects

Python Inheritance

Python Iterators

Python Scope

Python Modules

Sets

A set is a set of characters inside a pair of square brackets `[]` with a special meaning:

Set	Description	Try it
<code>[arn]</code>	Returns a match where one of the specified characters ( <code>a</code> , <code>r</code> , or <code>n</code> ) are present	<a href="#">Try it »</a>
<code>[a-n]</code>	Returns a match for any lower case character, alphabetically between <code>a</code> and <code>n</code>	<a href="#">Try it »</a>
<code>[^arn]</code>	Returns a match for any character EXCEPT <code>a</code> , <code>r</code> , and <code>n</code>	<a href="#">Try it »</a>
<code>[0123]</code>	Returns a match where any of the specified digits ( <code>0</code> , <code>1</code> , <code>2</code> , or <code>3</code> ) are present	<a href="#">Try it »</a>
<code>[0-9]</code>	Returns a match for any digit between <code>0</code> and <code>9</code>	<a href="#">Try it »</a>
<code>[0-5][0-9]</code>	Returns a match for any two-digit numbers from <code>00</code> and <code>59</code>	<a href="#">Try it »</a>
<code>[a-zA-Z]</code>	Returns a match for any character alphabetically between <code>a</code> and <code>z</code> , lower case OR upper case	<a href="#">Try it »</a>
<code>[+]</code>	In sets, <code>+</code> , <code>*</code> , <code>.</code> , <code> </code> , <code>()</code> , <code>\$</code> , <code>{}</code> has no special meaning, so <code>[+]</code> means: return a match for any <code>+</code> character in the string	<a href="#">Try it »</a>



HTML

CSS

JAVASCRIPT

SQL

PYTHON

PHP

BOOTSTRAP

HOW TO



## Python Tutorial

Python HOME

Python Intro

Python Get Started

Python Syntax

Python Comments

Python Variables

Python Data Types

Python Numbers

Python Casting

Python Strings

Python Booleans

Python Operators

Python Lists

Python Tuples

Python Sets

Python Dictionaries

Python If...Else

Python While Loops

Python For Loops

Python Functions

Python Lambda

# The findall() Function

The `findall()` function returns a list containing all matches.

## Example

Print a list of all matches:

```
import re

txt = "The rain in Spain"
x = re.findall("ai", txt)
print(x)
```

[Try it Yourself »](#)

The list contains the matches in the order they are found.

If no matches are found, an empty list is returned:



HTML

CSS

JAVASCRIPT

SQL

PYTHON

PHP

BOOTSTRAP

HOW TO



## Python Tutorial

Python HOME

Python Intro

Python Get Started

Python Syntax

Python Comments

Python Variables

Python Data Types

Python Numbers

Python Casting

Python Strings

Python Booleans

Python Operators

Python Lists

Python Tuples

Python Sets

Python Dictionaries

Python If...Else

Python While Loops

Python For Loops

Python Functions

Python Lambda

Python Arrays

Python Classes/Objects

## The search() Function

The `search()` function searches the string for a match, and returns a Match object if there is a match.

If there is more than one match, only the first occurrence of the match will be returned:

### Example

Search for the first white-space character in the string:

```
import re

txt = "The rain in Spain"
x = re.search("\s", txt)

print("The first white-space character is located in position:",
      x.start())
```

[Try it Yourself »](#)

If no matches are found, the value `None` is returned:



HTML

CSS

JAVASCRIPT

SQL

PYTHON

PHP

BOOTSTRAP

HOW TO



## Python Tutorial

Python HOME

Python Intro

Python Get Started

Python Syntax

Python Comments

Python Variables

Python Data Types

Python Numbers

Python Casting

Python Strings

Python Booleans

Python Operators

Python Lists

Python Tuples

Python Sets

Python Dictionaries

Python If...Else

Python While Loops

Python For Loops

Python Functions

Python Lambda

## The split() Function

The `split()` function returns a list where the string has been split at each match:

### Example

Split at each white-space character:

```
import re

txt = "The rain in Spain"
x = re.split("\s", txt)
print(x)
```

[Try it Yourself »](#)

You can control the number of occurrences by specifying the `maxsplit` parameter:



HTML

CSS

JAVASCRIPT

SQL

PYTHON

PHP

BOOTSTRAP

HOW TO



## Python Tutorial

Python HOME

Python Intro

Python Get Started

Python Syntax

Python Comments

Python Variables

Python Data Types

Python Numbers

Python Casting

Python Strings

Python Booleans

Python Operators

Python Lists

Python Tuples

Python Sets

Python Dictionaries

Python If...Else

Python While Loops

Python For Loops

Python Functions

Python Lambda

Try it Yourself »

## The sub() Function

The `sub()` function replaces the matches with the text of your choice:

### Example

Replace every white-space character with the number 9:

```
import re

txt = "The rain in Spain"
x = re.sub("\s", "9", txt)
print(x)
```

[Try it Yourself »](#)

```
def clean(text):
    # convert html escapes like & to characters.
    text = html.unescape(text)
    # tags like <tab>
    text = re.sub(r'<[^>]*>', ' ', text)
    # markdown URLs like [Some text](https://....)
    text = re.sub(r'\([([^\[]*)\]\([^\(\\)]*\)', r'\1', text)
    # text or code in brackets like [0]
    text = re.sub(r'\([([^\[]*)\]', ' ', text)
    # standalone sequences of specials, matches &# but not #cool
    text = re.sub(r'(?:^(|s)[&#<>{}\\[]+|\\:-]{1,})(?:\s|$)', ' ', text)
    # standalone sequences of hyphens like --- or ==
    text = re.sub(r'(?:^(|s)[\-=\+]{2,})(?:\s|$)', ' ', text)
    # sequences of white spaces
    text = re.sub(r'\s+', ' ', text)
    return text.strip()
```



FileEditViewHistoryBookmarksToolsHelp

wordcloud.WordCloud — wordcloud X

←→↻🏠🔒https://amueller.github.io/word\_cloud/generated/wordcloud.WordCloud.html📄☆📧📧📧☰

🏠 wordcloud  
1.8.1

Search docs

USER DOCUMENTATION

☰ API Reference

wordcloud.WordCloud

wordcloud.ImageColorGenerator

wordcloud.random\_color\_func

wordcloud.get\_single\_color\_func

Command Line Interface

Gallery of Examples

Changelog

CONTRIBUTOR DOCUMENTATION

Making a release

🏠 » API Reference » wordcloud.WordCloudView page source

# wordcloud.WordCloud

```
class wordcloud.WordCloud(font_path=None, width=400, height=200, margin=2, ranks_only=None, prefer_horizontal=0.9, mask=None, scale=1, color_func=None, max_words=200, min_font_size=4, stopwords=None, random_state=None, background_color='black', max_font_size=None, font_step=1, mode='RGB', relative_scaling='auto', regexp=None, collocations=True, colormap=None, normalize_plurals=True, contour_width=0, contour_color='black', repeat=False, include_numbers=False, min_word_length=0, collocation_threshold=30)
```

[source]

Word cloud object for generating and drawing.

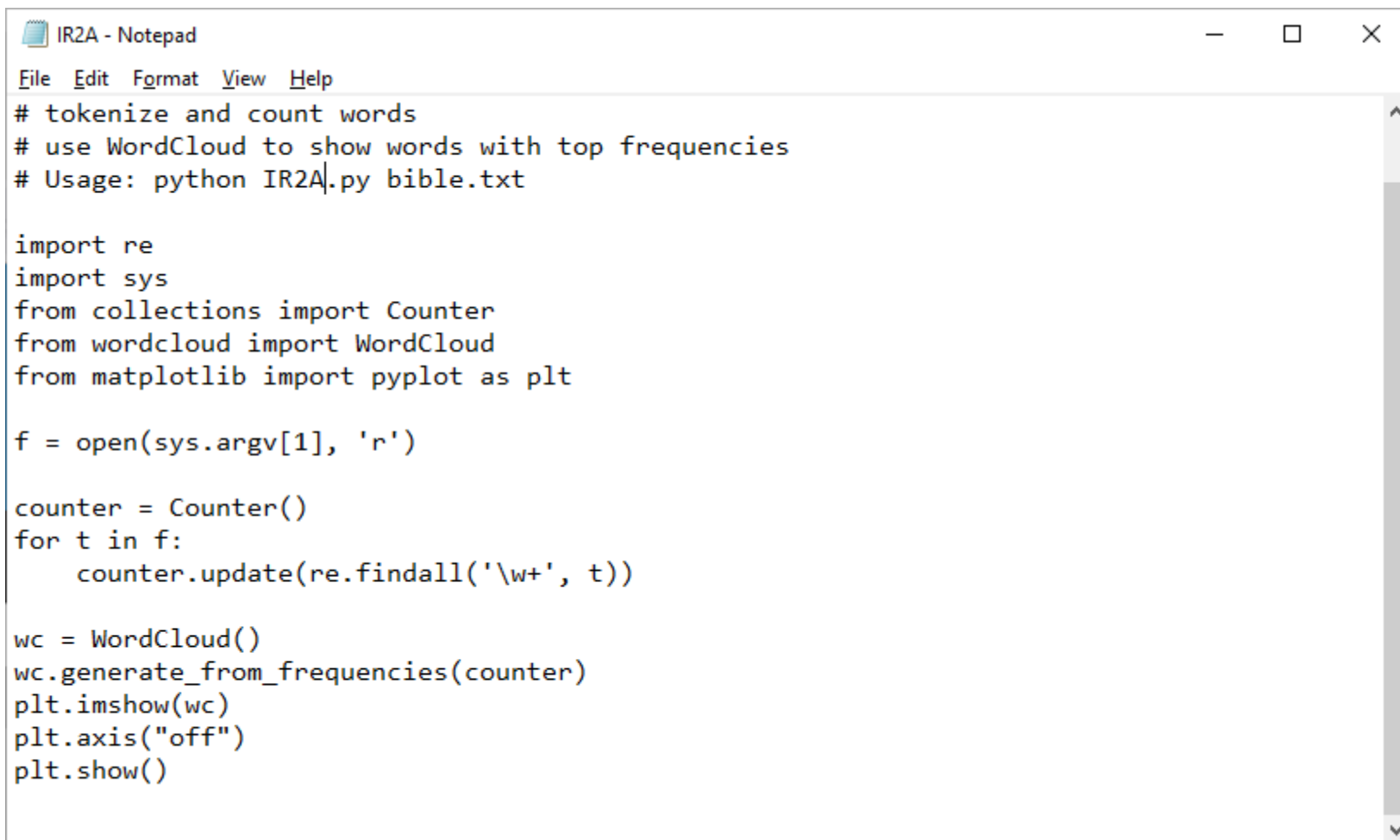
Parameters:

font\_path : string

Font path to the font that will be used (OTF or TTF). Defaults to DroidSansMono path on a Linux machine. If you are on another OS or don't have this font, you need to adjust this path.

width : int (default=400)

Width of the canvas.



```
IR2A - Notepad
File Edit Format View Help
# tokenize and count words
# use WordCloud to show words with top frequencies
# Usage: python IR2A.py bible.txt

import re
import sys
from collections import Counter
from wordcloud import WordCloud
from matplotlib import pyplot as plt

f = open(sys.argv[1], 'r')

counter = Counter()
for t in f:
    counter.update(re.findall('\w+', t))

wc = WordCloud()
wc.generate_from_frequencies(counter)
plt.imshow(wc)
plt.axis("off")
plt.show()
```

FileEditViewHistoryBookmarksToolsHelp

wordcloud.WordCloud — wordcloud X

←→↻🏠🔒https://amueller.github.io/word\_cloud/generated/wordcloud.WordCloud.html📄☆📧📧📧☰

🏠 wordcloud

1.8.1

Search docs

USER DOCUMENTATION

☰ API Reference

wordcloud.WordCloud

wordcloud.ImageColorGenerator

wordcloud.random\_color\_func

wordcloud.get\_single\_color\_func

Command Line Interface

Gallery of Examples

Changelog

CONTRIBUTOR DOCUMENTATION

Making a release

color\_func : callable, default=None

Callable with parameters word, font\_size, position, orientation, font\_path, random\_state that returns a PIL color for each word. Overwrites "colormap". See colormap for specifying a matplotlib colormap instead. To create a word cloud with a single color, use `color_func=lambda *args, **kwargs: "white"`. The single color can also be specified using RGB code. For example `color_func=lambda *args, **kwargs: (255,0,0)` sets color to red.

regexp : string or None (optional)

Regular expression to split the input text into tokens in process\_text. If None is specified, `r"\w[\w']+"` is used. Ignored if using generate\_from\_frequencies.

collocations : bool, default=True

Whether to include collocations (bigrams) of two words. Ignored if using generate\_from\_frequencies.

colormap : string or matplotlib colormap, default="viridis"

Matplotlib colormap to randomly draw colors from for each word. Ignored if "color\_func" is specified.

FileEditViewHistoryBookmarksToolsHelp

Choosing Colormaps in Matplotlib

←→↻🏠🔒https://matplotlib.org/stable/tutorials/colors/colormaps.html📄☆📧📺📺📺☰

matplotlib

Plot typesExamplesTutorialsReferenceUser guideDevelopRelease notes

📺🗨️📺🐦

🔍 Search the docs ...

Introductory▼Intermediate▼Advanced▼Colors^Specifying ColorsCustomized Colorbars TutorialCreating Colormaps in MatplotlibColormap NormalizationChoosing Colormaps in MatplotlibProvisional▼Text▼

# Choosing Colormaps in Matplotlib

Matplotlib has a number of built-in colormaps accessible via `matplotlib.cm.get_cmap`. There are also external libraries that have many extra colormaps, which can be viewed in the [Third-party colormaps](#) section of the Matplotlib documentation. Here we briefly discuss how to choose between the many options. For help on creating your own colormaps, see [Creating Colormaps in Matplotlib](#).

## Overview

The idea behind choosing a good colormap is to find a good representation in 3D colorspace for your data set. The best colormap for any given data set depends on many things including:

- Whether representing form or metric data ([Ware])
- Your knowledge of the data set (*e.g.*, is there a critical value from which the other values deviate?)

File Edit View History Bookmarks Tools Help





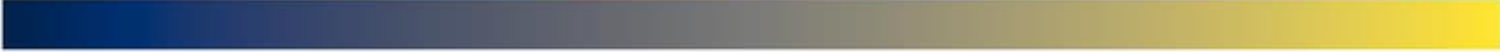
Choosing Colormaps in Matplotlib

https://matplotlib.org/stable/tutorials/colors/colormaps.html 150%

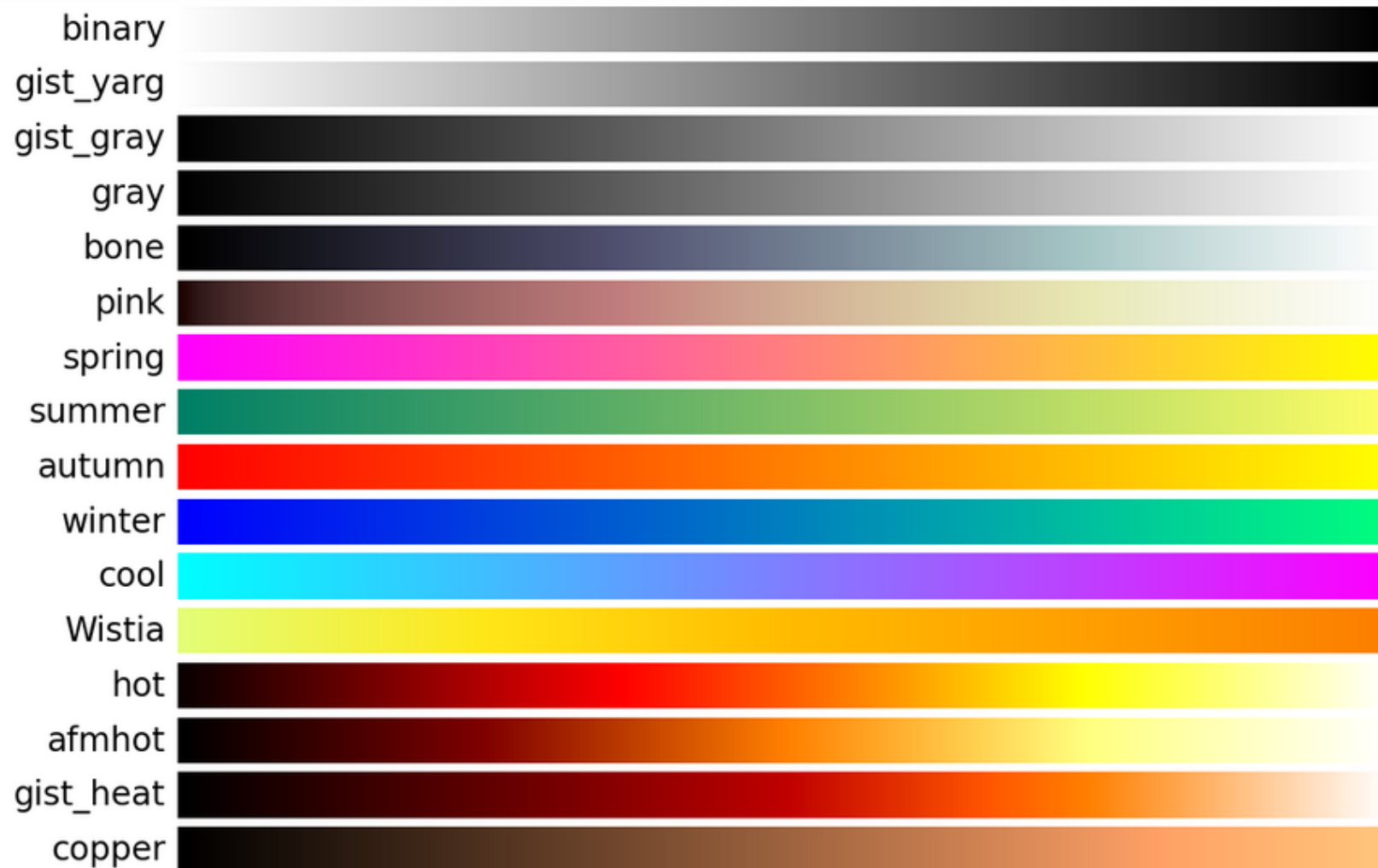
# matplotlib

```
plot_color_gradients('Perceptually Uniform Sequential',  
                    ['viridis', 'plasma', 'inferno', 'magma', 'cividis'])
```

## Perceptually Uniform Sequential colormaps

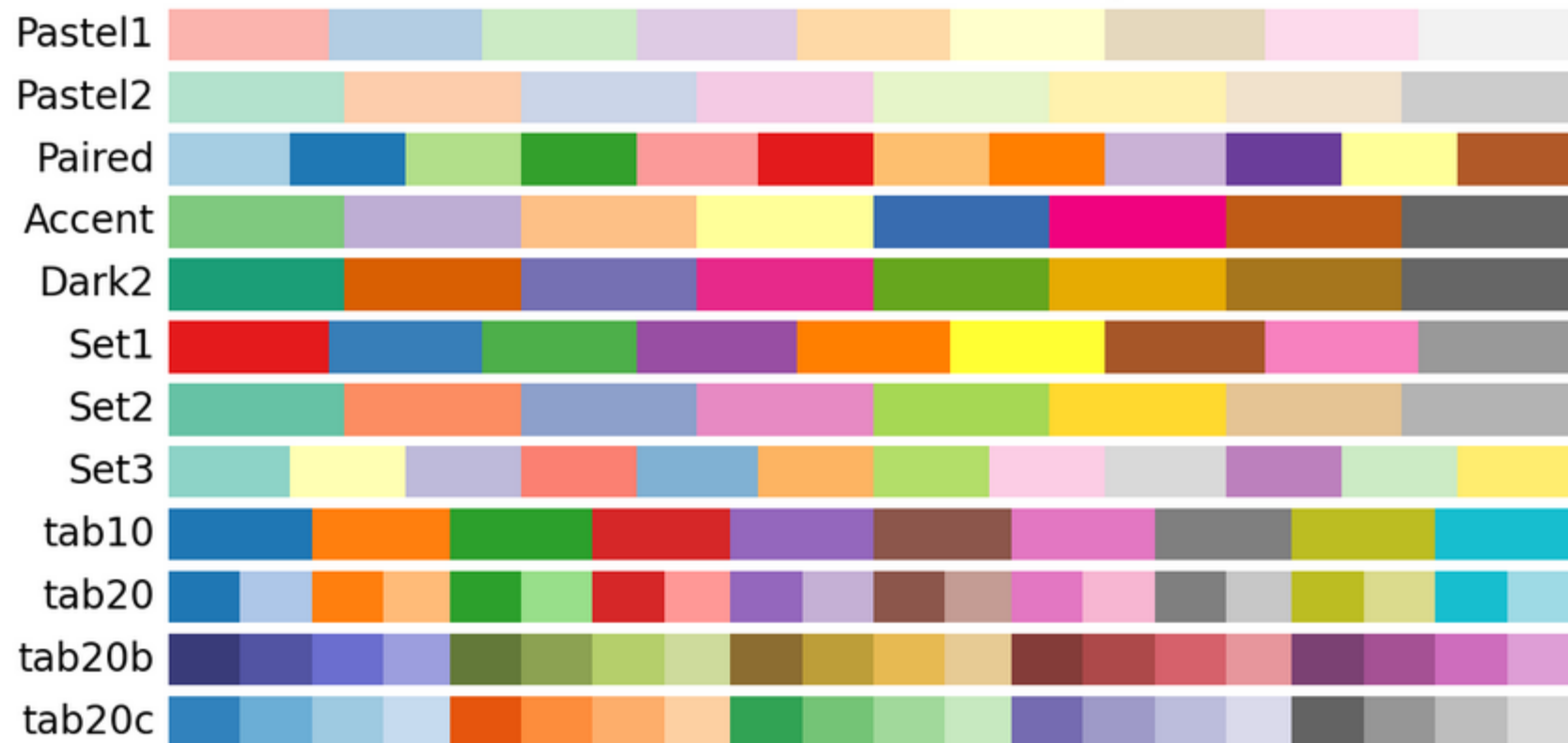
viridis	
plasma	
inferno	
magma	
cividis	

# matplotlib



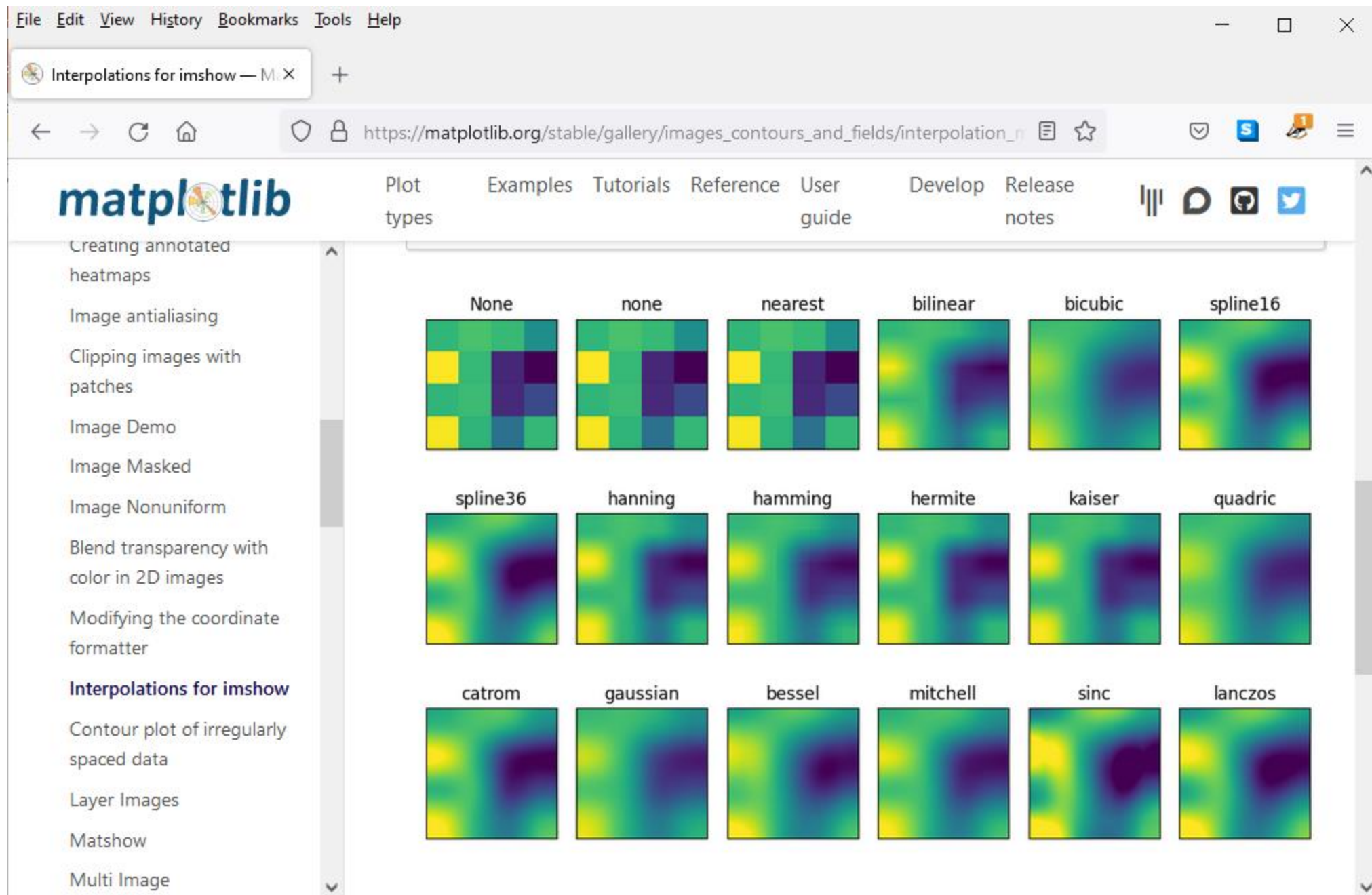


## Qualitative colormaps



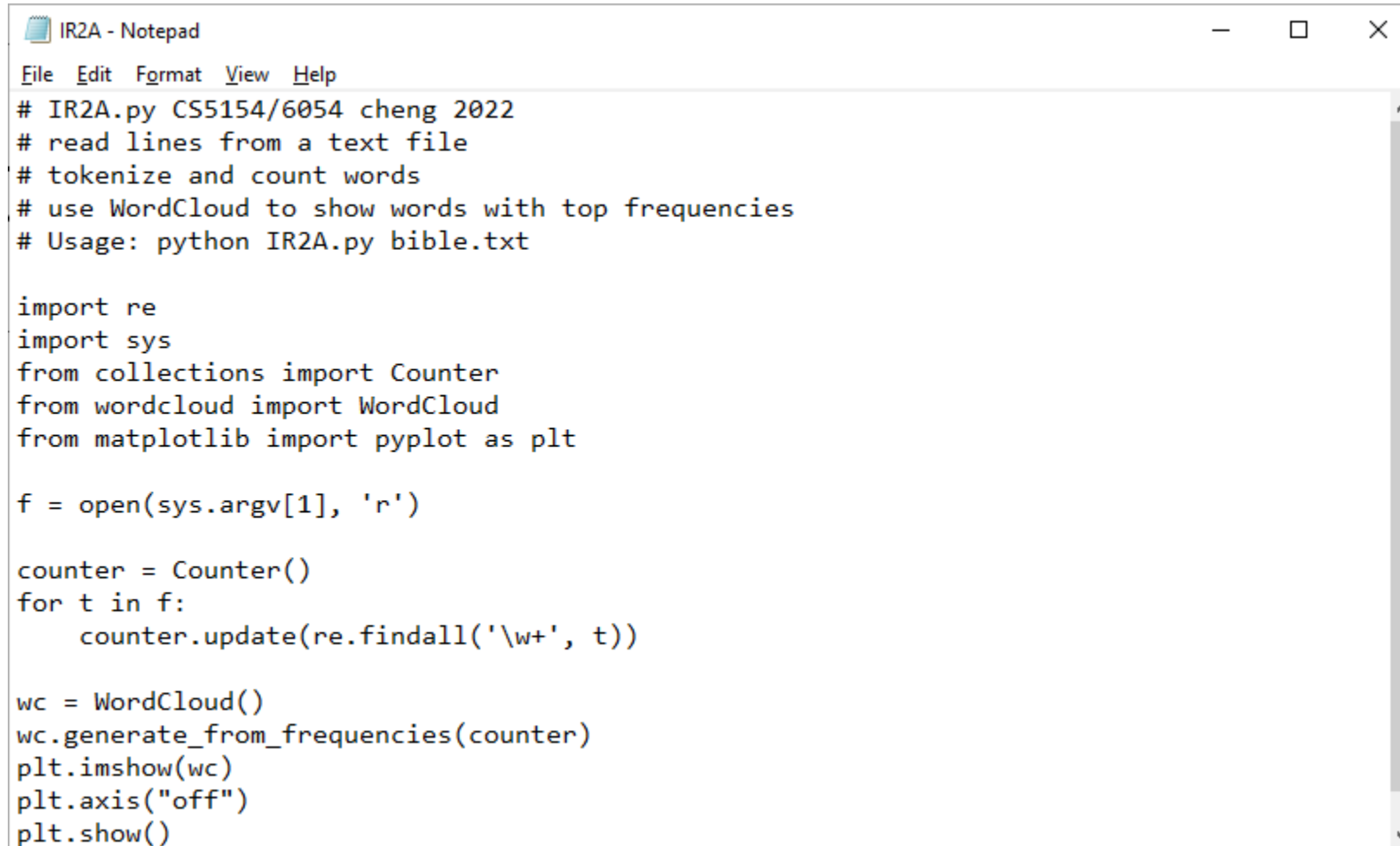


# imshow interpolation





# re.findall(), Counter, and WordCloud



```
IR2A - Notepad
File Edit Format View Help
# IR2A.py CS5154/6054 cheng 2022
# read lines from a text file
# tokenize and count words
# use WordCloud to show words with top frequencies
# Usage: python IR2A.py bible.txt

import re
import sys
from collections import Counter
from wordcloud import WordCloud
from matplotlib import pyplot as plt

f = open(sys.argv[1], 'r')

counter = Counter()
for t in f:
    counter.update(re.findall('\w+', t))

wc = WordCloud()
wc.generate_from_frequencies(counter)
plt.imshow(wc)
plt.axis("off")
plt.show()
```

