

1. Recap fra i går
2. Python primer
 1. Værktøjer
 2. Quiz / Titanic
3. Pentagon Papers

```
# -*- coding: utf-8 -*-
"""
Spyder Editor

This is a temporary script file.

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sn
from sentida import Sentida

df=pd.read_csv("data/home.csv", encoding="ISO-8859-1")

df.info()
df.columns
df['fn']=df['name'].str.split(' ')
df.loc[3,'fn']
df['clfn']=df['fn'].apply(lambda x: x[0])
df['Navn']=df['clfn'].str.replace('[^a-zA-ZæøåÄÅ -]', '', regex=True)

#
dfd=pd.read_excel("data/drenge.xlsx")
dfp=pd.read_excel("data/piger.xlsx")
dfgp=pd.merge(df, dfp, how="left", on="Navn")
dfgp.drop('Drengenavn', axis=1, inplace=True)

#hvor mange nan?
sum(dfgp['Pigenavn'].isnull())
dfgpn=pd.merge(dfgp, dfd, how="left", on="Navn")
dfgpn.drop('Pigenavn_y', axis=1, inplace=True)

dfgpn['gender']=dfgpn['Drengenavn'].apply(lambda x: "M" if x=="Ja" else "F" )

# drop irrelevant columns
dfgpn.drop('clfn', axis=1, inplace=True)

# count and lix
dfgpn['length']=df['content'].str.len()

dfgpn['lix']=df['content'].apply(computeLix)
idx=dfgpn['lix'].idxmin()
idx=dfgpn['lix'].idxmax()
dfgpn.loc[idx]['content']

# plot

sn.histplot(data=dfgpn, x='length', bins=30)
sn.histplot(data=dfgpn, x='lix', hue='gender', binwidth=1, kde=False)
→ sn.countplot(x='gender', data=dfgpn);

# sentida score → install sentida

dfgpn['sscore']=df['content'].apply(lambda x: Sentida().sentida(x,output="mean", normal=False)
# plot

sn.histplot(data=dfgpn, x='sscore', bins=40)
sn.boxplot(data=dfgpn, x='gender', y='sscore') ← bring ♀ flag layer over and ♂
```



TIDYTEXT

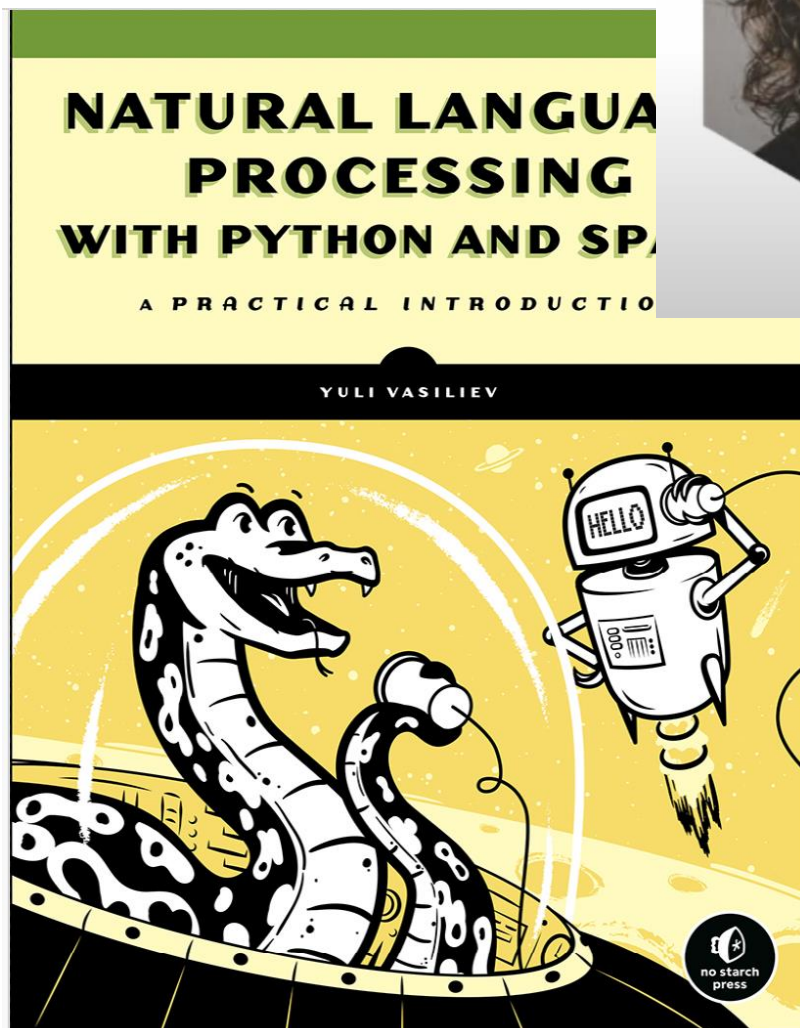
HELLO

I'm Julia Silge

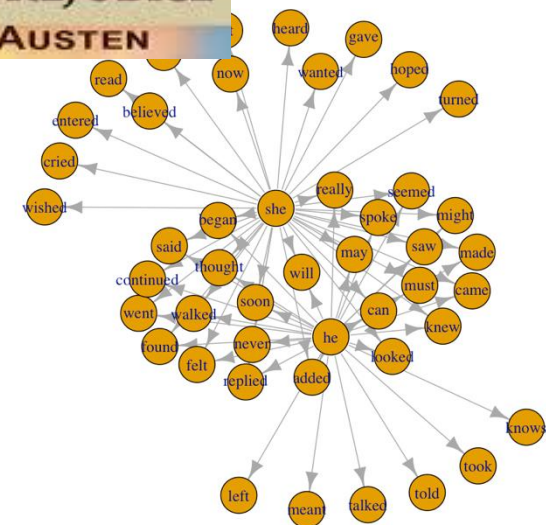
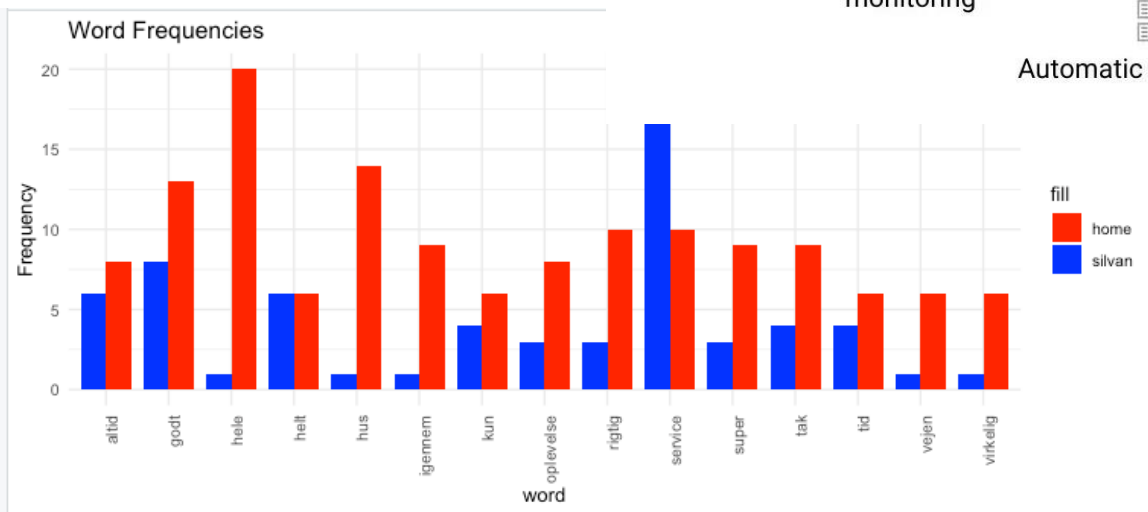
Software Engineer at RStudio

[@juliasilge](https://twitter.com/juliasilge)

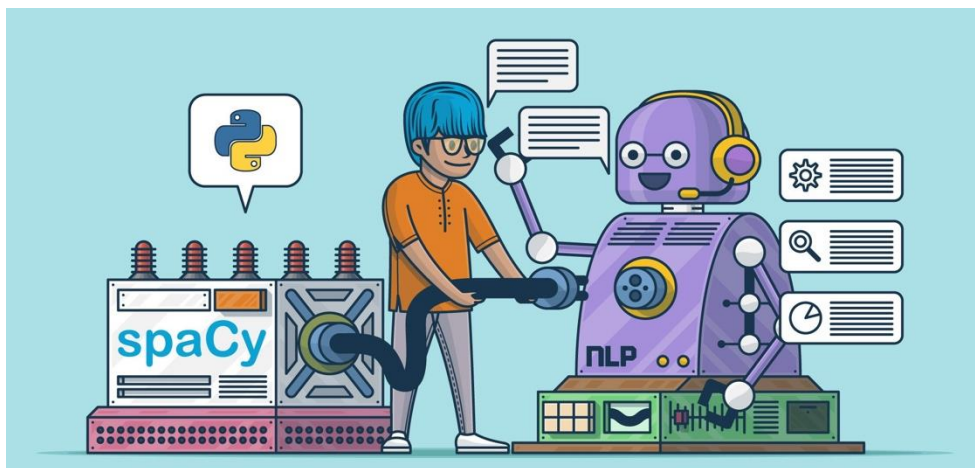
<https://juliasilge.com/>



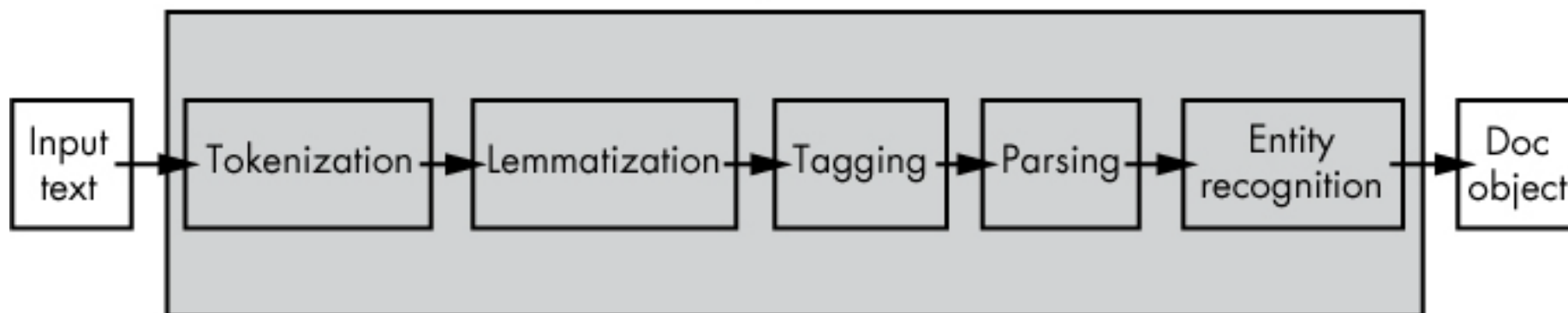
Planen



Værktøj



Pipeline



Advarsel



Hvad skete der med Biblioteket i Alexandria?

©Getty Images

Det Store Bibliotek i Alexandria var en af de mest ambitiøse intellektuelle bedrifter i menneskets historie – et sted hvor viden ingen grænser havde, og jagten på visdom blev betragtet som hellig. Oprettet i det 3. århundrede f.v.t. i den blomstrende by Alexandria i Egypten, var det langt mere end blot et bibliotek. Det var et centrum for læring, et kulturelt samlingspunkt og et fyrtårn for oplysning, der tiltrak datidens største tænkere, videnskabsfolk og filosoffer fra hele den kendte verden.

En vision om fremtiden

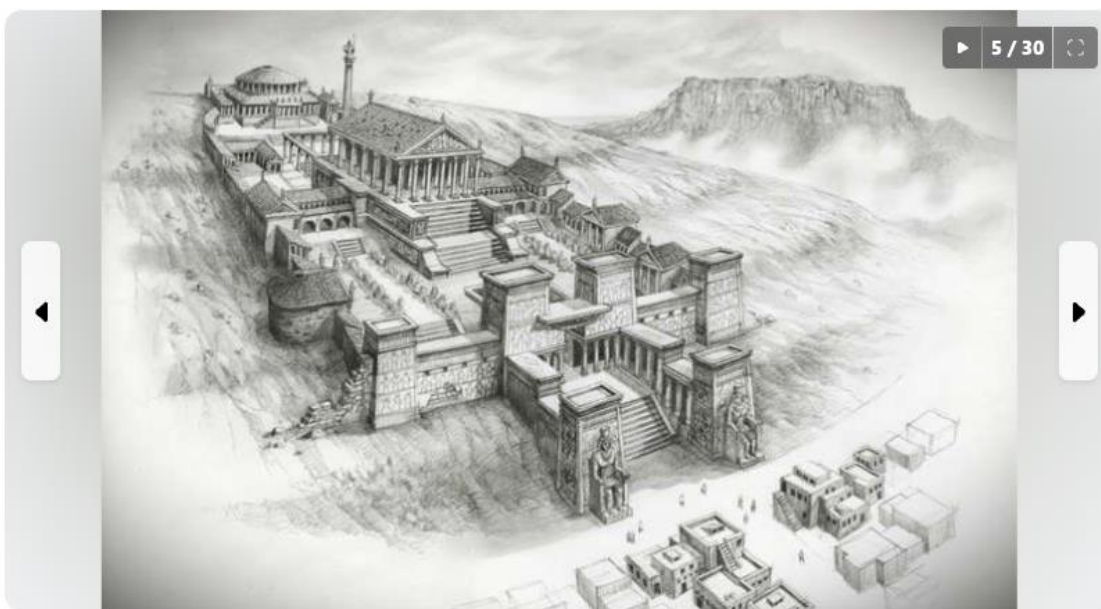
©Getty Images

Med udgangspunkt i Alexander den Stores drøm om en verden forenet gennem viden, forestillede Ptolemaios I og senere Ptolemaios II sig et storslået sted, hvor al menneskelig viden kunne samles, bevares og udvikles af tidens skarpeste hjerner.

Bøger ved kajen

©Getty Images

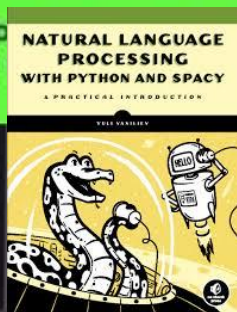
Under Ptolemaios III blev det lov, at alle skibe, der lagde til i Alexandrias havn, skulle aflevere eventuelle bøger. Skrivere lavede kopier, originalerne blev beholdt, og kopierne røg tilbage ombord. Sådan voksede samlingen konstant.



Musernes hellige hjem

©Public Domain

Biblioteket var en del af Mouseion – opkaldt efter muserne, de ni gudinder for kunst og viden. Med op mod en halv million skrifter blev det et magnetisk centrum for lærde fra hele den antikke verden.



Planen Tid

Uger	Dato	Fag	Tidspunkt	Indhold
Uge 18				
	30/4	NLP	Formiddag	Intro til valgfaget - 8.30 til 10.00
	1/5	NLP - Python	Formiddag	Intro til Python + regulære udtryk
	2/5	NLP - Python	Formiddag	"NLP" kapitel 2
Uge 19	5/5	NLP - Python	Eftermiddag	
	6/5	NLP - Python	Eftermiddag	
	8/5	NLP - R	Eftermiddag	
	9/5	NLP - R	Formiddag (OL)	
Uge 20	12/5		Formiddag	
	13/5		Formiddag	
	15/5		Formiddag	
	16/5		Formiddag	
Uge 21	19/5		Eftermiddag	Bigrams
	20/5		Eftermiddag	Bigrams
	22/5		Eftermiddag	
	23/5		Formiddag (OL)	
Uge 22	26/5	Deep Learning	Formiddag	
	27/5	Deep Learning	Formiddag	
	28/5			
	29/5	Fri		
	30/5			
Uge 23	2/6		Eftermiddag	
	3/6		Eftermiddag	
	4/6			
	5/6		Eftermiddag	
	6/6		Formiddag	
Uge 24	9/6	Pinse		
	10/6	Eksamen		
	11/6	Eksamen		

Planen Indhold

2

THE TEXT-PROCESSING PIPELINE

Setting Up Your Working Environment

Installing Statistical Models for spaCy

Basic NLP Operations with spaCy

Tokenization

Lemmatization

Applying Lemmatization for Meaning Recognition

Part-of-Speech Tagging

Using Part-of-Speech Tags to Find Relevant Verbs

Context Is Important

Syntactic Relations

Try This

Named Entity Recognition

Summary

1. The Tidy Text Format

Contrasting Tidy Text with Other Data Structures

The `unnest_tokens` Function

Tidying the Works of Jane Austen

The `gutenbergr` Package

Word Frequencies

Summary

2. Sentiment Analysis with Tidy Data

The sentiments Dataset

Sentiment Analysis with Inner Join

Comparing the Three Sentiment Dictionaries

Most Common Positive and Negative Words

Wordclouds

Looking at Units Beyond Just Words

Summary

3. Analyzing Word and Document Frequency: tf-idf

Term Frequency in Jane Austen's Novels

Zipf's Law

The `bind_tf_idf` Function

A Corpus of Physics Texts

Summary

4. Relationships Between Words: N-grams and Correlations

Tokenizing by N-gram

Counting and Filtering N-grams

Analyzing Bigrams

Using Bigrams to Provide Context in Sentiment Analysis

Visualizing a Network of Bigrams with `ggraph`

Visualizing Bigrams in Other Texts

Counting and Correlating Pairs of Words with the `widyr` Package

Counting and Correlating Among Sections

Examining Pairwise Correlation

Summary

5. Converting to and from Nontidy Formats

Tidying a Document-Term Matrix

Tidying `DocumentTermMatrix` Objects

Tidying `dfm` Objects

Casting Tidy Text Data into a Matrix

Tidying Corpus Objects with Metadata

Example: Mining Financial Articles

Summary

6. Topic Modeling

Latent Dirichlet Allocation

Word-Topic Probabilities

Document-Topic Probabilities

Example: The Great Library Heist

LDA on Chapters

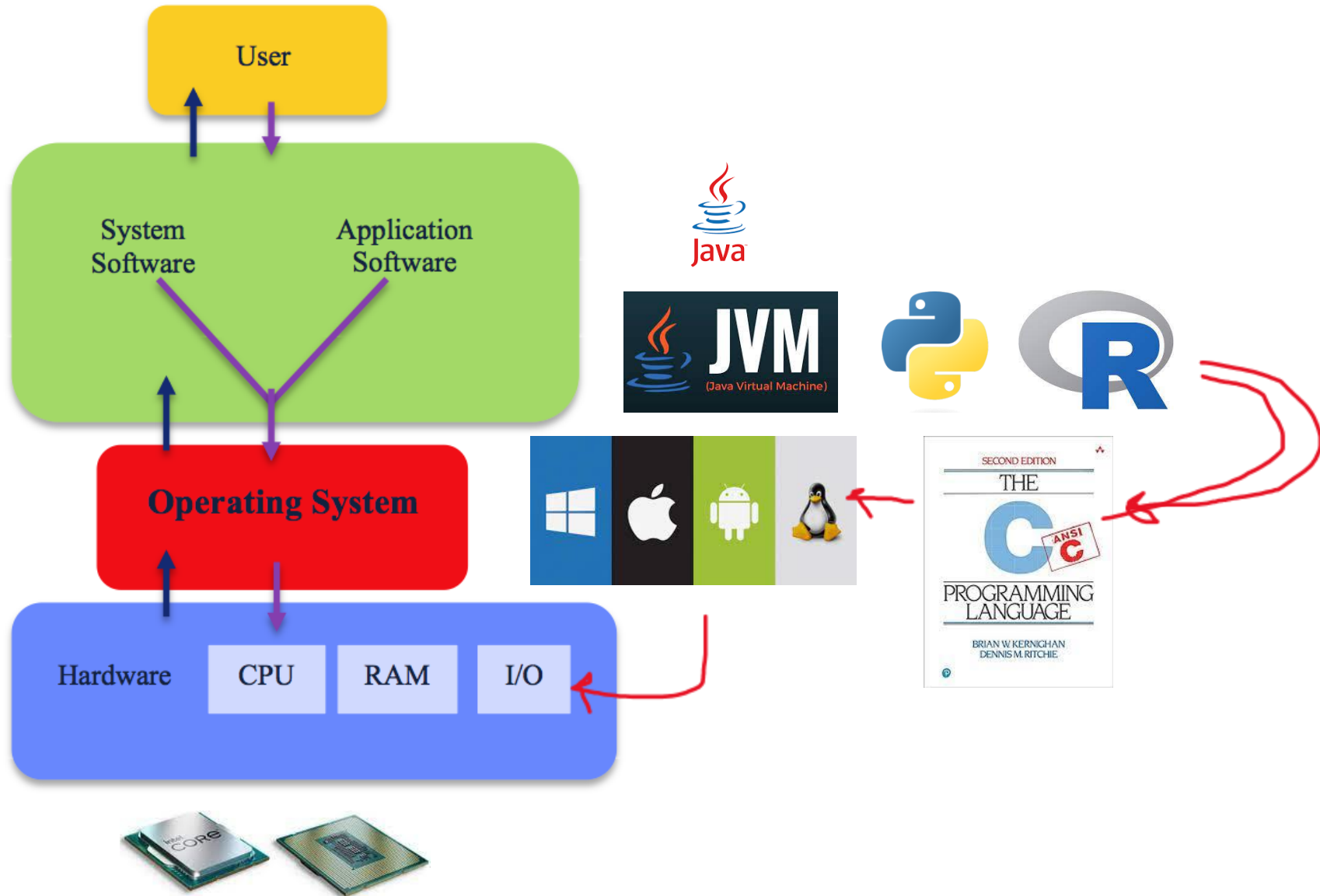
Per-Document Classification

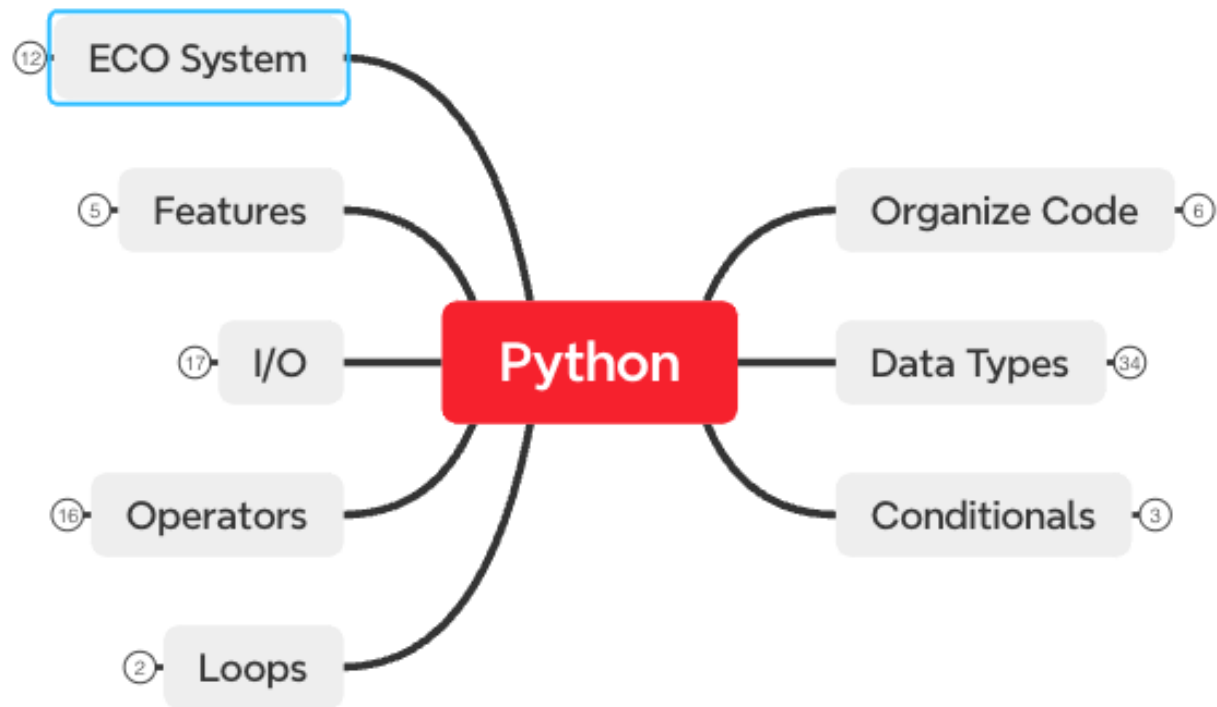
By-Word Assignments: `augment`

Alternative LDA Implementations

Summary

Python



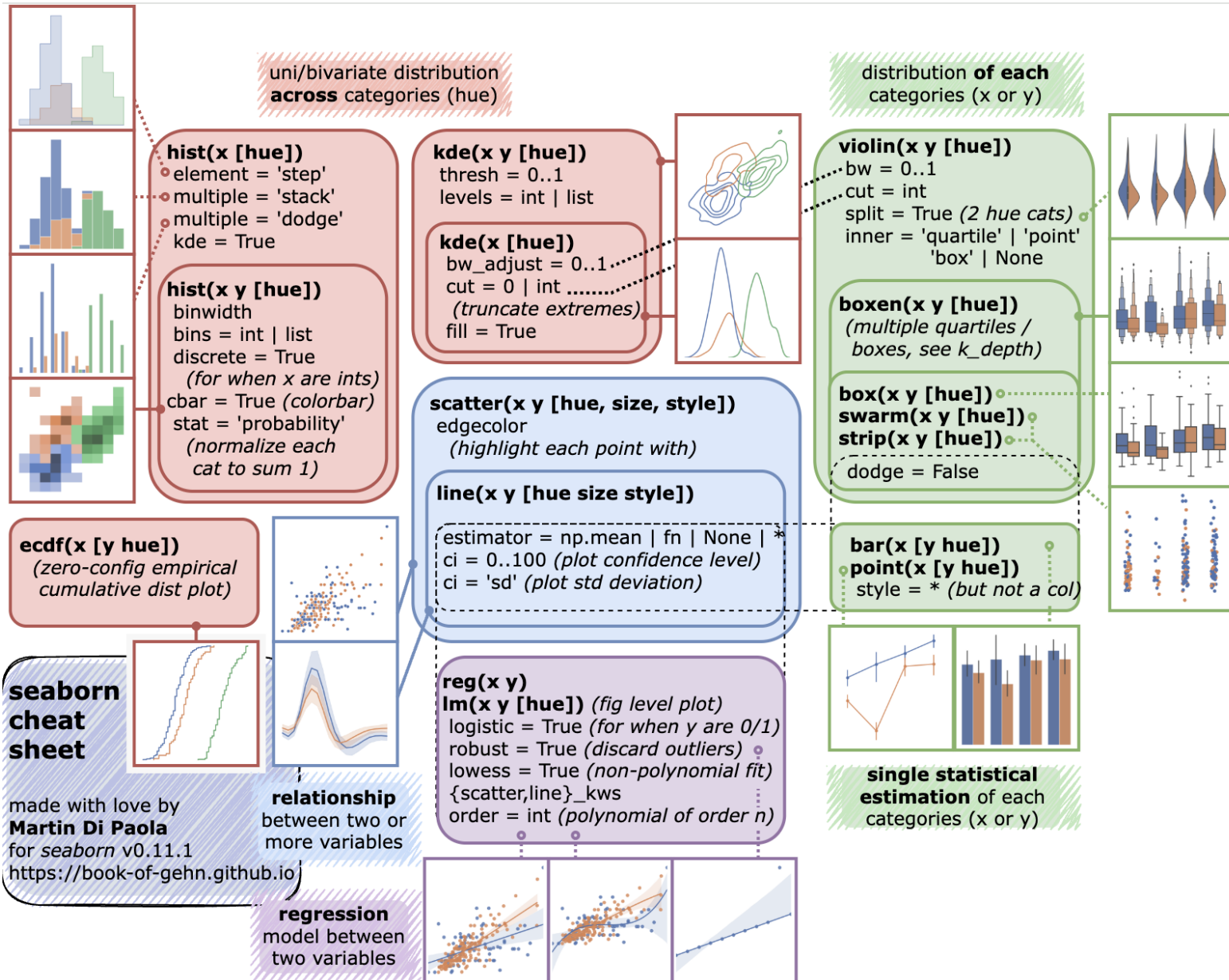


and	A logical operator
as	To create an alias
assert	For debugging
break	To break out of a loop
class	To define a class
continue	To go to the next iteration of a loop
def	To define a function
del	To delete an object
elif	A conditional statements, like else if
else	A conditional statements
except	Used with exceptions, what to do when an exception occurs
False	Boolean value
finally	Used with exceptions, will be executed no matter if there is an exception or not
for	To create a for loop
from	To import specific parts of a module
global	To declare a global variable
if	To make a conditional statement
import	To import a module
in	To check if a value is in a list, tuple
is	To test if two variables are equal
lambda	To create an anonymous function
None	Represents a null value
nonlocal	To declare a non-local variable
not	A logical operator
or	A logical operator
pass	A statement that will do nothing (null)
raise	To raise an exception
return	To exit a function and return a value
True	Boolean value
try	To make a try...except statement
while	To create a while loop
with	Used to simplify exception handling
yield	To end a function, returns a generator

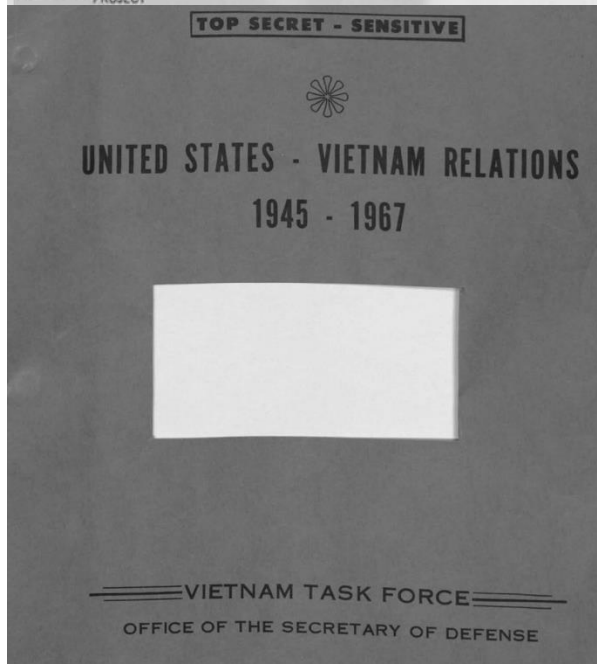
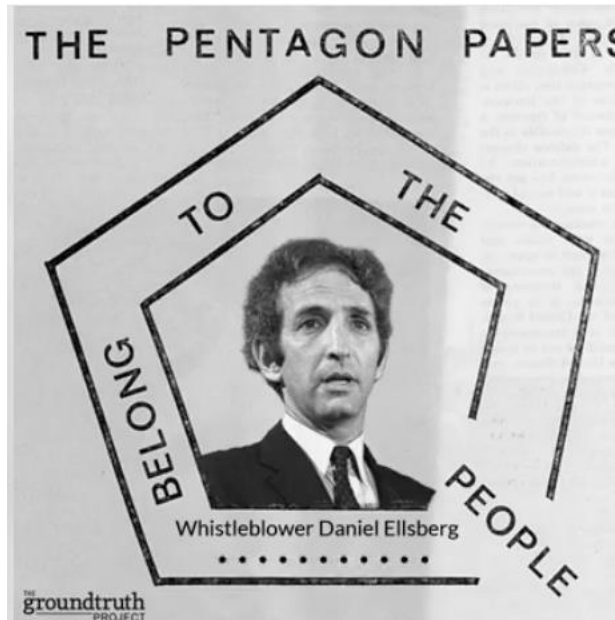
		Built-in Functions		
<code>abs()</code>	<code>divmod()</code>	<code>input()</code>	<code>open()</code>	<code>staticmethod()</code>
<code>all()</code>	<code>enumerate()</code>	<code>int()</code>	<code>ord()</code>	<code>str()</code>
<code>any()</code>	<code>eval()</code>	<code>isinstance()</code>	<code>pow()</code>	<code>sum()</code>
<code>basestring()</code>	<code>execfile()</code>	<code>issubclass()</code>	<code>print()</code>	<code>super()</code>
<code>bin()</code>	<code>file()</code>	<code>iter()</code>	<code>property()</code>	<code>tuple()</code>
<code>bool()</code>	<code>filter()</code>	<code>len()</code>	<code>range()</code>	<code>type()</code>
<code>bytearray()</code>	<code>float()</code>	<code>list()</code>	<code>raw_input()</code>	<code>unichr()</code>
<code>callable()</code>	<code>format()</code>	<code>locals()</code>	<code>reduce()</code>	<code>unicode()</code>
<code>chr()</code>	<code>frozenset()</code>	<code>long()</code>	<code>reload()</code>	<code>vars()</code>
<code>classmethod()</code>	<code>getattr()</code>	<code>map()</code>	<code>repr()</code>	<code>xrange()</code>
<code>cmp()</code>	<code>globals()</code>	<code>max()</code>	<code>reversed()</code>	<code>zip()</code>
<code>compile()</code>	<code>hasattr()</code>	<code>memoryview()</code>	<code>round()</code>	<code>__import__()</code>
<code>complex()</code>	<code>hash()</code>	<code>min()</code>	<code>set()</code>	
<code>delattr()</code>	<code>help()</code>	<code>next()</code>	<code>setattr()</code>	
<code>dict()</code>	<code>hex()</code>	<code>object()</code>	<code>slice()</code>	
<code>dir()</code>	<code>id()</code>	<code>oct()</code>	<code>sorted()</code>	

Function	Description	Syntax
compile	Returns a Regex pattern object	re.compile(pattern)
findall	Returns a list containing all the matches	re.findall(pattern, text)
match	Returns a match object if there is a match at the 0th position	re.match(pattern, text)
search	Returns a match object if there is a match anywhere in the string	re.search(pattern, text)
split	Returns a split where the string has been split at each match	re.split("separator", text)
sub	Replaces one or many matches with a string	re.sub(old_value, new_value, text)

Method / Property	Description	Example
<code>m.group()</code>	Returns the matched string	<code>'o'</code>
<code>m.start()</code>	Start index of the match	<code>0</code>
<code>m.end()</code>	End index of the match	<code>1</code>
<code>m.span()</code>	Tuple of (start, end) indices	<code>(0, 1)</code>
<code>m.re</code>	The compiled regular expression object	<code>re.compile('^o')</code>
<code>m.string</code>	The original string	<code>'otto'</code>



Projektet



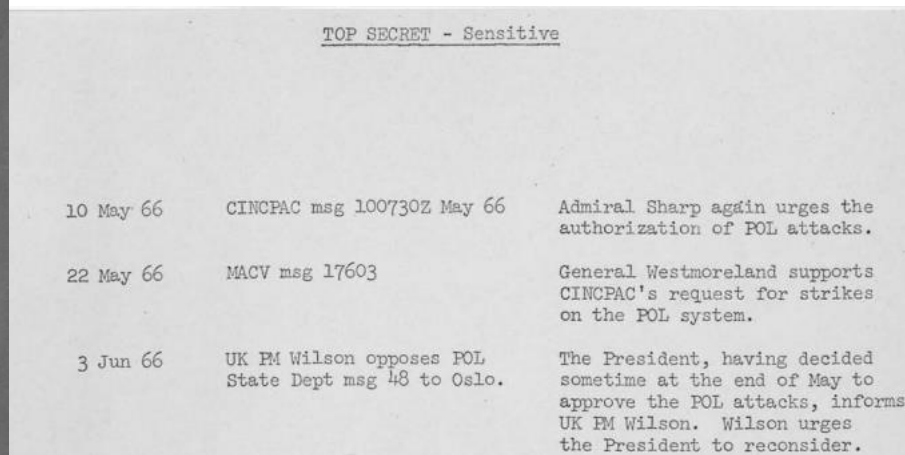
- <https://www.archives.gov/research/pentagon-papers>
- Download en pdf

TERMINALEN

- pdfseparate Pentagon-Papers-Part-IV-C-7-b.pdf page-%d.pdf
- ls *.pdf | while read x; do convert -density 300 \$x \$x.png; done
- ls *.png | while read x; do tesseract \$x \$x; done
- **LØS SORTERINGSPROBLEMET vha Python**
- ls *.txt | while read x; do cat \$x >> out.txt; done

SPYDER

- Indlæs out.txt i Spyder



LØS SORTERINGSPROBLEMET vha Python

```
ls -l | sort -n -k 5  
page-1.pdf  
page-10.pdf  
page-100.pdf  
page-101.pdf  
page-102.pdf  
page-103.pdf
```