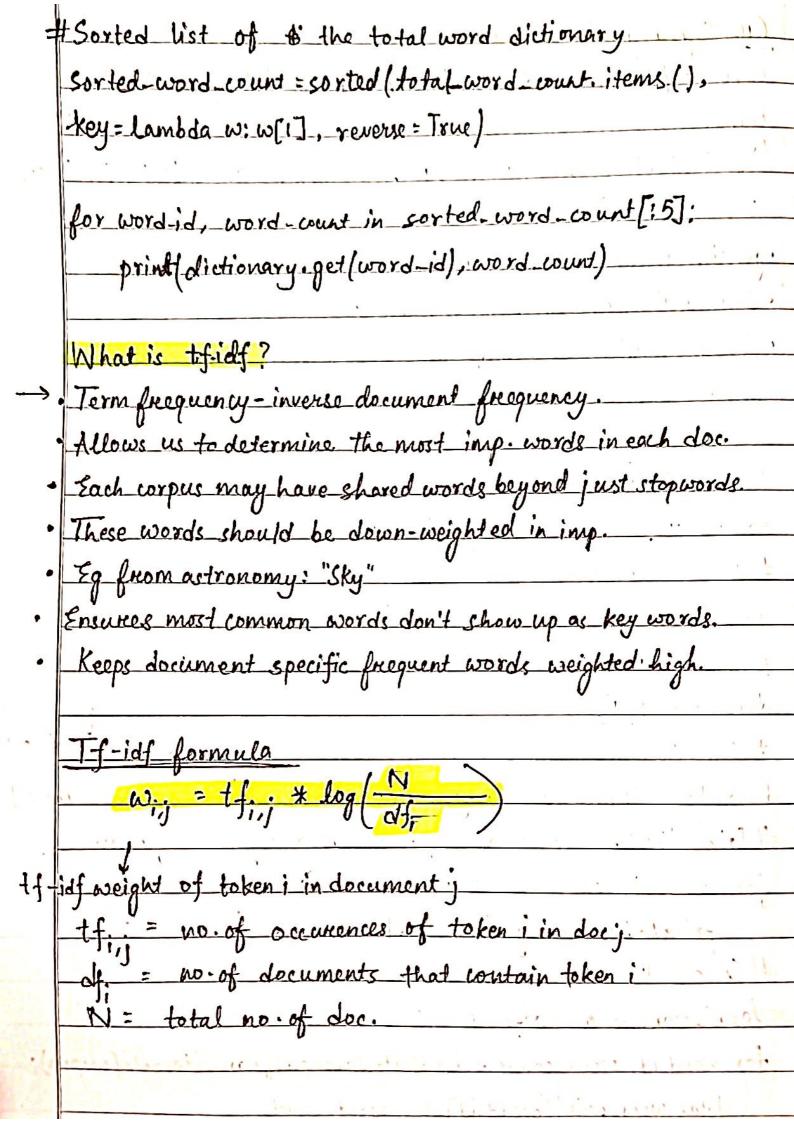
	In	troduction to Nodu	unal Languag	Page	,
09/20	20	troduction to Nodu Processing by Dat	acamp	rago	
	Regular express			1 11 11 11	
	Strings with a			white the	
	Allows us to m	ratch patterns in othe	n strings.	qui inte	
	Applications:		\ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \	- Mening	
		links in a document		- ini	
		addresses, remove		anted characte	us.
	import re		\ \ \	0	
		', 'abrdef') spai	1:(0,3), match	='abc'	
	word_regex =			1	
	, X	d_regex, 'hi thene!')	9		
	L> span=1	0,2) motch='hi'	1 4		
1			the set to see to		
	Pattern	Modelies Example	nple		
			lagic'		
	14	digit)		
3	10	space.	1	10	1
	*	1	ername74'	.1	
	+ or *		,		,
	10	greedy moth ac			
	<i>[</i> -, 7]		-spaces'		
	[a-7]	lowercase a	bedefg	U I'M U	
	~ "\ r *\"		·	1	7
	γ"\[.*]"	Search anything in squa		TILTI CON [ANTO	ripal
	7"[\S]+:"		auader:		

Python's re module	
re module	
split: split a string on regex	
findall: find all patterns in ast	ring
search: search for a pattern.	
match: match an entire string or.	substring based on a pattern.
Eearch for first occurence of 'Ass	harigaa' in
stringl = 'An lel us introduce you to	Antaripa Saha. She is a B. Teu
student. She belongs to Saha famil	- 21
Saha'.	
model = re-search ('cocono Saha', str	ingl)
print(match.start(), match.end())	-> will print the index
, ,	The state of the s
	24 31
pol	ર્ચ 4 31
pol	24 31
Regex groups	24 31
Regex groups OR is represented using !	
Regex groups OR is represented using ! We can define a group using!	
Regex groups OR is represented using !	
Regex groups OR is represented using ! We can define a group using ! We can define explicit character import re	ranger using []
Regex groups OR is represented using ! We can define a group using ! We can define explicit character import re match-digits—and-words = ('())	rangee using [] + \w+)')
Regex groups OR is represented using I We can define a group using I We can define explicit character import re match-digits-and-words = ('()) re find all (match-digits-and-word)	ranger using [] + \w+)') ords, 'He has 11 cats')
Regex groups OR is represented using ! We can define a group using ! We can define explicit character import re match-digits—and-words = ('())	ranger using [] + \w+)') ords, 'He has 11 cats')

use escape de gri	Regex ranges and groups Pattern Matches example [A-Za-z]+ 'ABCDghza' [0-9] [A-Za-z]-\.]+ A-Z, a-z, -,. 'My-Website.com' (a-z) [a-z]
escapeacteril who explains add a char add a char special char	(x+1, -1) spaces or a comma $(x+1, -1)$
	r'(\ω+ #\d \? !)'
	import ultk. tokenize import regexp-tokenize
	Pattern to find hashtags
	hastogs_pattern= r'#/w+'
	print(regenp tokenize (tweets[0], patternt) hastage pattern)
	Pattern = r'[@#]\w+)'
	Bag of words CODE
	from nelk. tokenize import word tokenize from collections import (ounter
	Counter (coord-tokenize (""The cat is in the box. The cal likes the box.
	The bon is over the cet .""")). most-con L) 5::: 3, 'The': 3, 'He:3 the': 2}

Lemmatize CODE
from nitk. stem import WordNetlemmatizer
from nltk. stem import WordNetlemmatizer alpha-only = [+ for t in lower clokers if t. is alpha()]
no-stops = [- for & in alpha-only if & not in english-stops
wordnet_lemmatizer = WordNetlemmatizer()
lemmatized = [wordnet_lemmatizer.lemmatize(+) for t in no-stops]
Create bag of words
bow = (ounter (lemmatized)
print (600. most_common (10)) -> 10 most common tokens.
T. L.
Introduction to gensim Spopular open-source NLP library
Use for building document or word vectors
Performing topic identification and document comparison
Cheating a gension dictionary:
CODE
from genim corpora dictionary import Dictionary
from nltk. tokenize import word-tokenize
my documents = [· · · · · · · · · · · · · · · · · ·
tokenized-docs = [word tokenize(doc.lower()) for doc in my-document distinguised docs)
dictionary = Dictionary (tokenized does) dictionary token 2id #Taking a look at tokens and their ids
Ly {'1': 11, '.': 7, 'a': 2, 'aboud': 14}

	Creating a gencin corpus:
	corpus = [dictionary. doc2bow(doc) for doc in tokenized-dous]
	corpus yist doc
	$\Gamma(0,1), (1,1), (2,1),$
word_id e	· · ·] · · · count
	·gensim models can be easily saved, updated, and reused
	· Our diction ary can also be updated
	· This more advanced and feature sich bag-of-words can be used in
	future exercises.
	CODE
	# Select id for word computer.
	computer id = dictionary get token 2 id get ("computer")
	# Use that id to get the word
	print (dictionary. get (computer_id))
	# Save the 2nd document
	# Sort the doc for frequency sorting the frequency sorting is telement in the list.
	# Sort the doc for frequency 1st element in the U.S.
	bow_doc = sorted (doc, key to = lambda w: w[1], reverse=True)
	# Print the top 5 words of the document along with frequency.
	for word id, word count in bow dow [: 5]:
	print (dictionary get (word id), word Court)
	# (neale the defaultdict
	total_word_count = defaultdid(int)
	for word id, word count in ind itertools chain from iterable (corpus):
	total_word_count[word_id] += word_cound.
	L) Dictionary of total word counts for all documents



	If-idf with gensim
	The state of the s
	from gensim models. Hiftfidfmodel import Tfidf Model
	tfidf = Tfidf Nordal (corpus)
	tfidf[corpus[1]]
	1
	[(0,0.174629), (1,0.174629), (9,6.298531),(10,0.771693),]
	What is Named Enlity Recognition?
•	NLP task to identify important named entities in the text
ti)	People, places, organizations
(ii:	Dates, states, works of art, etc.
•	Can be used alongside topic identification
•	Who? What? Where?
	Stanford Core NLP library supports:
•	Iwlegrated with Rython via Bulth
•	Java based core-ference
•	Support for NER as well as conference and dependency trees.
	CODE
	import nltk
	sentence = "In New York, I like to ride "17
	tokenized_sent = nltk.word_tokenize(sentence)
	tagged cent = nltk.pos_tago (tokenized_sent)
	tagged-sent [:3] - [('In', 'IN'), ('New', 'NNP), ('York', NNP)

	A COMMENT
# To print in tree form	
print (nuk. no-churk (tagged-sent))	
named entity	
NER in hltk	
CODE	
sentences = nllk, sent_tokenize (article)	
tokeni_sentenses = [nlb.word_tokenize(t) for [m	sentences)
# lag each sentence into part of species	
pos sentences = [nltk.pos_tag (sent) for sent in token	i-seniences
# Creale name entity chunks	
chunked_sentences = nltk.ne_chunk_sents (pos_senten	ces, binary=184
# Test for slems of the -live with 'NE' togs	
for sent in chunked-sentences;	
for chunk in sent:	
if hasaltr (churk, 'label') and churk. label ():	= 'NE':
privi(chunk)	
priva (opt word)	
"nasalla() to determine if each chunk has a 'label' and	9
then simply use the chunk's label () method as	
the dictionary key.	
0 0	

× 12

CODE

. 1	Egories:
import matphollib.pyplot as plt	egories:
ner_categories = defaultdict (int)	(1) 1.7
for sent in chunked-sentences:	N. L.
for chunk in sent:	
if hasalt's (chunk, 'label'):	
ner_categories[chunk.label()] +=1	No. 111
# Create a list from the dictionary keys for the	charts labels
labels = list(nex_categories.keys())	•
# Create a list of the values	1.
Values = [new categories get(v) for v in labels]	
# Create the pie chart	1
pt-pie (values, labels=labels, autopet="/.1.15/	1.1, startangle=140
plt-show()	100
	1
What is Spacy?	
NLP library cimilar to gensin, with diff. impleme	ntations.
Focus on excating NLP pipelines to generate model	
Open-source, with extra libraries and tools	6 17 ×
L-> Displacy	1. 1
1	. 1' 1' :
Entity recognition visualizer	· · · · · · · · · · · · · · · · · · ·

`

Inspecting the vectors	
CODE	and the state of t
# Create the CountVectorizer Datatra	me
MAINT OUT DO DO IN IN I COME	wintrainett, columns -
COL	und vectorizer get-feature nomes ()
I weate the I fid vectorizer Dala Frame	r kfr
tfidf-df: pd. Dataframe (tfidf-train. A	. Columns = Ifidf-vectorizer.
	got-features-names ()
print (count_df.head())	
print (stidf-of-head ())	The state of the s
# Calculdedifference in columns	The second secon
difference = set (count_df.columns) - se	et (Ifidf-df.columns)
# Check whether the dataframes are	equel
prival count of equals (A fidf df))	
	1
Naive Bayes with scikit-learn	
LODE	works well fwith count
from sklearn naive bayes import M	ultinomial NR expects integer in
from sklearn import metrics	for multilabel
Nb classifier = MultinomialNB()	· Jassi fication
nb classifier fit (count train, y train	n)
pred = nb_classifier-predict(count_te	
metrics, accureacy_score(y_test, pre	α.)
# CONFUSION matrix	
metrics.confusion_matrix (y-test, p	red, labels = [o, [])
	Change them acc. to labels
	like 'FAKE' and 'PEAL'
	for fate news classifier

		1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
	# If we use alpha in the classifier for improving of mb nbu classifier = Multinomial NB (alpha)	accuracy
	It if we use alpha in the carry	· + 1
	not non-classifier = Multinomiak N.B (alpha)	100
		0 130
d	Inspecting the model	
	J. Specing	
	CODE	
		100
	# het the class labels	
	class_labols = nb_classifier.classes_	150
	Caess_mores = 10=Ce is i grant	1
	# feature names	
	feature_names = tfidf-vectorizer.get-feature_names() It zip both features names together with coefficient ar	7/1/2
	the the	ray and
	It zip both features names together with coefficient	
	off CAVI Luca Diolett	100
	# sort by weights	llower is a world
	frat-with_weights = sorted (zip (nb-classifier.coef-[0], feat	ma-varies)
	8	100