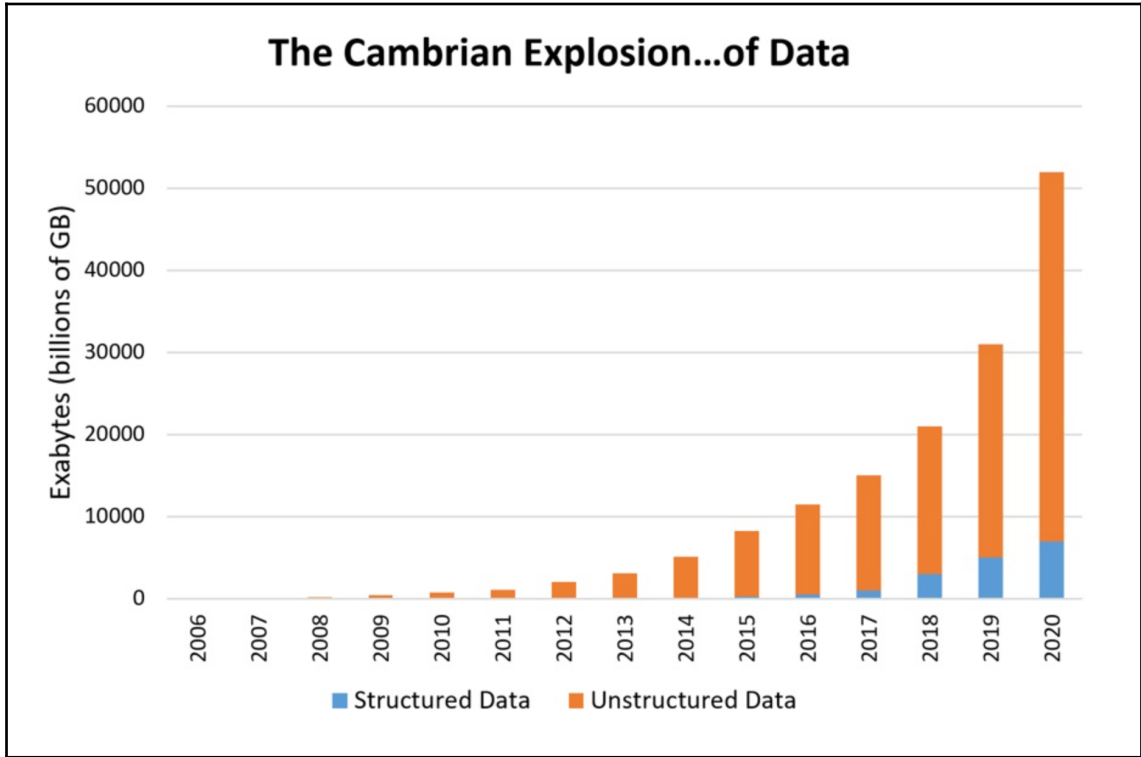
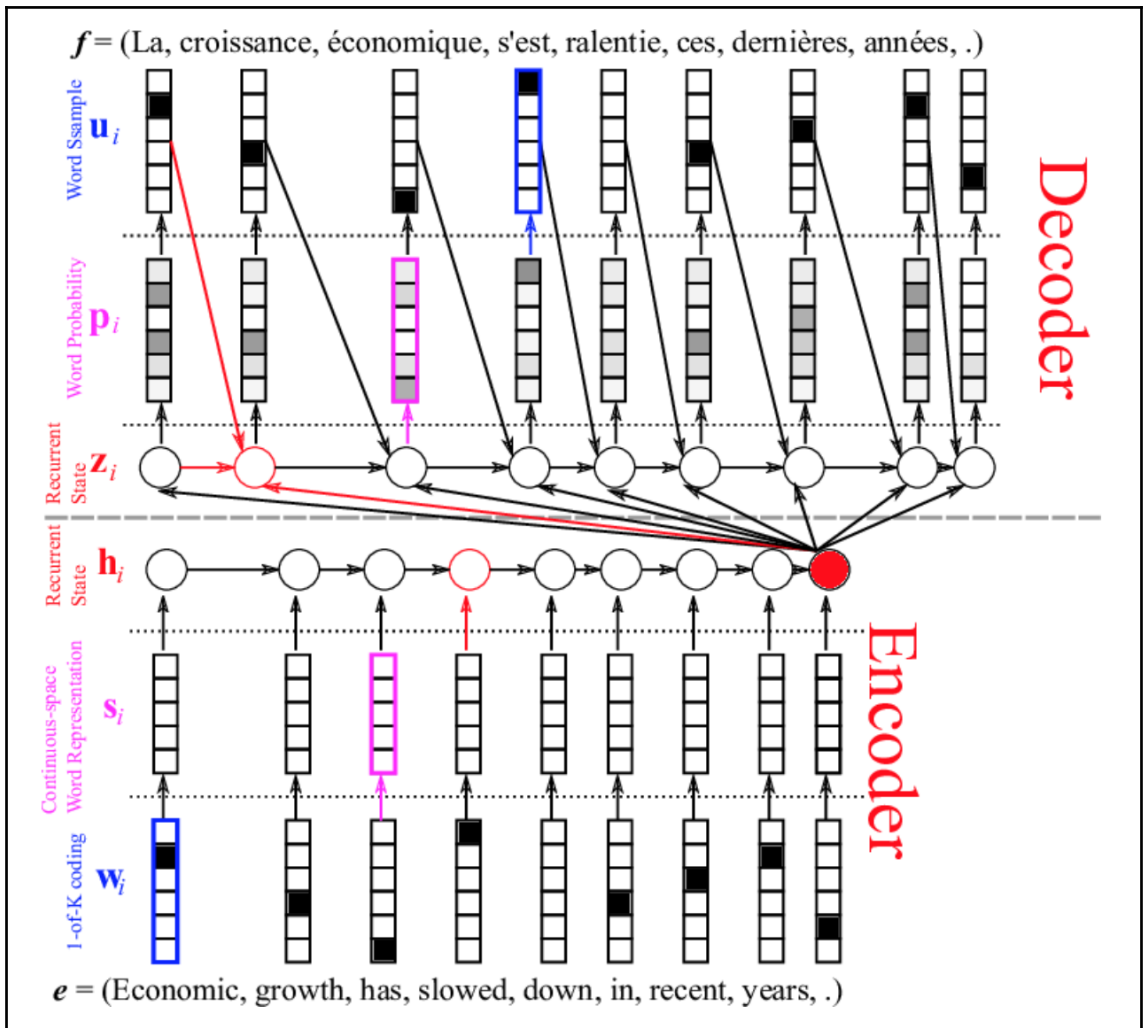
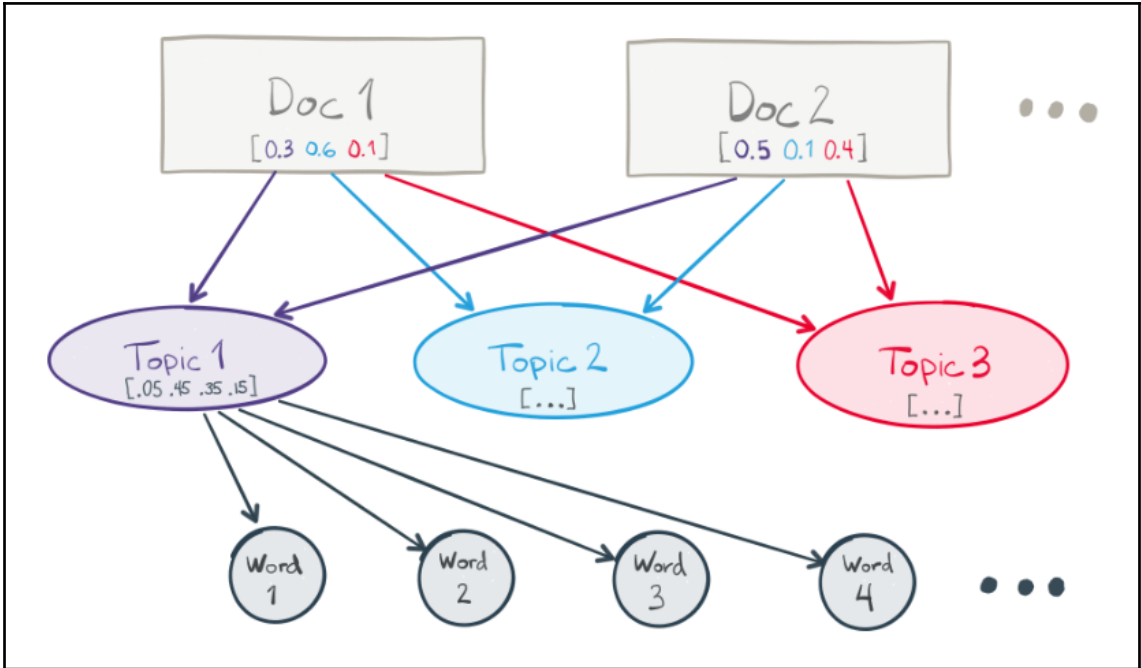


# Chapter 1: What is Text Analysis?



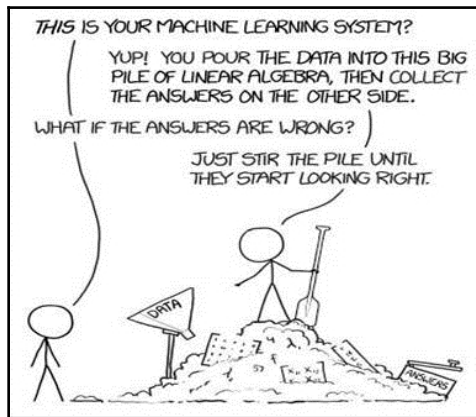




Dataset Name	Brief description	Preprocessing	Instances	Format	Default Task	Created (updated)	Reference	Creator
Amazon reviews	US product reviews from <a href="http://Amazon.com">Amazon.com</a> .	None.	~ 82M	Text	Classification, sentiment analysis	2015	[131]	McAuley et al.
OpinRank Review Dataset	Reviews of cars and hotels from <a href="http://Edmunds.com">Edmunds.com</a> and <a href="http://TripAdvisor">TripAdvisor</a> respectively.	None.	42,230 / ~259,000 respectively	Text	Sentiment analysis, clustering	2011	[132][133]	K. Ganesan et al.
MovieLens	22,000,000 ratings and 580,000 tags applied to 33,000 movies by 240,000 users.	None.	~ 22M	Text	Regression, clustering, classification	2016	[134]	<a href="http://GroupLens.org">GroupLens Research</a>
Yahoo! Music User Ratings of Musical Artists	Over 10M ratings of artists by Yahoo users.	None described.	~ 10M	Text	Clustering, regression	2004	[135][136]	Yahoo!
Car Evaluation Data Set	Car properties and their overall acceptability.	Six categorical features given.	1728	Text	Classification	1997	[137][138]	M. Bohanec
YouTube Comedy Slam Preference Dataset	User vote data for pairs of videos shown on YouTube. Users voted on funnier videos.	Video metadata given.	1,138,562	Text	Classification	2012	[139][140]	Google
Skytrax User Reviews Dataset	User reviews of airlines, airports, seats, and lounges from Skytrax.	Ratings are fine-grain and include many aspects of airport experience.	41396	Text	Classification, regression	2015	[141]	Q. Nguyen
Teaching Assistant Evaluation Dataset	Teaching assistant reviews.	Features of each instance such as class, class size, and instructor are given.	151	Text	Classification	1997	[142][143]	W. Loh et al.

**Twitter and tweets** [ edit ]

Dataset Name ↕	Brief description ↕	Preprocessing ↕	Instances ↕	Format ↕	Default Task ↕	Created (updated) ↕	Reference ↕	Creator ↕
Sentiment140	Tweet data from 2009 including original text, time stamp, user and sentiment.	Classified using distant supervision from presence of emoticon in tweet.	1,578,627	Tweets, comma, separated values	Sentiment analysis	2009	[161][162]	A. Go et al.
ASU Twitter Dataset	Twitter network data, not actual tweets. Shows connections between a large number of users.	None.	11,316,811 users, 85,331,846 connections	Text	Clustering, graph analysis	2009	[163][164]	R. Zafarani et al.
SNAP Social Circles: Twitter Database	Large twitter network data.	Node features, circles, and ego networks.	1,768,149	Text	Clustering, graph analysis	2012	[165][166]	J. McAuley et al.
Twitter Dataset for Arabic Sentiment Analysis	Arabic tweets.	Samples hand-labeled as positive or negative.	2000	Text	Classification	2014	[167][168]	N. Abdulla
Buzz in Social Media Dataset	Data from Twitter and Tom's Hardware. This dataset focuses on specific buzz topics being discussed on those sites.	Data is windowed so that the user can attempt to predict the events leading up to social media buzz.	140,000	Text	Regression, Classification	2013	[169][170]	F. Kawala et al.
Paraphrase and Semantic Similarity in Twitter (PIT)	This dataset focuses on whether tweets have (almost) same meaning/information or not. Manually labeled.	tokenization, part-of-speech and named entity tagging	18,762	Text	Regression, Classification	2015	[171][172]	Xu et al.

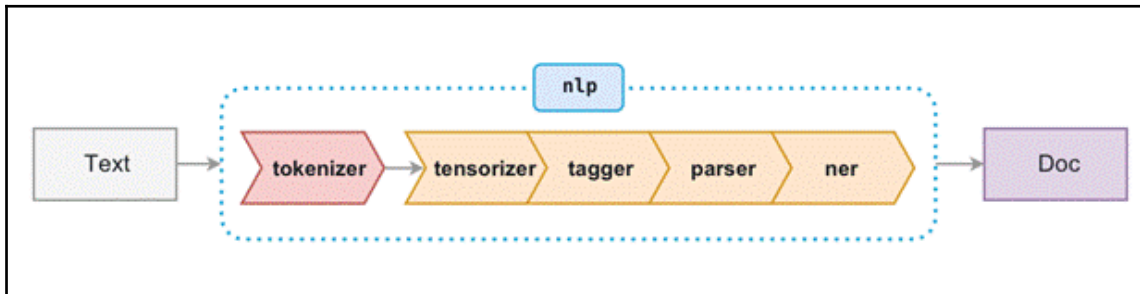


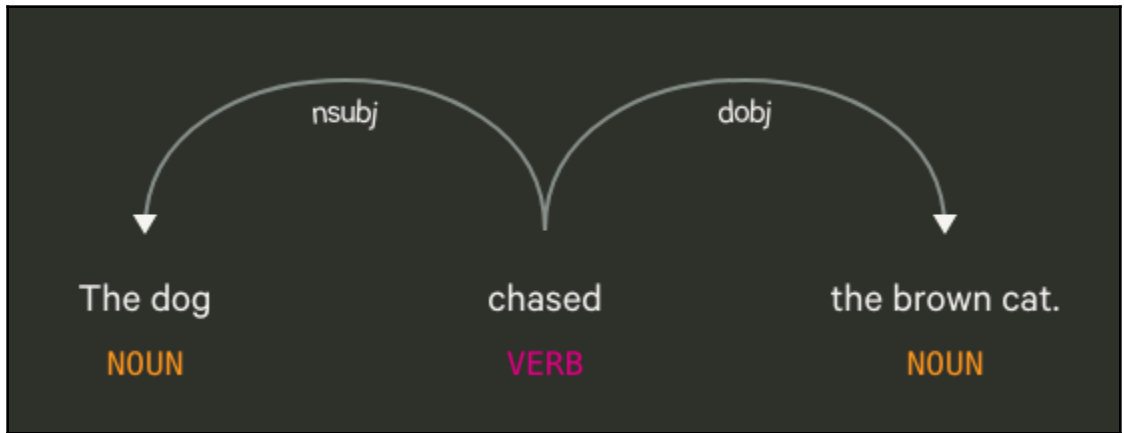
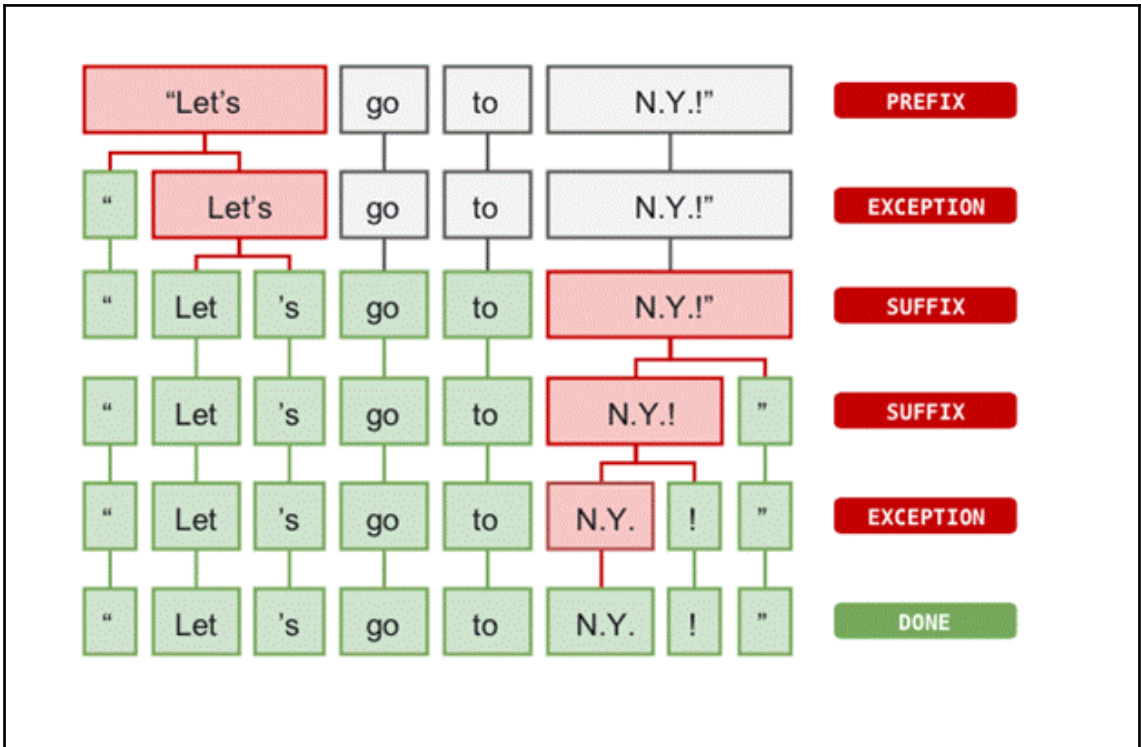
# Chapter 3: spaCy's Language Models

## Feature comparison

Here's a quick comparison of the functionalities offered by spaCy, [SyntaxNet](#), [NLTK](#) and [CoreNLP](#).

	SPACY	SYNTAXNET	NLTK	CORENLP
Programming language	Python	C++	Python	Java
Neural network models	●	●	●	●
Integrated word vectors	●	●	●	●
Multi-language support	●	●	●	●
Tokenization	●	●	●	●
Part-of-speech tagging	●	●	●	●
Sentence segmentation	●	●	●	●
Dependency parsing	●	●	●	●
Entity recognition	●	●	●	●
Coreference resolution	●	●	●	●

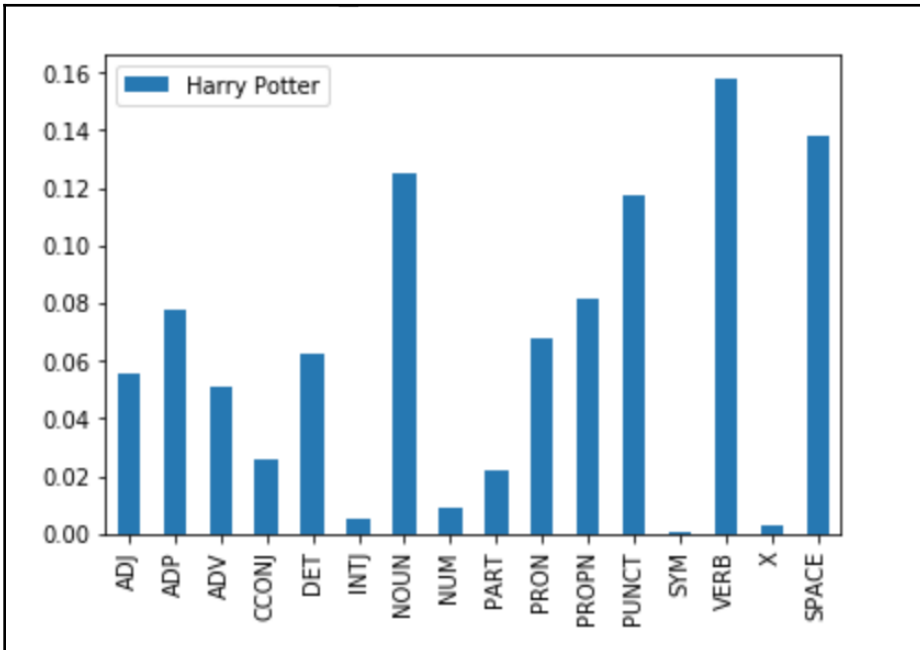
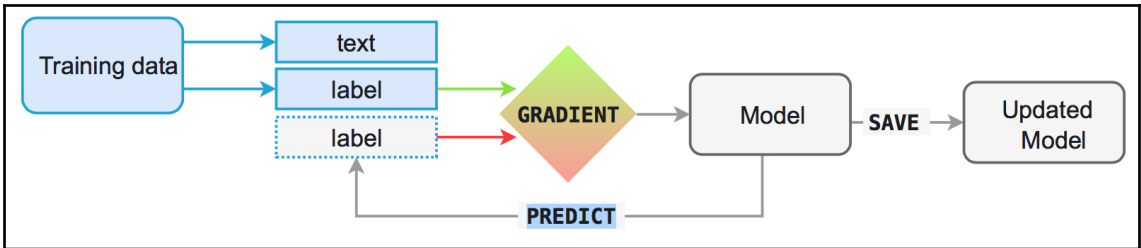
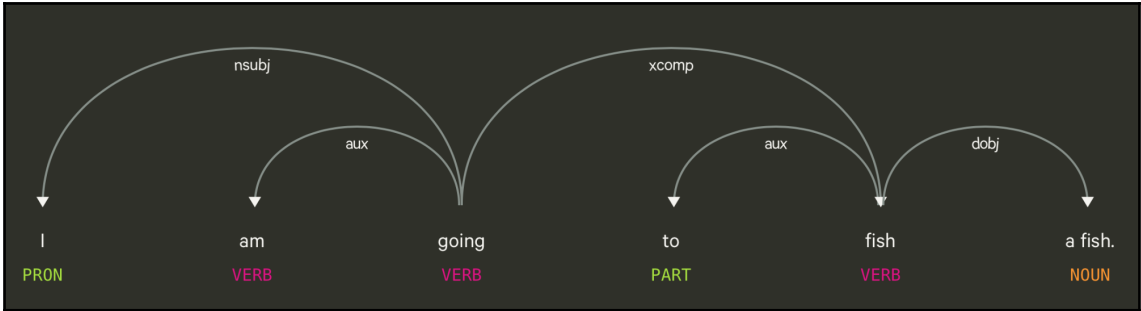




# Chapter 5: POS-Tagging and Its Applications

POS	DESCRIPTION	EXAMPLES
ADJ	adjective	<i>big, old, green, incomprehensible, first</i>
ADP	adposition	<i>in, to, during</i>
ADV	adverb	<i>very, tomorrow, down, where, there</i>
AUX	auxiliary	<i>is, has (done), will (do), should (do)</i>
CONJ	conjunction	<i>and, or, but</i>
CCONJ	coordinating conjunction	<i>and, or, but</i>
DET	determiner	<i>a, an, the</i>
INTJ	interjection	<i>psst, ouch, bravo, hello</i>
NOUN	noun	<i>girl, cat, tree, air, beauty</i>
NUM	numeral	<i>1, 2017, one, seventy-seven, IV, MMXIV</i>
PART	particle	<i>'s, not,</i>
PRON	pronoun	<i>I, you, he, she, myself, themselves, somebody</i>
PROPN	proper noun	<i>Mary, John, London, NATO, HBO</i>
PUNCT	punctuation	<i>., (, ), ?</i>
SCONJ	subordinating conjunction	<i>if, while, that</i>
SYM	symbol	<i>\$, %, \$, ©, +, -, x, ÷, =, :, 😊</i>
VERB	verb	<i>run, runs, running, eat, ate, eating</i>
X	other	<i>sfpkdspsxmsa</i>
SPACE	space	





---

# Chapter 6: NER-Tagging and Its Applications

TYPE	DESCRIPTION
PER	Named person or family.
LOC	Name of politically or geographically defined location (cities, provinces, countries, international regions, bodies of water, mountains).
ORG	Named corporate, governmental, or other organizational entity.
MISC	Miscellaneous entities, e.g. events, nationalities, products or works of art.

TYPE	DESCRIPTION
PERSON	People, including fictional.
NORP	Nationalities or religious or political groups.
FACILITY	Buildings, airports, highways, bridges, etc.
ORG	Companies, agencies, institutions, etc.
GPE	Countries, cities, states.
LOC	Non-GPE locations, mountain ranges, bodies of water.
PRODUCT	Objects, vehicles, foods, etc. (Not services.)
EVENT	Named hurricanes, battles, wars, sports events, etc.
WORK_OF_ART	Titles of books, songs, etc.
LAW	Named documents made into laws.
LANGUAGE	Any named language.
DATE	Absolute or relative dates or periods.
TIME	Times smaller than a day.
PERCENT	Percentage, including "%".
MONEY	Monetary values, including unit.
QUANTITY	Measurements, as of weight or distance.
ORDINAL	"first", "second", etc.
CARDINAL	Numerals that do not fall under another type.

TAG	DESCRIPTION
BEGIN	The first token of a multi-token entity.
IN	An inner token of a multi-token entity.
LAST	The final token of a multi-token entity.
UNIT	A single-token entity.
OUT	A non-entity token.

---

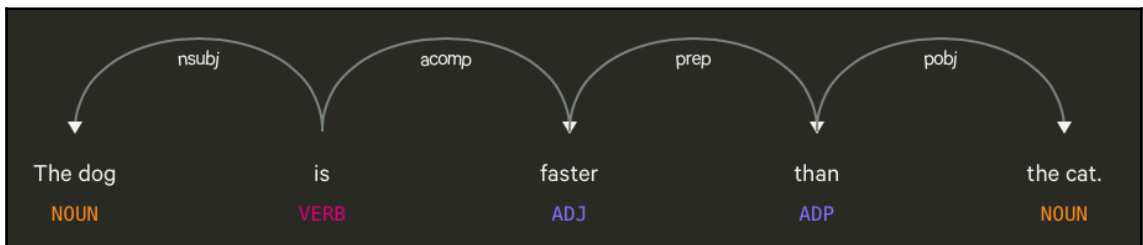
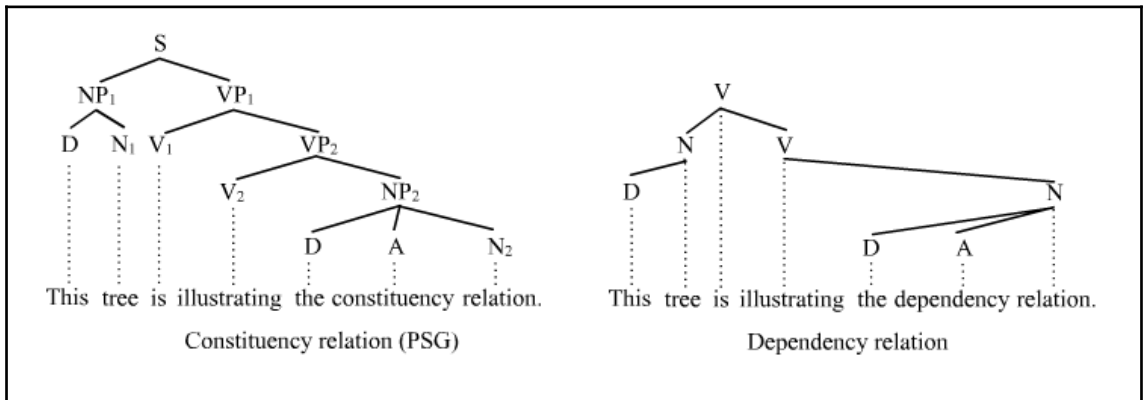
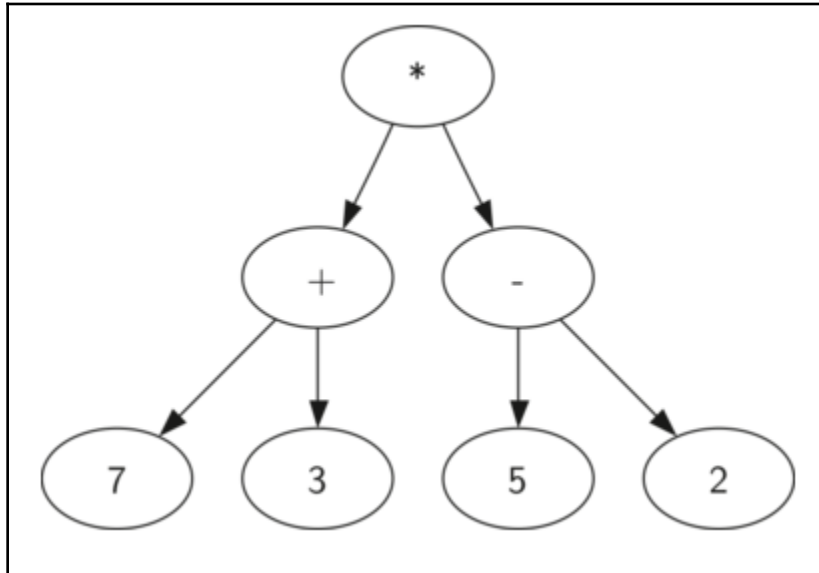
Elon Musk PERSON apparently wasn't aware that his company SpaceX had a Facebook ORG page. The SpaceX and Tesla PRODUCT CEO has responded to a comment on Twitter GPE calling for him to take down the SpaceX, Tesla and Elon Musk ORG official pages in support of the #deletefacebook movement by first ORDINAL acknowledging he didn't know one existed, and then following up with promises that he would indeed take them down.

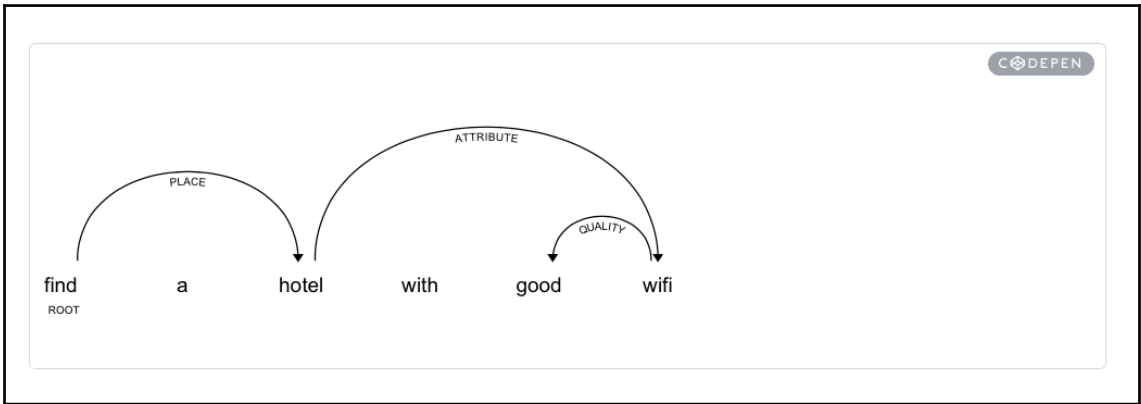
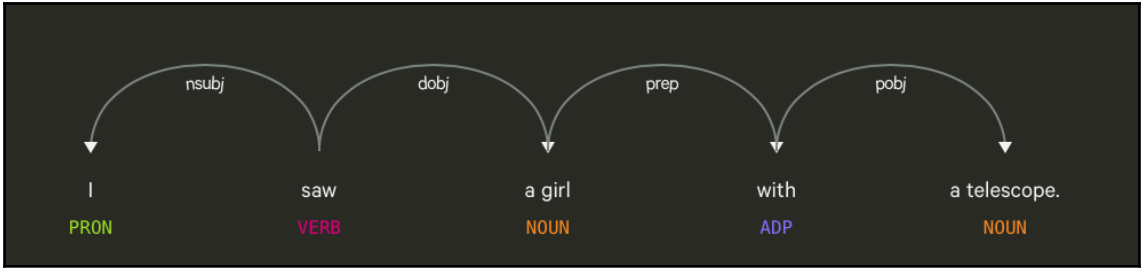
He's done just that, as the SpaceX NORP Facebook page is now gone, after having been live earlier today DATE (as you can see from the screenshot included taken at around 12:10 PM ET TIME).

Emmanuel Jean-Michel Frédéric Macron PERSON is a French NORP politician serving as President of France GPE and ex officio Co-Prince of Andorra LOC since 14 May 2017 DATE.

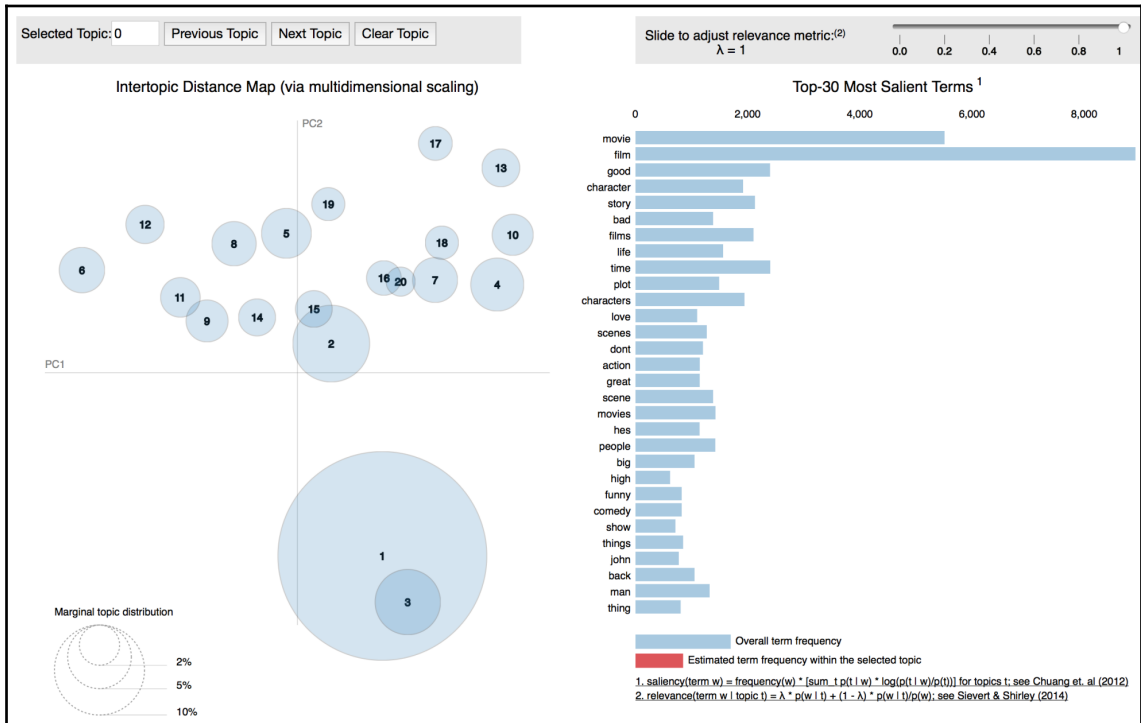
Before entering politics, he was a senior civil servant and investment banker. He studied philosophy at Paris Nanterre University ORG, completed a Master's of Public Affairs ORG at Sciences Po, and graduated from the École nationale d'administration ( PRODUCT ÉNA ORG ) in 2004 DATE. He worked at the Inspectorate General of Finances ORG, and later became an investment banker at Rothschild & Cie Banque ORG.

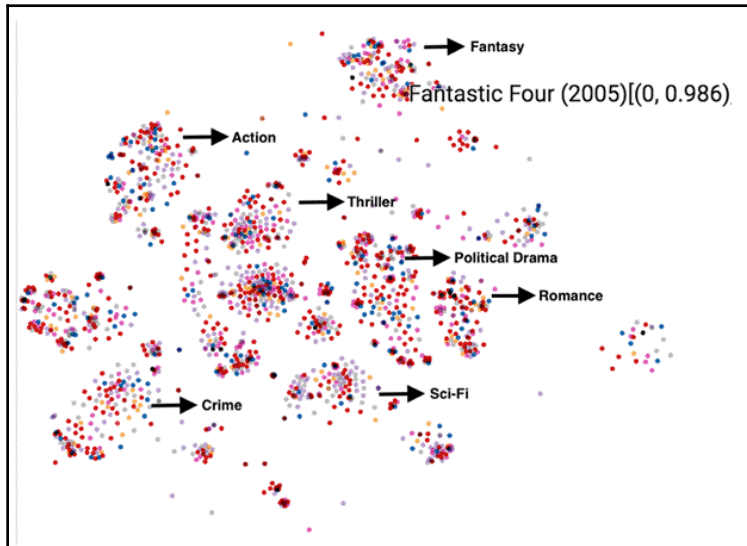
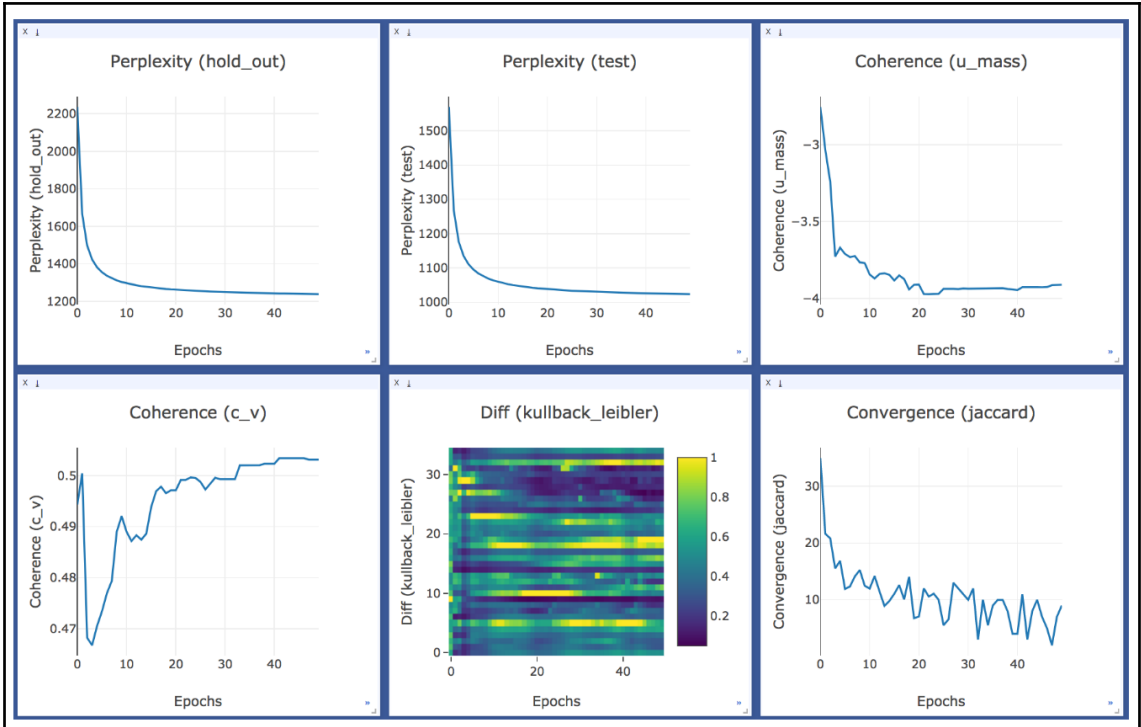
# Chapter 7: Dependency Parsing



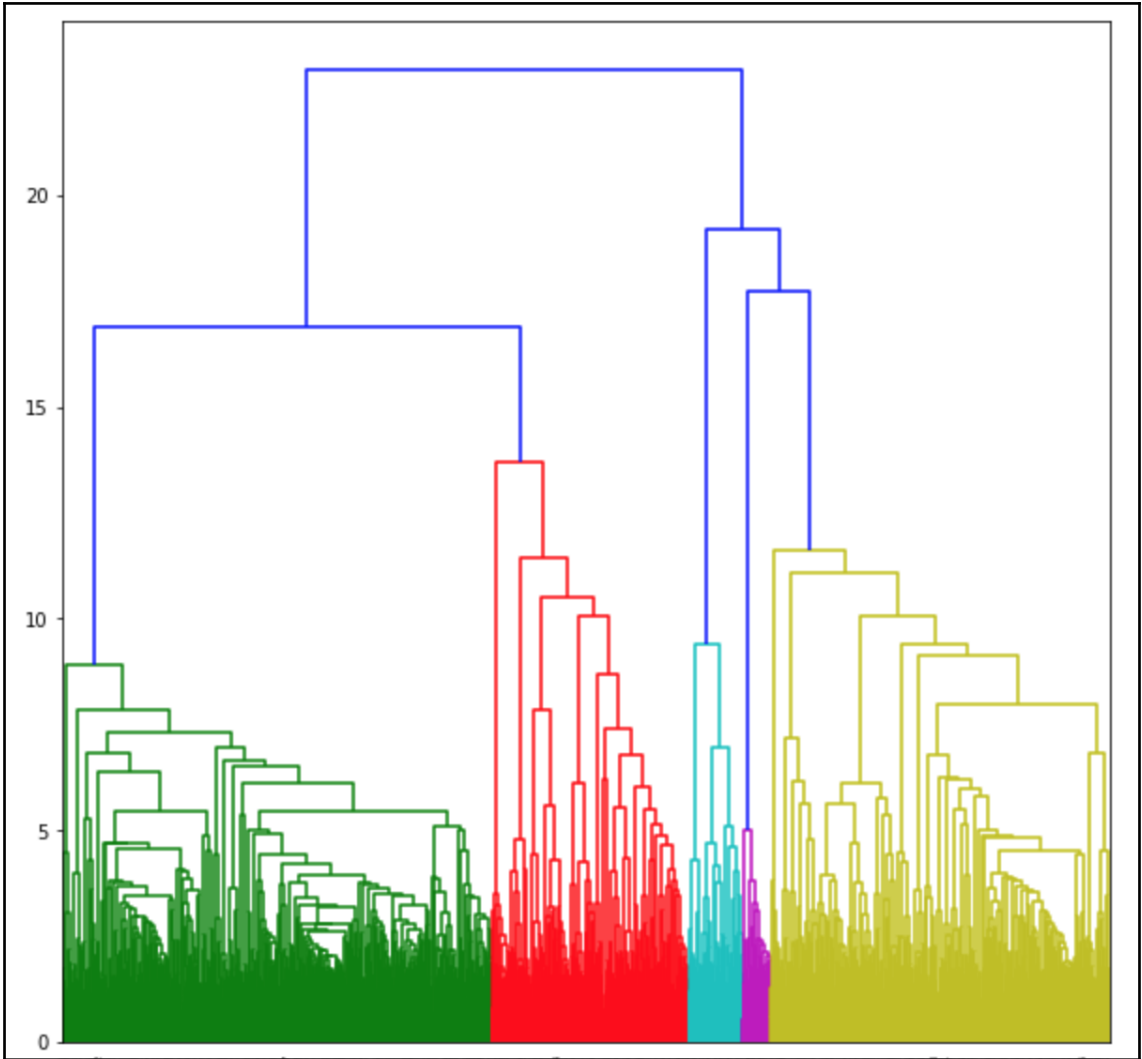


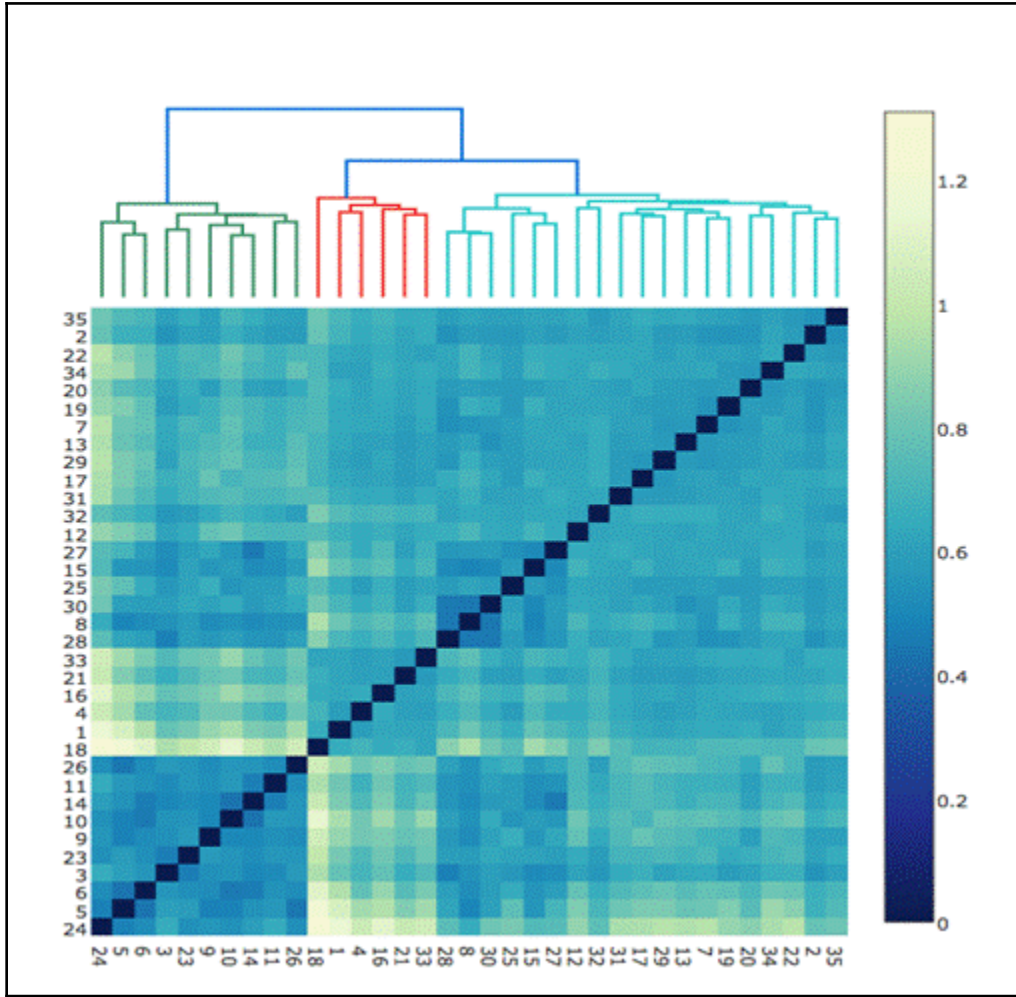
# Chapter 9: Advanced Topic Modeling











---

## Chapter 10: Clustering and Classifying Text

### Visualising Dataset

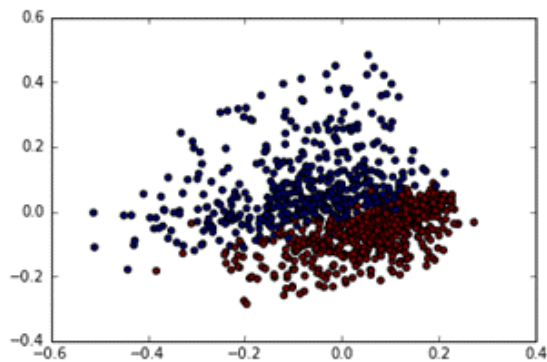
```
In [9]: from sklearn.decomposition import PCA

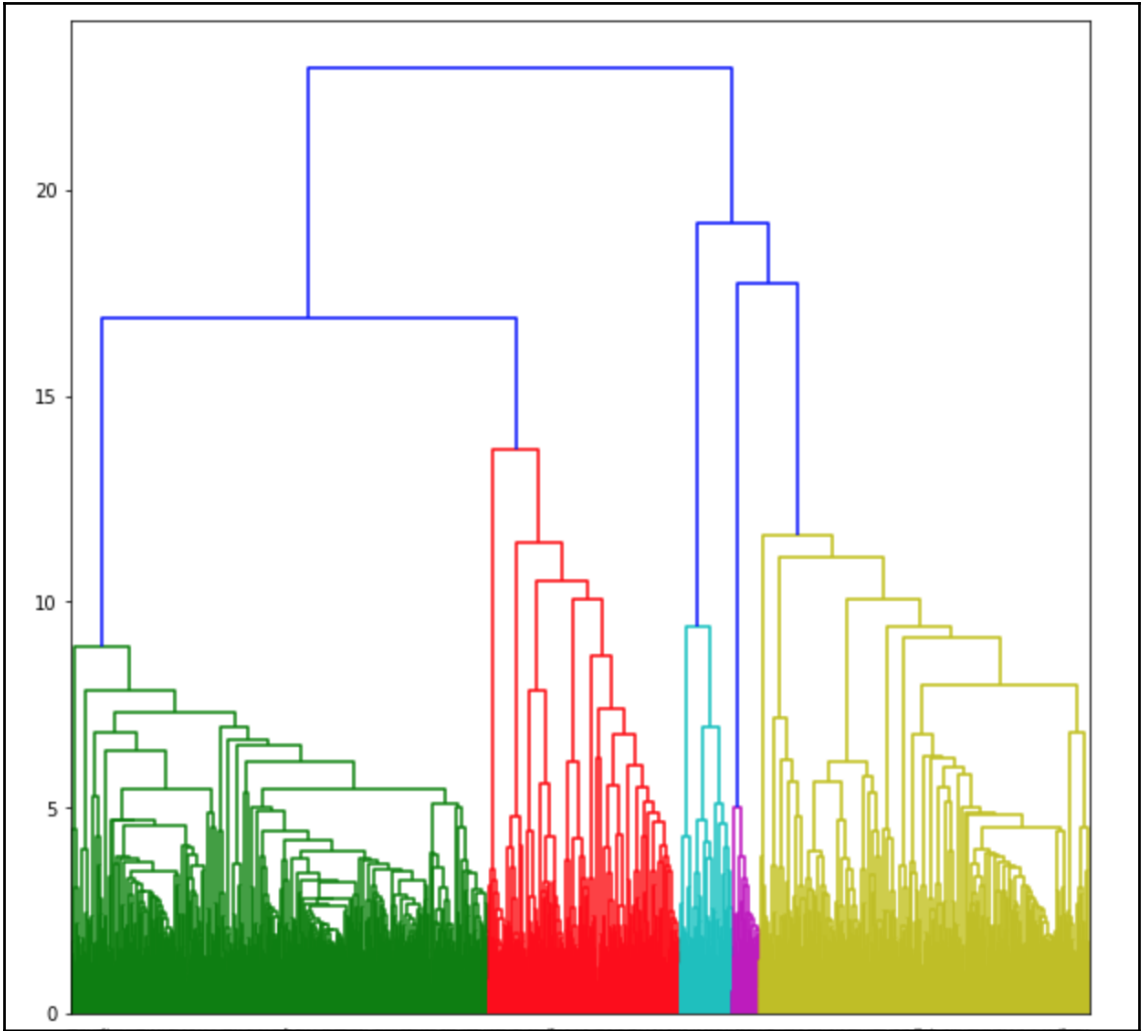
newsgroups_train = fetch_20newsgroups(subset='train',
                                      categories=['alt.atheism', 'sci.space'])

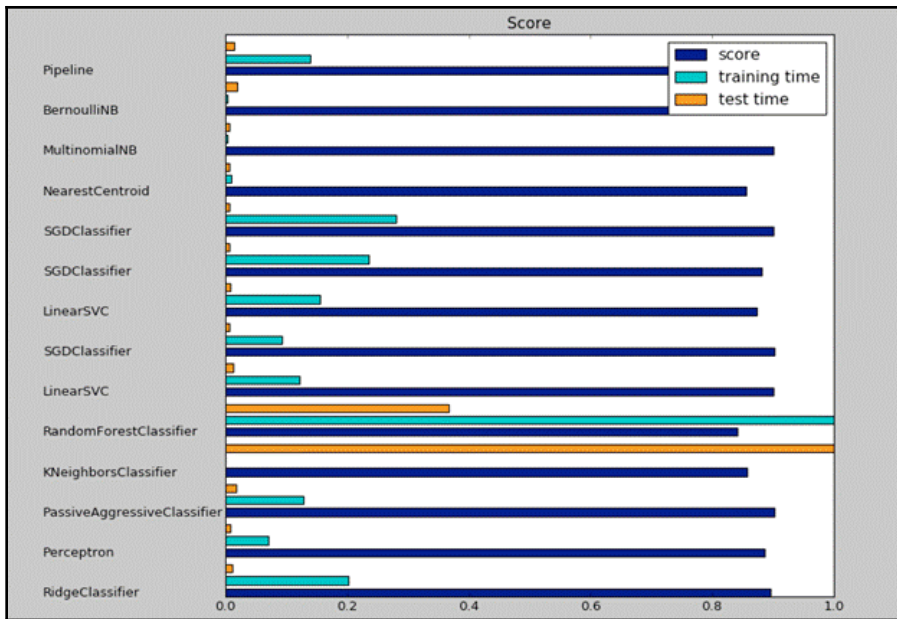
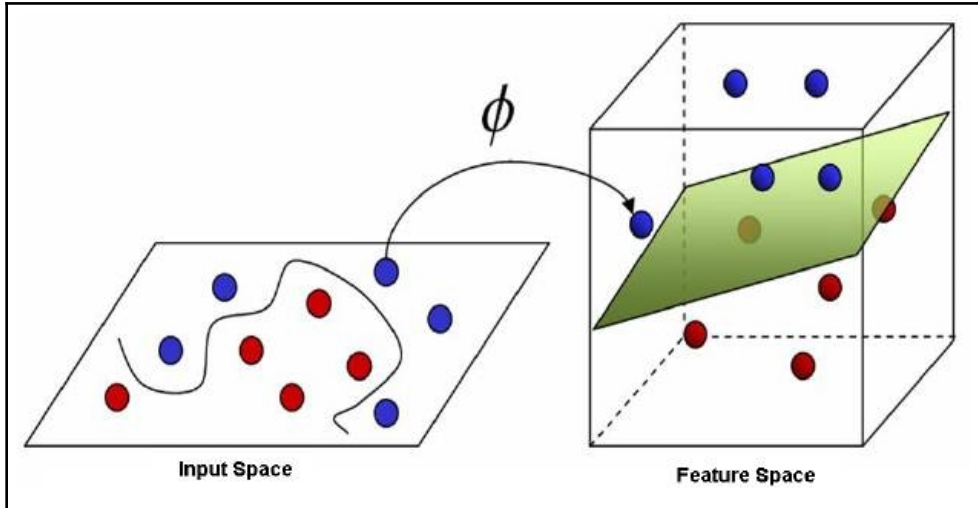
pipeline = Pipeline([
    ('vect', CountVectorizer()),
    ('tfidf', TfidfTransformer()),
])
X = pipeline.fit_transform(newsgroups_train.data).todense()

pca = PCA(n_components=2).fit(X)
data2D = pca.transform(X)
plt.scatter(data2D[:,0], data2D[:,1], c=newsgroups_train.target)
```

Out[9]: <matplotlib.collections.PathCollection at 0x10a7bae90>







---

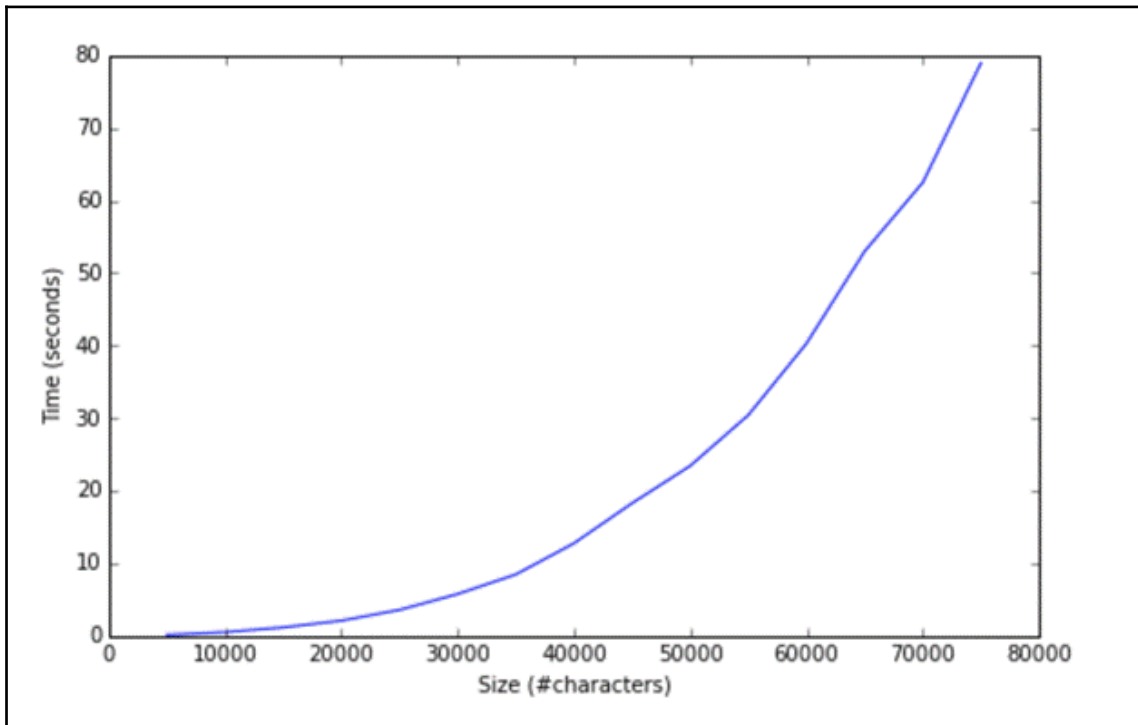
# Chapter 11: Similarity Queries and Summarization

A metric on a set  $X$  is a function (called the *distance function* or simply **distance**)

$$d : X \times X \rightarrow [0, \infty),$$

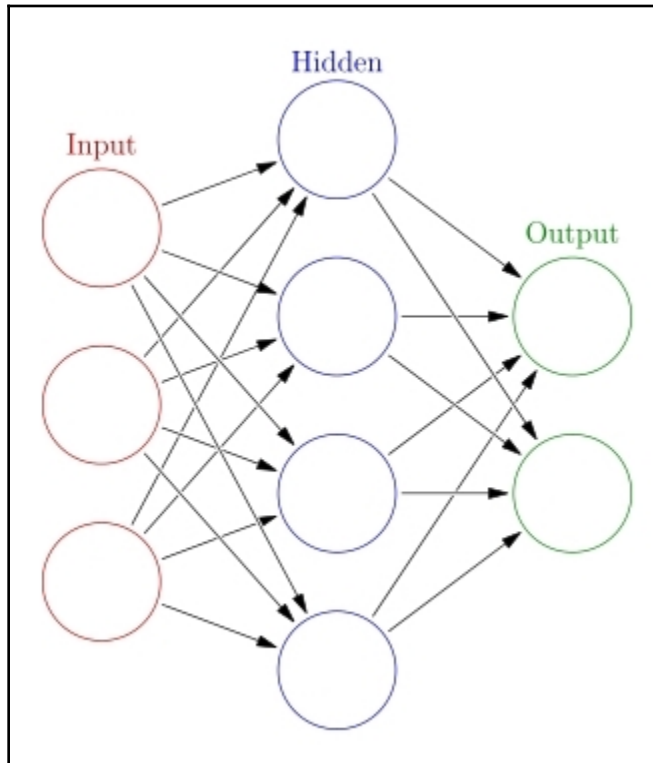
where  $[0, \infty)$  is the set of non-negative **real numbers** and for all  $x, y, z \in X$ , the following conditions are satisfied:

1.  $d(x, y) \geq 0$  non-negativity or separation axiom
2.  $d(x, y) = 0 \Leftrightarrow x = y$  identity of indiscernibles
3.  $d(x, y) = d(y, x)$  symmetry
4.  $d(x, z) \leq d(x, y) + d(y, z)$  subadditivity or triangle inequality



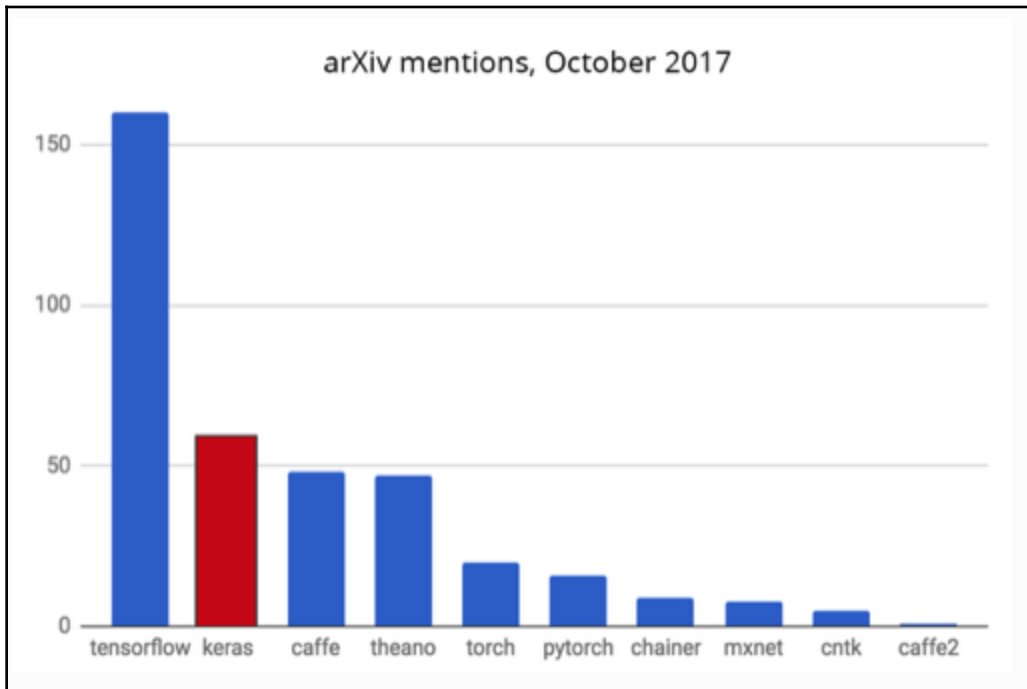
---

## Chapter 13: Deep Learning for Text



---

## Chapter 14: Keras and spaCy for Deep Learning





---

# Chapter 15: Sentiment Analysis and ChatBots

