

# Elements Of Data Science - F2020

## Week 10: NLP, Sentiment Analysis and Topic Modeling

11/23/2020

# TODOs

- Readings:
  - PDSH 5.11 k-Means
  - [Recommended] PML Chapter 11: Working with Unlabeled Data - Clustering Analysis except for last section on DBScan
  - [Optional] Data Science From Scratch Chap 22: Recommender Systems
- HW3, Due Friday Dec 4th 11:59pm
- Answer and submit Quiz 10, **Sunday Nov 29th, 11:59pm ET**

# Today

- Pipelines
- NLP
- Sentiment Analysis
- Topic Modeling

# Questions?

# Pipelines in sklearn

- Pipelines are wrappers used to string together transformers and estimators

```
In [2]: # Example from PML - scaling > feature extraction > classification
from sklearn.datasets import load_breast_cancer
from sklearn.model_selection import train_test_split
bc = load_breast_cancer()
X,y = bc['data'],bc['target']
X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.2,stratify=y,random_state=123)
```

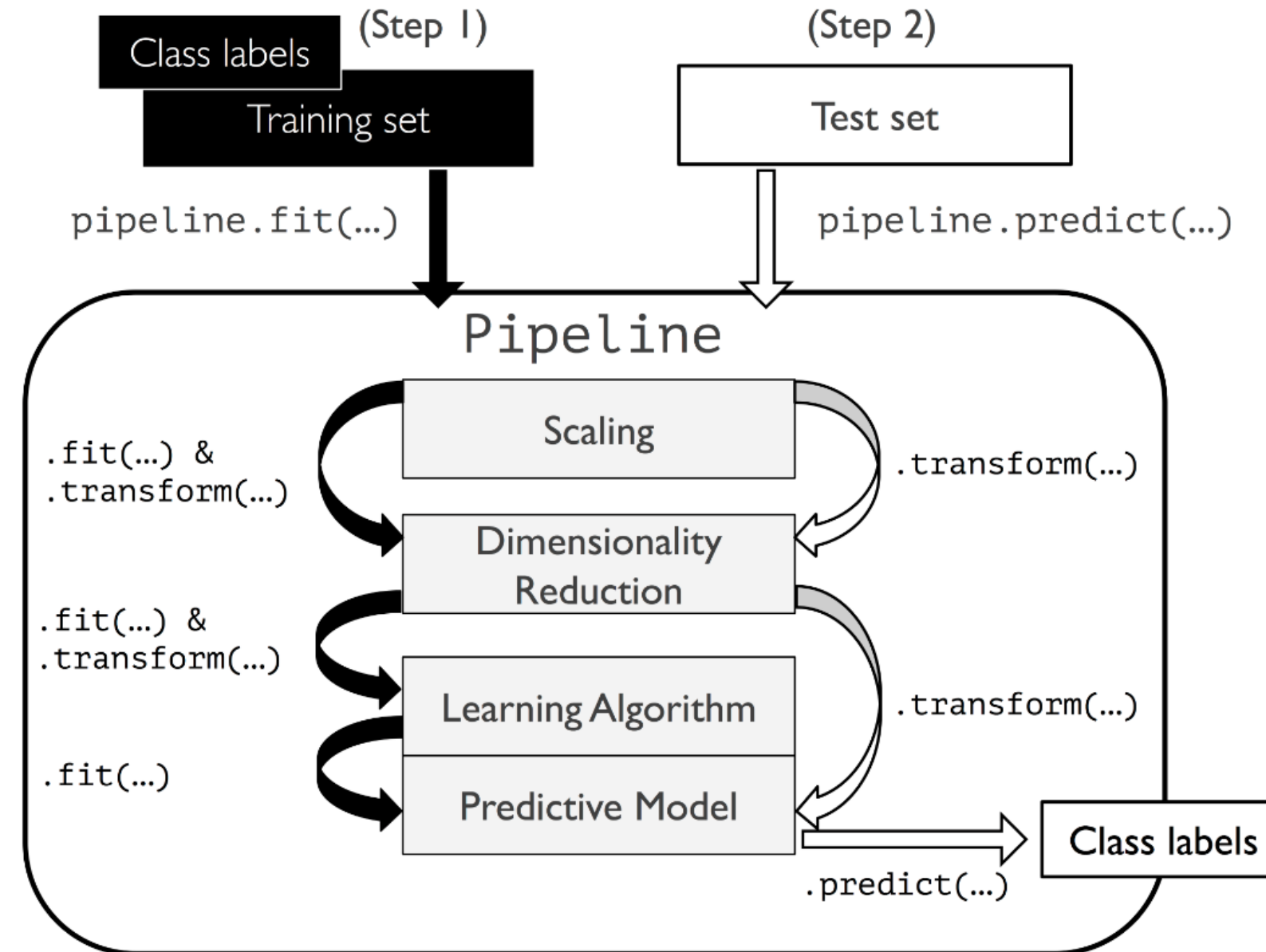
```
In [3]: from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import make_pipeline

# make_pipeline: arguments in order of how they should be applied
pipe_bc = make_pipeline(StandardScaler(),           # center and scale data
                        PCA(n_components=2),        # extract 2 dimensions
                        LogisticRegression(random_state=123) # classify using logistic regression
                        )
pipe_bc.fit(X_train,y_train)

score = pipe_bc.score(X_test,y_test)
print(f'test set accuracy: {score:0.2f}')

test set accuracy: 0.96
```

# Pipelines in sklearn



From PML

# Pipelines in sklearn: Named Steps

```
In [4]: from sklearn.pipeline import Pipeline

# Pipeline: list of (name,object) pairs
pipe_bc = Pipeline([('scale',StandardScaler()),
                    ('pca',PCA(n_components=2)),
                    ('lr',LogisticRegression(random_state=123)),
                    ])

pipe_bc.fit(X_train,y_train)

score = pipe_bc.score(X_test,y_test)
print(f'test set accuracy: {score:0.3f}')
```

test set accuracy: 0.956

```
In [5]: # access pipeline components by name like a dictionary
pipe_bc['lr'].coef_
```

```
Out[5]: array([[ -2.0068728 ,  1.12126495]])
```

```
In [6]: pipe_bc['pca'].components_[0]
```

```
Out[6]: array([0.21777854, 0.08876361, 0.22663097, 0.22043131, 0.14913361,
               0.23954684, 0.25974993, 0.26277752, 0.14518851, 0.06537618,
               0.20775303, 0.0074925 , 0.21143104, 0.2018041 , 0.0165253 ,
               0.17152404, 0.14891828, 0.18380569, 0.03639995, 0.09860293,
               0.22726391, 0.09186544, 0.23623194, 0.22416772, 0.13445762,
               0.21075345, 0.22996838, 0.25138607, 0.12409848, 0.13331693])
```

# Pipelines in sklearn: GridSearch with Pipelines

```
In [7]: from sklearn.model_selection import GridSearchCV

# separate step-names and argument-names with double-underscore '__'
params = {'pca__n_components':[2,10,100],
          'lr__penalty':['none','l1','l2'],
          'lr__C': [.01,1,10,100]}

gscv = GridSearchCV(pipe_bc, params, cv=3, n_jobs=-1).fit(X_train,y_train)

gscv.best_params_
```

```
Out[7]: {'lr__C': 10, 'lr__penalty': 'l2', 'pca__n_components': 10}
```

```
In [8]: score = gscv.score(X_test,y_test)
print(f'test set accuracy: {score:0.3f}')
```

```
test set accuracy: 0.965
```



# Natural Language Processing (NLP)

- Analyzing and interacting with natural language
- Python Libraries
  - **sklearn**
  - nltk
  - **spaCy**
  - gensim
  - ...

# Natural Language Processing (NLP)

- Many NLP Tasks
  - **sentiment analysis**
  - **topic modeling**
  - entity detection
  - machine translation
  - natural language generation
  - question answering
  - relationship extraction
  - automatic summarization
  - ...

# Aside: Python Builtin String Functions

```
In [9]: doc = "D.S. is fun!"  
doc
```

```
Out[9]: 'D.S. is fun!'
```

```
In [10]: doc.lower(), doc.upper()      # change capitalization
```

```
Out[10]: ('d.s. is fun!', 'D.S. IS FUN!')
```

```
In [11]: doc.split() , doc.split('.')  # split a string into parts (default is whitespace)
```

```
Out[11]: (['D.S.', 'is', 'fun!'], ['D', 'S', ' ' is fun!'])
```

```
In [12]: '|'.join(['ab', 'c', 'd'])     # join items in a list together
```

```
Out[12]: 'ab|c|d'
```

```
In [13]: '|'.join(doc[:5])              # a string itself is treated like a list of characters
```

```
Out[13]: 'D|. |S|. | '
```

```
In [14]: '  test  '.strip()             # remove whitespace from the beginning and end of a string
```

```
Out[14]: 'test'
```

- and many more, see <https://docs.python.org/3.8/library/string.html>

# NLP: The Corpus

- **corpus:** collection of documents
  - books
  - articles
  - reviews
  - tweets
  - resumes
  - sentences?
  - ...

# NLP: Doc Representation

- Documents usually represented as strings
  - string: a sequence (list) of unicode characters

```
In [15]: doc = "D.S. is fun!\nIt's  true."  
print(doc)
```

```
D.S. is fun!  
It's  true.
```

```
In [16]: '|'.join(doc)
```

```
Out[16]: "D|.|S|.| |i|s| |f|u|n|!|\n|I|t|'|s| | |t|r|u|e|."
```

- Need to split this up into parts (**tokens**)
- Good job for **Regular Expressions**

# Regular Expressions

- Strings that define search patterns over text
- Useful for finding/replacing/grouping
- python `re` library (others available)

```
In [17]: print(doc)
```

```
D.S. is fun!  
It's  true.
```

```
In [18]: import re  
# Find all of the whitespaces in doc  
# '\s+' means "one or more whitespace characters"  
re.findall(r'\s+',doc)
```

```
Out[18]: [' ', ' ', '\n', ' ']
```

# Regular Expressions

Just some of the special character definitions:

- `.` : any single character except newline (`r'.'` matches `'x'`)
- `*` : match 0 or more repetitions (`r'x*' matches 'x','xx',''`)
- `+` : match 1 or more repetitions (`r'x+' matches 'x','xx'`)
- `?` : match 0 or 1 repetitions (`r'x?' matches 'x' or ''`)
- `^` : beginning of string (`r'^D' matches 'D.S.'`)
- `$` : end of string (`r'fun!$' matches 'DS is fun!'`)

# Regular Expression Cont.

- `[]` : a set of characters (^ as first element = not)
- `\s` : whitespace character (Ex: `[\t\n\r\f\v]`)
- `\S` : non-whitespace character (Ex: `[^\t\n\r\f\v]`)
- `\w` : word character (Ex: `[a-zA-Z0-9_]`)
- `\W` : non-word character
- `\b` : boundary between `\w` and `\W`
- and many more!
- See [regex101.com](https://regex101.com) for examples and testing



# NLP: Tokenization

- **tokens:** strings that make up a document ('the', 'cat',...)
- **tokenization:** convert a document into tokens
- **vocabulary:** set of unique tokens (terms) in corpus

```
In [19]: # split on whitespace
re.split(r'\s+', doc)
```

```
Out[19]: ['D.S.', 'is', 'fun!', "It's", 'true.']
```

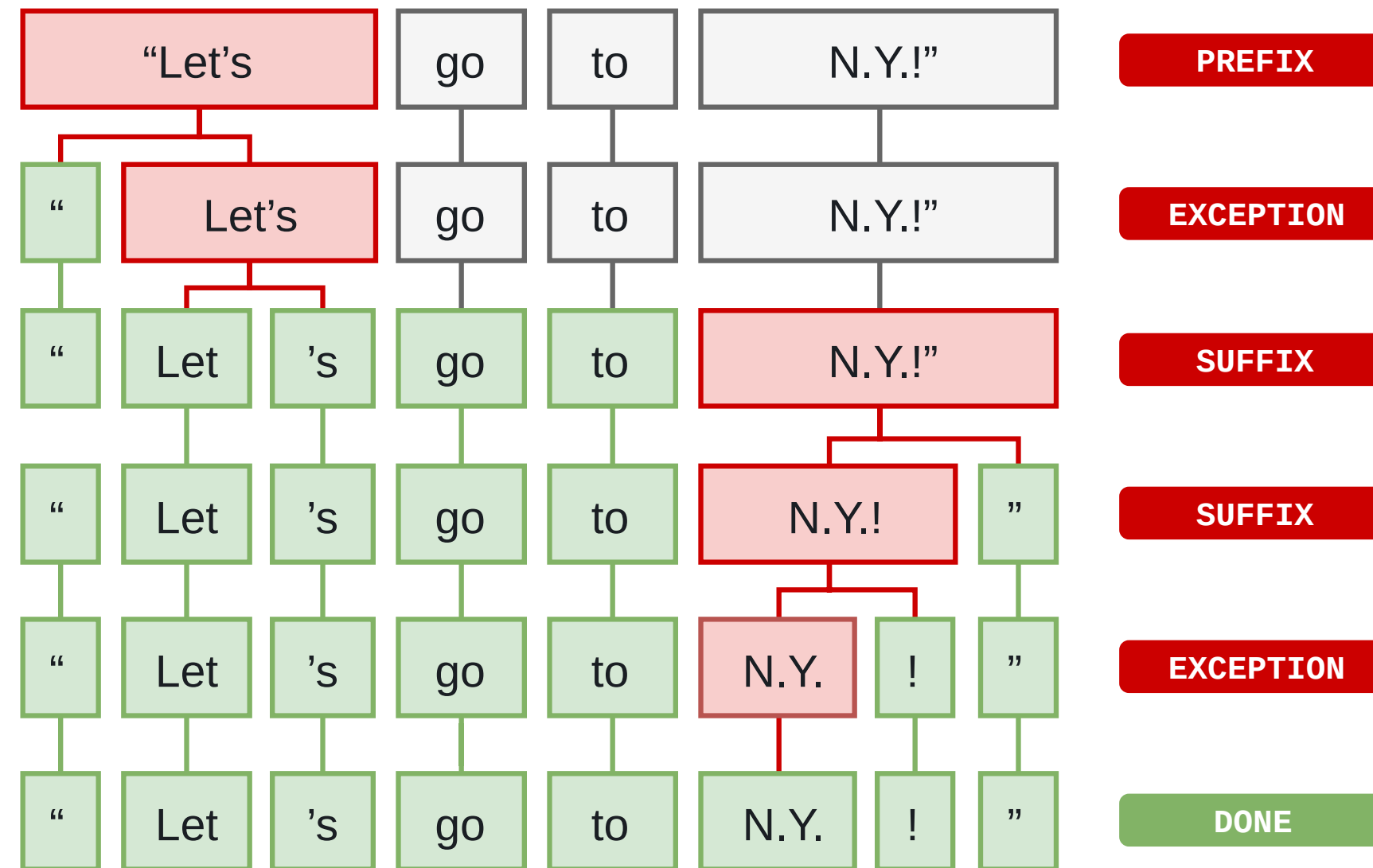
```
In [20]: # find tokens of length 2+ word characters
re.split('\b\w\w+\b', doc)
['is', 'fun', 'It', 'true']
```

```
Out[20]: ['is', 'fun', 'It', 'true']
```

```
In [21]: # find tokens of length 2+ non-space characters
re.findall(r"\b\S\S+\b", doc)
```

```
Out[21]: ['D.S', 'is', 'fun', "It's", 'true']
```

# NLP:Tokenization



From <https://spacy.io/usage/linguistic-features>

# NLP: Other Preprocessing

- lowercase
- remove special characters
- add , tags
- stemming: cut off beginning or ending of word
  - 'studies' becomes 'studi'
  - 'studying' becomes 'study'
- lemmatization: perform morphological analysis
  - 'studies' becomes 'study'
  - 'studying' becomes 'study'

# NLP: Bag of Words

- BOW representation: ignore token order

```
In [22]: sorted(re.findall(r'\b\S+\b', doc.lower()))
```

```
Out[22]: ['d.s', 'fun', 'is', "it's", 'true']
```

# NLP: n-Grams

- Unigram: single token
- Bigram: combination of two ordered tokens
- n-Gram: combination of n ordered tokens
- The larger n is, the larger the vocabulary

```
In [23]: # Bigram example:
tokens = '<start> data science is fun <end>'.split()
[tokens[i]+'_'+tokens[i+1] for i in range(len(tokens)-1)]

Out[23]: ['<start>_data', 'data_science', 'science_is', 'is_fun', 'fun_<end>']
```

# NLP: TF and DF

- **Term Frequency:** number of times term is seen per document

```
In [24]: corpus = ['red green blue', 'red blue blue']

#Vocabulary
vocab = sorted(set(' '.join(corpus).split()))
vocab
```

```
Out[24]: ['blue', 'green', 'red']
```

```
In [25]: #TF
from collections import Counter
tf = np.zeros((len(corpus), len(vocab)))
for i, doc in enumerate(corpus):
    for j, term in enumerate(vocab):
        tf[i, j] = Counter(doc.split())[term]
tf = pd.DataFrame(tf, index=['doc1', 'doc2'], columns=vocab)
tf
```

```
Out[25]:
```

	blue	green	red
doc1	1.0	1.0	1.0
doc2	2.0	0.0	1.0

# NLP: TF and DF

- **Document Frequency:** number of documents containing each term

```
In [26]: #DF  
tf.astype(bool).sum(axis=0)
```

```
Out[26]: blue      2  
         green     1  
         red       2  
         dtype: int64
```

# NLP: Stopwords

- terms that have high (or very low) DF and aren't informative
  - common english terms (ex: 'a', 'the','in',...)
  - domain specific (ex, in class slides: 'data\_science')
  - often removed prior to analysis
  - in sklearn
    - `min_df`, an integer  $> 0$ , keep terms that occur in at least  $n$  documents
    - `max_df`, a float in  $(0,1]$ , keep terms that occur in less than  $f\%$  of total documents



# NLP: CountVectorizer in sklearn

```
In [27]: corpus = ['blue green red', 'blue green green']

from sklearn.feature_extraction.text import CountVectorizer

cvect = CountVectorizer(lowercase=True,      # default, transform all docs to lowercase
                        ngram_range=(1,1),  # default, only unigrams
                        min_df=1,           # default, keep all terms
                        max_df=1.0,         # default, keep all terms
                        )
X_cv = cvect.fit_transform(corpus)
X_cv.shape
```

Out[27]: (2, 3)

```
In [28]: # learned vocabulary, sorted by column mapping id
sorted(cvect.vocabulary_.items(),key=lambda x: x[1])
```

Out[28]: [('blue', 0), ('green', 1), ('red', 2)]

```
In [29]: # term frequencies
X_cv.todense()
```

Out[29]: matrix([[1, 1, 1],
 [1, 2, 0]])

```
In [30]: # mapping back to terms via vocabulary mapping
cvect.inverse_transform(X_cv)
```

Out[30]: [array(['blue', 'green', 'red'], dtype='<U5'),
 array(['blue', 'green'], dtype='<U5')]

# NLP: Tfidf

- What if some terms are still uninformative?
- Can we downweight terms that occur in many documents?
- **Term Frequency \* Inverse Document Frequency (tf-idf)**
  - $\text{tf-idf}(t, d) = \text{tf}(t, d) \times \text{idf}(t)$
  - $\text{idf}(t) = \log \frac{1+n}{1+\text{df}(t)} + 1$

```
In [31]: from sklearn.feature_extraction.text import TfidfVectorizer

tfidfvect = TfidfVectorizer(norm='l2') # by default, also doing l2 normalization

X_tfidf = tfidfvect.fit_transform(corpus)
sorted(tfidfvect.vocabulary_.items(), key=lambda x: x[1])
```

```
Out[31]: [('blue', 0), ('green', 1), ('red', 2)]
```

```
In [32]: X_tfidf.todense()
```

```
Out[32]: matrix([[0.50154891, 0.50154891, 0.70490949],
                 [0.4472136 , 0.89442719, 0.          ]])
```

```
In [33]: # can also use to get term frequencies by setting use_idf to False and norm to none
TfidfVectorizer(use_idf=False, norm=None).fit_transform(corpus).todense()
```

```
Out[33]: matrix([[1., 1., 1.],
                 [1., 2., 0.]])
```

# NLP: Classification Example

```
In [34]: from sklearn.datasets import fetch_20newsgroups

ngs = fetch_20newsgroups(categories=['rec.sport.baseball', 'rec.sport.hockey']) # dataset has 20 categories, only get two

docs_ngs = ngs['data'] # get documents (emails)
y_ngs = ngs['target'] # get targets ([0,1])
target_names_ngs = ngs['target_names'] # get target names (['rec.autos', 'sci.space'])

print(y_ngs[0], target_names_ngs[y_ngs[0]]) # print target int and target name
print('-'*50) # print a string of 50 dashes
print(docs_ngs[0].strip()[:600]) # print beginning characters of first doc, after stripping whitespace
```

```
0 rec.sport.baseball
```

```
-----
From: dougb@comm.mot.com (Doug Bank)
Subject: Re: Info needed for Cleveland tickets
Reply-To: dougb@ecs.comm.mot.com
Organization: Motorola Land Mobile Products Sector
Distribution: usa
Nntp-Posting-Host: 145.1.146.35
Lines: 17
```

In article <1993Apr1.234031.4950@leland.Stanford.EDU>, bohnert@leland.Stanford.EDU (matthew bohnert) writes:

```
|> I'm going to be in Cleveland Thursday, April 15 to Sunday, April 18.
|> Does anybody know if the Tribe will be in town on those dates, and
|> if so, who're they playing and if tickets are available?
```

```
The tribe will be in town from April 16 to the 19th.
There
```

# NLP Example: Transform Docs

```
In [35]: from sklearn.model_selection import train_test_split
docs_ngs_train, docs_ngs_test, y_ngs_train, y_ngs_test = train_test_split(docs_ngs, y_ngs)

vect = TfidfVectorizer(lowercase=True,
                        min_df=5,          # occur in at least 5 documents
                        max_df=0.8,        # occur in at most 80% of documents
                        token_pattern='\\b\\S\\S+\\b', # tokens of at least 2 non-space characters
                        ngram_range=(1,1), # only unigrams
                        use_idf=False,     # term frequency counts instead of tf-idf
                        norm=None          # do not normalize
                        )
X_ngs_train = vect.fit_transform(docs_ngs_train)
X_ngs_train.shape
```

Out[35]: (897, 3676)

```
In [36]: # first few terms in learned vocabulary
list(vect.vocabulary_.items())[:5]
```

Out[36]: [('university', 3428),  
('computing', 870),  
('center', 744),  
('just', 1841),  
('distribution', 1082)]

```
In [37]: # first few terms in learned stopword list
list(vect.stop_words_)[:5]
```

Out[37]: ['upcomming', 'pony.1993apr15.223040.8733', 'gidp', 'alignment', 'memorable']

# NLP Example: Train and Evaluate Classifier

```
In [39]: from sklearn.model_selection import cross_val_score
from sklearn.linear_model import LogisticRegression
from sklearn.dummy import DummyClassifier

scores_dummy = cross_val_score(DummyClassifier(strategy='most_frequent'), X_ngs_train, y_ngs_train)
scores_lr = cross_val_score(LogisticRegression(), X_ngs_train, y_ngs_train)

print(f'dummy cv accuracy: {scores_dummy.mean():0.2f} +- {scores_dummy.std():0.2f}')
print(f'lr      cv accuracy: {scores_lr.mean():0.2f} +- {scores_lr.std():0.2f}')
```

dummy cv accuracy: 0.51 +- 0.00  
lr cv accuracy: 0.95 +- 0.01

# NLP Example: Using Pipeline

```
In [40]: from sklearn.pipeline import Pipeline

# use Pipeline instead of make_pipeline to add names to the steps
# (name,object) tuple pairs for each step
ngs_pipe = Pipeline([('vect',TfidfVectorizer(lowercase=True,
                                              min_df=5,
                                              max_df=0.8,
                                              token_pattern='\\b\\S\\S+\\b',
                                              ngram_range=(1,1),
                                              use_idf=False,
                                              norm=None ),
                      ('lr',LogisticRegression())
                    ])

ngs_pipe.fit(docs_ngs_train,y_ngs_train) # pass in docs, not transformed X

score_ngs = ngs_pipe.score(docs_ngs_train,y_ngs_train)
print(f'pipeline accuracy on training set: {score_ngs:0.2f}')
```

pipeline accuracy on training set: 1.00

```
In [41]: scores_pipe = cross_val_score(ngs_pipe,docs_ngs_train,y_ngs_train)
print(f'pipe cv accuracy: {scores_pipe.mean():0.2f} +- {scores_pipe.std():0.2f}')
```

pipe cv accuracy: 0.94 +- 0.01

```
In [42]: list(ngs_pipe['vect'].vocabulary_.items())[:3]
```

```
Out[42]: [('university', 3428), ('computing', 870), ('center', 744)]
```

# NLP Example: Add Feature Selection

```
In [43]: from sklearn.feature_selection import SelectFromModel

ngs_pipe = Pipeline([('vect', TfidfVectorizer(lowercase=True,
                                              min_df=5,
                                              max_df=0.8,
                                              token_pattern='\\b\\S\\S+\\b',
                                              ngram_range=(1,1),
                                              use_idf=False,
                                              norm=None ),
                      ('fs', SelectFromModel(estimator=LogisticRegression(C=.1, penalty='l1', solver='liblinear'))),
                      ('lr', LogisticRegression())
                     ])
ngs_pipe.fit(docs_ngs_train, y_ngs_train)
print(f'pipeline accuracy on training set: {ngs_pipe.score(docs_ngs_train, y_ngs_train):0.2f}')
scores_pipe = cross_val_score(ngs_pipe, docs_ngs_train, y_ngs_train)
print(f'pipe cv accuracy: {scores_pipe.mean():0.2f} +- {scores_pipe.std():0.2f}')
```

pipeline accuracy on training set: 0.97  
pipe cv accuracy: 0.92 +- 0.01

# NLP Example: Grid Search with Feature Selection

```
In [44]: params = {'vect__use_idf':[True,False],
                  'vect__norm':['l1','l2',None],
                  'fs__estimator__C': [.01,.1,1,10],
                  'lr__C': [.01,.1,1,10]}

gscv = GridSearchCV(ngs_pipe, params, cv=3, n_jobs=-1).fit(docs_ngs_train,y_ngs_train)

print(gscv.best_params_)
print(f'gscsv best cv accuracy : {gscv.best_score_:0.2f}')
print(f'gscsv test set accuracy: {gscv.score(docs_ngs_test,y_ngs_test):0.2f}')

{'fs__estimator__C': 0.1, 'lr__C': 0.1, 'vect__norm': None, 'vect__use_idf': True}
gscsv best cv accuracy : 0.95
gscsv test set accuracy: 0.95
```



# Sentiment Analysis and sklearn

- determine sentiment/opinion from unstructured text
  - usually positive/negative, but is domain specific
  - can be treated as a classification task (with a target, using all of the tools we know)
  - can also be treated as a linguistic task (sentence parsing)
- 
- Example: determine sentiment of movie reviews
  - see `sentiment_analysis_example.ipynb`

# Topic Modeling

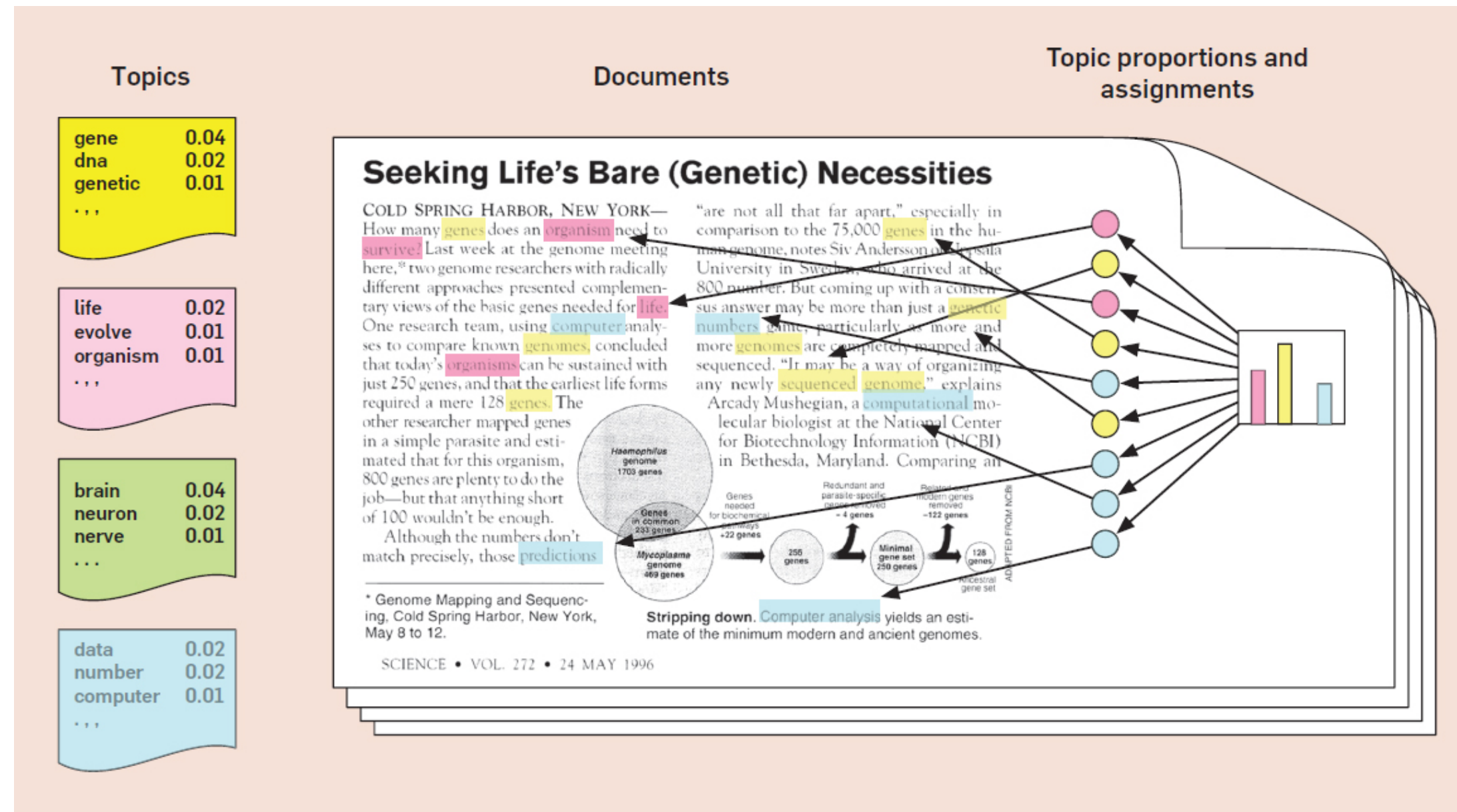
- What topics are our documents composed of?
  - How much of each topic does each document contain?
  - Can we represent documents using topic weights? (dimensionality reduction)
- 
- What is topic modeling?
  - How does Latent Dirichlet Allocation (LDA) work?
  - How to train and use LDA with sklearn?

# What is Topic Modeling?

- **topic:** a collection of related words
- A document can be composed of several topics
- Given a collection of documents, we can ask:
  - **What terms make up each topic?** (per topic term distribution)
  - **What topics make up each document?** (per document topic distribution)

# Topic Modeling with Latent Dirichlet Allocation (LDA)

- Unsupervised method for determining topics and topic assignments



# Topic Modeling: Example

- Guessing some **topics** (per topic term distribution  $\phi$ )

```
In [45]: vocab = ['baseball', 'cat', 'dog', 'pet', 'played', 'tennis']

V = len(vocab) # size of vocabulary

K = 2 # number of topics

# the probability of each term given topic 1 (high for sports terms)
topic_1 = [.33, 0, 0, 0, .33, .33]

# the probability of each term given topic 2 (high for pet terms)
topic_2 = [0, .25, .25, .25, .25, 0]

# per topic term distributions
phi = [topic_1, topic_2]

print(np.array(phi).shape) # K x V (number of topics x size of vocabulary)

(2, 6)
```

# Topic Modeling: Example

- Guessing the per document topic distributions  $\theta$  given the **topics**

```
In [46]: # recall
vocab = ['baseball', 'cat', 'dog', 'pet', 'played', 'tennis']

phi = [[.33, 0, 0, 0, .33, .33],
       [ 0, .25, .25, .25, .25, 0]]
```

```
In [47]: corpus = ['the dog and cat played tennis',
                   'tennis and baseball are sports',
                   'a dog or a cat can be a pet']

# per document topic distributions
theta = [[.50, .50],
         [.99, .01],
         [.01, .99]]

print(np.array(theta).shape) # M x K (number of documents x number of topics)

(3, 2)
```

# Topic Modeling With LDA

- Given
  - a set of documents
  - a number of topics  $k$
- Learn
  - the **per topic term distributions**  $\varphi$  (phi), size:  $k \times V$
  - the **per document topic distributions**  $\theta$  (theta), size:  $n \times k$
- How to learn  $\phi$  and  $\theta$ :
  - Latent Dirichlet Allocation (LDA)
  - generative statistical model
  - Blei, D., Ng, A., Jordan, M. Latent Dirichlet allocation. J. Mach. Learn. Res. 3 (Jan 2003)

# Topic Modeling With LDA

- Uses for  $\varphi$  (phi), the per topic word distributions:
  - inferring labels for topics
  - word clouds
- Uses for  $\theta$  (theta), the per document topic distributions:
  - dimensionality reduction
  - clustering
  - similarity



# LDA with sklearn

```
In [48]: # load data from all 20 newsgroups
newsgroups = fetch_20newsgroups()
ngs_all = newsgroups.data
len(ngs_all)
```

Out[48]: 11314

```
In [49]: # transform documents using tf-idf
tfidf = TfidfVectorizer(token_pattern=r'\b[a-zA-Z0-9-][a-zA-Z0-9-]+\b', min_df=50, max_df=.2)
X_tfidf = tfidf.fit_transform(ngs_all)
X_tfidf.shape
```

Out[49]: (11314, 4256)

```
In [50]: feature_names = tfidf.get_feature_names()
print(feature_names[:10])
print(feature_names[-10:])

['00', '000', '01', '02', '03', '04', '05', '06', '07', '08']
['yours', 'yourself', 'ysu', 'zealand', 'zero', 'zeus', 'zip', 'zone', 'zoo', 'zuma']
```

# LDA with sklearn Cont.

```
In [51]: from sklearn.decomposition import LatentDirichletAllocation

# create model with 20 topics
lda = LatentDirichletAllocation(n_components=20, # the number of topics
                               n_jobs=-1,      # use all cpus
                               random_state=123) # for reproducibility

# learn phi (lda.components_) and theta (X_lda)
# this will take a while!
X_lda = lda.fit_transform(X_tfidf)
```

```
In [52]: ngs_all[100][:100]
```

```
Out[52]: 'From: tchen@magnus.acs.ohio-state.edu (Tsung-Kun Chen)\nSubject: ** Software forsale (lots) **\nNntp-P'
```

```
In [53]: np.round(X_lda[100],2) # lda representation of document_100
```

```
Out[53]: array([0.01, 0.01, 0.01, 0.01, 0.1 , 0.01, 0.01, 0.01, 0.01, 0.01, 0.01,
                0.01, 0.01, 0.01, 0.38, 0.01, 0.14, 0.01, 0.01, 0.28])
```

```
In [54]: # Note: since this is unsupervised, these numbers may change
np.argsort(X_lda[100])[:, -1][:3] # the top topics of document_100
```

```
Out[54]: array([14, 19, 16])
```

# LDA: Per Topic Term Distributions

```
In [56]: print_top_words(lda, feature_names, 5)
```

```
Topic 0: uga ai georgia covington mcovingt
Topic 1: digex access turkish armenian armenians
Topic 2: god jesus bible christians christian
Topic 3: values objective frank morality ap
Topic 4: ohio-state magnus acs ohio cis
Topic 5: caltech keith sandvik livesey sgi
Topic 6: stratus msg usc indiana sw
Topic 7: alaska uci aurora colostate nsmca
Topic 8: wpi radar psu psuvm detector
Topic 9: columbia utexas gatech cc prism
Topic 10: scsi upenn simms ide bus
Topic 11: nhl team mit players hockey
Topic 12: lehigh duke jewish adobe ns1
Topic 13: henry toronto zoo ti dseg
Topic 14: sale card thanks please mac
Topic 15: virginia joel hall doug douglas
Topic 16: ca his new cs should
Topic 17: cleveland cwru freenet cramer ins
Topic 18: pitt gordon geb banks cs
Topic 19: windows file window files thanks
```

# LDA Review

- What did we learn?
  - per document topic distributions
  - per topic term distributions
- What can we use this for?
  - Dimensionality Reduction/Feature Extraction!
  - investigate topics (much like PCA components)

# Other NLP Features

- Part of Speech tags
- Dependency Parsing
- Entity Detection
- Word Vectors
- See spaCy!

# Using spaCy

```
In [57]: import spacy

#first run
#%run -m spacy download en_core_web_sm
try:

    nlp = spacy.load("en_core_web_sm")

except OSError as e:
    print('Need to run the following line in a new cell:')
    print('%run -m spacy download en_core_web_sm')
    print('or the following line from the commandline with eods-f20 activated:')
    print('python -m spacy download en_core_web_sm')

parsed = nlp("N.Y.C. isn't in New Jersey.")
'|'.join([token.text for token in parsed])
```

Out[57]: "N.Y.C.|is|n't|in|New|Jersey|."

# spaCy: Part of Speech Tagging

```
In [58]: doc = nlp("Apple is looking at buying U.K. startup for $1 billion.")

print(f"{'text':7s} {'lemma':7s} {'pos':5s} {'is_stop'}")
print('-'*30)
for token in doc:
    print(f"{'token.text':7s} {'token.lemma_':7s} {'token.pos_':5s} {'token.is_stop'}")
```

text	lemma	pos	is_stop
Apple	Apple	PROPN	False
is	be	AUX	True
looking	look	VERB	False
at	at	ADP	True
buying	buy	VERB	False
U.K.	U.K.	PROPN	False
startup	startup	NOUN	False
for	for	ADP	True
\$	\$	SYM	False
1	1	NUM	False
billion	billion	NUM	False
.	.	PUNCT	False

# spaCy: Part of Speech Tagging

```
In [59]: from spacy import displacy  
displacy.render(doc, style="dep")
```

Apple PROPN is AUX looking VERB at ADP buying VERB U.K. PROPN startup NOUN for ADP \$ SYM 1 NUM billion. NUM nsubj aux prep pcomp compound dobj  
prep quantmod compound pobj



# spaCy: Entity Detection

```
In [60]: [(ent.text, ent.label_) for ent in doc.ents]
```

```
Out[60]: [('Apple', 'ORG'), ('U.K.', 'GPE'), ('$1 billion', 'MONEY')]
```

```
In [61]: displacy.render(doc, style="ent")
```

Apple ORG is looking at buying U.K. GPE startup for \$1 billion MONEY .

# spaCy: Word Vectors

- word2vec
- shallow neural net
- predict a word given the surrounding context (SkipGram or CBOW)
- words used in similar context should have similar vectors

# spaCy: Word Vectors

- word2vec
- shallow neural net
- predict a word given the surrounding context (SkipGram or CBOW)
- words used in similar context should have similar vectors

```
In [62]: # Need either the _md or _lg models to get vector information
# Note: this takes a while!
%run -m spacy download en_core_web_md

✓ Download and installation successful
You can now load the model via spacy.load('en_core_web_md')
```

# spaCy: Word Vectors

- word2vec
- shallow neural net
- predict a word given the surrounding context (SkipGram or CBOW)
- words used in similar context should have similar vectors

```
In [62]: # Need either the _md or _lg models to get vector information
# Note: this takes a while!
%run -m spacy download en_core_web_md
```

✓ Download and installation successful  
You can now load the model via `spacy.load('en_core_web_md')`

```
In [63]: nlp = spacy.load('en_core_web_md') # _lg has a larger vocabulary

doc = nlp('Baseball is played on a diamond.')
doc[0].text, doc[0].vector.shape, list(doc[0].vector[:3])
```

```
Out[63]: ('Baseball', (300,), [0.55838, 0.42791, -0.11687])
```

# spaCy: Multiple Documents

```
In [64]: # Use nlp.pipe to transform multiple docs at once
docs = list(nlp.pipe(['Baseball is played on a diamond.',
                      'Hockey is played on ice.',
                      'Diamonds are clear as ice.']))
```

```
In [65]: # using average of token vectors for each document.
np.array([[ '{:.2f}'.format(docs[i].similarity(docs[j])) for j in range(3)]
         for i in range(3)])
```

```
Out[65]: array([[ '1.00', '0.85', '0.76'],
               [ '0.85', '1.00', '0.77'],
               [ '0.76', '0.77', '1.00']], dtype='<U4')
```

# Learning Sequences

- Hidden Markov Models
- Conditional Random Fields
- Recurrent Neural Networks
- LSTM
- BERT

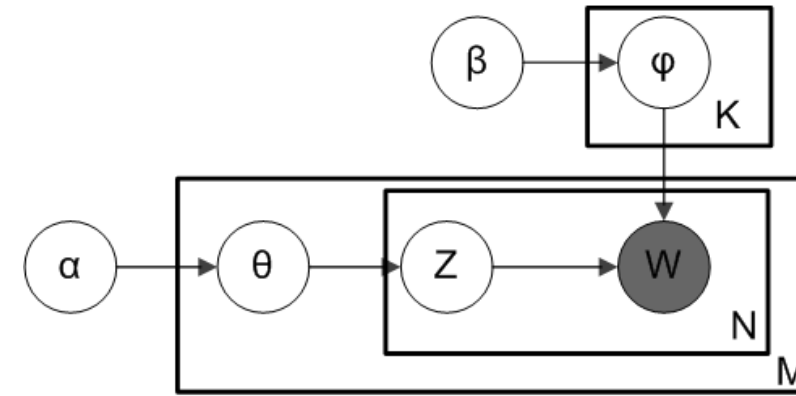
# NLP Review

- corpus, tokens, vocabulary, terms, n-grams, stopwords
- tokenization
- term frequency (TF), document frequency (DF)
- TF vs TF-IDF
- sentiment analysis
- topic modeling
- POS
- Dependency Parsing
- Entity Extraction
- Word Vectors

# Questions?



# LDA Plate Diagram



$K$  : number of topics

$\varphi$  : per topic term distributions

$\beta$  : parameters for word distribution die factory, length =  $V$  (size of vocab)

$M$  : number of documents

$N$  : number of words/tokens in each document

$\theta$  : per document topic distributions

$\alpha$  : parameters for topic die factory, length =  $K$  (number of topics)

$z$  : topic indexes

$w$  : observed tokens