

Analysis of unstructured data

Lecture 10 - text classification

Janusz Szwabiński

Outlook:

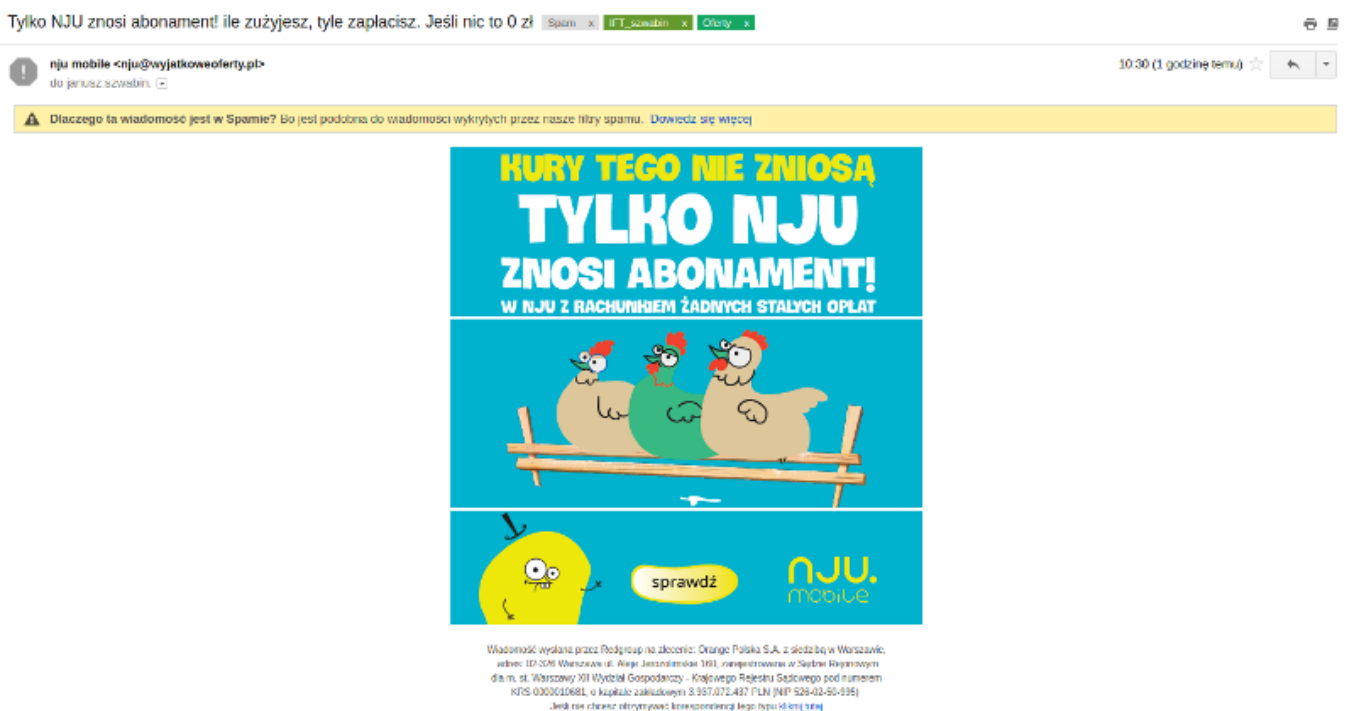
- Motivation
- Text classification
- Naive Bayes classifier
- Decision trees
- Maximum Entropy classifier
- How to improve accuracy of a text classifier?

References:

1. Ch. Manning, H. Schütze, *Foundations of Statistical Natural Language Processing*
2. P. Jackson, I. Moulinier, *Natural Language Processing for Online Applications*
3. Dan Jurafsky's presentation, *Text Classification and Naive Bayes*
4. NLTK Book, <http://www.nltk.org/book/> (<http://www.nltk.org/book/>)
5. <http://blog.datumbox.com/machine-learning-tutorial-the-max-entropy-text-classifier/> (<http://blog.datumbox.com/machine-learning-tutorial-the-max-entropy-text-classifier/>)

Motivation

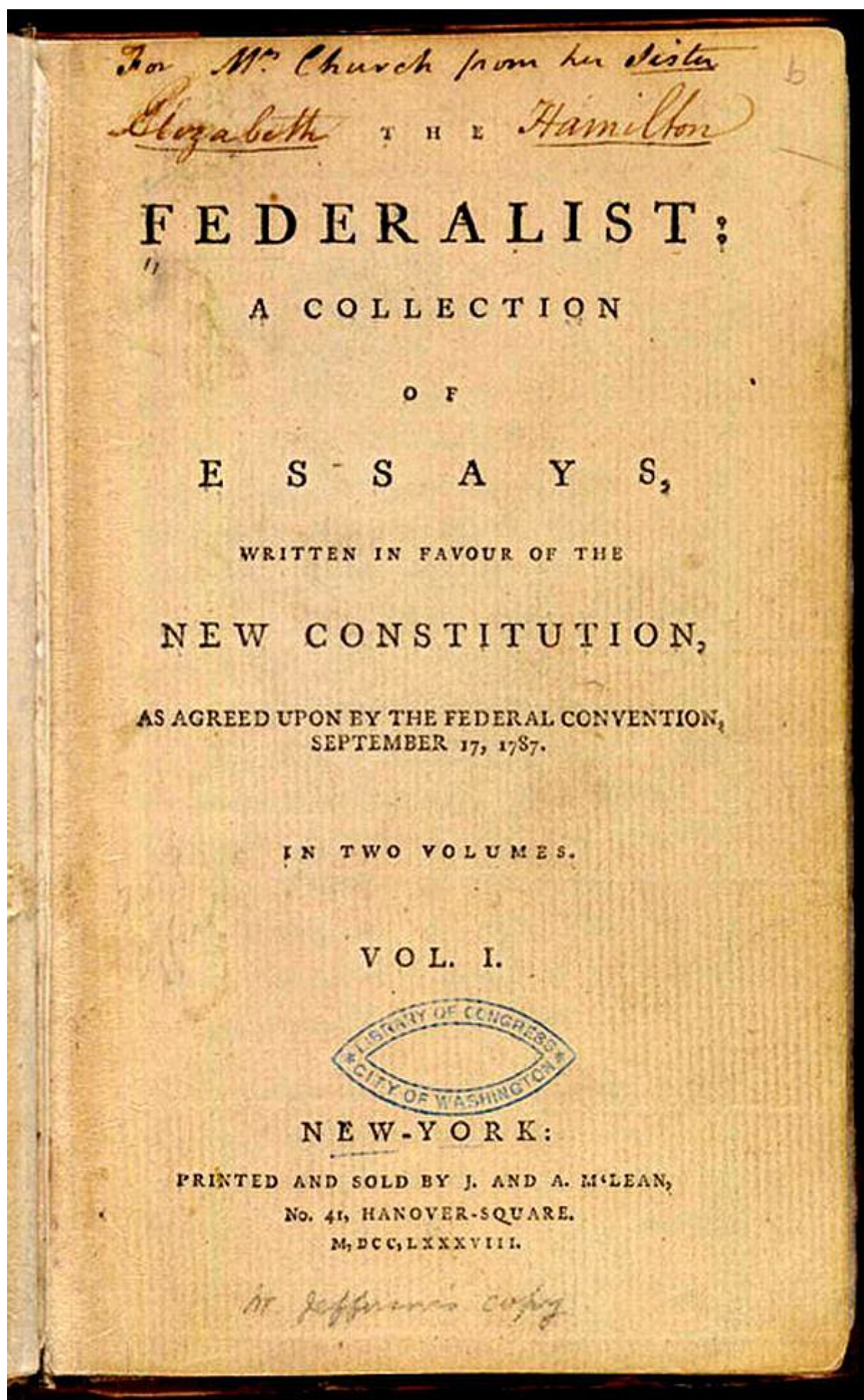
Is that spam?



The Federalist (or The Federalist Papers)

- 1787-1788

- a collection of 85 articles and essays written under the pseudonym "Publius" to promote the ratification of the United States Constitution
- at the time of publication the authorship of the articles was a closely guarded secret
- authors: Alexander Hamilton, James Madison, and John Jay
 - according to Hamilton, he was the author of 63 articles, Jay wrote 5 of them, Madison – 14, 3 were written by Hamilton and Madison
 - in 1818 Madison claimed the authorship of 26 essays
 - the scholarly detective work of Douglass Adair in 1944 postulated the following assignments of authorship:
 - Alexander Hamilton (51 articles: No. 1, 6–9, 11–13, 15–17, 21–36, 59–61, and 65–85)
 - James Madison (29 articles: No. 10, 14, 18–20, [10] 37–58 and 62–63)
 - John Jay (5 articles: No. 2–5 and 64).
 - in 1963 Mosteller and Wallace corroborated Adair's results by a computer aided text classification



Female or male author?

1925 present-day Vietnam was divided into three parts under French colonial rule. The southern region embracing Saigon and the Mekong delta was the colony of Cochin China; the central area with its imperial capital at Hue was the protectorate of Annam, and the northern region, Tongking, was also a separate ...

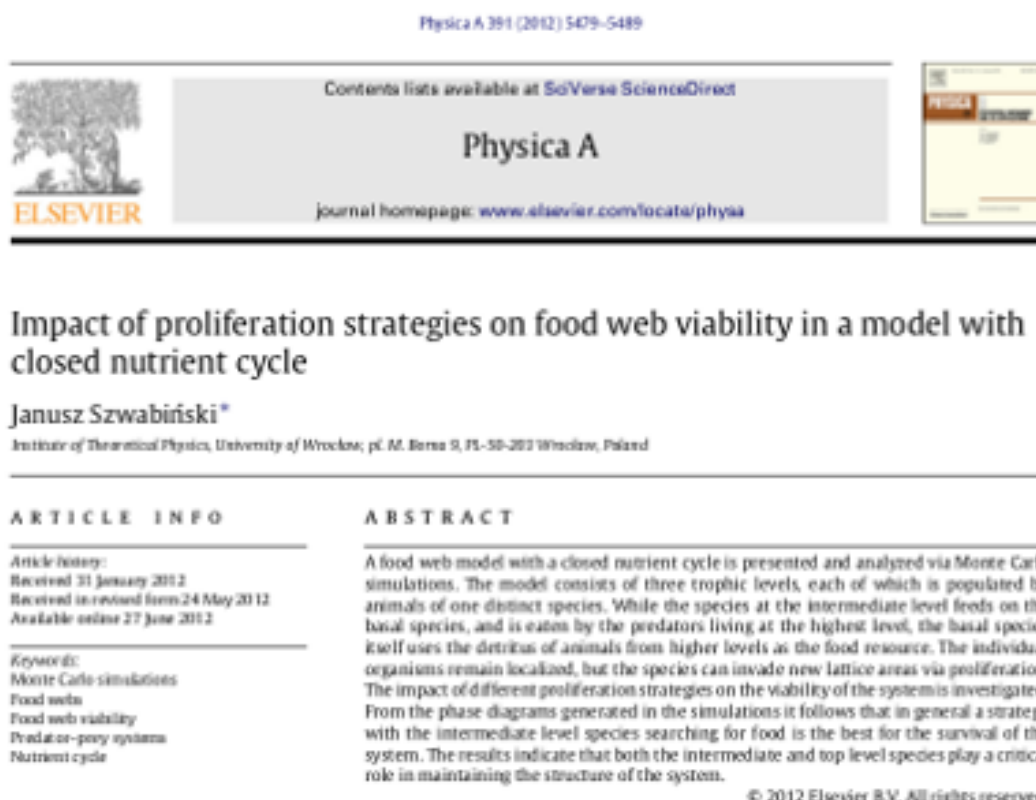
Clara never failed to be astonished by the extraordinary felicity of her own name. She found it hard to trust herself to the mercy of fate, which had managed over the years to convert her greatest shame into one of her greatest assets...

- S. Argamon, M. Koppel, J. Fine, A. R. Shimoni, 2003. "Gender, Genre, and Writing Style in Formal Written Texts", Text 23, pp. 321–346

Movie reviews

- Unbelievably disappointing
- Full of zany characters and richly applied satire, and some great plot twists
- This is the greatest screwball comedy ever filmed
- It was pathetic. The worst part about it was the boxing scenes

What is this paper about?



'Quantitative Biology' categories on arXiv:

- Quantitative Biology
- Biomolecules
- Cell Behavior
- Genomics
- Molecular Networks
- Neurons and Cognition
- Other Quantitative Biology

- Populations and Evolution
- Quantitative Methods
- Subcellular Processes
- Tissues and Organs

Text classification

From Wikipedia:

Document classification or document categorization is a problem in library science, information science and computer science. The task is to assign a document to one or more classes or categories. This may be done “manually” (or “intellectually”) or algorithmically. The intellectual classification of documents has mostly been the province of library science, while the algorithmic classification of documents is used mainly in information science and computer science. The problems are overlapping, however, and there is therefore also interdisciplinary research on document classification.

Formal definition

- input data:
 - document d
 - set of labels (classes, categories) $C = \{c_1, c_2, \dots, c_J\}$
- output data:
 - predicted class c of document d

Applications

- document cataloguing
- message filtering
- spam filtering
- language identification
- sentiment analysis
- genre identification (movies/books/music)
- authorship identification (gender, age)
- ... and many more

Methods

Hand-coded rules

- rules are usually expressed by a combination of words or other features
- example of a spam filtering rule:

```
black-list-address OR("dollars" AND "have been selected")
```

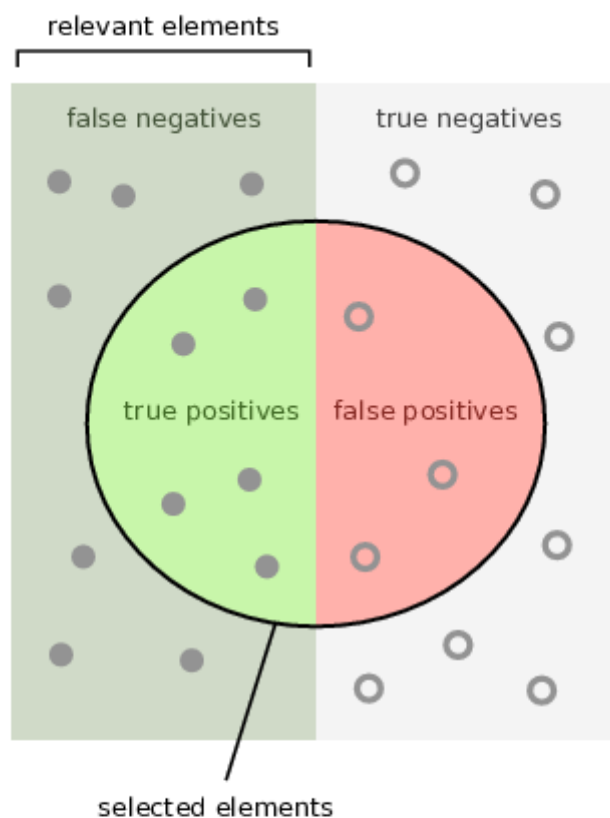
- high accuracy possible
- require continuous updates (by experts)

Supervised machine learning

- input data:
 - document d
 - set of labels (classes) $C = \{c_1, c_2, \dots, c_J\}$
 - training data, i.e. set of documents labeled manually, $(d_1, c_1), \dots, (d_m, c_m)$
- output data:

- learned classifier $\gamma : \mathcal{d} \rightarrow \mathcal{c}$
- possible classifiers:
 - naive Bayes
 - logistic regression
 - support vector machine
 - k nearest neighbors
 - maximum entropy
 - decision trees
 - decision rules (lists)
 - ...

Assessing a classifier



How many selected items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

The performance of the classifier can be also presented in form of a contingency matrix:

	N (predicted)	P (predicted)
N (actual)	True negatives, TN	False positives, fp
P (actual)	False negatives, fn	True positives, TP

Then:

$$Precision = \frac{TP}{TP + fp}$$
$$Recall = \frac{TP}{TP + fn}$$

Other measures are possible:

$$Accuracy = \frac{TP + TN}{N}$$
$$Fallout = \frac{fp}{fp + fn}$$
$$NumOfErrors = fp + fn$$

Naive Bayes classifier

Bayes' theorem

Definition - conditional probability

Conditional probability is a measure of the probability of an event given that (by assumption, presumption, assertion or evidence) another event has occurred,

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Bayes' theorem

Let A_i be a series of events such that $A_i \cap A_j \neq \emptyset$ for $i \neq j$, $\sum_i P(A_i) = 1$ and $P(B) > 0$. Then:

$$P(A_i|B) = \frac{P(B|A_i)P(A_i)}{P(B)}$$

where

$$P(B) = \sum_j P(B|A_j)P(A_j)$$

Example

Suppose there are two bowls of cookies. Bowl 1 contains 30 vanilla cookies and 10 chocolate cookies. Bowl 2 contains 20 of each.

Now suppose you choose one of the bowls at random and, without looking, select a cookie at random. The cookie is vanilla. What is the probability that it came from Bowl 1?

We are looking for conditional probability

$$P(box1|van)$$

It is not obvious how to compute it. On the other hand, the probability of a vanilla cookie given Bowl 1 is much easier to calculate:

$$P(van|box1) = \frac{30}{40} = \frac{3}{4}$$

Now we can use the Bayes' theorem:

$$P(box1|van) = \frac{P(van|box1)P(box1)}{P(van)}$$

In our example we have:

$$P(box1) = \frac{1}{2}$$

$$P(van) = \frac{50}{80} = \frac{5}{8}$$

Thus we get:

$$P(box1|van) = \frac{(1/2)(3/4)}{5/8} = \frac{3}{5}$$

A (slightly) different interpretation

Given a hypothesis H and evidence E , Bayes' theorem states that the relationship between the probability of the hypothesis before getting the evidence and the probability of the hypothesis after getting the evidence is

$$P(H|E) = \frac{P(E|H)P(H)}{P(E)}$$

where:

- $P(H)$ - prior (probability that H is true before knowing E)
- $P(H|E)$ - posterior (probability that H is true after knowing E)
- $P(E)$ - normalization constant (probability of getting E for any H)
- $P(E|H)$ - likelihood (probability of getting E if H)

Naive Bayes classifier

- a classifier based on the Bayes' theorem
- appropriate to multidimensional problems
- "naive" because of the assumption of independence between features
 - very often an oversimplification
 - nevertheless, the performance of the classifier is usually better than expected

Example with structured data

In the following table a training set with customer data of a computer shop is presented:

Id	Age	Income	Graduated	Credit rating	Buyer
1	<30	high	no	good	no
2	<30	high	no	excellent	no
3	[30;40]	high	no	good	yes
4	>40	medium	no	good	yes
5	>40	low	yes	good	yes
6	>40	low	yes	excellent	no
7	[30;40]	low	yes	excellent	yes
8	<30	medium	no	good	no
9	<30	low	yes	good	yes
10	>40	medium	yes	good	yes
11	<30	medium	yes	excellent	yes
12	[30;40]	medium	no	excellent	yes
13	[30;40]	high	yes	good	yes
14	>40	medium	no	excellent	no

Our goal is to check if person X (age under 30, medium income, graduated and with a good credit score) will buy a computer in the shop.

To this end we have to check for which value of i the product

$$P(X|C_i) * P(C_i)$$

has its maximum. Here:

- C_i , $i = 1, 2$ - two categories: will buy or not
- $P(C_i)$ - prior describing the probability of being a member of class C_i

From the training set we have:

$$P(C_1) = \frac{9}{14} = 0,643$$

$$P(C_2) = \frac{5}{14} = 0,357$$

We calculate the conditional probabilities for each attribute values:

$$P(X|C_1) = P(\text{age} = "< 30" | C_1) * P(\text{income} = "medium" | C_1) * P(\text{graduated} = "ye$$

$$P(X|C_2) = P(\text{age} = "< 30" | C_2) * P(\text{income} = "medium" | C_2) * P(\text{graduated} = "ye$$

Based on the data in the training set we get:

$$P(\text{age} = "< 30" | C_1) = \frac{2}{9}, \quad P(\text{income} = "medium" | C_1) = \frac{4}{9}, \quad P(\text{graduated} = "yes$$

$$P(\text{credit} = "good" | C_1) = \frac{6}{9}$$

$$P(\text{age} = "< 30" | C_2) = \frac{3}{5}, \quad P(\text{income} = "medium" | C_2) = \frac{2}{5}, \quad P(\text{graduated} = "yes$$

$$P(\text{income} = "good" | C_2) = \frac{2}{5}$$

Hence

$$P(X|C_1) = 0,044, \quad P(X|C_2) = 0,019$$

and

$$P(X|C_1) * P(C_1) = 0,028$$

$$P(X|C_2) * P(C_2) = 0,007$$

Since the probability is higher for C_1 , we classify the person X as a buyer.

Text classification

- bag of word as text representation
- use reduced representation!!!
- additional assumption --> position of a word in a document does not matter

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet.



I **love** this movie! It's **sweet**, but with **satirical** humor. The dialogue is **great** and the adventure scenes are **fun**... It manages to be **whimsical** and **romantic** while **laughing** at the conventions of the fairy tale genre. I would **recommend** it to just about anyone. I've seen it **several** times, and I'm always **happy** to see it **again** whenever I have a friend who hasn't seen it yet.



x love xxxxxxxxxxxxxxxxxxxx sweet
xxxxxxxx satirical xxxxxxxxxxxx
xxxxxxxxxxxxxx great xxxxxxxx
xxxxxxxxxxxxxxxxxxxxxxxxxxxx fun xxxx
xxxxxxxxxxxxxxxxxxxx whimsical xxxx
romantic xxxx laughing
xx
xxxxxxxxxxxxxxxxxxxxxx recommend xxxxxx
xx
xx several xxxxxxxxxxxxxxxxxxxxxxxxxxxx
xxxxx happy xxxxxxxxxxxx again
xx
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx



great	2
love	2
recommend	1
laugh	1
happy	1
...	...

According to the Bayes' theorem, for a document d and a category c we have

$$P(c|d) = \frac{P(d|c)P(c)}{P(d)}$$

Classification of a document is nothing but finding the category with highest probability (*maximum a posteriori*, MAP):

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(c|d) = \operatorname{argmax}_{c \in C} \frac{P(d|c)P(c)}{P(d)}$$

For the sake of simplicity let us omit the normalization constant:

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(d|c)P(c)$$

A document is simply a set of features, thus:

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n|c)P(c)$$

Problems:

- $O(|X|^n \times |C|)$ parameters
- probability can be estimated only for a large training set

"Solution":

- we assume that for every class c the features x_i are independent, i.e.

$$P(x_1, x_2, \dots, x_n|c) = P(x_1|c) * P(x_2|c) * \dots * P(x_n|c)$$

$$c_{NB} = \operatorname{argmax}_{c \in C} P(c) \prod_{x \in X} P(x|c)$$

- in the first approximation we can take word frequencies as the probabilities of word occurrences

$$\hat{P}(c_j) = \frac{\text{doccount}(C = c_j)}{N_{\text{doc}}}$$

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

Is it enough?

No! It may happen that there was no word 'fantastic' in the training set. Then

$$\hat{P}(\text{"fantastic"} | \text{"positive"}) = \frac{\text{count}(\text{"fantastic"}, \text{"positive"})}{\sum_{w \in V} \text{count}(w, \text{"positive"})} = 0$$

which gives

$$c_{NB} = \operatorname{argmax}_{c \in C} \hat{P}(c) \prod_{x \in X} \hat{P}(x | c) = 0$$

Laplace smoothing

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j) + 1}{\sum_{w \in V} (\text{count}(w, c_j) + 1)}$$

Naive Bayes in NLTK

In [1]:

```
import nltk
```

Name classification

Let us come back once again to the example from the previous lecture - we want to classify names according to gender. First, we read in the corpus:

In [2]:

```
from nltk.corpus import names
names = [(name, 'male') for name in names.words('male.txt')] + [(name, 'female') for name in names.words('female.txt')]
```

In [3]:

```
names[:10]
```

Out[3]:

```
[('Aamir', 'male'),
 ('Aaron', 'male'),
 ('Abbey', 'male'),
 ('Abbie', 'male'),
 ('Abbot', 'male'),
 ('Abbott', 'male'),
 ('Abby', 'male'),
 ('Abdel', 'male'),
 ('Abdul', 'male'),
 ('Abdulkarim', 'male')]
```

In [4]:

```
import random
random.shuffle(names)
names[:10]
```

Out[4]:

```
[('Matthew', 'male'),
 ('Burgess', 'male'),
 ('Mendel', 'male'),
 ('Cordy', 'female'),
 ('Katharyn', 'female'),
 ('Faythe', 'female'),
 ('Sigfried', 'male'),
 ('Mindy', 'female'),
 ('Elora', 'female'),
 ('Shay', 'male')]
```

In order to build a classifier, we need a feature extractor. Suppose the last letter of the name will be used as the feature for classification:

In [5]:

```
def gender_features(word):
    return {'last_letter': word[-1]}
gender_features('Shrek')
```

Out[5]:

```
{'last_letter': 'k'}
```

We will use the extractor function to process the input data and train the classifier:

In [6]:

```
featuresets = [(gender_features(n), gender) for (n, gender) in names]
train_set, test_set = featuresets[500:], featuresets[:500]
print(train_set[:10])
```

```
[({'last_letter': 'a'}, 'female'), ({'last_letter': 'a'}, 'female'),
 ({'last_letter': 'e'}, 'female'), ({'last_letter': 'l'}, 'male'),
 ({'last_letter': 'd'}, 'male'), ({'last_letter': 'e'}, 'female'),
 ({'last_letter': 'd'}, 'male'), ({'last_letter': 'a'}, 'female'),
 ({'last_letter': 'l'}, 'female'), ({'last_letter': 'a'}, 'female')]
```

In [7]:

```
classifier = nltk.NaiveBayesClassifier.train(train_set)
```

Let us check how the classifier works on names that are not in the corpus:

In [8]:

```
classifier.classify(gender_features('Neo'))
```

Out[8]:

'male'

In [9]:

```
classifier.classify(gender_features('Trinity'))
```

Out[9]:

'female'

We may use the test set to check the accuracy of the classifier:

In [10]:

```
print(nltk.classify.accuracy(classifier, test_set))
```

0.766

Checking of the most important features is possible as well:

In [11]:

```
classifier.show_most_informative_features(5)
```

Most Informative Features

5 : 1.0	last_letter = 'a'	female : male =	35.
6 : 1.0	last_letter = 'k'	male : female =	28.
0 : 1.0	last_letter = 'f'	male : female =	14.
0 : 1.0	last_letter = 'p'	male : female =	12.
7 : 1.0	last_letter = 'd'	male : female =	10.

Let us do some serious testing:

In [12]:

```
from nltk.metrics.scores import precision, recall
```

In [13]:

```
import collections
from nltk.metrics.scores import precision, recall
from nltk.classify import NaiveBayesClassifier

refsets = collections.defaultdict(set)
testsets = collections.defaultdict(set)

for i, (feats, label) in enumerate(test_set):
    refsets[label].add(i)
    observed = classifier.classify(feats)
    testsets[observed].add(i)

print('male precision:', precision(refsets['male'], testsets['male']))
print('male recall:', recall(refsets['male'], testsets['male']))
print('female precision:', precision(refsets['female'], testsets['female']))
print('female recall:', recall(refsets['female'], testsets['female']))

male precision: 0.7248677248677249
male recall: 0.6782178217821783
female precision: 0.7909967845659164
female recall: 0.825503355704698
```

If we are not happy with the performance, we can extract a different feature set:

In [14]:

```
def gender_features2(name):
    features = {}
    features["first_letter"] = name[0].lower()
    features["last_letter"] = name[-1].lower()
    for letter in 'abcdefghijklmnopqrstuvwxyz':
        features["count({})".format(letter)] = name.lower().count(letter)
        features["has({})".format(letter)] = (letter in name.lower())
    return features
```

In [15]:

```
featuresets = [(gender_features2(n), gender) for (n, gender) in names]
train_set, test_set = featuresets[500:], featuresets[:500]
classifier = nltk.NaiveBayesClassifier.train(train_set)
print(nltk.classify.accuracy(classifier, test_set))
```

0.746

It may happen that it decreases after adding new features, especially in case of small and noisy dataset (*overfitting*). The choice of the feature set is crucial and usually one uses trial and error to find the optimal features.

Document classification

Let us consider now the movie review corpus included in NLTK:

In [16]:

```
from nltk.corpus import movie_reviews
documents = [(list(movie_reviews.words(fileid)), category)
              for category in movie_reviews.categories()
              for fileid in movie_reviews.fileids(category)]
random.shuffle(documents)
print(documents[0])
```

(['this', 'is', 'one', 'of', 'the', 'worst', 'big', '-', 'screen', 'film', 'experiences', 'i', '"', 've', 'had', 'for', 'a', 'while', '.', 'with', 'this', 'film', ',', 'plus', '`', 'showgirls', '"', 'a nd', '`', 'basic', 'instinct', '"', ',', 'paul', 'verhoeven', 'has', 'stamped', 'himself', 'as', 'currently', 'one', 'of', 'the', 'worst', 'blockbuster', 'directors', '.', 'his', 'celebrated', 'film', '`', 'total', 'recall', '"', 'was', '?', 'i', 'admit', '?', 'successfully', 'scripted', ',', 'but', 'it', 'nonetheless', 'contained', 'directorial', 'flaws', '.', 'obviously', 'nobody', 'wanted', 'to', 'invest', 'too', 'much', 'money', 'in', 'a', 'production', 'from', 'someone', 'like', 'verhoeven', ',', 'the', 'result', 'being', 'that', 'much', 'of', 'the', 'special', 'effects', 'in', '`', 'starship', '"', 'seemed', 'fake', '.', 'but', 'not', 'everything', 'bad', 'in', 'the', 'film', 'was', 'the', 'director', '"', 's', 'fault', ',', 'even', 'though', 'he', 'was', 'one', 'of', 'the', 'guys', 'who', 'employed', 'the', 'actors', '.', 'it', 'is', 'surprising', 'that', 'none', 'of', 'the', 'actors', 'received', 'nominations', 'for', 'the', 'razzie', 'awards', '(', 'i', 'expected', 'five', 'for', 'the', 'acting', 'categories', ')', '.', 'casserole', 'vanity', 'devoid', ',', 'dense', 'ribald', ',', 'dingy', 'miasma', ',', 'and', 'jackass', 'bushy', 'are', 'in', 'serious', 'need', 'of', 'acting', 'school', '.', 'no', ',', 'they', 'have', 'to', 'pass', 'primary', 'school', 'drama', 'classes', 'first', '.', 'while', '`', 'total', '"', 'was', 'written', 'well', ',', 'starship', '"', 'is', 'purely', 'pathetic', '.', 'all', 'right', ',', 'it', 'is', 'supposed', 'to', 'be', 'a', 'fast', '-', 'paced', 'entertainment', 'film', ',', 'and', 'you', '"', 're', 'supposed', 'to', 'turn', 'off', 'your', 'intellect', '(', 'completely', ')', 'and', 'enjoy', 'the', 'action', 'sequences', 'and', 'special', 'effects', '(', 'that', 'is', ',', 'guts', 'and', 'gore', ')', '.', 'as', 'a', 'matter', 'of', 'fact', ',', 'i', 'found', 'the', 'activity', 'incredibly', 'boring', ',', 'a', 'complete', 'waste', 'of', 'more', 'than', 'two', 'hours', '.', 'half', 'of', 'the', 'film', 'was', 'a', 'bad', 'episode', 'of', '`', 'beverly', 'hills', '90210', '"', '(', 'dina', 'meyer', 'was', 'in', '"', 'beverly', 'hills', '"', ')', ',', 'while', 'another', 'quarter', 'was', 'simply', 'nothing', '(', 'things', 'like', 'presenting', 'irrelevant', 'information', 'in', 'an', 'irritating', 'way', 'on', 'the', 'web', ')', ',', 'and', 'the', 'rest', 'was', 'a', 'display', 'of', 'humans', 'fighting', 'computer', '-', 'generated', 'images', '.', 'the', 'battles', 'were', 'all', 'the', 'same', '?', 'jumping', 'around', ',', 'shoot', 'or', 'get', 'stabbed', '?', 'and', 'on', 'barren', 'planets', 'that', 'only', 'had', 'giant', 'insects', '.', 'there', 'weren', '"', 't', 'even', 'any', 'stunts', ',', 'which', 'i', 'consider', 'slightly', 'more', 'exciting', 'than', 'pictures', 'running', 'around', '.', 'i', 'wonder', 'what', 'the', 'insects', 'eat', ',', 'if', 'there', '"', 's', 'nothing', 'but', 'them', 'on', 'the', 'planets', '?', 'there', 'is', 'so', 'much', 'laughable', 'treatment', 'in', 'this', 'film', ',', 'and', 'it', 'is', 'frankly', 'not', 'amusing', 'when', 'jokes', 'are', 'intended', '.', 'this', 'type', 'of', 'story', 'is', 'obviously', 'aimed', 'at', '10', '-', 'year', '-', 'olds', ',', 'who', 'can', '"', 't', 'see', 'it', 'anyway', 'because', 'of', 'the', 'violence', 'and', 'some', 'sexuality', '.', 'but', 'then', ',', 'there', 'are', 'always', '16', '-', 'year', '-', 'olds', 'who', 'have', 'that', 'frame', 'of', 'mind', '.', 'the', 'pointless', 'plot', 'begins', 'when', 'johnny', '"', 's', '(', 'vanity', 'devoid', ')', 'girlfriend', 'carmen', '(', 'richards', ')', 'decides', 'that', 'she', 'wants', 'to', 'join', 'the', 'troopers', 'to', 'fight', 'the', 'insects', 'who', 'are', 'throwing', 'asteroids', 'at', 'earth', '.', 'johnny', 'then', 'signs', 'up', 'as', 'a', 'trooper', 'also', ',', 'after', 'an',

'overacted', 'argument', 'with', 'his', 'parents', '.', 'but', 'there', 'is', 'another', 'girl', ',', 'dizzy', '(', 'meyer', ')',
 ',', 'who', 'likes', 'johnny', 'and', 'then', 'there', 'is', 'another', 'boy', 'who', 'likes', 'carmen', ',', 'which', 'results', 'in', 'a', 'love', 'quadrangle', ',', 'which', 'isn', '"', 't', 'better', ',', 'because', 'it', 'means', 'augmented', 'worse', '-', 'than', '-', 'stereotyped', 'soap', 'opera', ',', 'increased', 'bitchiness', ',', 'and', 'more', 'bad', 'beverly', 'hills', '+', 'melrose', '.', 'and', 'the', 'result', 'of', 'this', 'love', 'quadrangle', 'at', 'the', 'end', 'is', 'also', 'rather', 'stupid', '.', 'anyway', ',', 'getting', 'back', 'to', 'the', 'thing', 'you', 'might', 'call', 'plot', ',', 'johnny', 'is', 'too', 'stupid', 'to', 'be', 'a', 'pilot', 'and', 'has', 'to', 'join', 'the', 'infantry', ',', 'while', 'his', 'girl', 'and', 'the', 'other', 'dude', 'are', 'in', 'the', 'same', 'league', '.', 'dizzy', 'comes', 'chasing', 'johnny', 'and', 'joins', 'the', 'infantry', 'also', '.', 'they', 'then', 'start', 'training', ',', 'which', 'contains', 'what', 'roger', 'hebert', 'calls', 'ips', '(', 'idiot', 'plot', 'syndrome', ',', 'moments', 'when', 'only', 'an', 'idiot', 'would', 'have', 'made', 'such', 'obvious', 'mistakes', ')', ',', 'then', 'real', 'combat', '.', 'and', 'guess', 'what', '?', 'that', '"', 's', 'about', 'as', 'complex', 'as', 'it', 'gets', '.', 'oh', ',', 'and', 'one', 'of', 'their', 'friends', ',', 'carl', '(', 'neil', 'patrick', 'harris', 'a', '.', 'k', '.', 'a', '.', 'doogie', 'howser', 'm', '.', 'd', ',', ')', ',', 'becomes', 'involved', 'in', 'war', 'intelligence', ',', 'and', 'his', 'abilities', 'at', 'the', 'end', 'are', 'really', 'corny', 'and', 'make', 'me', 'want', 'to', 'spray', 'insecticide', 'on', 'someone', 'for', 'it', '.', 'he', '"', 's', 'my', 'fifth', 'nom', 'for', 'a', 'razzie', '.', 'the', 'troopers', 'fight', ',', 'fall', 'in', 'love', ',', 'die', ',', 'kill', ',', 'and', 'try', 'to', 'act', '.', 'naturally', ',', 'they', 'win', ',', 'or', 'sort', 'of', 'half', '-', 'win', '.', 'of', 'course', ',', 'in', 'between', '(', 'and', 'at', 'the', 'end', ')', 'there', 'are', 'soldiers', 'chatting', 'and', 'smiling', 'while', 'carrying', 'grievous', 'wounds', 'caused', 'by', 'bug', 'legs', '.', 'the', 'bugs', 'also', 'suffer', 'from', 'ips', ':', 'why', 'would', 'you', 'release', 'hold', 'of', 'your', 'captive', 'before', 'killing', 'it', '?', 'isn', '"', 't', 'it', 'also', 'amazing', 'that', 'earthlings', 'haven', '"', 't', 'invented', 'better', 'hand', '-', 'held', 'weapons', 'by', 'then', '?', 'the', 'only', 'question', 'that', 'remains', 'is', 'why', 'i', 'gave', 'it', 'one', 'star', 'instead', 'of', 'zero', '.', 'well', ',', 'maybe', 'a', '-', 'quarter', '(', 'of', 'a', 'star', ')', 'for', 'the', 'originality', 'of', 'the', 'co', '-', 'sex', 'shower', 'scene', 'and', 'the', '(', 'very', ')', 'brief', 'moments', 'of', 'suspense', ',', 'another', 'quarter', 'for', 'copying', ',', 'zulu', '"', 'and', 'letting', 'the', 'good', 'guys', '(', 'the', 'bugs', ')', 'win', ',', 'and', 'half', 'a', 'star', 'for', 'the', 'sucking', '-', 'out', 'of', 'the', 'brain', 'of', 'one', 'of', 'those', 'people', 'who', 'call', 'themselves', 'actors', '(', 'but', 'there', 'should', 'have', 'been', 'more', ',', 'the', 'troopers', 'deserved', 'to', 'die', ')', '.,', 'ne g')

Again, we define a helper function to extract the features:

In [17]:

```
all_words = nltk.FreqDist(w.lower() for w in movie_reviews.words())
word_features = list(all_words)[:2000]

def document_features(document):
    document_words = set(document)
    features = {}
    for word in word_features:
        features['contains({})'.format(word)] = (word in document_words)
    return features
```

Let us check how this function works:

In [18]:

```
print(document_features(movie_reviews.words('pos/cv957_8737.txt')))
```

```
{'contains(six)': False, 'contains(8034)': False, 'contains(cityscape)': False, 'contains(eloquence)': False, 'contains(alain)': False, 'contains(cheeky)': False, 'contains(oil)': False, 'contains(milla)': False, 'contains(deficiency)': False, 'contains(dots)': False, 'contains(gorefest)': False, 'contains(sonorra)': False, 'contains(dusting)': False, 'contains(meaning)': False, 'contains(excitement)': False, 'contains(hides)': False, 'contains(mergers)': False, 'contains(bucket)': False, 'contains(chemo)': False, 'contains(fullfillment)': False, 'contains(separation)': False, 'contains(accomplishment)': False, 'contains(offence)': False, 'contains(wealth)': False, 'contains(shocker)': False, 'contains(mazzello)': False, 'contains(kitschiest)': False, 'contains(woman)': True, 'contains(immortal)': False, 'contains(casserole)': False, 'contains(kettle)': False, 'contains(lavatorial)': False, 'contains(confidence)': False, 'contains(apparent)': False, 'contains(nationalist)': False, 'contains(rake)': False, 'contains(patch)': False, 'contains(edgar)': False, 'contains(prioritized)': False, 'contains(premeditated)': False, 'contains(jaunts)': False, 'contains(unspecified)': False, 'contains(slayers)': False, 'contains(kewpies)': False, 'contains(decomposing)': False, 'contains(evidencing)': False, 'contains(scuffle)': False, 'contains(ploys)': False, 'contains(deane)': False, 'contains(drunks)': False, 'contains(problem)': False, 'contains(endemic)': False, 'contains(douglas)': False, 'contains(coloreds)': False, 'contains(hairline)': False, 'contains(actosta)': False, 'contains(stairwells)': False, 'contains(timm)': False, 'contains(galas)': False, 'contains(doctors)': False, 'contains(misappropriating)': False, 'contains(starship)': False, 'contains(spouses)': False, 'contains(traumatic)': False, 'contains(prescribed)': False, 'contains(shwarzenegger)': False, 'contains(tsi)': False, 'contains(snipe)': False, 'contains(fashioning)': False, 'contains(disliking)': False, 'contains(handwaving)': False, 'contains(deed)': False, 'contains(flot sam)': False, 'contains(jimkendrick)': False, 'contains(penn)': False, 'contains(rave)': False, 'contains(loeb)': False, 'contains(whalberg)': False, 'contains(emannuelle)': False, 'contains(operational)': False, 'contains(motivate)': False, 'contains(imitation)': False, 'contains(tripp)': False, 'contains(smartass)': False, 'contains(thereus)': False, 'contains(goggins)': False, 'contains(folded)': False, 'contains(humbled)': False, 'contains(direct)': False, 'contains(godess)': False, 'contains(midland)': False, 'contains(verging)': False, 'contains(dreamworld)': False, 'contains(uncharacteristic)': False, 'contains(reptilian)': False, 'contains(unreal)': False, 'contains(fudd)': False, 'contains(uninterrupted)': False, 'contains(exacted)': False, 'contains(awakens)': False, 'contains(complimented)': False, 'contains(indentity)': False, 'contains(hairbrush)': False, 'contains(naoki)': False, 'contains(enthusiastic)': False, 'contains(shard)': False, 'contains(inanely)': False, 'contains(patronize)': False, 'contains(birdy)': False, 'contains(affordable)': False, 'contains(blurts)': False, 'contains(tragically)': False, 'contains(holly)': False, 'contains(understatement)': False, 'contains(tonya)': False, 'contains(narrowed)': False, 'contains(asked)': False, 'contains(bash)': False, 'contains(dancefloor)': False, 'contains(stooped)': False, 'contains(coutroom)': False, 'contains(vaughns)': False, 'contains(ferraris)': False, 'contains(applesauce)': False, 'contains(mesa)': False, 'contains(classifies)': False, 'contains(elspeth)': False, 'contains(pagniacchi)': False, 'contains(cate)': False, 'contains(swampy)': False, 'contains(rhythm)': False, 'contains(hollywood)': False, 'contains(remaining)': False, 'contains(godard)': False, 'contains(solemnly)': False, 'contains(kershner)': False, 'contains(organizing)': False, 'contains(siblings)': False, 'contains(arcs)': False, 'contains(runs)': False, 'contains(career)': False, 'contains(izabella)': False, 'contains(mat)': False}
```

se, 'contains(errupts)': False, 'contains(corresponding)': False, 'contains(nero)': False, 'contains(depends)': False, 'contains(mal kovich)': False, 'contains(sanctimonious)': False, 'contains(guitar)': False, 'contains(girlfriend)': True, 'contains(staged)': False, 'contains(ezsterhas)': False, 'contains(etch)': False, 'contains(counterprogramming)': False, 'contains(distrusts)': False, 'contains(encapsulates)': False, 'contains(parodizing)': False, 'contains(agent)': False, 'contains(raider)': False, 'contains(faintly)': False, 'contains(shoah)': False, 'contains(insurrection)': False, 'contains(affair)': False, 'contains(funky)': False, 'contains(final s)': False, 'contains(herbie)': False, 'contains(grindhouse)': False, 'contains(oilrig)': False, 'contains(chemisty)': False, 'contains(houseboat)': False, 'contains(belinda)': False, 'contains(disarming)': False, 'contains(forbes)': False, 'contains(gratuitous)': False, 'contains(carlotta)': False, 'contains(affront)': False, 'contains(rhythm)': False, 'contains(skarsgard)': False, 'contains(xerox)': False, 'contains(synonyms)': False, 'contains(emotional)': False, 'contains(baseman)': False, 'contains(stirred)': False, 'contains(riot)': False, 'contains(shinohara)': False, 'contains(vaginal)': False, 'contains(manga)': False, 'contains(newfound)': False, 'contains(ramses)': False, 'contains(blinding)': False, 'contains(occult)': False, 'contains(jester)': False, 'contains(preteen)': False, 'contains(brigitte)': False, 'contains(credit)': False, 'contains(mother)': False, 'contains(flood)': False, 'contains(picturesque)': False, 'contains(thrusts)': False, 'contains(positronic)': False, 'contains(crush)': False, 'contains(peculiarities)': False, 'contains(recommendation)': False, 'contains(bathroom)': False, 'contains(giles)': False, 'contains(competent)': False, 'contains(presided)': False, 'contains(tenuous)': False, 'contains(schellenberg)': False, 'contains(lancaster)': False, 'contains(performances)': False, 'contains(kilmer)': False, 'contains(dah)': False, 'contains(channing)': False, 'contains(conference)': False, 'contains(1847)': False, 'contains(healy)': False, 'contains(usurper)': False, 'contains(cane)': False, 'contains(circles)': False, 'contains(zanin)': False, 'contains(fake)': False, 'contains(burying)': False, 'contains(reserve)': False, 'contains(harmonizing)': False, 'contains(interviewers)': False, 'contains(flabby)': False, 'contains(thoughtfulness)': False, 'contains(effects)': False, 'contains(skeletal)': False, 'contains(Animator)': False, 'contains(spineless)': False, 'contains(bogus)': False, 'contains(gaius)': False, 'contains(timberland)': False, 'contains(hunk)': False, 'contains(confrontatory)': False, 'contains(meting)': False, 'contains(smacking)': False, 'contains(footprints)': False, 'contains(lifeforce)': False, 'contains(prefecture)': False, 'contains(ambiguities)': False, 'contains(pyrotechnic)': False, 'contains(kewpie)': False, 'contains(zookeepers)': False, 'contains(wetmore)': False, 'contains(pushy)': False, 'contains(salami)': False, 'contains(journalists)': False, 'contains(respond)': False, 'contains(gently)': False, 'contains(uppington)': False, 'contains(&#)': False, 'contains(aspirins)': False, 'contains(sedated)': False, 'contains(knott)': False, 'contains(enchanted)': False, 'contains(ship)': False, 'contains(cotterill)': False, 'contains(lollipop)': False, 'contains(lockers)': False, 'contains(violently)': False, 'contains(turf)': False, 'contains(wookie)': False, 'contains(summarizing)': False, 'contains(menage)': False, 'contains(prague)': False, 'contains(delaurentiis)': False, 'contains(insinuations)': False, 'contains(knees)': False, 'contains(goldberg)': False, 'contains(interlinked)': False, 'contains(lulu)': False, 'contains(dga)': False, 'contains(speed2)': False, 'contains(ackland)': False, 'contains(yawner)': False, 'contains(cking)': False, 'contains(cyanide)': False, 'contains(ballyhooed)': False, 'contains(morale)': False, 'contains(wrongdoings)': False, 'contains(rattigan)': F

also, 'contains(remaking)': False, 'contains(pairing)': False, 'contains(bonkers)': False, 'contains(duarte)': False, 'contains(facing)': False, 'contains(gifford)': False, 'contains(superfluous)': False, 'contains(renfro)': False, 'contains(sing)': False, 'contains(prescribed)': False, 'contains(dell)': False, 'contains(hostile)': False, 'contains(contemplating)': False, 'contains(unearths)': False, 'contains(unbeknownst)': False, 'contains(befall)': False, 'contains(garlic)': False, 'contains(attending)': False, 'contains(maybe)': False, 'contains(phillippe)': False, 'contains(complain)': False, 'contains(nuptials)': False, 'contains(rewatch)': False, 'contains(blamed)': False, 'contains(unexpected)': False, 'contains(breeze)': False, 'contains(deception)': False, 'contains(stayin)': False, 'contains(incongruities)': False, 'contains(fictionalized)': False, 'contains(hermit)': False, 'contains(giveaways)': False, 'contains(hypnotic)': False, 'contains(hathaway)': False, 'contains(upholds)': False, 'contains(rabal)': False, 'contains(ask)': False, 'contains(website)': False, 'contains(surrounding)': False, 'contains(constipated)': False, 'contains(retread)': False, 'contains(dietz)': False, 'contains(stinks)': False, 'contains(ze)': False, 'contains(waters)': False, 'contains(gooders)': False, 'contains(talkie)': False, 'contains(smeared)': False, 'contains(octopussy)': False, 'contains(heart)': False, 'contains(cohesiveness)': False, 'contains(worded)': False, 'contains(gators)': False, 'contains(robob)': False, 'contains(protruding)': False, 'contains(sire)': False, 'contains(mouthful)': False, 'contains(infection)': False, 'contains(dong)': False, 'contains(brainard)': False, 'contains(lisbon)': False, 'contains(tally)': False, 'contains(plops)': False, 'contains(predilection)': False, 'contains(costumed)': False, 'contains(gershwin)': False, 'contains(davidovitch)': False, 'contains(bender)': False, 'contains(conceivably)': False, 'contains(arrested)': False, 'contains(offender)': False, 'contains(glee)': False, 'contains(puffing)': False, 'contains(moviereviews)': False, 'contains(stucker)': False, 'contains(schmaltzfest)': False, 'contains(booked)': False, 'contains(1773)': False, 'contains(skepticism)': False, 'contains(love)': False, 'contains(recipes)': False, 'contains(juliet)': False, 'contains(evidenced)': False, 'contains(tel)': False, 'contains(callisto)': False, 'contains(mordantly)': False, 'contains(meets)': False, 'contains(niagara)': False, 'contains(dashing)': False, 'contains(streetwalker)': False, 'contains(lagpacan)': False, 'contains(pollsters)': False, 'contains(snyder)': False, 'contains(inexpensive)': False, 'contains(ideology)': False, 'contains(connects)': False, 'contains(pausing)': False, 'contains(attainable)': False, 'contains(glen)': False, 'contains(saim)': False, 'contains(kasa)': False, 'contains(permissiveness)': False, 'contains(punchline)': False, 'contains(file)': False, 'contains(desensitization)': False, 'contains(leathers)': False, 'contains(slaughtering)': False, 'contains(tucci)': False, 'contains(spraying)': False, 'contains(conjunctions)': False, 'contains(ore)': False, 'contains(beltrami)': False, 'contains(collectively)': False, 'contains(jemmons)': False, 'contains(iris)': False, 'contains(impregnating)': False, 'contains(vaccinate)': False, 'contains(ahh)': False, 'contains(evacuation)': False, 'contains(copulate)': False, 'contains(archetypal)': False, 'contains(scrawny)': False, 'contains(ballstein)': False, 'contains(gasp)': True, 'contains(heavy)': False, 'contains(plains)': False, 'contains(diametric)': False, 'contains(gehrig)': False, 'contains(braveheart)': False, 'contains(smile)': False, 'contains(kristen)': False, 'contains(derelect)': False, 'contains(lee)': False, 'contains(intellectual)': False, 'contains(twitching)': False, 'contains(tasha)': False, 'contains(fighting)': False, 'contains(sadler)': False, 'contains(snoozing)': False, 'contains(pedophile)': False, 'contains(obtuse)': False, 'contains(dirtier)': False, 'contains(excepti

ons)': False, 'contains(curiosities)': False, 'contains(schweethaa
t)': False, 'contains(singular)': False, 'contains(deloreans)': Fal
se, 'contains(penniless)': False, 'contains(normandy)': False, 'con
tains(deadlier)': False, 'contains(singers)': False, 'contains(tod
d)': False, 'contains(vore)': False, 'contains(pollster)': False,
'contains(twirling)': False, 'contains(erupt)': False, 'contains(m
ickelberry)': False, 'contains(reminicent)': False, 'contains(geeki
sh)': False, 'contains(envies)': False, 'contains(doodles)': False,
'contains(investigation)': False, 'contains(manage)': False, 'conta
ins(ditches)': False, 'contains(yong)': False, 'contains(lustre)':
False, 'contains(brighter)': False, 'contains(nhu)': False, 'conta
ins(faze)': False, 'contains(extracurricular)': False, 'contains(ex
perience)': False, 'contains(rejoiced)': False, 'contains(preferabl
y)': False, 'contains(distorts)': False, 'contains(remington)': Fal
se, 'contains(rivals)': False, 'contains(films)': False, 'contains
(funnyman)': False, 'contains(restauranteur)': False, 'contains(pro
jectionist)': False, 'contains(goop)': False, 'contains(jena)': Fal
se, 'contains(redneck)': False, 'contains(writes)': False, 'contain
s(depreciating)': False, 'contains(trenchcoat)': False, 'contains(m
ambos)': False, 'contains(science)': False, 'contains(pysche)': Fal
se, 'contains(lowdown)': False, 'contains(occupied)': False, 'conta
ins(crabs)': False, 'contains(philes)': False, 'contains(crushed)':
False, 'contains(congregation)': False, 'contains(proportions)': Fa
lse, 'contains(odin)': False, 'contains(entitlement)': False, 'cont
ains(tarnish)': False, 'contains(crooning)': False, 'contains(abeya
nce)': False, 'contains(traditionally)': False, 'contains(aristocra
tic)': False, 'contains(striding)': False, 'contains(transvestite
s)': False, 'contains(degree)': False, 'contains(forays)': False,
'contains(retrospectives)': False, 'contains(computerized)': Fals
e, 'contains(failure)': False, 'contains(besson)': False, 'contains
(appreciative)': False, 'contains(misquote)': False, 'contains(vaug
hn)': False, 'contains(minbo)': False, 'contains(oversee)': False,
'contains(thermal)': False, 'contains(dip)': False, 'contains(dozie
r)': False, 'contains(freakish)': False, 'contains(screwball)': Fa
lse, 'contains(grrrl)': False, 'contains(beaulieu)': False, 'contai
ns(glean)': False, 'contains(iconography)': False, 'contains(rela
y)': False, 'contains(manic)': False, 'contains(drab)': False, 'con
tains(specifically)': False, 'contains(bigwigs)': False, 'contains
(souza)': False, 'contains(rainbows)': False, 'contains(mommies)':
False, 'contains(flexibility)': False, 'contains(bedrock)': False,
'contains(intermediary)': False, 'contains(warburton)': False, 'con
tains(sequelitis)': False, 'contains(untold)': False, 'contains(isl
amic)': False, 'contains(demean)': False, 'contains(aiming)': Fals
e, 'contains(melaine)': False, 'contains(transylvania)': False, 'co
ntains(contributions)': False, 'contains(skimming)': False, 'contai
ns(thus)': False, 'contains(academia)': False, 'contains(gascogn
e)': False, 'contains(hackers)': False, 'contains(favreau)': False,
'contains(francois)': False, 'contains(darlene)': False, 'contains
(theft)': False, 'contains(allied)': False, 'contains(h2k)': False,
'contains(amass)': False, 'contains(underlining)': False, 'contains
(sticking)': False, 'contains(fourth)': False, 'contains(atrice)':
False, 'contains(abdomen)': False, 'contains(ariane)': False, 'con
tains(massacre)': False, 'contains(avenging)': False, 'contains(mac
ht)': False, 'contains(tokens)': False, 'contains(_disturbing_behav
ior_)': False, 'contains(boggling)': False, 'contains(go)': False,
'contains(revealing)': False, 'contains(wrights)': False, 'contain
s(biting)': False, 'contains(lovlier)': False, 'contains(renfo)': F
alse, 'contains(unwraps)': False, 'contains(deems)': False, 'contai
ns(motivation)': False, 'contains(benben)': False, 'contains(state
s)': False, 'contains(rancorously)': False, 'contains(febre)': Fals
e, 'contains(prude)': False, 'contains(relieves)': False, 'contains

(diluting)': False, 'contains(dredge)': False, 'contains(sickboy)': False, 'contains(reason)': False, 'contains(allah)': False, 'contains(underscored)': False, 'contains(naunton)': False, 'contains(bangkok)': False, 'contains(pattern)': False, 'contains(antisocial)': False, 'contains(advancements)': False, 'contains(thankless)': False, 'contains(commercialist)': False, 'contains(joel)': False, 'contains(pubescent)': False, 'contains(abbott)': False, 'contains(internalizing)': False, 'contains(neckbraces)': False, 'contains(virtuosity)': False, 'contains(pacing)': False, 'contains(consumerism)': False, 'contains(swann)': False, 'contains(cassette)': False, 'contains(polly)': False, 'contains(blood)': False, 'contains(possum)': False, 'contains(enlightenment)': False, 'contains(wagnerian)': False, 'contains(grainy)': False, 'contains(laconically)': False, 'contains(horseshoes)': False, 'contains(fatass)': False, 'contains(glossed)': False, 'contains(wilson)': False, 'contains(kaaya)': False, 'contains(terminating)': False, 'contains(glo)': False, 'contains(walk)': False, 'contains(suing)': False, 'contains(mil)': False, 'contains(god)': False, 'contains(stroll)': False, 'contains(chronicles)': False, 'contains(94)': False, 'contains(fedoras)': False, 'contains(modernized)': False, 'contains(props)': False, 'contains(concerning)': False, 'contains(gunmen)': False, 'contains(labeling)': False, 'contains(batarangs)': False, 'contains(augusts)': False, 'contains(trained)': False, 'contains(miscalculation)': False, 'contains(posits)': False, 'contains(henslowe)': False, 'contains(suppression)': False, 'contains(blueprints)': False, 'contains(attic)': False, 'contains(uninhabited)': False, 'contains(galactica)': False, 'contains(paperbacks)': False, 'contains(quoc)': False, 'contains(miniscule)': False, 'contains(domestically)': False, 'contains(smolder)': False, 'contains(professors)': False, 'contains(easier)': False, 'contains(protective)': False, 'contains(slot)': False, 'contains(scrubbers)': False, 'contains(earthquakes)': False, 'contains(craig)': False, 'contains(characteristically)': False, 'contains(timidity)': False, 'contains(inadequate)': False, 'contains(pitching)': False, 'contains(pets)': False, 'contains(dedication)': False, 'contains(analyzation)': False, 'contains(gone)': False, 'contains(_huge_)': False, 'contains(mistake)': False, 'contains(slots)': False, 'contains(remix)': False, 'contains(vinegar)': False, 'contains(joker)': False, 'contains(mckee)': False, 'contains(doffs)': False, 'contains(hostility)': False, 'contains(_it)': False, 'contains(vcrs)': False, 'contains(precursor)': False, 'contains(boobies)': False, 'contains(infantry)': False, 'contains(episodes)': False, 'contains(macneill)': False, 'contains(chewed)': False, 'contains(permeated)': False, 'contains(vastly)': False, 'contains(jarring)': False, 'contains(stymied)': False, 'contains(ugly)': False, 'contains(korner)': False, 'contains(bela)': False, 'contains(september)': False, 'contains(vamps)': False, 'contains(trade)': False, 'contains(plex)': False, 'contains(lurks)': False, 'contains(kidlets)': False, 'contains(skyjacked)': False, 'contains(reading)': False, 'contains(zookeeper)': False, 'contains(glimcher)': False, 'contains(uproiously)': False, 'contains(perplexed)': False, 'contains(feigns)': False, 'contains(heartbreak)': False, 'contains(untapped)': False, 'contains(emergency)': False, 'contains(dashiell)': False, 'contains(pertinent)': False, 'contains(aladdin)': False, 'contains(sermon)': False, 'contains(fridge)': False, 'contains(veronica)': False, 'contains(gloat)': False, 'contains(geisha)': False, 'contains(murray)': False, 'contains(interchange)': False, 'contains(apprehensively)': False, 'contains(transpired)': False, 'contains(paltry)': False, 'contains(flabbergastingly)': False, 'contains(excitable)': False, 'contains(recreation)': False, 'contains(settles)': False, 'contains(foredoomed)': False, 'contains(textures)': False, 'contains(quote)': False, 'contains(prototype)': False, 'contains(men_)': False

se, 'contains(stead)': False, 'contains(masturbator)': False, 'contains(orgasms)': False, 'contains(nba)': False, 'contains(trustworthy)': False, 'contains(decadenet)': False, 'contains(debatable)': False, 'contains(raving)': False, 'contains(eflman)': False, 'contains(smug)': False, 'contains(taunts)': False, 'contains(gadfly)': False, 'contains(unspool)': False, 'contains(seann)': False, 'contains(mailing)': False, 'contains(faulted)': False, 'contains(gather)': False, 'contains(rubbish)': False, 'contains(heartland)': False, 'contains(escorted)': False, 'contains(insistance)': False, 'contains(gazillionaire)': False, 'contains(pual)': False, 'contains(evelyn)': False, 'contains(containers)': False, 'contains(ragtag)': False, 'contains(shannon)': False, 'contains(schmeer)': False, 'contains(pinon)': False, 'contains(levine)': False, 'contains(darryl)': False, 'contains(invaders)': False, 'contains(institutionalized)': False, 'contains(bodyguards)': False, 'contains(cousine)': False, 'contains(254)': False, 'contains(permanant)': False, 'contains(redeveloped)': False, 'contains(adaptable)': False, 'contains(ditz)': False, 'contains(nephew)': False, 'contains(enlightening)': False, 'contains(propose)': False, 'contains(indignity)': False, 'contains(stodge)': False, 'contains(burkettsville)': False, 'contains(charecter)': False, 'contains(hideousness)': False, 'contains(dishonor)': False, 'contains(flack)': False, 'contains(waging)': False, 'contains(mandated)': False, 'contains(encased)': False, 'contains(blossoming)': False, 'contains(tsunami)': False, 'contains(loathed)': False, 'contains(mute)': False, 'contains(byplay)': False, 'contains(hologado)': False, 'contains(mon)': False, 'contains(membership)': False, 'contains(ruth)': False, 'contains(persuade)': False, 'contains(terrorize)': False, 'contains(realize)': False, 'contains(gilligan)': False, 'contains(languishing)': False, 'contains(partition)': False, 'contains(motormouth)': False, 'contains(insight)': False, 'contains(prevue)': False, 'contains(mail)': False, 'contains(mishandled)': False, 'contains(powered)': False, 'contains(intermingling)': False, 'contains(cgi)': False, 'contains(dellaplane)': False, 'contains(lrighly)': False, 'contains(porcelain)': False, 'contains(ultimo)': False, 'contains(hateful)': False, 'contains(bridget)': False, 'contains(1985)': False, 'contains(overjoyed)': False, 'contains(grouper)': False, 'contains(asswhole)': False, 'contains(overlight)': False, 'contains(mesmerize)': False, 'contains(swordsman)': False, 'contains(loophole)': False, 'contains(weil)': False, 'contains(nationals)': False, 'contains(yuji)': False, 'contains(muddies)': False, 'contains(compulsive)': False, 'contains(uh)': False, 'contains(installs)': False, 'contains(nullified)': False, 'contains(brewing)': False, 'contains(prologue)': False, 'contains(droppe d)': False, 'contains(guitarist)': False, 'contains(divorcee)': False, 'contains(rang)': False, 'contains(1990s)': False, 'contains(dopey)': False, 'contains(frumpy)': False, 'contains(dizzy)': False, 'contains(stint)': False, 'contains(furrowed)': False, 'contains(toll)': False, 'contains(weighty)': False, 'contains(insulin)': False, 'contains(plow)': False, 'contains(junichiro)': False, 'contains(knowingly)': False, 'contains(amyls)': False, 'contains(curt)': False, 'contains(commoner)': False, 'contains(dishes)': False, 'contains(fortresses)': False, 'contains(dabbling)': False, 'contains(overachieve)': False, 'contains(hallways)': False, 'contains(cycle)': False, 'contains(wildly)': False, 'contains(taxes)': False, 'contains(patrolman)': False, 'contains(underworld)': False, 'contains(newsgroup)': False, 'contains(similar)': False, 'contains(garfield)': False, 'contains(grimmest)': False, 'contains(ads)': False, 'contains(sites)': False, 'contains(dispair)': False, 'contains(designers)': False, 'contains(arangements)': False, 'contains(leguizimo)': False, 'contains(23rd)': False, 'contains(chord)': False, 'contains(disappointed)': False, 'contains(part)': False, 'contains(entrail

s)': False, 'contains(refrains)': False, 'contains(stephane)': False, 'contains(perry)': False, 'contains(homages)': False, 'contains(hurts)': False, 'contains(entrusting)': False, 'contains(autumnal)': False, 'contains(approachment)': False, 'contains(pans)': False, 'contains(bones)': False, 'contains(paquin)': False, 'contains(sniggering)': False, 'contains(whereabouts)': False, 'contains(discerning)': False, 'contains(gilsig)': False, 'contains(degeneration)': False, 'contains(tricks)': False, 'contains(worms)': False, 'contains(touchstone)': False, 'contains(whut)': False, 'contains(retro)': False, 'contains(donning)': False, 'contains(flabbergasted)': False, 'contains(ripping)': False, 'contains(forgiveness)': False, 'contains(skies)': False, 'contains(irreversible)': False, 'contains(yon)': False, 'contains(differs)': False, 'contains(shandling)': False, 'contains(cates)': False, 'contains(tm)': False, 'contains(potboilers)': False, 'contains(gulp)': False, 'contains(nerd)': False, 'contains(stivaletti)': False, 'contains(giacomo)': False, 'contains(stempel)': False, 'contains(therese)': False, 'contains(walker)': False, 'contains(anabella)': False, 'contains(milo)': False, 'contains(nada)': False, 'contains(cuddy)': False, 'contains(fifth)': False, 'contains(planned)': False, 'contains(mm)': False, 'contains(outfit)': False, 'contains(omnipresence)': False, 'contains(hutton)': False, 'contains(weaving)': False, 'contains(wore)': False, 'contains(humans)': False, 'contains(vaudeville)': False, 'contains(honorary)': False, 'contains(wai)': False, 'contains(traffic)': False, 'contains(tiering)': False, 'contains(treachery)': False, 'contains(ang)': False, 'contains(hijacker)': False, 'contains(certainly)': False, 'contains(whammy)': False, 'contains(rebs)': False, 'contains(prevalence)': False, 'contains(moviestar)': False, 'contains(oscar)': False, 'contains(san)': False, 'contains(shadyac)': False, 'contains(grunts)': False, 'contains(calahan)': False, 'contains(utmost)': False, 'contains(bombed)': False, 'contains(hagerty)': False, 'contains(bare)': False, 'contains(assists)': False, 'contains(fleder)': False, 'contains(reform)': False, 'contains(perpetually)': False, 'contains(bopper)': False, 'contains(roselyn)': False, 'contains(amadeus)': False, 'contains(monroth)': False, 'contains(sabara)': False, 'contains(saying)': False, 'contains(frightens)': False, 'contains(60s)': False, 'contains(menno)': False, 'contains(h20)': False, 'contains(bishop)': False, 'contains(pumped)': False, 'contains(impressiveness)': False, 'contains(bart)': False, 'contains(boosts)': False, 'contains(keenan)': False, 'contains(docu)': False, 'contains(pit)': False, 'contains(merrill)': False, 'contains(store)': False, 'contains(captor)': False, 'contains(clause)': False, 'contains(triumphing)': False, 'contains(terseness)': False, 'contains(loyalty)': False, 'contains(abounds)': False, 'contains(claustrophobia)': False, 'contains(hanako)': False, 'contains(provides)': False, 'contains(defuse)': False, 'contains(peeled)': False, 'contains(intruding)': False, 'contains(_mafia_)': False, 'contains(strengthening)': False, 'contains(donnie)': False, 'contains(mcmillan)': False, 'contains(vigilante)': False, 'contains(picking)': False, 'contains(tiring)': False, 'contains(medicine)': False, 'contains(sargent)': False, 'contains(windon)': False, 'contains(shone)': False, 'contains(helming)': False, 'contains(peep)': False, 'contains(fraternalizing)': False, 'contains(transexual)': False, 'contains(tykes)': False, 'contains(lowell)': False, 'contains(ditching)': False, 'contains(armstrong)': False, 'contains(enlightens)': False, 'contains(spitting)': False, 'contains(prehensile)': False, 'contains(pictorial)': False, 'contains(rehabilitation)': False, 'contains(page)': False, 'contains(robotically)': False, 'contains(accumulate)': False, 'contains(aniston)': False, 'contains(sharply)': False, 'contains(appendages)': False, 'contains(wen)': False, 'contains(patiently)': False, 'contains(slicing)': False, 'contains

(commiserate)': False, 'contains(whirl)': False, 'contains(temper)': False, 'contains(unsettling)': False, 'contains(ursula)': False, 'contains(lopsided)': False, 'contains(overtakes)': False, 'contains(slacker)': False, 'contains(glutton)': False, 'contains(mentioned)': False, 'contains(lightness)': False, 'contains(bogg)': False, 'contains(reprograms)': False, 'contains(splicing)': False, 'contains(moranis)': False, 'contains(lautrec)': False, 'contains(toughs)': False, 'contains(bridal)': False, 'contains(151)': False, 'contains(emotion)': False, 'contains(dumb)': False, 'contains(sears)': False, 'contains(blips)': False, 'contains(mongolian)': False, 'contains(lasers)': False, 'contains(goode)': False, 'contains(172)': False, 'contains(solos)': False, 'contains(expletitive)': False, 'contains(increased)': False, 'contains(pharaoh)': False, 'contains(lukewarm)': False, 'contains(curing)': False, 'contains(discoverer)': False, 'contains(18s)': False, 'contains(shoot)': False, 'contains(superlative)': False, 'contains(escapism)': False, 'contains(barbara)': False, 'contains(unparalleled)': False, 'contains(1950s)': False, 'contains(affiliations)': False, 'contains(fulfilling)': False, 'contains(delightfully)': False, 'contains(buzzcocks)': False, 'contains(carpenters)': False, 'contains(1521)': False, 'contains(cozart)': False, 'contains(connundrum)': False, 'contains(created)': False, 'contains(maimed)': False, 'contains(stolz)': False, 'contains(craftsmanship)': False, 'contains(cough)': False, 'contains(scorsese)': False, 'contains(consultation)': False, 'contains(abetting)': False, 'contains(landscapes)': False, 'contains(clockwatchers)': False, 'contains(stank)': False, 'contains(palate)': False, 'contains(repertoire)': False, 'contains(foulmouthed)': False, 'contains(collars)': False, 'contains(gazing)': False, 'contains(emittd)': False, 'contains(decreed)': False, 'contains(carried)': False, 'contains(poon)': False, 'contains(thrash)': False, 'contains(innocuously)': False, 'contains(kissed)': False, 'contains(cynically)': False, 'contains(adherents)': False, 'contains(laurene)': False, 'contains(ripley)': False, 'contains(undirected)': False, 'contains(deguerin)': False, 'contains(conventional)': False, 'contains(edu)': False, 'contains(elderly)': False, 'contains(adjustment)': False, 'contains(ignorant)': False, 'contains(trough)': False, 'contains(puppets)': False, 'contains(arc)': False, 'contains(essentially)': False, 'contains(kafka)': False, 'contains(wackier)': False, 'contains(goldsmith)': False, 'contains(crewman)': False, 'contains(nouri)': False, 'contains(rehabilitate)': False, 'contains(arizona)': False, 'contains(flaquer)': False, 'contains(filmgoers)': False, 'contains(schreiber)': False, 'contains(crucifixion)': False, 'contains(hipness)': False, 'contains(flyboy)': False, 'contains(salt)': False, 'contains(xtdl)': False, 'contains(secure)': False, 'contains(calvinist)': False, 'contains(photograph)': False, 'contains(webcams)': False, 'contains(revisiting)': False, 'contains(burkhart)': False, 'contains(flashcards)': False, 'contains(ominously)': False, 'contains(chaired)': False, 'contains(outgrown)': False, 'contains(machette)': False, 'contains(restrain)': False, 'contains(screenplay)': False, 'contains(dodger)': False, 'contains(wringing)': False, 'contains(constrained)': False, 'contains(typist)': False, 'contains(earthward)': False, 'contains(preposterousness)': False, 'contains(bestselling)': False, 'contains(sneeze)': False, 'contains(screen)': True, 'contains(plan)': False, 'contains(opus)': False, 'contains(scaffolding)': False, 'contains(evangelical)': False, 'contains(gruelling)': False, 'contains(deadline)': False, 'contains(twosomes)': False, 'contains(decipher)': False, 'contains(unaffected)': False, 'contains(okeday)': False, 'contains(productive)': False, 'contains(ricky)': False, 'contains(danced)': False, 'contains(snowfall)': False, 'contains(skarsgaard)': False, 'contains(leftover)': False, 'contains(fibber)': False, 'contains(pagan)': F

```

also, 'contains(distributors)': False, 'contains(vantages)': False,
'contains(religios)': False, 'contains(nyah)': False, 'contains(ber
tolini)': False, 'contains(rectangles)': False, 'contains(fluoro)':
False, 'contains(kathy)': False, 'contains()': False, 'contains(mu
ltifaceted)': False, 'contains(hunts)': False, 'contains(guccion
e)': False, 'contains(yessssss)': False, 'contains(expanded)': Fals
e, 'contains(remembers)': False, 'contains(meteorite)': False, 'con
tains(translation)': False, 'contains(copacabana)': False, 'contain
s(swoops)': False, 'contains(outweigh)': False, 'contains(dreams)':
False, 'contains(confusingly)': False, 'contains(scrolled)': False,
'contains(rube)': False, 'contains(deserved)': False, 'contains(ela
stica)': False, 'contains(tape)': False, 'contains(knowledgeable)':
False, 'contains(basil)': False, 'contains(fracture)': False, 'cont
ains(puts)': False, 'contains(overdose)': False, 'contains(overemph
asizing)': False, 'contains(16th)': False, 'contains(overpopulate
d)': False, 'contains(standup)': False, 'contains(analysis)': Fals
e, 'contains(hitting)': False, 'contains(bologna)': False, 'contain
s(wheeler)': False, 'contains(literacy)': False, 'contains(ki)': Fa
lse, 'contains(hypocritical)': False, 'contains(vaults)': False, 'c
ontains(knights)': False, 'contains(fleeting)': False, 'contains(sw
at)': False, 'contains(apprenticeship)': False, 'contains(exam)': F
alse, 'contains(wattage)': False, 'contains(declan)': False, 'conta
ins(rejuvenated)': False, 'contains(respect)': False, 'contains(pre
ening)': False, 'contains(faired)': False, 'contains(hecklers)': Fa
lse, 'contains(overview)': False, 'contains(shave)': False, 'contai
ns(bujold)': False, 'contains(jam)': False, 'contains(betcha)': Fal
se, 'contains(toughened)': False, 'contains(zeffirelli)': False, 'c
ontains(gratuities)': False, 'contains(migraines)': False, 'contain
s(clash)': False, 'contains(rampage)': False, 'contains(brought)':
False, 'contains(dunn)': False, 'contains(connoisseurs)': False,
'contains(donned)': False, 'contains(amnesiac)': False, 'contains
(misrepresenting)': False, 'contains(sequoia)': False, 'contains(st
accatto)': False, 'contains(homcoming)': False, 'contains(fully)':
False, 'contains(dwaine)': False, 'contains(alteration)': False,
'contains(uncorking)': False, 'contains(maps)': False, 'contains(w
etter)': False, 'contains(imported)': False, 'contains(tossable)':
False, 'contains(merhi)': False, 'contains(interview)': False, 'co
ntains(jefferson)': False, 'contains(ellie)': False, 'contains(kiss
ing)': False, 'contains(duties)': False, 'contains(fassbinder)': Fa
lse, 'contains(hops)': False, 'contains(imprisoning)': False, 'cont
ains(hatable)': False, 'contains(spit)': False, 'contains(vices)':
False, 'contains(ancestors)': False, 'contains(demi)': False, 'con
tains(emilio)': False, 'contains(mayoral)': False, 'contains(640)':
False, 'contains(meowth)': False, 'contains(hire)': False, 'contain
s(touring)': False, 'contains(interiors)': False, 'contains(submi
t)': False, 'contains(woodsboro)': False, 'contains(consists)': Fal
se, 'contains(tatooine)': False, 'contains(someone)': False, 'conta
ins(blossom)': False, 'contains(lacking)': False, 'contains(frill
s)': False, 'contains(loafer)': False, 'contains(sinclair)': False,
'contains(unsuspenseful)': False, 'contains(hear)': False, 'contain
s(predictably)': False, 'contains(teaspoons)': False, 'contains(not
ifying)': False, 'contains(merciful)': False, 'contains(sassines
s)': False, 'contains(climactic)': False, 'contains(home)': False,
'contains(protgaonist)': False, 'contains(decently)': False, 'cont
ains(gaby)': False, 'contains(revolutionized)': False, 'contains(me
ister)': False, 'contains(detection)': False, 'contains(rein)': Fal
se, 'contains(bingo)': False, 'contains(curmudgeons)': False, 'cont
ains(tantric)': False, 'contains(creepers)': False, 'contains(brani
ff)': False, 'contains(payload)': False, 'contains(invests)': Fals
e, 'contains(wrenched)': False, 'contains(layout)': False, 'contain
s(und)': False, 'contains(-)': False, 'contains(culls)': False, 'c

```

ontains(whooshed)': False, 'contains(field)': True, 'contains(sige
 l)': False, 'contains(applicants)': False, 'contains(transistor)':
 False, 'contains(cattle)': False, 'contains(flick)': False, 'conta
 ins(dustbuster)': False, 'contains(extracted)': False, 'contains(ab
 ba)': False, 'contains(horrifying)': False, 'contains(scriptiing)':
 False, 'contains(lyon)': False, 'contains(sociology)': False, 'cont
 ains(koch)': False, 'contains(pidgeonhole)': False, 'contains(origi
 nated)': False, 'contains(stereotypical)': False, 'contains(safel
 y)': False, 'contains(sharks)': False, 'contains(reveal)': False,
 'contains(reprehensible)': False, 'contains(logistical)': False,
 'contains(felix)': False, 'contains(fond)': False, 'contains(dream
 quest)': False, 'contains(assuredness)': False, 'contains(boaz)': F
 alse, 'contains(hoggett)': False, 'contains(storm)': False, 'contai
 ns(brute)': False, 'contains(rochelle)': False, 'contains(daredevil
 s)': False, 'contains(survey)': False, 'contains(tagalong)': False,
 'contains(exacting)': False, 'contains(appearing)': False, 'contain
 s(quake)': False, 'contains(steering)': False, 'contains(farting)':
 False, 'contains(muff)': False, 'contains(dum)': False, 'contains(d
 ewey)': False, 'contains(min)': False, 'contains(fatalities)': Fals
 e, 'contains(macabre)': False, 'contains(courage)': False, 'contain
 s(allayah)': False, 'contains(parenthetically)': False, 'contains(f
 acade)': False, 'contains(masseur)': False, 'contains(emotionles
 s)': False, 'contains(actor)': False, 'contains(majesty)': False,
 'contains(alienated)': False, 'contains(confided)': False, 'contai
 ns(active)': False, 'contains(shiftlessness)': False, 'contains(win
 ningham)': False, 'contains(007)': False, 'contains(weight)': Fals
 e, 'contains(undermined)': False, 'contains(rap)': False, 'contains
 (ply)': False, 'contains(potente)': False, 'contains(attend)': Fals
 e, 'contains(foreman)': False, 'contains(advisor)': False, 'contain
 s(vulcan)': False, 'contains(exploded)': False, 'contains(planne
 r)': False, 'contains(atreus)': False, 'contains(vitriol)': False,
 'contains(preparations)': False, 'contains(neater)': False, 'conta
 ins(mannered)': False, 'contains(sculpting)': False, 'contains(perm
 anetly)': False, 'contains(gain)': False, 'contains(presnell)': Fal
 se, 'contains(copping)': True, 'contains(marley)': False, 'contains
 (coca)': False, 'contains(unzipped)': False, 'contains(nighttime)':
 False, 'contains(vibrates)': False, 'contains(drink)': False, 'cont
 ains(rubbery)': False, 'contains(ethos)': False, 'contains(jovia
 l)': False, 'contains(major)': False, 'contains(footloose)': False,
 'contains(input)': False, 'contains(bondian)': False, 'contains(kre
 mp)': False, 'contains(ferocious)': False, 'contains(assets)': Fals
 e, 'contains(sicker)': False, 'contains(francine)': False, 'contain
 s(paleontology)': False, 'contains(noraruth)': False, 'contains(sa
 b)': False, 'contains(depote)': False, 'contains(randomly)': False,
 'contains(wisely)': False, 'contains(inroad)': False, 'contains(lax
 atives)': False, 'contains(hydraulic)': False, 'contains(directio
 n)': False, 'contains(soul)': False, 'contains(distance)': False,
 'contains(relentlessly)': False, 'contains(gq)': False, 'contains
 (headfirst)': False, 'contains(procol)': False, 'contains(&)': Fals
 e, 'contains(neverland)': False, 'contains(smirk)': False, 'contain
 s(ex)': False, 'contains(misogyny)': False, 'contains(seague)': Fal
 se, 'contains(excerpts)': False, 'contains(fine)': False, 'contains
 (steelworks)': False, 'contains(donation)': False, 'contains(munr
 o)': False, 'contains(glamorising)': False, 'contains(washout)': Fa
 lse, 'contains(vichy)': False, 'contains(lend)': False, 'contains(e
 xecuting)': False, 'contains(thinks)': False, 'contains(whimsica
 l)': False, 'contains(garfunkel)': False, 'contains(subtler)': Fals
 e, 'contains(onlookers)': False, 'contains(submerged)': False, 'con
 tains(sprocket)': False, 'contains(hardy)': False, 'contains(shell
 s)': False, 'contains(snuggling)': False, 'contains(_monster_movie
 _)': False, 'contains(lemay)': False, 'contains(eloquent)': False,

'contains(burlesque)': False, 'contains(eruting)': False, 'contains(articulately)': False, 'contains(horseshit)': False, 'contains(calvin)': False, 'contains(elliott)': False, 'contains(elise)': False, 'contains(interaction)': False, 'contains(dostoevski)': False, 'contains(snub)': False, 'contains(tunes)': False, 'contains(prancing)': False, 'contains(1995)': False, 'contains(ultraconservative)': False, 'contains(feat)': False, 'contains(indiglo)': False, 'contains(concocts)': False, 'contains(figaro)': False, 'contains(dalben)': False, 'contains(lubezki)': False, 'contains(beaker)': False, 'contains(bribe)': False, 'contains(confusions)': False, 'contains(joplin)': False, 'contains(miscarriage)': False, 'contains(dope d)': False, 'contains(contented)': False, 'contains(characteristics)': False, 'contains(cremer)': False, 'contains(wyatt)': False, 'contains(turrets)': False, 'contains(timeline)': False, 'contains(asiduously)': False, 'contains(publicist)': False, 'contains(mauna u)': False, 'contains(keyzer)': False, 'contains(identifiable)': False, 'contains(parallel)': False, 'contains(slathers)': False, 'contains(garber)': False, 'contains(arsed)': False, 'contains(reinventions)': False, 'contains(malle)': False, 'contains(punished)': False, 'contains(dignities)': False, 'contains(chili)': False, 'contains(stealth)': False, 'contains(pages)': False, 'contains(feelgood)': False, 'contains(couteney)': False, 'contains(straightest)': False, 'contains(exonerated)': False, 'contains(egos)': False, 'contains(witnessing)': False, 'contains(bragging)': False, 'contains(toes)': False, 'contains(fantasizing)': False, 'contains(chung)': False, 'contains(asleep)': False, 'contains(thierry)': False, 'contains(conglomerates)': False, 'contains(gravedigging)': False, 'contains(outlet)': False, 'contains(sheryl)': False, 'contains(remakes)': False, 'contains(partment)': False, 'contains(jacqueline)': False, 'contains(stupefied)': False, 'contains(event)': False, 'contains(terrestrials)': False, 'contains(jia)': False, 'contains(bothersome)': False, 'contains(serial)': False, 'contains(twistet)': False, 'contains(richman)': False, 'contains(abundance)': False, 'contains(intuits)': False, 'contains(subtitled)': False, 'contains(reach)': False, 'contains(irvin)': False, 'contains(boorman)': False, 'contains(misconceptions)': False, 'contains(sheepish)': False, 'contains(cleaners)': False, 'contains(loans)': False, 'contains(replete)': False, 'contains(arqua)': False, 'contains(laconic)': False, 'contains(penitentiary)': False, 'contains(proposing)': False, 'contains(thing_not)': False, 'contains(zinnia)': False, 'contains(primo)': False, 'contains(swigert)': False, 'contains(breezy)': False, 'contains(moore)': False, 'contains(disoriented)': False, 'contains(charlize)': False, 'contains(related)': True, 'contains(uptift)': False, 'contains(adaptions)': False, 'contains(maxim)': False, 'contains(wresting)': False, 'contains(headlights)': False, 'contains(accounts)': False, 'contains(stretch)': False, 'contains(malibu)': False, 'contains(puppeteering)': False, 'contains(heed)': False, 'contains(filmming)': False, 'contains(robs)': False, 'contains(rhett)': False, 'contains(kleiser)': False, 'contains(dreads)': False, 'contains(fujioka)': False, 'contains(offended)': False, 'contains(threatens)': False, 'contains(boyishly)': False, 'contains(paul l)': False, 'contains(dolenz)': False, 'contains(supplying)': False, 'contains(intolerable)': False, 'contains(nick)': False, 'contains(target)': False, 'contains(wanderlust)': False, 'contains(drix)': False, 'contains(sweden)': False, 'contains(seriously)': False, 'contains(reserves)': False, 'contains(counterparts)': False, 'contains(ratted)': False, 'contains(seng)': False, 'contains(performance)': False, 'contains(tournaments)': False, 'contains(brutish)': False, 'contains(nominee)': False, 'contains(boatman)': False, 'contains(5000)': False, 'contains(penetrate)': False, 'contains(skaro)': False, 'contains(sprayed)': False, 'contains(bresson)': False, 'contains(munchi

e)': False, 'contains(thanksgiving)': False, 'contains(horseman)': False, 'contains(immediately)': False, 'contains(butchers)': False, 'contains(forty)': False, 'contains(professionalism)': False, 'contains(nighthawks)': False, 'contains(unfetching)': False, 'contains(accelerate)': False, 'contains(icebergs)': False, 'contains(tiaras)': False, 'contains(spurs)': False, 'contains(nutter)': False, 'contains(scarlett)': False, 'contains(promote)': False, 'contains(dapper)': False, 'contains(hirschfeld)': False, 'contains(purging)': False, 'contains(workout)': False, 'contains(mccracken)': False, 'contains(vertigo)': False, 'contains(examination)': False, 'contains(graciously)': False, 'contains(stoker)': False, 'contains(daryl)': False, 'contains(potatohead)': False, 'contains(burglary)': False, 'contains(piven)': False, 'contains(apart)': False, 'contains(cannery)': False, 'contains(rubell)': False, 'contains(apologetically)': False, 'contains(restless)': False, 'contains(bogeyman)': False, 'contains(derek)': False, 'contains(sequels)': False, 'contains(turgidson)': False, 'contains(bleaches)': False, 'contains(than)': False, 'contains(sketching)': False, 'contains(witty)': False, 'contains(brauvara)': False, 'contains(_ferris)': False, 'contains(akiko)': False, 'contains(jing)': False, 'contains(autism)': False, 'contains(goodly)': False, 'contains(soprana)': False, 'contains(publicity)': False, 'contains(whoopee)': False, 'contains(expect)': False, 'contains(bistro)': False, 'contains(avalon)': False, 'contains(_breakfast_)': False, 'contains(necessary)': False, 'contains(feebly)': False, 'contains(buttocks)': False, 'contains(farewell)': False, 'contains(hymn)': False, 'contains(goes)': False, 'contains(mixer)': False, 'contains(approve)': False, 'contains(mcdowell)': False, 'contains(animated)': False, 'contains(duane)': False, 'contains(natacha)': False, 'contains(seattle)': False, 'contains(brinkford)': False, 'contains(flutters)': False, 'contains(til)': False, 'contains(trick)': False, 'contains(longshot)': False, 'contains(forsyth)': False, 'contains(martyrs)': False, 'contains(etchings)': False, 'contains(cunning)': False, 'contains(impossibilities)': False, 'contains(duck)': True, 'contains(mild)': False, 'contains(ultimatum)': False, 'contains(gleaming)': False, 'contains(resigning)': False, 'contains(celebrity)': False, 'contains(vowed)': False, 'contains(evoking)': False, 'contains(recommending)': False, 'contains(violated)': False, 'contains(requirement)': False, 'contains(sickness)': False, 'contains(response)': False, 'contains(slayer)': False, 'contains(bullet)': False, 'contains(ousting)': False, 'contains(wishes)': False, 'contains(blunts)': False, 'contains(svelte)': False, 'contains(blip)': False, 'contains(fertile)': False, 'contains(immune)': False, 'contains(bawls)': False, 'contains(clarisse)': False, 'contains(mayersberg)': False, 'contains(1996)': False, 'contains(feasibility)': False, 'contains(medgar)': False, 'contains(cooperation)': False, 'contains(persuasively)': False, 'contains(regain)': False, 'contains(issues)': False, 'contains(immediate)': False, 'contains(litmus)': False, 'contains(highschool)': False, 'contains(rangoon)': False, 'contains(unpardonable)': False, 'contains(1988)': False, 'contains(roughly)': False, 'contains(unengaging)': False, 'contains(artsier)': False, 'contains(gregarious)': False, 'contains(owl)': False, 'contains(rutger)': False, 'contains(pirate)': False, 'contains(kinfolk)': False, 'contains(khachaturian)': False, 'contains(conformist)': False, 'contains(hooligans)': False, 'contains(seaside)': False, 'contains(mutiny)': False, 'contains(smother)': False, 'contains(wager)': False, 'contains(obscenities)': False, 'contains(navigator)': False, 'contains(sponge)': False, 'contains(koop)': False, 'contains(pilgrims)': False, 'contains(walters)': False, 'contains(looker)': False, 'contains(2013)': False, 'contains(conspiracy)': False, 'contains(agutter)': False, 'contains(courtning)': False, 'contains(mafioso)': False, 'contains(canoeing)': False

e, 'contains(hubley)': False, 'contains(calhoun)': False, 'contains(jinn)': False, 'contains(benches)': False, 'contains(vapors)': False, 'contains(coyle)': False, 'contains(devious)': False, 'contains(opportunities)': False, 'contains(photogenic)': False, 'contains(hoary)': False, 'contains(sparing)': False, 'contains(kazantzakis)': False, 'contains(emotions)': False, 'contains(chipmunk)': False, 'contains(unsworth)': False, 'contains(grapevine)': False, 'contains(stamp)': False, 'contains(sadist)': False, 'contains(disgruntled)': False, 'contains(khe)': False, 'contains(acknowledges)': False, 'contains(tempers)': False, 'contains(depended)': False, 'contains(revving)': False, 'contains(released)': False, 'contains(telekinesis)': False, 'contains(makeover)': False, 'contains(drain)': False, 'contains(cruel)': False, 'contains(glamorized)': False, 'contains(amistad)': False, 'contains(circa)': False, 'contains(cheesiness)': False, 'contains(homemade)': False, 'contains(enterprising)': False, 'contains(cuddly)': False, 'contains(hazel)': False, 'contains(_murder_)': False, 'contains(rothchild)': False, 'contains(spicer)': False, 'contains(obstructed)': False, 'contains(jost)': False, 'contains(guessed)': False, 'contains(cultured)': False, 'contains(towards)': False, 'contains(clashed)': False, 'contains(ent)': False, 'contains(political)': False, 'contains(gobbling)': False, 'contains(had)': False, 'contains(struts)': False, 'contains(tamala)': False, 'contains(johner)': False, 'contains(irresponsible)': False, 'contains(shooters)': False, 'contains(subordination)': False, 'contains(war)': False, 'contains(orangutan)': False, 'contains(polarization)': False, 'contains(remark)': False, 'contains(alloy)': False, 'contains(distracted)': False, 'contains(figs)': False, 'contains(confounds)': False, 'contains(tentative)': False, 'contains(yanni)': False, 'contains(descending)': False, 'contains(quirkiness)': False, 'contains(quartet)': False, 'contains(supervising)': False, 'contains(request)': False, 'contains(paint)': False, 'contains(cliffs)': False, 'contains(sacrifices)': False, 'contains(gunning)': False, 'contains(anand)': False, 'contains(stifling)': False, 'contains(risks)': False, 'contains(spinozza)': False, 'contains(guinee)': False, 'contains(spiritualistic)': False, 'contains(obligation)': False, 'contains(drew)': False, 'contains(37th)': False, 'contains(fungus)': False, 'contains(strikes)': False, 'contains(bobo)': False, 'contains(declaration)': False, 'contains(topic)': False, 'contains(rin)': False, 'contains(push)': False, 'contains(unwise)': False, 'contains(heading)': False, 'contains(loquacious)': False, 'contains(markings)': False, 'contains(adlai)': False, 'contains(lela)': False, 'contains(corresponded)': False, 'contains(1982)': False, 'contains(pedlar)': False, 'contains(sexualized)': False, 'contains(gordy)': False, 'contains(pauses)': False, 'contains(balaban)': False, 'contains(overbearingly)': False, 'contains(fiercer)': False, 'contains(pare)': False, 'contains(articulated)': False, 'contains(delegates)': False, 'contains(stitch)': False, 'contains(bhatnagar)': False, 'contains(labors)': False, 'contains(coz)': False, 'contains(harlin)': False, 'contains(neglecting)': False, 'contains(pesudo)': False, 'contains(stresses)': False, 'contains(waisted)': False, 'contains(hurls)': False, 'contains(cellophane)': False, 'contains(fairfax)': False, 'contains(phew)': False, 'contains(rhod)': False, 'contains(bikers)': False, 'contains(mirren)': False, 'contains(cling)': False, 'contains(anakins)': False, 'contains(appointed)': False, 'contains(toiling)': False, 'contains(harvie)': False, 'contains(elegy)': False, 'contains(seething)': False, 'contains(avenue)': False, 'contains(desktop)': False, 'contains(climaxing)': False, 'contains(schuemacher)': False, 'contains(greg)': False, 'contains(letup)': False, 'contains(coffee)': False, 'contains(corral)': False, 'contains(brew)': False, 'contains(brevity)': False, 'contains(follows)': False, 'contains(avital)': False, 'contains(thawed)':

```

False, 'contains(lindberg)': False, 'contains(rottingham)': False,
'contains(gymnasium)': False, 'contains(sonar)': False, 'contains(h
allucinatory)': False, 'contains(internment)': False, 'contains(saf
ina)': False, 'contains(coagulate)': False, 'contains(columbus)': F
alse, 'contains(aretha)': False, 'contains(golden)': False, 'contai
ns(desparate)': False, 'contains(perlich)': False, 'contains(mainte
nance)': False, 'contains(putzes)': False, 'contains(anally)': Fals
e, 'contains(delacroix)': False, 'contains(stagecoach)': False, 'co
ntains(postlethwaite)': False, 'contains(conferring)': False, 'cont
ains(heighten)': False, 'contains(pain)': False, 'contains(parsle
y)': False, 'contains(partnered)': False, 'contains(bootstraps)': F
alse, 'contains(geiger)': False, 'contains(squashed)': False, 'cont
ains(meanders)': False, 'contains(tourfilm)': False, 'contains(inte
llectualizing)': False, 'contains(forbidding)': False, 'contains(su
o)': False, 'contains(featuring)': False, 'contains(coax)': False,
'contains(wrapping)': False, 'contains(niftiest)': False, 'contain
s(rudiments)': False, 'contains(deride)': False, 'contains(offshoo
t)': False, 'contains(romance)': False, 'contains(groundbreaking)':
False, 'contains(dreamt)': False, 'contains(marianne)': False, 'con
tains(and)': True, 'contains(bolster)': False, 'contains(base)': Fa
lse, 'contains(nuns)': False, 'contains(brilliant)': False, 'contai
ns(highlighting)': False, 'contains(duplicates)': False, 'contains
(trance)': False, "contains('--)": False, 'contains(psychic)': Fals
e, 'contains(zit)': False, 'contains(illiteracy)': False, 'contains
(commence)': False, 'contains(longshanks)': False, 'contains(beleiv
e)': False, 'contains(swaggering)': False, 'contains(9mm)': False,
'contains(ahem)': False, 'contains(crowned)': False, 'contains(mar
athon)': False, 'contains(marienbad)': False, 'contains(accountabl
e)': False, 'contains(shadowing)': False, 'contains(rather)': Fals
e, 'contains(alzhiemer)': False, 'contains(eligible)': False, 'cont
ains(dislocation)': False, 'contains(voiceovers)': False, 'contains
(affay)': False, 'contains(independence)': False, 'contains(mst3
k)': False, 'contains(sulking)': False, 'contains(cloying)': False,
'contains(skinheads)': False, 'contains(garp)': False, 'contains(ro
bin)': False, 'contains(suppose)': False, 'contains(confounding)':
False, 'contains(priority)': False, 'contains(word)': False, 'cont
ains(silliest)': False, 'contains(gapes)': False, 'contains(pointle
ssly)': False, 'contains(unflinching)': False, 'contains(intruder
s)': False, 'contains(deserving)': False, 'contains(quintessentia
l)': False, 'contains(jealously)': False, 'contains(regaled)': Fals
e, 'contains(emptying)': False, 'contains(outinen)': False, 'contai
ns(distorted)': False, 'contains(youths)': False, 'contains(tell)':
False, 'contains(atheist)': False, 'contains(departs)': False, 'con
tains(handgun)': False, 'contains(slaughter)': False, 'contains(pay
s)': False, 'contains(albert)': False, 'contains(damnation)': Fals
e, 'contains(hopper)': False, 'contains(recylcled)': False, 'contai
ns(truckloads)': False, 'contains(jnr)': False, 'contains(layout
s)': False, 'contains(coffins)': False, 'contains(slate)': False,
'contains(pessimistic)': False, 'contains(agile)': False, 'contain
s(ravera)': False, 'contains(cheif)': False, 'contains(masur)': Fal
se, 'contains(diners)': False, 'contains(lengthy)': False, 'contain
s(teriffic)': False, 'contains(kilgore)': False, 'contains(sinfull
y)': False, 'contains(cash)': False, 'contains(tippi)': False, 'con
tains(blissful)': False, 'contains(sequences)': False, 'contains(mc
gavin)': False, 'contains(talmud)': False, 'contains(apprehended)':
False, 'contains(batfans)': False, 'contains(families)': False, 'co
ntains(synthesizer)': False, 'contains(macy)': False, 'contains(nee
dles)': False, 'contains(pathos)': False, 'contains(dramatism)': Fa
lse, 'contains(weasely)': False, 'contains(sterling)': False, 'cont
ains(improvisation)': False, 'contains(rile)': False, 'contains(spi
nster)': False}

```

In [19]:

```
len(movie_reviews.words('pos/cv957_8737.txt'))
```

Out[19]:

597

We are ready to train the classifier:

In [20]:

```
featuresets = [(document_features(d), c) for (d,c) in documents]
train_set, test_set = featuresets[100:], featuresets[:100]
classifier = nltk.NaiveBayesClassifier.train(train_set)
```

Simple check:

In [21]:

```
classifier.show_most_informative_features(5)
```

Most Informative Features

3 : 1.0	contains(courage) = True	pos : neg	=	6.
2 : 1.0	contains(cunning) = True	pos : neg	=	6.
9 : 1.0	contains(stinks) = True	neg : pos	=	5.
5 : 1.0	contains(dedication) = True	pos : neg	=	5.
5 : 1.0	contains(elise) = True	pos : neg	=	5.

In [22]:

```
nltk.classify.accuracy(classifier,test_set)
```

Out[22]:

0.59

In [23]:

```
doc = "this movie was pathetic".split()
```

In [24]:

```
classifier.classify(document_features(doc))
```

Out[24]:

'neg'

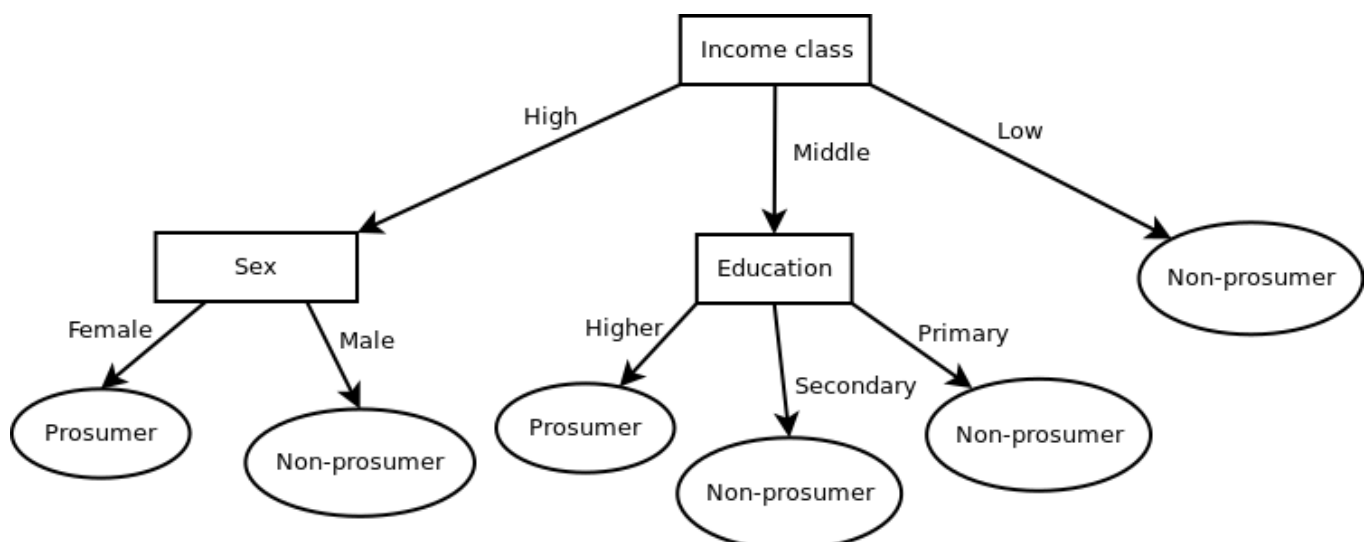
Decision trees

- a data mining method used very often for classification purposes
- other common applications:
 - variable selection, i.e. selection of the most relevant attributes that should be used to form a model
 - relative importance of variables
 - prediction, i.e. forecasting results for future data by using a tree built on historical data
- a decision tree depicts rules for dividing existing data into groups
 - the first rule splits the entire data set into some number of subsets
 - another rules may be applied to each of those subsets forming a next generation of them
 - procedure is repeated until the data in each subset forms a final group
- trees accept several types of variables (nominal, ordinal and interval ones)
- they can handle records with missing data and do not require to drop them

To give an understanding of the concept of decision trees for a non-expert reader let us assume that our goal is to predict if a person is a prosumer taking into consideration features like income, education and sex. A sample data set could look like the one presented in the following table (**training set**):

Income class	Sex	Education	Prosument
High	F	Secondary	YES
High	M	Higher	NO
High	F	Higher	YES
Middle	M	Higher	YES
Middle	F	Secondary	NO
Middle	F	Primary	NO
Low	F	Higher	NO
Low	M	Secondary	NO
Middle	F	Higher	YES
High	F	Higher	YES

A corresponding decision tree is shown below:



- a graph-like structure consisting of non-leaf nodes (depicted by rectangles) and leaves (ovals) connected with each other
- **non-leaf nodes** are attributes characterizing the items in the data set
- **leaves** represent predicted variable
- prediction is easy:
 - a new person - a female with high income and a secondary level education
 - going along appropriate branches will lead us to the conclusion that the person is a prosumer
- interpretation of the decision tree is straightforward and no expert knowledge

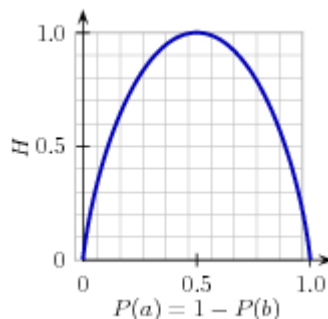
Building a tree

- finding a proper attribute to split the data is not a trivial task
- several algorithms available
- *Iterative Dichotomiser 3* (ID3) - an algorithm invented by Ross Quinlan (1986)
 - it two concepts from the information theory: information entropy and information gain

The entropy H (see C. Shannon, *A Mathematical Theory of Communication*, 1948) is understood as the average information contained in a message or needed to generate it. In our example the message would be simply the Prosumer or Not Prosumer classes returned by the tree. The entropy of a data set is calculated according to the following formula:

$$H(p, n) = -\frac{p}{p+n} \log \frac{p}{p+n} - \frac{n}{p+n} \log \frac{n}{p+n}$$

Here, p is the number of positive examples in the data set and n - the number of negative ones. Please note that for pure data samples, i.e. all records belonging to the same class, the entropy is equal to zero.



In our example we have 5 prosumers and 5 non-prosumers in the data, thus the initial entropy is

$$H(5, 5) = -\frac{5}{10} \log_2 \frac{5}{10} - \frac{5}{10} \log_2 \frac{5}{10} = -\log_2 \frac{1}{2} = 1$$

The information gain is the difference between the entropy before and after the split. Let us take the attribute Income class as a candidate for breaking down the data. After the split we get following entropies in each subset:

$$\begin{aligned} H_{High}(3, 1) &= -0.75 * \log_2(0.75) - 0.25 * \log_2(0.25) = 0.811 \\ H_{Middle}(2, 2) &= 1 \\ H_{Low}(0, 2) &= 0 \end{aligned}$$

For the average entropy after the split we have

$$\bar{H}(\text{Income}) = 0.4 * 0.811 + 0.4 * 1.0 + 0.2 * 0 = 0.7244.$$

The corresponding information gain is given by

$$I_G(\text{Income}) = H(5, 5) - \bar{H}(\text{Income}) = 0.2756$$

Similarly, for the other attributes we get

$$I_G(\textit{Sex}) = 0.03$$
$$I_G(\textit{Education}) = 0.16$$

Since `Income` class yields the largest information gain, we use it to split the data. Then we repeat the same procedure recursively on each subset, considering only attributes not selected before.

In [25]:

```
import math
def entropy(labels):
    freqdist = nltk.FreqDist(labels)
    probs = [freqdist.freq(l) for l in freqdist]
    return -sum(p * math.log(p,2) for p in probs)
```

In [26]:

```
print(entropy(['male', 'male', 'male', 'male']))
```

-0.0

In [27]:

```
print(entropy(['male', 'female', 'male', 'male']))
```

0.8112781244591328

In [28]:

```
print(entropy(['female', 'male', 'female', 'male']))
```

1.0

In [29]:

```
print(entropy(['female', 'female', 'male', 'female']))
```

0.8112781244591328

In [30]:

```
print(entropy(['female', 'female', 'female', 'female']))
```

-0.0

Decision trees in NLTK

We will work with the movie reviews from the previous example. Let us define first a couple of helper functions:

In [31]:

```
def bag_of_words(words):
    return dict([(word, True) for word in words])
```

In [32]:

```
def bag_of_words_not_in_set(words, badwords):  
    return bag_of_words(set(words) - set(badwords))
```

In [33]:

```
import collections  
def label_feats_from_corpus(corp, feature_detector=bag_of_words):  
    label_feats = collections.defaultdict(list)  
    for label in corp.categories():  
        for fileid in corp.fileids(categories=[label]):  
            feats = feature_detector(corp.words(fileids=[fileid]))  
            label_feats[label].append(feats)  
    return label_feats
```

In [34]:

```
def split_label_feats(lfeats, split=0.75):  
    train_feats = []  
    test_feats = []  
    for label, feats in lfeats.items():  
        cutoff = int(len(feats) * split)  
        train_feats.extend([(feat, label) for feat in feats[:cutoff]])  
        test_feats.extend([(feat, label) for feat in feats[cutoff:]])  
    return train_feats, test_feats
```

In [35]:

```
from nltk.corpus import movie_reviews  
print(movie_reviews.categories())
```

```
['neg', 'pos']
```

In [36]:

```
lfeats = label_feats_from_corpus(movie_reviews)  
print(lfeats.keys())
```

```
dict_keys(['neg', 'pos'])
```

In [37]:

```
train_feats, test_feats = split_label_feats(lfeats)  
print(len(train_feats))  
print(len(test_feats))
```

```
1500
```

```
500
```


In [38]:

```
print(test_feats[0])
```

```
{('s': True, 'idea': True, 'esteem': True, 'gets': True, 'cute': True, 'enjoyable': True, 'twist': True, 'last': True, '-': True, 'it': True, 'quite': True, 'own': True, 'when': True, 'funny': True, 'from': True, 'pet': True, 'both': True, 'predictability': True, 'zegers': True, '"': True, 'alone': True, 'if': True, 'feat': True, 'e', 'wags': True, 'true': True, 'credits': True, 'not': True, '.': True, '!': True, 'snively': True, ',': True, 'used': True, 'see': True, 'actual': True, 'can': True, 'pairs': True, 'yeller': True, 'two': True, 'michael': True, ')': True, 'been': True, 'a': True, 'are': True, 'opens': True, 'movie': True, 'just': True, 'school': True, 'comedy': True, 'death': True, 'straight': True, 'light': True, 'marches': True, 'very': True, '(' : True, 'norm': True, 'block': True, 'had': True, 'number': True, 'faux': True, 'face': True, 'stupid': True, 'shots': True, 'surfaces': True, ':': True, 'buried': True, 'cope': True, 'big': True, 'lie': True, 'moment': True, 'slapstick': True, 'on': True, 'cans': True, 'fernwell': True, 'would': True, 'accomplished': True, 'reluctant': True, 'the': True, 'and': True, 'comeuppance': True, 'i': True, 'its': True, 'calculated': True, 'court': True, 'letterman': True, 'matter': True, 'david': True, 'abusive': True, 'dog': True, 'formula': True, 'no': True, 'paint': True, 'mopey': True, 'chain': True, 'climax': True, 'send': True, 'may': True, 'least': True, 'well': True, 'once': True, 'heavy': True, 'rockets': True, 'old': True, 'end': True, 'contracts': True, 'events': True, '--': True, 'animal': True, 'asked': True, 'an': True, 'musical': True, 'himself': True, 'k9': True, '?': True, 'motion': True, 'where': True, 'whatever': True, 'style': True, 'sight': True, 'clown': True, 'have': True, 'clad': True, 'win': True, 'fades': True, 'finals': True, 'buddy': True, 'there': True, 'boy': True, 'one': True, 'else': True, 'although': True, 'is': True, 'visual': True, 'friend': True, 'note': True, 'gloom': True, 'all': True, 'for': True, 'game': True, 'proves': True, 'father': True, 'okay': True, 'escapes': True, 'exist': True, 'up': True, 'special': True, 'basketball': True, 'seem': True, 'current': True, 'newspapers': True, 'name': True, 'with': True, 'approach': True, 'everything': True, 'absurdity': True, 'picture': True, 'possibly': True, 'moments': True, 'back': True, 'surprisingly': True, 'move': True, 'off': True, 'quicker': True, 'trying': True, 'splish': True, 'think': True, 'tricks': True, 'ends': True, 'his': True, 'this': True, 'or': True, 'effects': True, 'cleaned': True, 'save': True, 't': True, 'owner': True, 'he': True, 'as': True, 'connection': True, 'washington': True, 'must': True, 'self': True, 'executed': True, 'of': True, 'even': True, 'occasional': True, 'be': True, 'sequences': True, 'mine': True, 'yeah': True, 'solemn': True, 'begin': True, 'here': True, 'josh': True, 'we': True, 'courtroom': True, 'prowess': True, 'reclaim': True, 'jersey': True, 'appropriate': True, 'actually': True, 'then': True, 'splash': True, 'in': True, 'at': True, 'tells': True, 'insists': True, 'new': True, 'than': True, 'plays': True, 'cannot': True, 'kevin': True, 'air': True, 'out': True, 'their': True, 'sneakers': True, 'team': True, 'that': True, 'spilled': True, 'follows': True, 'player': True, 'realized': True, 'make': True, '"': True, 'could': True, 'successful': True, 'sink': True, 'places': True, 'to': True, 'hire': True, 'anything': True, 're': True, 'before': True, 'were': True, 'trick': True, 'tale': True, 'segment': True, 'while': True, 'recent': True, 'interested': True, 'family': True, 'doesn': True, 'forced': True, 'joke': True, 'disney': True, 'jeter': True, 'rabies': True, 'montage': True, 'bud': True, 'pooch': True, 'mascot': True, 'golden': True, 'saw': True, 'they': True, 'retriever': True, 'kid': True, 'more': True, 'story': True, 'but': True, 'granted': True, 'better': True, 'cheer': True}, 'neg')
```

Now we can train and test the classifier:

In [39]:

```
from nltk.classify import DecisionTreeClassifier  
#dt_classifier = DecisionTreeClassifier.train(train_feats, binary=True, entropy_  
cutoff=0.8, depth_cutoff=5, support_cutoff=30)  
dt_classifier = DecisionTreeClassifier.train(train_feats, depth_cutoff=5)
```

In [40]:

```
print(nltk.classify.accuracy(dt_classifier, test_feats))
```

0.682

In [41]:

```
for i in test_feats:  
    print("Human: ",i[-1]," | Machine: ", dt_classifier.classify(i[0]))
```

Human:	neg		Machine:	pos
Human:	neg		Machine:	neg
Human:	neg		Machine:	neg
Human:	neg		Machine:	pos
Human:	neg		Machine:	neg
Human:	neg		Machine:	neg
Human:	neg		Machine:	pos
Human:	neg		Machine:	neg
Human:	neg		Machine:	neg
Human:	neg		Machine:	neg
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	neg
Human:	neg		Machine:	neg
Human:	neg		Machine:	neg
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	neg
Human:	neg		Machine:	neg
Human:	neg		Machine:	neg
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	neg
Human:	neg		Machine:	neg
Human:	neg		Machine:	pos
Human:	neg		Machine:	neg
Human:	neg		Machine:	neg
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	neg
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	neg
Human:	neg		Machine:	pos
Human:	neg		Machine:	neg
Human:	neg		Machine:	pos
Human:	neg		Machine:	neg
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	neg
Human:	neg		Machine:	pos
Human:	neg		Machine:	pos
Human:	neg		Machine:	neg
Human:	neg		Machine:	pos
Human:	neg		Machine:	neg

Human:	neg	Machine:	neg
Human:	neg	Machine:	pos
Human:	neg	Machine:	neg
Human:	neg	Machine:	pos
Human:	neg	Machine:	pos
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	pos
Human:	neg	Machine:	neg
Human:	neg	Machine:	pos
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	pos
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	pos
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	pos
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	pos
Human:	neg	Machine:	pos
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	pos
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	pos
Human:	neg	Machine:	pos
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	pos
Human:	neg	Machine:	neg
Human:	neg	Machine:	pos
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg
Human:	neg	Machine:	neg

Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg

Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: neg
Human: neg		Machine: pos
Human: neg		Machine: pos
Human: neg		Machine: neg

[illegible]

[illegible]

[illegible]

[illegible]

Human:	pos		Machine:	neg
Human:	pos		Machine:	neg
Human:	pos		Machine:	pos
Human:	pos		Machine:	pos
Human:	pos		Machine:	pos
Human:	pos		Machine:	pos
Human:	pos		Machine:	pos
Human:	pos		Machine:	pos
Human:	pos		Machine:	pos
Human:	pos		Machine:	pos
Human:	pos		Machine:	pos
Human:	pos		Machine:	pos

In [44]:

```
doc = bag_of_words("this movie was bad".split())  
print(doc)
```

```
{'movie': True, 'was': True, 'this': True, 'bad': True}
```

In [45]:

```
dt_classifier.classify(doc)
```

Out[45]:

```
'neg'
```

In [46]:

```
doc = bag_of_words("this movie was pathetic".split())  
print(doc)
```

```
{'movie': True, 'was': True, 'this': True, 'pathetic': True}
```

In [47]:

```
dt_classifier.classify(doc)
```

Out[47]:

```
'pos'
```

Maximum Entropy Classifier

- a probabilistic classifier which belongs to the class of exponential models
- independence of features not required
- based on the principle of Maximum Entropy - from all the models that fit the training data it selects the one which has the largest entropy

From Wikipedia:

The principle of maximum entropy states that the probability distribution which best represents the current state of knowledge is the one with largest entropy, in the context of precisely stated prior data (such as a proposition that expresses testable information).

Another way of stating this: Take precisely stated prior data or testable information about a probability distribution function. Consider the set of all trial probability distributions that would encode the prior data. According to this principle, the distribution with maximal information entropy is the proper one.

The principle was first expounded by E. T. Jaynes in two papers in 1957, where he emphasized a natural correspondence between statistical mechanics and information theory. In particular, Jaynes offered a new and very general rationale why the Gibbsian method of statistical mechanics works. He argued that the entropy of statistical mechanics and the information entropy of information theory are basically the same thing. Consequently, statistical mechanics should be seen just as a particular application of a general tool of logical inference and information theory.

- can be used to solve a large variety of text classification problems:
 - language detection
 - topic classification
 - sentiment analysis
- implementing MaxEnt classifier in a standard programming language is non-trivial primarily due to the numerical optimization problem that one should solve in order to estimate the weights of the model

When to use?

- when we don't know anything about the prior distributions and when it is unsafe to make any such assumptions
- when we can't assume the conditional independence of the features (particularly true in Text Classification problems where the features are usually words which obviously are not independent)

Theoretical Background

Our target is to use the contextual information of the document (unigrams, bigrams, other characteristics within the text) in order to categorize it to a given class (positive/neutral/negative, objective/subjective etc). Following the standard bag-of-words framework that is commonly used in natural language processing and information retrieval, let $\{w_1, \dots, w_m\}$ be the m words that can appear in a document. Then each document is represented by a sparse array with 1s and 0s that indicate whether a particular word w_i exists or not in the context of the document.

In order to construct a stochastic model, the first step is to collect a large number of training data which consists of samples in the following format: (x_i, y_i) where the x_i includes the contextual information of the document (the sparse array) and y_i its class.

The second step is to summarize the training sample in terms of its empirical probability distribution:

$$\tilde{p}(x, y) \equiv \frac{1}{N} \times \text{number of times that } (x, y) \text{ occurs in the sample}$$

Here, N is the size of the training dataset.

The above empirical distribution will be used to construct the statistical model of the random process which assigns texts to a particular class by taking into account their contextual information.

We introduce the following indicator function, i.e. **feature**:

$$f_j(x, y) = \begin{cases} 1 & \text{if } y = c_i \text{ and } x \text{ contains } w_k \\ 0 & \text{otherwise} \end{cases}$$

This binary valued indicator function returns 1 only when the class of a particular document is c_i and the document contains the word w_k .

We express any statistic of the training dataset as the expected value of the appropriate binary-valued indicator function f_j . Thus the expected value of feature f_j with respect to the empirical distribution $\tilde{p}(x, y)$ is equal to:

$$\tilde{p}(f_j) \equiv \sum_{x, y} \tilde{p}(x, y) f_j(x, y)$$

If each training sample (x, y) occurs once in training dataset then $\tilde{p}(x, y)$ is equal to $1/N$.

When a particular statistic is useful to our classification, we require our model to accord with it. To do so, we constrain the expected value that the model assigns to the expected value of the feature function f_j , i.e.

$$p(f_j) \equiv \sum_{x, y} \tilde{p}(x) p(y | x) f_j(x, y)$$

Here, $\tilde{p}(x)$ is the empirical distribution of x in the training dataset and it is usually set equal to $1/N$.

By constraining the expected value to be the equal to the empirical value and from the above equations we have that:

$$\sum_{x, y} \tilde{p}(x) p(y | x) f_j(x, y) = \sum_{x, y} \tilde{p}(x, y) f_j(x, y)$$

The last equation is called **constrain** and we have as many constrains as the number of feature functions j .

The above constrains can be satisfied by an infinite number of models. In order to build our model, we need to select the best candidate based on a specific criterion. According to the principle of **Maximum Entropy**, we should select the model that is as close as possible to uniform. In other words, we should select the model p^* with Maximum Entropy:

$$p^* = \arg \max_{p \in \mathcal{C}} \left(- \sum_{x, y} \tilde{p}(x) p(y | x) \log p(y | x) \right)$$

We have:

1. $p(y|x) \geq 0$ for all x and y
2. $\sum_y p(y|x) = 1$ for all x
3. $\sum_{x,y} \tilde{p}(x) p(y|x) f_j(x, y) = \sum_{x,y} \tilde{p}(x, y) f_j(x, y)$ for $j \in \{1, 2, \dots, n\}$

To solve the above optimization problem we introduce the Lagrangian multipliers and estimate them with the maximum likelihood method.

It can be proven that if we find the multipliers $\{\lambda_1, \dots, \lambda_n\}$ that maximize the dual problem, the probability of a document x to be classified as y is equal to:

$$p^*(y|x) = \frac{\exp\left(\sum_i \lambda_i f_i(x, y)\right)}{\sum_y \exp\left(\sum_i \lambda_i f_i(x, y)\right)}$$

Thus, once the lamda parameters are found, all we need to do in order to classify a new document is to use the **maximum a posteriori** decision rule and select the category with the highest probability.

Estimating the lamda parameters requires using an **iterative scaling algorithm** such as the GIS (Generalized Iterative Scaling) or the IIS (Improved Iterative Scaling):

Input: Feature functions f_1, \dots, f_n , empirical distribution $\tilde{p}(x)$

Output: Optimal parameter values λ_i , optimal model p^*

Steps:

1. Start with $\lambda_i = 0$ for all $i \in \{1, 2, \dots, n\}$
2. For each $i \in \{1, 2, \dots, n\}$ do:
 - A. Let $\Delta\lambda_i$ be the solution to

$$\sum_{x,y} \tilde{p}(x) p(y|x) f_j(x, y) \exp(\Delta\lambda_i f^s(x, y)) = \tilde{p}(f_i)$$
 with $f^s(x, y) = \sum_{i=1}^n f_i(x, y)$
 - B. Update the value of λ_i according to $\lambda_i \leftarrow \lambda_i + \Delta\lambda_i$
3. Go to step 2 if not all parameters λ_i have converged

The $f^s(x, y)$ is the total number of features which are active for a particular (x, y) pair. If this number is constant for all documents then $\Delta\lambda_i$ can be calculated in closed-form:

$$\Delta\lambda_i = \frac{1}{C} \log \frac{\tilde{p}(f_i)}{p(f_i)}$$

where $C = f^s(x, y)$. The assumption of $f^s(x, y)$ being constant is rarely realistic in practice.

Fortunately, it is only necessary that C is bounded by $f^s(x, y)$ (Goodman 2002, Ratnaparkhi 1997). Thus in some versions of ISS one selects simply

$$C = C_{max} = \arg \max_{x,y} f^s(x, y)$$

MaxEnt classifier in NLTK

- GIS(Generalized Iterative Scaling), IIS(Improved Iterative Scaling), and LM-BFGS (limited memory Broyden-Fletcher-Goldfarb-Shanno) training methods
- GIS and ISS are implemented within NLTK module
 - seem very slow
 - cost large memory

- LM-BFGS supports external libraries like MEGAM (MEGA Model Optimization Package)

In [25]:

```
from nltk import MaxentClassifier  
me_classifier = MaxentClassifier.train(train_set,algorithm='gis')
```

==> Training (100 iterations)

Iteration	Log Likelihood	Accuracy
1	-0.69315	0.505
2	-0.69248	0.505
3	-0.69193	0.505
4	-0.69138	0.505
5	-0.69083	0.505
6	-0.69028	0.507
7	-0.68973	0.513
8	-0.68919	0.526
9	-0.68864	0.543
10	-0.68810	0.558
11	-0.68755	0.580
12	-0.68701	0.599
13	-0.68647	0.621
14	-0.68593	0.643
15	-0.68538	0.662
16	-0.68485	0.673
17	-0.68431	0.691
18	-0.68377	0.702
19	-0.68323	0.713
20	-0.68270	0.725
21	-0.68216	0.738
22	-0.68163	0.750
23	-0.68110	0.761
24	-0.68057	0.767
25	-0.68004	0.773
26	-0.67951	0.777
27	-0.67898	0.783
28	-0.67845	0.782
29	-0.67793	0.790
30	-0.67740	0.796
31	-0.67688	0.801
32	-0.67635	0.804
33	-0.67583	0.807
34	-0.67531	0.808
35	-0.67479	0.812
36	-0.67427	0.813
37	-0.67375	0.814
38	-0.67323	0.815
39	-0.67272	0.815
40	-0.67220	0.814
41	-0.67168	0.817
42	-0.67117	0.819
43	-0.67066	0.819
44	-0.67015	0.821
45	-0.66963	0.822
46	-0.66912	0.824
47	-0.66862	0.827
48	-0.66811	0.828
49	-0.66760	0.828
50	-0.66709	0.829
51	-0.66659	0.831
52	-0.66608	0.832
53	-0.66558	0.832
54	-0.66508	0.832
55	-0.66458	0.831
56	-0.66408	0.829
57	-0.66358	0.830

58	-0.66308	0.831
59	-0.66258	0.831
60	-0.66208	0.831
61	-0.66159	0.831
62	-0.66109	0.831
63	-0.66060	0.830
64	-0.66010	0.830
65	-0.65961	0.831
66	-0.65912	0.832
67	-0.65863	0.832
68	-0.65814	0.832
69	-0.65765	0.832
70	-0.65716	0.832
71	-0.65667	0.833
72	-0.65619	0.833
73	-0.65570	0.834
74	-0.65522	0.834
75	-0.65474	0.834
76	-0.65425	0.834
77	-0.65377	0.834
78	-0.65329	0.835
79	-0.65281	0.835
80	-0.65233	0.835
81	-0.65186	0.835
82	-0.65138	0.835
83	-0.65090	0.835
84	-0.65043	0.835
85	-0.64995	0.836
86	-0.64948	0.835
87	-0.64901	0.835
88	-0.64853	0.835
89	-0.64806	0.835
90	-0.64759	0.836
91	-0.64712	0.836
92	-0.64666	0.836
93	-0.64619	0.836
94	-0.64572	0.836
95	-0.64526	0.837
96	-0.64479	0.837
97	-0.64433	0.838
98	-0.64387	0.838
99	-0.64340	0.839
Final	-0.64294	0.839

In [26]:

```
me_classifier.show_most_informative_features(5)
```

```
-0.110 contains(cunning)==True and label is 'neg'
-0.102 contains(skarsgard)==True and label is 'neg'
-0.102 contains(dedication)==True and label is 'neg'
-0.101 contains(tatooine)==True and label is 'neg'
-0.100 contains(elise)==True and label is 'neg'
```

In [27]:

```
nltk.classify.accuracy(classifier,test_set)
```

Out[27]:

```
0.59
```

How to improve accuracy of a text classifier?

- eliminate low quality features (words)
 - L1 regularization
- recursively grow your stopwords list (by analyzing the top features):
 - frequently used words
 - countries
 - cities
 - adjectives
 - temporal words (Tuesday, tomorrow, January)
- look beyond unigrams
 - considering bigrams often boosts performance
- diversify your corpus
 - it helps dilute word features that are specific to one particular corpus
 - context of the classification problem is important while selecting the corpus
 - example: social media texts classifier - it is not enough to include only Twitter corpus
- tweak precision and recall
 - tweak the system so when it fails, it does so in a manner that is more tolerable
 - shifting False Positive (or Precision) under-performance to False Negative (or recall) under-performance
 - or viceversa
- eliminate low quality predictions
 - instead, return "class unknown" if necessary
- canonicalize words through lemma reduction
 - reduced probability space of the algorithm
- normalize exaggerations
 - haaaaapppyyyy --> happy
- eliminate numerals, punctuation and corpus specific text
 - if you are dealing with a dataset from Twitter, you might want to eliminate (using a RegEx perhaps) any usernames (of the format @user) because they do not contribute towards the classification problem you have at hand
- try a different classification algorithm
- lowering all text is not always smart
 - if you are classifying **Intensity** or **Mood**, then capitalization might be an important feature that contributes positively to the accuracy of the predictions
 - if you are trying to classify text into topics or categories, then lowering all text might have a very healthy impact on overall accuracy
- manual curation of corpus data
 - of particular importance if the training set involves human entry, or people trying to game the system for their own benefit