

# 07 - Sentiment-Analysis (movies)

April 24, 2020

## 1 Sentiment Analysis on movie reviews

Analysis of Social Media Contents

Alessandro Ortis - University of Catania

Based on *Exercise B: Sentiment Analysis* in [Natural Language Processing with Python/NLTK](#) by Luciano M. Guasco

### 1.1 1. Exploring the movie\_reviews corpus

The `nltk.corpus` package defines a collection of corpus reader classes, which can be used to access the contents of a diverse set of corpora. The list of available corpora is given at:

[http://www.nltk.org/nltk\\_data/](http://www.nltk.org/nltk_data/)

Each corpus reader class is specialized to handle a specific corpus format. In addition, the `nltk.corpus` package automatically creates a set of corpus reader instances that can be used to access the corpora in the NLTK data package.

```
In [1]: from nltk.corpus import movie_reviews # These are movie reviews already separated as p
        #movie_reviews.readme().replace('\n', ' ').replace('\t', '').replace('`', '').replac
```

Most corpora consist of a set of files, each containing a document (or other pieces of text). A list of identifiers for these files is accessed via the `fileids()` method of the corpus reader.

```
In [2]: movie_reviews.fileids()
```

```
Out[2]: ['neg/cv000_29416.txt',
        'neg/cv001_19502.txt',
        'neg/cv002_17424.txt',
        'neg/cv003_12683.txt',
        'neg/cv004_12641.txt',
        'neg/cv005_29357.txt',
        'neg/cv006_17022.txt',
        'neg/cv007_4992.txt',
        'neg/cv008_29326.txt',
        'neg/cv009_29417.txt',
        'neg/cv010_29063.txt',
        'neg/cv011_13044.txt',
        'neg/cv012_29411.txt',
```

'neg/cv013\_10494.txt',  
'neg/cv014\_15600.txt',  
'neg/cv015\_29356.txt',  
'neg/cv016\_4348.txt',  
'neg/cv017\_23487.txt',  
'neg/cv018\_21672.txt',  
'neg/cv019\_16117.txt',  
'neg/cv020\_9234.txt',  
'neg/cv021\_17313.txt',  
'neg/cv022\_14227.txt',  
'neg/cv023\_13847.txt',  
'neg/cv024\_7033.txt',  
'neg/cv025\_29825.txt',  
'neg/cv026\_29229.txt',  
'neg/cv027\_26270.txt',  
'neg/cv028\_26964.txt',  
'neg/cv029\_19943.txt',  
'neg/cv030\_22893.txt',  
'neg/cv031\_19540.txt',  
'neg/cv032\_23718.txt',  
'neg/cv033\_25680.txt',  
'neg/cv034\_29446.txt',  
'neg/cv035\_3343.txt',  
'neg/cv036\_18385.txt',  
'neg/cv037\_19798.txt',  
'neg/cv038\_9781.txt',  
'neg/cv039\_5963.txt',  
'neg/cv040\_8829.txt',  
'neg/cv041\_22364.txt',  
'neg/cv042\_11927.txt',  
'neg/cv043\_16808.txt',  
'neg/cv044\_18429.txt',  
'neg/cv045\_25077.txt',  
'neg/cv046\_10613.txt',  
'neg/cv047\_18725.txt',  
'neg/cv048\_18380.txt',  
'neg/cv049\_21917.txt',  
'neg/cv050\_12128.txt',  
'neg/cv051\_10751.txt',  
'neg/cv052\_29318.txt',  
'neg/cv053\_23117.txt',  
'neg/cv054\_4101.txt',  
'neg/cv055\_8926.txt',  
'neg/cv056\_14663.txt',  
'neg/cv057\_7962.txt',  
'neg/cv058\_8469.txt',  
'neg/cv059\_28723.txt',  
'neg/cv060\_11754.txt',

'neg/cv061\_9321.txt',  
'neg/cv062\_24556.txt',  
'neg/cv063\_28852.txt',  
'neg/cv064\_25842.txt',  
'neg/cv065\_16909.txt',  
'neg/cv066\_11668.txt',  
'neg/cv067\_21192.txt',  
'neg/cv068\_14810.txt',  
'neg/cv069\_11613.txt',  
'neg/cv070\_13249.txt',  
'neg/cv071\_12969.txt',  
'neg/cv072\_5928.txt',  
'neg/cv073\_23039.txt',  
'neg/cv074\_7188.txt',  
'neg/cv075\_6250.txt',  
'neg/cv076\_26009.txt',  
'neg/cv077\_23172.txt',  
'neg/cv078\_16506.txt',  
'neg/cv079\_12766.txt',  
'neg/cv080\_14899.txt',  
'neg/cv081\_18241.txt',  
'neg/cv082\_11979.txt',  
'neg/cv083\_25491.txt',  
'neg/cv084\_15183.txt',  
'neg/cv085\_15286.txt',  
'neg/cv086\_19488.txt',  
'neg/cv087\_2145.txt',  
'neg/cv088\_25274.txt',  
'neg/cv089\_12222.txt',  
'neg/cv090\_0049.txt',  
'neg/cv091\_7899.txt',  
'neg/cv092\_27987.txt',  
'neg/cv093\_15606.txt',  
'neg/cv094\_27868.txt',  
'neg/cv095\_28730.txt',  
'neg/cv096\_12262.txt',  
'neg/cv097\_26081.txt',  
'neg/cv098\_17021.txt',  
'neg/cv099\_11189.txt',  
'neg/cv100\_12406.txt',  
'neg/cv101\_10537.txt',  
'neg/cv102\_8306.txt',  
'neg/cv103\_11943.txt',  
'neg/cv104\_19176.txt',  
'neg/cv105\_19135.txt',  
'neg/cv106\_18379.txt',  
'neg/cv107\_25639.txt',  
'neg/cv108\_17064.txt',

'neg/cv109\_22599.txt',  
'neg/cv110\_27832.txt',  
'neg/cv111\_12253.txt',  
'neg/cv112\_12178.txt',  
'neg/cv113\_24354.txt',  
'neg/cv114\_19501.txt',  
'neg/cv115\_26443.txt',  
'neg/cv116\_28734.txt',  
'neg/cv117\_25625.txt',  
'neg/cv118\_28837.txt',  
'neg/cv119\_9909.txt',  
'neg/cv120\_3793.txt',  
'neg/cv121\_18621.txt',  
'neg/cv122\_7891.txt',  
'neg/cv123\_12165.txt',  
'neg/cv124\_3903.txt',  
'neg/cv125\_9636.txt',  
'neg/cv126\_28821.txt',  
'neg/cv127\_16451.txt',  
'neg/cv128\_29444.txt',  
'neg/cv129\_18373.txt',  
'neg/cv130\_18521.txt',  
'neg/cv131\_11568.txt',  
'neg/cv132\_5423.txt',  
'neg/cv133\_18065.txt',  
'neg/cv134\_23300.txt',  
'neg/cv135\_12506.txt',  
'neg/cv136\_12384.txt',  
'neg/cv137\_17020.txt',  
'neg/cv138\_13903.txt',  
'neg/cv139\_14236.txt',  
'neg/cv140\_7963.txt',  
'neg/cv141\_17179.txt',  
'neg/cv142\_23657.txt',  
'neg/cv143\_21158.txt',  
'neg/cv144\_5010.txt',  
'neg/cv145\_12239.txt',  
'neg/cv146\_19587.txt',  
'neg/cv147\_22625.txt',  
'neg/cv148\_18084.txt',  
'neg/cv149\_17084.txt',  
'neg/cv150\_14279.txt',  
'neg/cv151\_17231.txt',  
'neg/cv152\_9052.txt',  
'neg/cv153\_11607.txt',  
'neg/cv154\_9562.txt',  
'neg/cv155\_7845.txt',  
'neg/cv156\_11119.txt',

'neg/cv157\_29302.txt',  
'neg/cv158\_10914.txt',  
'neg/cv159\_29374.txt',  
'neg/cv160\_10848.txt',  
'neg/cv161\_12224.txt',  
'neg/cv162\_10977.txt',  
'neg/cv163\_10110.txt',  
'neg/cv164\_23451.txt',  
'neg/cv165\_2389.txt',  
'neg/cv166\_11959.txt',  
'neg/cv167\_18094.txt',  
'neg/cv168\_7435.txt',  
'neg/cv169\_24973.txt',  
'neg/cv170\_29808.txt',  
'neg/cv171\_15164.txt',  
'neg/cv172\_12037.txt',  
'neg/cv173\_4295.txt',  
'neg/cv174\_9735.txt',  
'neg/cv175\_7375.txt',  
'neg/cv176\_14196.txt',  
'neg/cv177\_10904.txt',  
'neg/cv178\_14380.txt',  
'neg/cv179\_9533.txt',  
'neg/cv180\_17823.txt',  
'neg/cv181\_16083.txt',  
'neg/cv182\_7791.txt',  
'neg/cv183\_19826.txt',  
'neg/cv184\_26935.txt',  
'neg/cv185\_28372.txt',  
'neg/cv186\_2396.txt',  
'neg/cv187\_14112.txt',  
'neg/cv188\_20687.txt',  
'neg/cv189\_24248.txt',  
'neg/cv190\_27176.txt',  
'neg/cv191\_29539.txt',  
'neg/cv192\_16079.txt',  
'neg/cv193\_5393.txt',  
'neg/cv194\_12855.txt',  
'neg/cv195\_16146.txt',  
'neg/cv196\_28898.txt',  
'neg/cv197\_29271.txt',  
'neg/cv198\_19313.txt',  
'neg/cv199\_9721.txt',  
'neg/cv200\_29006.txt',  
'neg/cv201\_7421.txt',  
'neg/cv202\_11382.txt',  
'neg/cv203\_19052.txt',  
'neg/cv204\_8930.txt',

'neg/cv205\_9676.txt',  
'neg/cv206\_15893.txt',  
'neg/cv207\_29141.txt',  
'neg/cv208\_9475.txt',  
'neg/cv209\_28973.txt',  
'neg/cv210\_9557.txt',  
'neg/cv211\_9955.txt',  
'neg/cv212\_10054.txt',  
'neg/cv213\_20300.txt',  
'neg/cv214\_13285.txt',  
'neg/cv215\_23246.txt',  
'neg/cv216\_20165.txt',  
'neg/cv217\_28707.txt',  
'neg/cv218\_25651.txt',  
'neg/cv219\_19874.txt',  
'neg/cv220\_28906.txt',  
'neg/cv221\_27081.txt',  
'neg/cv222\_18720.txt',  
'neg/cv223\_28923.txt',  
'neg/cv224\_18875.txt',  
'neg/cv225\_29083.txt',  
'neg/cv226\_26692.txt',  
'neg/cv227\_25406.txt',  
'neg/cv228\_5644.txt',  
'neg/cv229\_15200.txt',  
'neg/cv230\_7913.txt',  
'neg/cv231\_11028.txt',  
'neg/cv232\_16768.txt',  
'neg/cv233\_17614.txt',  
'neg/cv234\_22123.txt',  
'neg/cv235\_10704.txt',  
'neg/cv236\_12427.txt',  
'neg/cv237\_20635.txt',  
'neg/cv238\_14285.txt',  
'neg/cv239\_29828.txt',  
'neg/cv240\_15948.txt',  
'neg/cv241\_24602.txt',  
'neg/cv242\_11354.txt',  
'neg/cv243\_22164.txt',  
'neg/cv244\_22935.txt',  
'neg/cv245\_8938.txt',  
'neg/cv246\_28668.txt',  
'neg/cv247\_14668.txt',  
'neg/cv248\_15672.txt',  
'neg/cv249\_12674.txt',  
'neg/cv250\_26462.txt',  
'neg/cv251\_23901.txt',  
'neg/cv252\_24974.txt',

'neg/cv253\_10190.txt',  
'neg/cv254\_5870.txt',  
'neg/cv255\_15267.txt',  
'neg/cv256\_16529.txt',  
'neg/cv257\_11856.txt',  
'neg/cv258\_5627.txt',  
'neg/cv259\_11827.txt',  
'neg/cv260\_15652.txt',  
'neg/cv261\_11855.txt',  
'neg/cv262\_13812.txt',  
'neg/cv263\_20693.txt',  
'neg/cv264\_14108.txt',  
'neg/cv265\_11625.txt',  
'neg/cv266\_26644.txt',  
'neg/cv267\_16618.txt',  
'neg/cv268\_20288.txt',  
'neg/cv269\_23018.txt',  
'neg/cv270\_5873.txt',  
'neg/cv271\_15364.txt',  
'neg/cv272\_20313.txt',  
'neg/cv273\_28961.txt',  
'neg/cv274\_26379.txt',  
'neg/cv275\_28725.txt',  
'neg/cv276\_17126.txt',  
'neg/cv277\_20467.txt',  
'neg/cv278\_14533.txt',  
'neg/cv279\_19452.txt',  
'neg/cv280\_8651.txt',  
'neg/cv281\_24711.txt',  
'neg/cv282\_6833.txt',  
'neg/cv283\_11963.txt',  
'neg/cv284\_20530.txt',  
'neg/cv285\_18186.txt',  
'neg/cv286\_26156.txt',  
'neg/cv287\_17410.txt',  
'neg/cv288\_20212.txt',  
'neg/cv289\_6239.txt',  
'neg/cv290\_11981.txt',  
'neg/cv291\_26844.txt',  
'neg/cv292\_7804.txt',  
'neg/cv293\_29731.txt',  
'neg/cv294\_12695.txt',  
'neg/cv295\_17060.txt',  
'neg/cv296\_13146.txt',  
'neg/cv297\_10104.txt',  
'neg/cv298\_24487.txt',  
'neg/cv299\_17950.txt',  
'neg/cv300\_23302.txt',

'neg/cv301\_13010.txt',  
'neg/cv302\_26481.txt',  
'neg/cv303\_27366.txt',  
'neg/cv304\_28489.txt',  
'neg/cv305\_9937.txt',  
'neg/cv306\_10859.txt',  
'neg/cv307\_26382.txt',  
'neg/cv308\_5079.txt',  
'neg/cv309\_23737.txt',  
'neg/cv310\_14568.txt',  
'neg/cv311\_17708.txt',  
'neg/cv312\_29308.txt',  
'neg/cv313\_19337.txt',  
'neg/cv314\_16095.txt',  
'neg/cv315\_12638.txt',  
'neg/cv316\_5972.txt',  
'neg/cv317\_25111.txt',  
'neg/cv318\_11146.txt',  
'neg/cv319\_16459.txt',  
'neg/cv320\_9693.txt',  
'neg/cv321\_14191.txt',  
'neg/cv322\_21820.txt',  
'neg/cv323\_29633.txt',  
'neg/cv324\_7502.txt',  
'neg/cv325\_18330.txt',  
'neg/cv326\_14777.txt',  
'neg/cv327\_21743.txt',  
'neg/cv328\_10908.txt',  
'neg/cv329\_29293.txt',  
'neg/cv330\_29675.txt',  
'neg/cv331\_8656.txt',  
'neg/cv332\_17997.txt',  
'neg/cv333\_9443.txt',  
'neg/cv334\_0074.txt',  
'neg/cv335\_16299.txt',  
'neg/cv336\_10363.txt',  
'neg/cv337\_29061.txt',  
'neg/cv338\_9183.txt',  
'neg/cv339\_22452.txt',  
'neg/cv340\_14776.txt',  
'neg/cv341\_25667.txt',  
'neg/cv342\_20917.txt',  
'neg/cv343\_10906.txt',  
'neg/cv344\_5376.txt',  
'neg/cv345\_9966.txt',  
'neg/cv346\_19198.txt',  
'neg/cv347\_14722.txt',  
'neg/cv348\_19207.txt',



'neg/cv349\_15032.txt',  
'neg/cv350\_22139.txt',  
'neg/cv351\_17029.txt',  
'neg/cv352\_5414.txt',  
'neg/cv353\_19197.txt',  
'neg/cv354\_8573.txt',  
'neg/cv355\_18174.txt',  
'neg/cv356\_26170.txt',  
'neg/cv357\_14710.txt',  
'neg/cv358\_11557.txt',  
'neg/cv359\_6751.txt',  
'neg/cv360\_8927.txt',  
'neg/cv361\_28738.txt',  
'neg/cv362\_16985.txt',  
'neg/cv363\_29273.txt',  
'neg/cv364\_14254.txt',  
'neg/cv365\_12442.txt',  
'neg/cv366\_10709.txt',  
'neg/cv367\_24065.txt',  
'neg/cv368\_11090.txt',  
'neg/cv369\_14245.txt',  
'neg/cv370\_5338.txt',  
'neg/cv371\_8197.txt',  
'neg/cv372\_6654.txt',  
'neg/cv373\_21872.txt',  
'neg/cv374\_26455.txt',  
'neg/cv375\_9932.txt',  
'neg/cv376\_20883.txt',  
'neg/cv377\_8440.txt',  
'neg/cv378\_21982.txt',  
'neg/cv379\_23167.txt',  
'neg/cv380\_8164.txt',  
'neg/cv381\_21673.txt',  
'neg/cv382\_8393.txt',  
'neg/cv383\_14662.txt',  
'neg/cv384\_18536.txt',  
'neg/cv385\_29621.txt',  
'neg/cv386\_10229.txt',  
'neg/cv387\_12391.txt',  
'neg/cv388\_12810.txt',  
'neg/cv389\_9611.txt',  
'neg/cv390\_12187.txt',  
'neg/cv391\_11615.txt',  
'neg/cv392\_12238.txt',  
'neg/cv393\_29234.txt',  
'neg/cv394\_5311.txt',  
'neg/cv395\_11761.txt',  
'neg/cv396\_19127.txt',

'neg/cv397\_28890.txt',  
'neg/cv398\_17047.txt',  
'neg/cv399\_28593.txt',  
'neg/cv400\_20631.txt',  
'neg/cv401\_13758.txt',  
'neg/cv402\_16097.txt',  
'neg/cv403\_6721.txt',  
'neg/cv404\_21805.txt',  
'neg/cv405\_21868.txt',  
'neg/cv406\_22199.txt',  
'neg/cv407\_23928.txt',  
'neg/cv408\_5367.txt',  
'neg/cv409\_29625.txt',  
'neg/cv410\_25624.txt',  
'neg/cv411\_16799.txt',  
'neg/cv412\_25254.txt',  
'neg/cv413\_7893.txt',  
'neg/cv414\_11161.txt',  
'neg/cv415\_23674.txt',  
'neg/cv416\_12048.txt',  
'neg/cv417\_14653.txt',  
'neg/cv418\_16562.txt',  
'neg/cv419\_14799.txt',  
'neg/cv420\_28631.txt',  
'neg/cv421\_9752.txt',  
'neg/cv422\_9632.txt',  
'neg/cv423\_12089.txt',  
'neg/cv424\_9268.txt',  
'neg/cv425\_8603.txt',  
'neg/cv426\_10976.txt',  
'neg/cv427\_11693.txt',  
'neg/cv428\_12202.txt',  
'neg/cv429\_7937.txt',  
'neg/cv430\_18662.txt',  
'neg/cv431\_7538.txt',  
'neg/cv432\_15873.txt',  
'neg/cv433\_10443.txt',  
'neg/cv434\_5641.txt',  
'neg/cv435\_24355.txt',  
'neg/cv436\_20564.txt',  
'neg/cv437\_24070.txt',  
'neg/cv438\_8500.txt',  
'neg/cv439\_17633.txt',  
'neg/cv440\_16891.txt',  
'neg/cv441\_15276.txt',  
'neg/cv442\_15499.txt',  
'neg/cv443\_22367.txt',  
'neg/cv444\_9975.txt',

'neg/cv445\_26683.txt',  
'neg/cv446\_12209.txt',  
'neg/cv447\_27334.txt',  
'neg/cv448\_16409.txt',  
'neg/cv449\_9126.txt',  
'neg/cv450\_8319.txt',  
'neg/cv451\_11502.txt',  
'neg/cv452\_5179.txt',  
'neg/cv453\_10911.txt',  
'neg/cv454\_21961.txt',  
'neg/cv455\_28866.txt',  
'neg/cv456\_20370.txt',  
'neg/cv457\_19546.txt',  
'neg/cv458\_9000.txt',  
'neg/cv459\_21834.txt',  
'neg/cv460\_11723.txt',  
'neg/cv461\_21124.txt',  
'neg/cv462\_20788.txt',  
'neg/cv463\_10846.txt',  
'neg/cv464\_17076.txt',  
'neg/cv465\_23401.txt',  
'neg/cv466\_20092.txt',  
'neg/cv467\_26610.txt',  
'neg/cv468\_16844.txt',  
'neg/cv469\_21998.txt',  
'neg/cv470\_17444.txt',  
'neg/cv471\_18405.txt',  
'neg/cv472\_29140.txt',  
'neg/cv473\_7869.txt',  
'neg/cv474\_10682.txt',  
'neg/cv475\_22978.txt',  
'neg/cv476\_18402.txt',  
'neg/cv477\_23530.txt',  
'neg/cv478\_15921.txt',  
'neg/cv479\_5450.txt',  
'neg/cv480\_21195.txt',  
'neg/cv481\_7930.txt',  
'neg/cv482\_11233.txt',  
'neg/cv483\_18103.txt',  
'neg/cv484\_26169.txt',  
'neg/cv485\_26879.txt',  
'neg/cv486\_9788.txt',  
'neg/cv487\_11058.txt',  
'neg/cv488\_21453.txt',  
'neg/cv489\_19046.txt',  
'neg/cv490\_18986.txt',  
'neg/cv491\_12992.txt',  
'neg/cv492\_19370.txt',

'neg/cv493\_14135.txt',  
'neg/cv494\_18689.txt',  
'neg/cv495\_16121.txt',  
'neg/cv496\_11185.txt',  
'neg/cv497\_27086.txt',  
'neg/cv498\_9288.txt',  
'neg/cv499\_11407.txt',  
'neg/cv500\_10722.txt',  
'neg/cv501\_12675.txt',  
'neg/cv502\_10970.txt',  
'neg/cv503\_11196.txt',  
'neg/cv504\_29120.txt',  
'neg/cv505\_12926.txt',  
'neg/cv506\_17521.txt',  
'neg/cv507\_9509.txt',  
'neg/cv508\_17742.txt',  
'neg/cv509\_17354.txt',  
'neg/cv510\_24758.txt',  
'neg/cv511\_10360.txt',  
'neg/cv512\_17618.txt',  
'neg/cv513\_7236.txt',  
'neg/cv514\_12173.txt',  
'neg/cv515\_18484.txt',  
'neg/cv516\_12117.txt',  
'neg/cv517\_20616.txt',  
'neg/cv518\_14798.txt',  
'neg/cv519\_16239.txt',  
'neg/cv520\_13297.txt',  
'neg/cv521\_1730.txt',  
'neg/cv522\_5418.txt',  
'neg/cv523\_18285.txt',  
'neg/cv524\_24885.txt',  
'neg/cv525\_17930.txt',  
'neg/cv526\_12868.txt',  
'neg/cv527\_10338.txt',  
'neg/cv528\_11669.txt',  
'neg/cv529\_10972.txt',  
'neg/cv530\_17949.txt',  
'neg/cv531\_26838.txt',  
'neg/cv532\_6495.txt',  
'neg/cv533\_9843.txt',  
'neg/cv534\_15683.txt',  
'neg/cv535\_21183.txt',  
'neg/cv536\_27221.txt',  
'neg/cv537\_13516.txt',  
'neg/cv538\_28485.txt',  
'neg/cv539\_21865.txt',  
'neg/cv540\_3092.txt',

'neg/cv541\_28683.txt',  
'neg/cv542\_20359.txt',  
'neg/cv543\_5107.txt',  
'neg/cv544\_5301.txt',  
'neg/cv545\_12848.txt',  
'neg/cv546\_12723.txt',  
'neg/cv547\_18043.txt',  
'neg/cv548\_18944.txt',  
'neg/cv549\_22771.txt',  
'neg/cv550\_23226.txt',  
'neg/cv551\_11214.txt',  
'neg/cv552\_0150.txt',  
'neg/cv553\_26965.txt',  
'neg/cv554\_14678.txt',  
'neg/cv555\_25047.txt',  
'neg/cv556\_16563.txt',  
'neg/cv557\_12237.txt',  
'neg/cv558\_29376.txt',  
'neg/cv559\_0057.txt',  
'neg/cv560\_18608.txt',  
'neg/cv561\_9484.txt',  
'neg/cv562\_10847.txt',  
'neg/cv563\_18610.txt',  
'neg/cv564\_12011.txt',  
'neg/cv565\_29403.txt',  
'neg/cv566\_8967.txt',  
'neg/cv567\_29420.txt',  
'neg/cv568\_17065.txt',  
'neg/cv569\_26750.txt',  
'neg/cv570\_28960.txt',  
'neg/cv571\_29292.txt',  
'neg/cv572\_20053.txt',  
'neg/cv573\_29384.txt',  
'neg/cv574\_23191.txt',  
'neg/cv575\_22598.txt',  
'neg/cv576\_15688.txt',  
'neg/cv577\_28220.txt',  
'neg/cv578\_16825.txt',  
'neg/cv579\_12542.txt',  
'neg/cv580\_15681.txt',  
'neg/cv581\_20790.txt',  
'neg/cv582\_6678.txt',  
'neg/cv583\_29465.txt',  
'neg/cv584\_29549.txt',  
'neg/cv585\_23576.txt',  
'neg/cv586\_8048.txt',  
'neg/cv587\_20532.txt',  
'neg/cv588\_14467.txt',

'neg/cv589\_12853.txt',  
'neg/cv590\_20712.txt',  
'neg/cv591\_24887.txt',  
'neg/cv592\_23391.txt',  
'neg/cv593\_11931.txt',  
'neg/cv594\_11945.txt',  
'neg/cv595\_26420.txt',  
'neg/cv596\_4367.txt',  
'neg/cv597\_26744.txt',  
'neg/cv598\_18184.txt',  
'neg/cv599\_22197.txt',  
'neg/cv600\_25043.txt',  
'neg/cv601\_24759.txt',  
'neg/cv602\_8830.txt',  
'neg/cv603\_18885.txt',  
'neg/cv604\_23339.txt',  
'neg/cv605\_12730.txt',  
'neg/cv606\_17672.txt',  
'neg/cv607\_8235.txt',  
'neg/cv608\_24647.txt',  
'neg/cv609\_25038.txt',  
'neg/cv610\_24153.txt',  
'neg/cv611\_2253.txt',  
'neg/cv612\_5396.txt',  
'neg/cv613\_23104.txt',  
'neg/cv614\_11320.txt',  
'neg/cv615\_15734.txt',  
'neg/cv616\_29187.txt',  
'neg/cv617\_9561.txt',  
'neg/cv618\_9469.txt',  
'neg/cv619\_13677.txt',  
'neg/cv620\_2556.txt',  
'neg/cv621\_15984.txt',  
'neg/cv622\_8583.txt',  
'neg/cv623\_16988.txt',  
'neg/cv624\_11601.txt',  
'neg/cv625\_13518.txt',  
'neg/cv626\_7907.txt',  
'neg/cv627\_12603.txt',  
'neg/cv628\_20758.txt',  
'neg/cv629\_16604.txt',  
'neg/cv630\_10152.txt',  
'neg/cv631\_4782.txt',  
'neg/cv632\_9704.txt',  
'neg/cv633\_29730.txt',  
'neg/cv634\_11989.txt',  
'neg/cv635\_0984.txt',  
'neg/cv636\_16954.txt',

'neg/cv637\_13682.txt',  
'neg/cv638\_29394.txt',  
'neg/cv639\_10797.txt',  
'neg/cv640\_5380.txt',  
'neg/cv641\_13412.txt',  
'neg/cv642\_29788.txt',  
'neg/cv643\_29282.txt',  
'neg/cv644\_18551.txt',  
'neg/cv645\_17078.txt',  
'neg/cv646\_16817.txt',  
'neg/cv647\_15275.txt',  
'neg/cv648\_17277.txt',  
'neg/cv649\_13947.txt',  
'neg/cv650\_15974.txt',  
'neg/cv651\_11120.txt',  
'neg/cv652\_15653.txt',  
'neg/cv653\_2107.txt',  
'neg/cv654\_19345.txt',  
'neg/cv655\_12055.txt',  
'neg/cv656\_25395.txt',  
'neg/cv657\_25835.txt',  
'neg/cv658\_11186.txt',  
'neg/cv659\_21483.txt',  
'neg/cv660\_23140.txt',  
'neg/cv661\_25780.txt',  
'neg/cv662\_14791.txt',  
'neg/cv663\_14484.txt',  
'neg/cv664\_4264.txt',  
'neg/cv665\_29386.txt',  
'neg/cv666\_20301.txt',  
'neg/cv667\_19672.txt',  
'neg/cv668\_18848.txt',  
'neg/cv669\_24318.txt',  
'neg/cv670\_2666.txt',  
'neg/cv671\_5164.txt',  
'neg/cv672\_27988.txt',  
'neg/cv673\_25874.txt',  
'neg/cv674\_11593.txt',  
'neg/cv675\_22871.txt',  
'neg/cv676\_22202.txt',  
'neg/cv677\_18938.txt',  
'neg/cv678\_14887.txt',  
'neg/cv679\_28221.txt',  
'neg/cv680\_10533.txt',  
'neg/cv681\_9744.txt',  
'neg/cv682\_17947.txt',  
'neg/cv683\_13047.txt',  
'neg/cv684\_12727.txt',

'neg/cv685\_5710.txt',  
'neg/cv686\_15553.txt',  
'neg/cv687\_22207.txt',  
'neg/cv688\_7884.txt',  
'neg/cv689\_13701.txt',  
'neg/cv690\_5425.txt',  
'neg/cv691\_5090.txt',  
'neg/cv692\_17026.txt',  
'neg/cv693\_19147.txt',  
'neg/cv694\_4526.txt',  
'neg/cv695\_22268.txt',  
'neg/cv696\_29619.txt',  
'neg/cv697\_12106.txt',  
'neg/cv698\_16930.txt',  
'neg/cv699\_7773.txt',  
'neg/cv700\_23163.txt',  
'neg/cv701\_15880.txt',  
'neg/cv702\_12371.txt',  
'neg/cv703\_17948.txt',  
'neg/cv704\_17622.txt',  
'neg/cv705\_11973.txt',  
'neg/cv706\_25883.txt',  
'neg/cv707\_11421.txt',  
'neg/cv708\_28539.txt',  
'neg/cv709\_11173.txt',  
'neg/cv710\_23745.txt',  
'neg/cv711\_12687.txt',  
'neg/cv712\_24217.txt',  
'neg/cv713\_29002.txt',  
'neg/cv714\_19704.txt',  
'neg/cv715\_19246.txt',  
'neg/cv716\_11153.txt',  
'neg/cv717\_17472.txt',  
'neg/cv718\_12227.txt',  
'neg/cv719\_5581.txt',  
'neg/cv720\_5383.txt',  
'neg/cv721\_28993.txt',  
'neg/cv722\_7571.txt',  
'neg/cv723\_9002.txt',  
'neg/cv724\_15265.txt',  
'neg/cv725\_10266.txt',  
'neg/cv726\_4365.txt',  
'neg/cv727\_5006.txt',  
'neg/cv728\_17931.txt',  
'neg/cv729\_10475.txt',  
'neg/cv730\_10729.txt',  
'neg/cv731\_3968.txt',  
'neg/cv732\_13092.txt',



'neg/cv733\_9891.txt',  
'neg/cv734\_22821.txt',  
'neg/cv735\_20218.txt',  
'neg/cv736\_24947.txt',  
'neg/cv737\_28733.txt',  
'neg/cv738\_10287.txt',  
'neg/cv739\_12179.txt',  
'neg/cv740\_13643.txt',  
'neg/cv741\_12765.txt',  
'neg/cv742\_8279.txt',  
'neg/cv743\_17023.txt',  
'neg/cv744\_10091.txt',  
'neg/cv745\_14009.txt',  
'neg/cv746\_10471.txt',  
'neg/cv747\_18189.txt',  
'neg/cv748\_14044.txt',  
'neg/cv749\_18960.txt',  
'neg/cv750\_10606.txt',  
'neg/cv751\_17208.txt',  
'neg/cv752\_25330.txt',  
'neg/cv753\_11812.txt',  
'neg/cv754\_7709.txt',  
'neg/cv755\_24881.txt',  
'neg/cv756\_23676.txt',  
'neg/cv757\_10668.txt',  
'neg/cv758\_9740.txt',  
'neg/cv759\_15091.txt',  
'neg/cv760\_8977.txt',  
'neg/cv761\_13769.txt',  
'neg/cv762\_15604.txt',  
'neg/cv763\_16486.txt',  
'neg/cv764\_12701.txt',  
'neg/cv765\_20429.txt',  
'neg/cv766\_7983.txt',  
'neg/cv767\_15673.txt',  
'neg/cv768\_12709.txt',  
'neg/cv769\_8565.txt',  
'neg/cv770\_11061.txt',  
'neg/cv771\_28466.txt',  
'neg/cv772\_12971.txt',  
'neg/cv773\_20264.txt',  
'neg/cv774\_15488.txt',  
'neg/cv775\_17966.txt',  
'neg/cv776\_21934.txt',  
'neg/cv777\_10247.txt',  
'neg/cv778\_18629.txt',  
'neg/cv779\_18989.txt',  
'neg/cv780\_8467.txt',

'neg/cv781\_5358.txt',  
'neg/cv782\_21078.txt',  
'neg/cv783\_14724.txt',  
'neg/cv784\_16077.txt',  
'neg/cv785\_23748.txt',  
'neg/cv786\_23608.txt',  
'neg/cv787\_15277.txt',  
'neg/cv788\_26409.txt',  
'neg/cv789\_12991.txt',  
'neg/cv790\_16202.txt',  
'neg/cv791\_17995.txt',  
'neg/cv792\_3257.txt',  
'neg/cv793\_15235.txt',  
'neg/cv794\_17353.txt',  
'neg/cv795\_10291.txt',  
'neg/cv796\_17243.txt',  
'neg/cv797\_7245.txt',  
'neg/cv798\_24779.txt',  
'neg/cv799\_19812.txt',  
'neg/cv800\_13494.txt',  
'neg/cv801\_26335.txt',  
'neg/cv802\_28381.txt',  
'neg/cv803\_8584.txt',  
'neg/cv804\_11763.txt',  
'neg/cv805\_21128.txt',  
'neg/cv806\_9405.txt',  
'neg/cv807\_23024.txt',  
'neg/cv808\_13773.txt',  
'neg/cv809\_5012.txt',  
'neg/cv810\_13660.txt',  
'neg/cv811\_22646.txt',  
'neg/cv812\_19051.txt',  
'neg/cv813\_6649.txt',  
'neg/cv814\_20316.txt',  
'neg/cv815\_23466.txt',  
'neg/cv816\_15257.txt',  
'neg/cv817\_3675.txt',  
'neg/cv818\_10698.txt',  
'neg/cv819\_9567.txt',  
'neg/cv820\_24157.txt',  
'neg/cv821\_29283.txt',  
'neg/cv822\_21545.txt',  
'neg/cv823\_17055.txt',  
'neg/cv824\_9335.txt',  
'neg/cv825\_5168.txt',  
'neg/cv826\_12761.txt',  
'neg/cv827\_19479.txt',  
'neg/cv828\_21392.txt',

'neg/cv829\_21725.txt',  
'neg/cv830\_5778.txt',  
'neg/cv831\_16325.txt',  
'neg/cv832\_24713.txt',  
'neg/cv833\_11961.txt',  
'neg/cv834\_23192.txt',  
'neg/cv835\_20531.txt',  
'neg/cv836\_14311.txt',  
'neg/cv837\_27232.txt',  
'neg/cv838\_25886.txt',  
'neg/cv839\_22807.txt',  
'neg/cv840\_18033.txt',  
'neg/cv841\_3367.txt',  
'neg/cv842\_5702.txt',  
'neg/cv843\_17054.txt',  
'neg/cv844\_13890.txt',  
'neg/cv845\_15886.txt',  
'neg/cv846\_29359.txt',  
'neg/cv847\_20855.txt',  
'neg/cv848\_10061.txt',  
'neg/cv849\_17215.txt',  
'neg/cv850\_18185.txt',  
'neg/cv851\_21895.txt',  
'neg/cv852\_27512.txt',  
'neg/cv853\_29119.txt',  
'neg/cv854\_18955.txt',  
'neg/cv855\_22134.txt',  
'neg/cv856\_28882.txt',  
'neg/cv857\_17527.txt',  
'neg/cv858\_20266.txt',  
'neg/cv859\_15689.txt',  
'neg/cv860\_15520.txt',  
'neg/cv861\_12809.txt',  
'neg/cv862\_15924.txt',  
'neg/cv863\_7912.txt',  
'neg/cv864\_3087.txt',  
'neg/cv865\_28796.txt',  
'neg/cv866\_29447.txt',  
'neg/cv867\_18362.txt',  
'neg/cv868\_12799.txt',  
'neg/cv869\_24782.txt',  
'neg/cv870\_18090.txt',  
'neg/cv871\_25971.txt',  
'neg/cv872\_13710.txt',  
'neg/cv873\_19937.txt',  
'neg/cv874\_12182.txt',  
'neg/cv875\_5622.txt',  
'neg/cv876\_9633.txt',

'neg/cv877\_29132.txt',  
'neg/cv878\_17204.txt',  
'neg/cv879\_16585.txt',  
'neg/cv880\_29629.txt',  
'neg/cv881\_14767.txt',  
'neg/cv882\_10042.txt',  
'neg/cv883\_27621.txt',  
'neg/cv884\_15230.txt',  
'neg/cv885\_13390.txt',  
'neg/cv886\_19210.txt',  
'neg/cv887\_5306.txt',  
'neg/cv888\_25678.txt',  
'neg/cv889\_22670.txt',  
'neg/cv890\_3515.txt',  
'neg/cv891\_6035.txt',  
'neg/cv892\_18788.txt',  
'neg/cv893\_26731.txt',  
'neg/cv894\_22140.txt',  
'neg/cv895\_22200.txt',  
'neg/cv896\_17819.txt',  
'neg/cv897\_11703.txt',  
'neg/cv898\_1576.txt',  
'neg/cv899\_17812.txt',  
'neg/cv900\_10800.txt',  
'neg/cv901\_11934.txt',  
'neg/cv902\_13217.txt',  
'neg/cv903\_18981.txt',  
'neg/cv904\_25663.txt',  
'neg/cv905\_28965.txt',  
'neg/cv906\_12332.txt',  
'neg/cv907\_3193.txt',  
'neg/cv908\_17779.txt',  
'neg/cv909\_9973.txt',  
'neg/cv910\_21930.txt',  
'neg/cv911\_21695.txt',  
'neg/cv912\_5562.txt',  
'neg/cv913\_29127.txt',  
'neg/cv914\_2856.txt',  
'neg/cv915\_9342.txt',  
'neg/cv916\_17034.txt',  
'neg/cv917\_29484.txt',  
'neg/cv918\_27080.txt',  
'neg/cv919\_18155.txt',  
'neg/cv920\_29423.txt',  
'neg/cv921\_13988.txt',  
'neg/cv922\_10185.txt',  
'neg/cv923\_11951.txt',  
'neg/cv924\_29397.txt',

'neg/cv925\_9459.txt',  
'neg/cv926\_18471.txt',  
'neg/cv927\_11471.txt',  
'neg/cv928\_9478.txt',  
'neg/cv929\_1841.txt',  
'neg/cv930\_14949.txt',  
'neg/cv931\_18783.txt',  
'neg/cv932\_14854.txt',  
'neg/cv933\_24953.txt',  
'neg/cv934\_20426.txt',  
'neg/cv935\_24977.txt',  
'neg/cv936\_17473.txt',  
'neg/cv937\_9816.txt',  
'neg/cv938\_10706.txt',  
'neg/cv939\_11247.txt',  
'neg/cv940\_18935.txt',  
'neg/cv941\_10718.txt',  
'neg/cv942\_18509.txt',  
'neg/cv943\_23547.txt',  
'neg/cv944\_15042.txt',  
'neg/cv945\_13012.txt',  
'neg/cv946\_20084.txt',  
'neg/cv947\_11316.txt',  
'neg/cv948\_25870.txt',  
'neg/cv949\_21565.txt',  
'neg/cv950\_13478.txt',  
'neg/cv951\_11816.txt',  
'neg/cv952\_26375.txt',  
'neg/cv953\_7078.txt',  
'neg/cv954\_19932.txt',  
'neg/cv955\_26154.txt',  
'neg/cv956\_12547.txt',  
'neg/cv957\_9059.txt',  
'neg/cv958\_13020.txt',  
'neg/cv959\_16218.txt',  
'neg/cv960\_28877.txt',  
'neg/cv961\_5578.txt',  
'neg/cv962\_9813.txt',  
'neg/cv963\_7208.txt',  
'neg/cv964\_5794.txt',  
'neg/cv965\_26688.txt',  
'neg/cv966\_28671.txt',  
'neg/cv967\_5626.txt',  
'neg/cv968\_25413.txt',  
'neg/cv969\_14760.txt',  
'neg/cv970\_19532.txt',  
'neg/cv971\_11790.txt',  
'neg/cv972\_26837.txt',

```
'neg/cv973_10171.txt',
'neg/cv974_24303.txt',
'neg/cv975_11920.txt',
'neg/cv976_10724.txt',
'neg/cv977_4776.txt',
'neg/cv978_22192.txt',
'neg/cv979_2029.txt',
'neg/cv980_11851.txt',
'neg/cv981_16679.txt',
'neg/cv982_22209.txt',
'neg/cv983_24219.txt',
'neg/cv984_14006.txt',
'neg/cv985_5964.txt',
'neg/cv986_15092.txt',
'neg/cv987_7394.txt',
'neg/cv988_20168.txt',
'neg/cv989_17297.txt',
'neg/cv990_12443.txt',
'neg/cv991_19973.txt',
'neg/cv992_12806.txt',
'neg/cv993_29565.txt',
'neg/cv994_13229.txt',
'neg/cv995_23113.txt',
'neg/cv996_12447.txt',
'neg/cv997_5152.txt',
'neg/cv998_15691.txt',
'neg/cv999_14636.txt',
...]
```

```
In [3]: len(movie_reviews.fileids())
```

```
Out[3]: 2000
```

Each corpus reader provides a variety of methods to read data from the corpus, depending on the format of the corpus. For example, plaintext corpora support methods to read the corpus as raw text, a list of words, a list of sentences, or a list of paragraphs.

```
In [4]: #movie_reviews.raw("neg/cv000_29416.txt").replace("\n", "").replace("'", "'').replace(
```

## 1.2 2. Building and testing the classifier

```
In [5]: from nltk.corpus import stopwords
```

```
stops = stopwords.words('english')
stops.extend('.', [','], (,), ;, /, -, \', ?, ", :, <, >, n\ 't, |, #, \ 's, \ ", \ 're, \ 've, \ 'll, \ 'd, \ 're'.sp
stops.extend(',', ')
stops
```

```
Out[5]: ['i',
         'me',
         'my',
         'myself',
         'we',
         'our',
         'ours',
         'ourselves',
         'you',
         "you're",
         "you've",
         "you'll",
         "you'd",
         'your',
         'yours',
         'yourself',
         'yourselves',
         'he',
         'him',
         'his',
         'himself',
         'she',
         "she's",
         'her',
         'hers',
         'herself',
         'it',
         "it's",
         'its',
         'itself',
         'they',
         'them',
         'their',
         'theirs',
         'themselves',
         'what',
         'which',
         'who',
         'whom',
         'this',
         'that',
         "that'll",
         'these',
         'those',
         'am',
         'is',
         'are',
         'was',
```

'were',  
'be',  
'been',  
'being',  
'have',  
'has',  
'had',  
'having',  
'do',  
'does',  
'did',  
'doing',  
'a',  
'an',  
'the',  
'and',  
'but',  
'if',  
'or',  
'because',  
'as',  
'until',  
'while',  
'of',  
'at',  
'by',  
'for',  
'with',  
'about',  
'against',  
'between',  
'into',  
'through',  
'during',  
'before',  
'after',  
'above',  
'below',  
'to',  
'from',  
'up',  
'down',  
'in',  
'out',  
'on',  
'off',  
'over',  
'under',



'again',  
'further',  
'then',  
'once',  
'here',  
'there',  
'when',  
'where',  
'why',  
'how',  
'all',  
'any',  
'both',  
'each',  
'few',  
'more',  
'most',  
'other',  
'some',  
'such',  
'no',  
'nor',  
'not',  
'only',  
'own',  
'same',  
'so',  
'than',  
'too',  
'very',  
's',  
't',  
'can',  
'will',  
'just',  
'don',  
'don't',  
'should',  
'should've',  
'now',  
'd',  
'll',  
'm',  
'o',  
're',  
've',  
'y',  
'ain',

'aren',  
"aren't",  
'couldn',  
"couldn't",  
'didn',  
"didn't",  
'doesn',  
"doesn't",  
'hadn',  
"hadn't",  
'hasn',  
"hasn't",  
'haven',  
"haven't",  
'isn',  
"isn't",  
'ma',  
'mightn',  
"mightn't",  
'mustn',  
"mustn't",  
'needn',  
"needn't",  
'shan',  
"shan't",  
'shouldn',  
"shouldn't",  
'wasn',  
"wasn't",  
'weren',  
"weren't",  
'won',  
"won't",  
'wouldn',  
"wouldn't",  
'.',  
'[',  
']',  
'(',  
)',  
';',  
'/',  
'-',  
"''",  
'?',  
"''",  
':',  
'<',

```
'>',
'n't",
'|',
'#',
's",
"',
're",
've",
'll",
'd",
're",
',']
```

```
In [6]: from nltk.classify import NaiveBayesClassifier
import nltk.classify.util # Utility functions and classes for classifiers. Contains fu

# Given a word, returns a dict {word: True}. This will be our feature in the classifier
def word_feats(words):
    return dict([(word, True) for word in words if word not in stops and word.isalpha()

pos_ids = movie_reviews.fileids('pos')
neg_ids = movie_reviews.fileids('neg')

len(pos_ids) + len(neg_ids)
```

Out[6]: 2000

```
In [7]: # We take the positive/negative words, create the feature for such words, and store it
pos_feats = [(word_feats(movie_reviews.words(fileids=[f])), 'pos') for f in pos_ids]
neg_feats = [(word_feats(movie_reviews.words(fileids=[f])), 'neg') for f in neg_ids]

#pos_feats
```

```
In [8]: # 3/4 of the features will be used for training.
pos_len_train = int(len(pos_feats) * 3 / 4)
neg_len_train = int(len(neg_feats) * 3 / 4)

pos_len_train
```

Out[8]: 750

```
In [9]: train_feats = neg_feats[:neg_len_train] + pos_feats[:pos_len_train]
test_feats = neg_feats[neg_len_train:] + pos_feats[pos_len_train:]

# Training a NaiveBayesClassifier with our training feature words.
classifier = NaiveBayesClassifier.train(train_feats)

print('Accuracy: ', nltk.classify.util.accuracy(classifier, test_feats))
```

Accuracy: 0.712

```
In [10]: # We can see which words fit best in each class.
         classifier.show_most_informative_features()
```

Most Informative Features

magnificent = True	pos : neg	=	15.0 : 1.0
outstanding = True	pos : neg	=	13.6 : 1.0
insulting = True	neg : pos	=	13.0 : 1.0
vulnerable = True	pos : neg	=	12.3 : 1.0
ludicrous = True	neg : pos	=	11.8 : 1.0
avoids = True	pos : neg	=	11.7 : 1.0
uninvolving = True	neg : pos	=	11.7 : 1.0
astounding = True	pos : neg	=	10.3 : 1.0
fascination = True	pos : neg	=	10.3 : 1.0
idiotic = True	neg : pos	=	9.8 : 1.0

### 1.3 3. Classifying new data

```
In [11]: from nltk import word_tokenize, pos_tag
```

```
        sentence = "I feel so miserable, it makes me amazing"
        tokens = [word for word in word_tokenize(sentence) if word not in stops]
        tokens
```

```
Out[11]: ['I', 'feel', 'miserable', 'makes', 'amazing']
```

```
In [12]: feats = word_feats(word for word in tokens)
        feats
```

```
Out[12]: {'I': True, 'feel': True, 'miserable': True, 'makes': True, 'amazing': True}
```

```
In [13]: classifier.classify(feats)
```

```
Out[13]: 'pos'
```

```
In [14]: sentence2 = "You are a pathetic fool, a terrible excuse for a human being."
        tokens2 = [word for word in word_tokenize(sentence2) if word not in stops]
        tokens2
```

```
Out[14]: ['You', 'pathetic', 'fool', 'terrible', 'excuse', 'human']
```

```
In [15]: pos_tags2 = [pos for pos in pos_tag(tokens2) if pos[1] == 'JJ']
        pos_tags2
```

```
Out[15]: [('pathetic', 'JJ'), ('terrible', 'JJ')]
```

```
In [16]: feats2 = word_feats([word for (word,_) in pos_tags2])
        feats2
```

```
Out[16]: {'pathetic': True, 'terrible': True}
```

```
In [17]: classifier.classify(feats2)
```

```
Out[17]: 'neg'
```

#### 1.4 4. Classifying News Documents in Categories (sport, humor, adventure, etc..)

```
In [18]: from nltk.corpus import brown
```

```
from random import shuffle
```

```
from nltk import NaiveBayesClassifier
```

```
from nltk import FreqDist
```

```
from nltk.classify import accuracy, apply_features
```

```
from nltk.corpus import stopwords
```

```
#given a document extract features (the presence or not of the best1500 top 1500 of f
```

```
def document_features(doc):
```

```
    doc_set_words = set(doc)
```

```
    features_dic = {} #features is a dictionary
```

```
    for word in best1500_words_corpora:
```

```
        features_dic['has(%s)' % word] = (word in doc_set_words)
```

```
    return features_dic
```

```
#Word frequency of words in corpora.
```

```
#Return the top 1500 most used words.
```

```
def init_corpora():
```

```
    words_in_corpora = FreqDist(w.lower() for w in brown.words())
```

```
    # Python 2
```

```
    best1500 = words_in_corpora.keys()[:1500]
```

```
    # Python 3
```

```
    best1500 = list(words_in_corpora.keys()[:1500])
```

```
    return nonstop(best1500)
```

```
#Receives a list of words, and remove those in stopwords.
```

```
def nonstop(listwords):
```

```
    return [word for word in listwords if word not in stopw]
```

```
stopw = stopwords.words('english')
```

```
best1500_words_corpora = init_corpora()
```

```
documents = [(list(nonstop(brown.words(fileid))), categoria) for categoria in brown.
```

```
shuffle(documents)
```

```
print(len(documents))
```

```
#print(documents)
```

```
doc_features_set = [(document_features(d),c) for (d,c) in documents]
```

```
categories = set([c for (_,c) in documents])
```

```
print(categories)
#print(len(doc_features_set))
```

500

```
{'religion', 'romance', 'adventure', 'learned', 'government', 'reviews', 'news', 'lore', 'edit
```

```
In [19]: train_set = doc_features_set[:350]
        #train_set = apply_features(document, features, doc_features_set[:200])
        classifier = NaiveBayesClassifier.train(train_set)
        classifier.show_most_informative_features(30)

        #We expect Adventura as 'cn02' contains words of a books of adventures
        print(classifier.classify(document_features(brown.words('cn02'))))
```

#### Most Informative Features

has(television) = True	humor : learne =	24.9 : 1.0
has(saw) = True	fictio : learne =	20.1 : 1.0
has(quiet) = True	scienc : learne =	19.3 : 1.0
has(afternoon) = True	humor : learne =	19.3 : 1.0
has(midnight) = True	scienc : learne =	19.3 : 1.0
has(announced) = True	news : learne =	19.3 : 1.0
has(pay) = True	govern : learne =	19.3 : 1.0
has(1961) = True	news : belles =	18.7 : 1.0
has(funds) = True	govern : belles =	18.7 : 1.0
has(wife) = True	humor : learne =	18.2 : 1.0
has(church) = True	religi : learne =	18.1 : 1.0
has(returned) = True	myster : learne =	18.0 : 1.0
has(blue) = True	romanc : learne =	17.7 : 1.0
has(race) = True	review : learne =	17.4 : 1.0
has(audience) = True	review : learne =	17.4 : 1.0
has(yes) = True	myster : learne =	17.4 : 1.0
has(notice) = True	myster : belles =	17.3 : 1.0
has(starts) = True	review : belles =	16.8 : 1.0
has(fiscal) = True	govern : belles =	16.7 : 1.0
has(board) = True	scienc : learne =	16.2 : 1.0
has(night) = True	fictio : hobbie =	15.8 : 1.0
has(stood) = True	myster : learne =	15.7 : 1.0
has(schools) = True	editor : learne =	15.7 : 1.0
has(says) = True	review : learne =	15.1 : 1.0
has(telephone) = True	humor : learne =	14.9 : 1.0
has(school) = True	humor : learne =	14.9 : 1.0
has(meeting) = True	humor : belles =	14.4 : 1.0
has(son) = True	romanc : learne =	14.3 : 1.0
has(went) = False	govern : romanc =	14.1 : 1.0
has(insisted) = True	humor : learne =	13.8 : 1.0

romance

```
In [20]: test_set = doc_features_set[150:]  
         print(accuracy(classifier, test_set))
```

0.7057142857142857