# Assignment 1: Exploring Word Vectors (23 Points)

### Estimated Time: ~3 hours

The objective of this assignment is to warm-up you of some python coding, and also get you to familarize with some NLP concepts.

```
In [1]:
```

```python
import sys
assert sys.version_info[0]==3
assert sys.version_info[1] >= 5

from gensim.models import KeyedVectors
from gensim.test.utils import datapath
import pprint
import matplotlib.pyplot as plt
plt.rcParams['figure.figsize'] = [6, 5]
import nltk
nltk.download('reuters')
from nltk.corpus import reuters
import numpy as np
import random
import scipy as sp
from sklearn.decomposition import TruncatedSVD
from sklearn.decomposition import PCA

START_TOKEN = '<START>'
END_TOKEN = '<END>'

np.random.seed(0)
random.seed(0)
```

```
[nltk_data] Downloading package reuters to /root/nltk_data...
```

## Word Vectors

Word Vectors are often used as a fundamental component for downstream NLP tasks, e.g. question answering, text generation, translation, etc., so it is important to build some intuitions as to their strengths and weaknesses. Here, you will explore two types of word vectors: those derived from *co-occurrence matrices*, and those derived via *GloVe*.

**Note on Terminology:** The terms "word vectors" and "word embeddings" are often used interchangeably. The term "embedding" refers to the fact that we are encoding aspects of a word's meaning in a lower dimensional space. As [Wikipedia](#) states, "*conceptually it involves a mathematical embedding from a space with one dimension per word to a continuous vector space with a much lower dimension*".

## Part 1: Count-Based Word Vectors (10 points)

Most word vector models start from the following idea:

*You shall know a word by the company it keeps ([Firth, J. R. 1957:11](#))*

Many word vector implementations are driven by the idea that similar words, i.e., (near) synonyms, will be used in similar contexts. As a result, similar words will often be spoken or written along with a shared subset of words, i.e., contexts. By examining these contexts, we can try to develop embeddings for our words. With this intuition in mind, many "old school" approaches to constructing word vectors relied on word counts. Here we elaborate upon one of those strategies, *co-occurrence matrices* (for more information, see [here](#) or [here](#)).

### Co-Occurrence

A co-occurrence matrix counts how often things co-occur in some environment. Given some word $w_i$ occurring in the document, we consider the *context window* surrounding $w_i$. Supposing our fixed window size is $n$, then this is the $n$ preceding and $n$ subsequent words in that document, i.e. words $w_{i-n} \dots w_{i-1}$ and $w_{i+1} \dots w_{i+n}$. We build a *co-occurrence matrix* $M$, which is a symmetric word-by-word matrix in which $M_{ij}$ is the number of times $w_j$ appears inside $w_i$'s window among all documents.

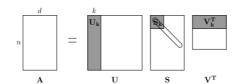**Example: Co-Occurrence with Fixed Window of n=1:**

Document 1: "all that glitters is not gold"

Document 2: "all is well that ends well"

| * | `<START>` | all | that | glitters | is | not | gold | well | ends | `<END>` |
|---|---|---|---|---|---|---|---|---|---|---|
| `<START>` | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| all | 2 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| that | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 |
| glitters | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| is | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 |
| not | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 |
| gold | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 |
| well | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 1 |
| ends | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| `<END>` | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 |

**Note:** In NLP, we often add `<START>` and `<END>` tokens to represent the beginning and end of sentences, paragraphs or documents. In thise case we imagine `<START>` and `<END>` tokens encapsulating each document, e.g., "`<START>` All that glitters is not gold `<END>`", and include these tokens in our co-occurrence counts.

The rows (or columns) of this matrix provide one type of word vectors (those based on word-word co-occurrence), but the vectors will be large in general (linear in the number of distinct words in a corpus). Thus, our next step is to run *dimensionality reduction*. In particular, we will run *SVD (Singular Value Decomposition)*, which is a kind of generalized *PCA (Principal Components Analysis)* to select the top $k$ principal components. Here's a visualization of dimensionality reduction with SVD. In this picture our co-occurrence matrix is $A$ with $n$ rows corresponding to $n$ words. We obtain a full matrix decomposition, with the singular values ordered in the diagonal $S$ matrix, and our new, shorter length-$k$ word vectors in $U_k$.



This reduced-dimensionality co-occurrence representation preserves semantic relationships between words, e.g. *doctor* and *hospital* will be closer than *doctor* and *dog*.

**Notes:** If you can barely remember what an eigenvalue is, here's [a slow, friendly introduction to SVD](#). Though, for the purpose of this class, you only need to know how to extract the k-dimensional embeddings by utilizing pre-programmed implementations of these algorithms from the numpy, scipy, or sklearn python packages. In practice, it is challenging to apply full SVD to large corpora because of the memory needed to perform PCA or SVD. However, if you only want the top $k$ vector components for relatively small $k$ — known as [Truncated SVD](#) — then there are reasonably scalable techniques to compute those iteratively.

### Plotting Co-Occurrence Word Embeddings

Here, we will be using the Reuters (business and financial news) corpus. If you haven't run the import cell at the top of this page, please run it now (click it and press SHIFT-RETURN). The corpus consists of 10,788 news documents totaling 1.3 million words. These documents span 90 categories and are split into train and test. For more details, please see [https://www.nltk.org/book/ch02.html](https://www.nltk.org/book/ch02.html). We provide a `read_corpus` function below that pulls out only articles from the "crude" (i.e. news articles about oil, gas, etc.) category. The function also adds `<START>` and `<END>` tokens to each of the documents, and lowercases words. You do **not** have to perform any other kind of pre-processing.

```
In [2]:
```

```python
def read_corpus(category="crude"):
    """ Read files from the specified Reuter's category.
        Params:
            category (string): category name
        Return:
            list of lists, with words from each of the processed files
    """
    files = reuters.fileids(category)
    return [[START_TOKEN] + [w.lower() for w in list(reuters.words(f))] + [END_TOKEN] for f in files]
```

```
In [5]:
```

```
!unzip /root/nltk_data/corpora/reuters.zip -d /root/nltk_data/corpora
```

```
Streaming output truncated to the last 5000 lines.
  inflating: /root/nltk_data/corpora/reuters/training/2231
  inflating: /root/nltk_data/corpora/reuters/training/2232
  inflating: /root/nltk_data/corpora/reuters/training/2234
  inflating: /root/nltk_data/corpora/reuters/training/2236
  inflating: /root/nltk_data/corpora/reuters/training/2237
  inflating: /root/nltk_data/corpora/reuters/training/2238
  inflating: /root/nltk_data/corpora/reuters/training/2239
  inflating: /root/nltk_data/corpora/reuters/training/2240
  inflating: /root/nltk_data/corpora/reuters/training/2244
  inflating: /root/nltk_data/corpora/reuters/training/2246
  inflating: /root/nltk_data/corpora/reuters/training/2247
  inflating: /root/nltk_data/corpora/reuters/training/2249
  inflating: /root/nltk_data/corpora/reuters/training/225
  inflating: /root/nltk_data/corpora/reuters/training/2251
  inflating: /root/nltk_data/corpora/reuters/training/2252
  inflating: /root/nltk_data/corpora/reuters/training/2253
  inflating: /root/nltk_data/corpora/reuters/training/2257
  inflating: /root/nltk_data/corpora/reuters/training/2259
  inflating: /root/nltk_data/corpora/reuters/training/2260
  inflating: /root/nltk_data/corpora/reuters/training/2264
  inflating: /root/nltk_data/corpora/reuters/training/2265
  inflating: /root/nltk_data/corpora/reuters/training/2266
  inflating: /root/nltk_data/corpora/reuters/training/2268
  inflating: /root/nltk_data/corpora/reuters/training/2269
  inflating: /root/nltk_data/corpora/reuters/training/227
  inflating: /root/nltk_data/corpora/reuters/training/2270
  inflating: /root/nltk_data/corpora/reuters/training/2271
  inflating: /root/nltk_data/corpora/reuters/training/2273
  inflating: /root/nltk_data/corpora/reuters/training/2275
  inflating: /root/nltk_data/corpora/reuters/training/2276
  inflating: /root/nltk_data/corpora/reuters/training/2278
  inflating: /root/nltk_data/corpora/reuters/training/2279
  inflating: /root/nltk_data/corpora/reuters/training/228
  inflating: /root/nltk_data/corpora/reuters/training/2281
  inflating: /root/nltk_data/corpora/reuters/training/2286
  inflating: /root/nltk_data/corpora/reuters/training/229
  inflating: /root/nltk_data/corpora/reuters/training/2290
  inflating: /root/nltk_data/corpora/reuters/training/2291
  inflating: /root/nltk_data/corpora/reuters/training/2297
  inflating: /root/nltk_data/corpora/reuters/training/2298
  inflating: /root/nltk_data/corpora/reuters/training/23
  inflating: /root/nltk_data/corpora/reuters/training/2301
  inflating: /root/nltk_data/corpora/reuters/training/2303
  inflating: /root/nltk_data/corpora/reuters/training/2304
  inflating: /root/nltk_data/corpora/reuters/training/2307
  inflating: /root/nltk_data/corpora/reuters/training/2308
  inflating: /root/nltk_data/corpora/reuters/training/2309
  inflating: /root/nltk_data/corpora/reuters/training/2310
  inflating: /root/nltk_data/corpora/reuters/training/2312
  inflating: /root/nltk_data/corpora/reuters/training/2316
  inflating: /root/nltk_data/corpora/reuters/training/232
  inflating: /root/nltk_data/corpora/reuters/training/2320
  inflating: /root/nltk_data/corpora/reuters/training/2322
  inflating: /root/nltk_data/corpora/reuters/training/2323
  inflating: /root/nltk_data/corpora/reuters/training/2325
  inflating: /root/nltk_data/corpora/reuters/training/2326
  inflating: /root/nltk_data/corpora/reuters/training/2327
  inflating: /root/nltk_data/corpora/reuters/training/233
  inflating: /root/nltk_data/corpora/reuters/training/2332
  inflating: /root/nltk_data/corpora/reuters/training/2333
  inflating: /root/nltk_data/corpora/reuters/training/2335
  inflating: /root/nltk_data/corpora/reuters/training/2336
  inflating: /root/nltk_data/corpora/reuters/training/2337
  inflating: /root/nltk_data/corpora/reuters/training/2338
  inflating: /root/nltk_data/corpora/reuters/training/2339
  inflating: /root/nltk_data/corpora/reuters/training/2340
  inflating: /root/nltk_data/corpora/reuters/training/2342
  inflating: /root/nltk_data/corpora/reuters/training/2344
  inflating: /root/nltk_data/corpora/reuters/training/2346
  inflating: /root/nltk_data/corpora/reuters/training/2347
  inflating: /root/nltk_data/corpora/reuters/training/2348
  inflating: /root/nltk_data/corpora/reuters/training/235
  inflating: /root/nltk_data/corpora/reuters/training/2350
  inflating: /root/nltk_data/corpora/reuters/training/2352
  inflating: /root/nltk_data/corpora/reuters/training/2353
  inflating: /root/nltk_data/corpora/reuters/training/2354
  inflating: /root/nltk_data/corpora/reuters/training/2357
  inflating: /root/nltk_data/corpora/reuters/training/2358
  inflating: /root/nltk_data/corpora/reuters/training/2359
  inflating: /root/nltk_data/corpora/reuters/training/236
  inflating: /root/nltk_data/corpora/reuters/training/2365
  inflating: /root/nltk_data/corpora/reuters/training/2366
  inflating: /root/nltk_data/corpora/reuters/training/2367
  inflating: /root/nltk_data/corpora/reuters/training/237
  inflating: /root/nltk_data/corpora/reuters/training/2370
  inflating: /root/nltk_data/corpora/reuters/training/2371
  inflating: /root/nltk_data/corpora/reuters/training/2374
  inflating: /root/nltk_data/corpora/reuters/training/2375
  inflating: /root/nltk_data/corpora/reuters/training/2376
  inflating: /root/nltk_data/corpora/reuters/training/2377
  inflating: /root/nltk_data/corpora/reuters/training/2379
  inflating: /root/nltk_data/corpora/reuters/training/2380
  inflating: /root/nltk_data/corpora/reuters/training/2382
  inflating: /root/nltk_data/corpora/reuters/training/2383
  inflating: /root/nltk_data/corpora/reuters/training/2385
  inflating: /root/nltk_data/corpora/reuters/training/2386
  inflating: /root/nltk_data/corpora/reuters/training/2389
  inflating: /root/nltk_data/corpora/reuters/training/2390
  inflating: /root/nltk_data/corpora/reuters/training/2391
  inflating: /root/nltk_data/corpora/reuters/training/2393
  inflating: /root/nltk_data/corpora/reuters/training/2394
  inflating: /root/nltk_data/corpora/reuters/training/2396
  inflating: /root/nltk_data/corpora/reuters/training/2397
  inflating: /root/nltk_data/corpora/reuters/training/2398
  inflating: /root/nltk_data/corpora/reuters/training/24
  inflating: /root/nltk_data/corpora/reuters/training/2402
  inflating: [/root/nltk_data/corpora/reuters/training/2403
  inflating: /root/nltk_data/corpora/reuters/training/2405
  inflating: /root/nltk_data/corpora/reuters/training/2406
  inflating: /root/nltk_data/corpora/reuters/training/2409
  inflating: /root/nltk_data/corpora/reuters/training/241
  inflating: /root/nltk_data/corpora/reuters/training/2410
  inflating: /root/nltk_data/corpora/reuters/training/2411
  inflating: /root/nltk_data/corpora/reuters/training/2417
  inflating: /root/nltk_data/corpora/reuters/training/2419
  inflating: /root/nltk_data/corpora/reuters/training/242
  inflating: /root/nltk_data/corpora/reuters/training/2420
```

```
inflating: /root/nltk_data/corpora/reuters/training/2421
inflating: /root/nltk_data/corpora/reuters/training/2425
inflating: /root/nltk_data/corpora/reuters/training/2428
inflating: /root/nltk_data/corpora/reuters/training/2429
inflating: /root/nltk_data/corpora/reuters/training/2433
inflating: /root/nltk_data/corpora/reuters/training/2434
inflating: /root/nltk_data/corpora/reuters/training/2435
inflating: /root/nltk_data/corpora/reuters/training/2436
inflating: /root/nltk_data/corpora/reuters/training/2437
inflating: /root/nltk_data/corpora/reuters/training/2440
inflating: /root/nltk_data/corpora/reuters/training/2441
inflating: /root/nltk_data/corpora/reuters/training/2442
inflating: /root/nltk_data/corpora/reuters/training/2447
inflating: /root/nltk_data/corpora/reuters/training/2448
inflating: /root/nltk_data/corpora/reuters/training/2449
inflating: /root/nltk_data/corpora/reuters/training/2452
inflating: /root/nltk_data/corpora/reuters/training/2456
inflating: /root/nltk_data/corpora/reuters/training/2458
inflating: /root/nltk_data/corpora/reuters/training/2459
inflating: /root/nltk_data/corpora/reuters/training/246
inflating: /root/nltk_data/corpora/reuters/training/2460
inflating: /root/nltk_data/corpora/reuters/training/2462
inflating: /root/nltk_data/corpora/reuters/training/2465
inflating: /root/nltk_data/corpora/reuters/training/2466
inflating: /root/nltk_data/corpora/reuters/training/2467
inflating: /root/nltk_data/corpora/reuters/training/2468
inflating: /root/nltk_data/corpora/reuters/training/247
inflating: /root/nltk_data/corpora/reuters/training/2470
inflating: /root/nltk_data/corpora/reuters/training/2471
inflating: /root/nltk_data/corpora/reuters/training/2474
inflating: /root/nltk_data/corpora/reuters/training/2475
inflating: /root/nltk_data/corpora/reuters/training/2476
inflating: /root/nltk_data/corpora/reuters/training/2478
inflating: /root/nltk_data/corpora/reuters/training/2479
inflating: /root/nltk_data/corpora/reuters/training/248
inflating: /root/nltk_data/corpora/reuters/training/2480
inflating: /root/nltk_data/corpora/reuters/training/2482
inflating: /root/nltk_data/corpora/reuters/training/2483
inflating: /root/nltk_data/corpora/reuters/training/2484
inflating: /root/nltk_data/corpora/reuters/training/2488
inflating: /root/nltk_data/corpora/reuters/training/249
inflating: /root/nltk_data/corpora/reuters/training/2490
inflating: /root/nltk_data/corpora/reuters/training/2491
inflating: /root/nltk_data/corpora/reuters/training/2492
inflating: /root/nltk_data/corpora/reuters/training/2494
inflating: /root/nltk_data/corpora/reuters/training/2495
inflating: /root/nltk_data/corpora/reuters/training/2496
inflating: /root/nltk_data/corpora/reuters/training/2497
inflating: /root/nltk_data/corpora/reuters/training/2499
inflating: /root/nltk_data/corpora/reuters/training/2500
inflating: /root/nltk_data/corpora/reuters/training/2504
inflating: /root/nltk_data/corpora/reuters/training/2505
inflating: /root/nltk_data/corpora/reuters/training/2508
inflating: /root/nltk_data/corpora/reuters/training/2511
inflating: /root/nltk_data/corpora/reuters/training/2512
inflating: /root/nltk_data/corpora/reuters/training/2513
inflating: /root/nltk_data/corpora/reuters/training/2515
inflating: /root/nltk_data/corpora/reuters/training/2517
inflating: /root/nltk_data/corpora/reuters/training/2518
inflating: /root/nltk_data/corpora/reuters/training/2519
inflating: /root/nltk_data/corpora/reuters/training/2520
inflating: /root/nltk_data/corpora/reuters/training/2521
inflating: /root/nltk_data/corpora/reuters/training/2522
inflating: /root/nltk_data/corpora/reuters/training/2524
inflating: /root/nltk_data/corpora/reuters/training/2526
inflating: /root/nltk_data/corpora/reuters/training/2528
inflating: /root/nltk_data/corpora/reuters/training/2529
inflating: /root/nltk_data/corpora/reuters/training/253
inflating: /root/nltk_data/corpora/reuters/training/2530
inflating: /root/nltk_data/corpora/reuters/training/2531
inflating: /root/nltk_data/corpora/reuters/training/2532
inflating: /root/nltk_data/corpora/reuters/training/2534
inflating: /root/nltk_data/corpora/reuters/training/2535
inflating: /root/nltk_data/corpora/reuters/training/2538
inflating: /root/nltk_data/corpora/reuters/training/2539
inflating: /root/nltk_data/corpora/reuters/training/254
inflating: /root/nltk_data/corpora/reuters/training/2540
inflating: /root/nltk_data/corpora/reuters/training/2542
inflating: /root/nltk_data/corpora/reuters/training/2543
inflating: /root/nltk_data/corpora/reuters/training/2544
inflating: /root/nltk_data/corpora/reuters/training/2546
inflating: /root/nltk_data/corpora/reuters/training/2548
inflating: /root/nltk_data/corpora/reuters/training/2549
inflating: /root/nltk_data/corpora/reuters/training/2550
inflating: /root/nltk_data/corpora/reuters/training/2551
inflating: /root/nltk_data/corpora/reuters/training/2552
inflating: /root/nltk_data/corpora/reuters/training/2553
inflating: /root/nltk_data/corpora/reuters/training/2554
inflating: /root/nltk_data/corpora/reuters/training/2555
inflating: /root/nltk_data/corpora/reuters/training/2557
inflating: /root/nltk_data/corpora/reuters/training/2558
inflating: /root/nltk_data/corpora/reuters/training/2559
inflating: /root/nltk_data/corpora/reuters/training/256
inflating: /root/nltk_data/corpora/reuters/training/2560
inflating: /root/nltk_data/corpora/reuters/training/2561
inflating: /root/nltk_data/corpora/reuters/training/2562
inflating: /root/nltk_data/corpora/reuters/training/2564
inflating: /root/nltk_data/corpora/reuters/training/2565
inflating: /root/nltk_data/corpora/reuters/training/2566
inflating: /root/nltk_data/corpora/reuters/training/2567
inflating: /root/nltk_data/corpora/reuters/training/2568
inflating: /root/nltk_data/corpora/reuters/training/2569
inflating: /root/nltk_data/corpora/reuters/training/2570
inflating: /root/nltk_data/corpora/reuters/training/2572
inflating: /root/nltk_data/corpora/reuters/training/2574
inflating: /root/nltk_data/corpora/reuters/training/2576
inflating: /root/nltk_data/corpora/reuters/training/2577
inflating: /root/nltk_data/corpora/reuters/training/2578
inflating: /root/nltk_data/corpora/reuters/training/2579
inflating: /root/nltk_data/corpora/reuters/training/2580
inflating: /root/nltk_data/corpora/reuters/training/2582
inflating: /root/nltk_data/corpora/reuters/training/2583
inflating: /root/nltk_data/corpora/reuters/training/2584
inflating: /root/nltk_data/corpora/reuters/training/2585
inflating: /root/nltk_data/corpora/reuters/training/2587
inflating: /root/nltk_data/corpora/reuters/training/259
inflating: /root/nltk_data/corpora/reuters/training/2590
inflating: /root/nltk_data/corpora/reuters/training/2591
inflating: /root/nltk_data/corpora/reuters/training/2593
inflating: /root/nltk_data/corpora/reuters/training/2594
inflating: /root/nltk_data/corpora/reuters/training/2595
inflating: /root/nltk_data/corpora/reuters/training/2596
inflating: /root/nltk_data/corpora/reuters/training/2599
inflating: /root/nltk_data/corpora/reuters/training/260
inflating: /root/nltk_data/corpora/reuters/training/2601
inflating: /root/nltk_data/corpora/reuters/training/2602
inflating: /root/nltk_data/corpora/reuters/training/2604
inflating: /root/nltk_data/corpora/reuters/training/2606
inflating: /root/nltk_data/corpora/reuters/training/2608
inflating: /root/nltk_data/corpora/reuters/training/2609
inflating: /root/nltk_data/corpora/reuters/training/2611
inflating: /root/nltk_data/corpora/reuters/training/2613
inflating: /root/nltk_data/corpora/reuters/training/2617
inflating: /root/nltk_data/corpora/reuters/training/2618
```

```
inflating: /root/nltk_data/corpora/reuters/training/2619
inflating: /root/nltk_data/corpora/reuters/training/2620
inflating: /root/nltk_data/corpora/reuters/training/2622
inflating: /root/nltk_data/corpora/reuters/training/2623
inflating: /root/nltk_data/corpora/reuters/training/2626
inflating: /root/nltk_data/corpora/reuters/training/2630
inflating: /root/nltk_data/corpora/reuters/training/2633
inflating: /root/nltk_data/corpora/reuters/training/2634
inflating: /root/nltk_data/corpora/reuters/training/2635
inflating: /root/nltk_data/corpora/reuters/training/2637
inflating: /root/nltk_data/corpora/reuters/training/2638
inflating: /root/nltk_data/corpora/reuters/training/2645
inflating: /root/nltk_data/corpora/reuters/training/2647
inflating: /root/nltk_data/corpora/reuters/training/2648
inflating: /root/nltk_data/corpora/reuters/training/2649
inflating: /root/nltk_data/corpora/reuters/training/2650
inflating: /root/nltk_data/corpora/reuters/training/2652
inflating: /root/nltk_data/corpora/reuters/training/2658
inflating: /root/nltk_data/corpora/reuters/training/2659
inflating: /root/nltk_data/corpora/reuters/training/2660
inflating: /root/nltk_data/corpora/reuters/training/2663
inflating: /root/nltk_data/corpora/reuters/training/2664
inflating: /root/nltk_data/corpora/reuters/training/2666
inflating: /root/nltk_data/corpora/reuters/training/2667
inflating: /root/nltk_data/corpora/reuters/training/267
inflating: /root/nltk_data/corpora/reuters/training/2671
inflating: /root/nltk_data/corpora/reuters/training/2672
inflating: /root/nltk_data/corpora/reuters/training/2674
inflating: /root/nltk_data/corpora/reuters/training/2677
inflating: /root/nltk_data/corpora/reuters/training/2678
inflating: /root/nltk_data/corpora/reuters/training/2681
inflating: /root/nltk_data/corpora/reuters/training/2683
inflating: /root/nltk_data/corpora/reuters/training/2686
inflating: /root/nltk_data/corpora/reuters/training/2687
inflating: /root/nltk_data/corpora/reuters/training/2688
inflating: /root/nltk_data/corpora/reuters/training/2691
inflating: /root/nltk_data/corpora/reuters/training/2693
inflating: /root/nltk_data/corpora/reuters/training/2695
inflating: /root/nltk_data/corpora/reuters/training/2696
inflating: /root/nltk_data/corpora/reuters/training/2697
inflating: /root/nltk_data/corpora/reuters/training/2699
inflating: /root/nltk_data/corpora/reuters/training/27
inflating: /root/nltk_data/corpora/reuters/training/2700
inflating: /root/nltk_data/corpora/reuters/training/2708
inflating: /root/nltk_data/corpora/reuters/training/2709
inflating: /root/nltk_data/corpora/reuters/training/271
inflating: /root/nltk_data/corpora/reuters/training/2711
inflating: /root/nltk_data/corpora/reuters/training/2714
inflating: /root/nltk_data/corpora/reuters/training/2718
inflating: /root/nltk_data/corpora/reuters/training/272
inflating: /root/nltk_data/corpora/reuters/training/2720
inflating: /root/nltk_data/corpora/reuters/training/2722
inflating: /root/nltk_data/corpora/reuters/training/2723
inflating: /root/nltk_data/corpora/reuters/training/2725
inflating: /root/nltk_data/corpora/reuters/training/2727
inflating: /root/nltk_data/corpora/reuters/training/2728
inflating: /root/nltk_data/corpora/reuters/training/2729
inflating: /root/nltk_data/corpora/reuters/training/273
inflating: /root/nltk_data/corpora/reuters/training/2737
inflating: /root/nltk_data/corpora/reuters/training/274
inflating: /root/nltk_data/corpora/reuters/training/2741
inflating: /root/nltk_data/corpora/reuters/training/2742
inflating: /root/nltk_data/corpora/reuters/training/2743
inflating: /root/nltk_data/corpora/reuters/training/2744
inflating: /root/nltk_data/corpora/reuters/training/2746
inflating: /root/nltk_data/corpora/reuters/training/2747
inflating: /root/nltk_data/corpora/reuters/training/2748
inflating: /root/nltk_data/corpora/reuters/training/2749
inflating: /root/nltk_data/corpora/reuters/training/275
inflating: /root/nltk_data/corpora/reuters/training/2750
inflating: /root/nltk_data/corpora/reuters/training/2751
inflating: /root/nltk_data/corpora/reuters/training/2752
inflating: /root/nltk_data/corpora/reuters/training/2754
inflating: /root/nltk_data/corpora/reuters/training/2756
inflating: /root/nltk_data/corpora/reuters/training/2757
inflating: /root/nltk_data/corpora/reuters/training/2760
inflating: /root/nltk_data/corpora/reuters/training/2761
inflating: /root/nltk_data/corpora/reuters/training/2762
inflating: /root/nltk_data/corpora/reuters/training/2763
inflating: /root/nltk_data/corpora/reuters/training/2764
inflating: /root/nltk_data/corpora/reuters/training/2765
inflating: /root/nltk_data/corpora/reuters/training/2767
inflating: /root/nltk_data/corpora/reuters/training/2768
inflating: /root/nltk_data/corpora/reuters/training/2769
inflating: /root/nltk_data/corpora/reuters/training/2771
inflating: /root/nltk_data/corpora/reuters/training/2772
inflating: /root/nltk_data/corpora/reuters/training/2773
inflating: /root/nltk_data/corpora/reuters/training/2775
inflating: /root/nltk_data/corpora/reuters/training/2777
inflating: /root/nltk_data/corpora/reuters/training/2779
inflating: /root/nltk_data/corpora/reuters/training/2781
inflating: /root/nltk_data/corpora/reuters/training/2782
inflating: /root/nltk_data/corpora/reuters/training/2783
inflating: /root/nltk_data/corpora/reuters/training/2784
inflating: /root/nltk_data/corpora/reuters/training/2785
inflating: /root/nltk_data/corpora/reuters/training/2787
inflating: /root/nltk_data/corpora/reuters/training/2789
inflating: /root/nltk_data/corpora/reuters/training/279
inflating: /root/nltk_data/corpora/reuters/training/2791
inflating: /root/nltk_data/corpora/reuters/training/2792
inflating: /root/nltk_data/corpora/reuters/training/2793
inflating: /root/nltk_data/corpora/reuters/training/2794
inflating: /root/nltk_data/corpora/reuters/training/2797
inflating: /root/nltk_data/corpora/reuters/training/2798
inflating: /root/nltk_data/corpora/reuters/training/2799
inflating: /root/nltk_data/corpora/reuters/training/2800
inflating: /root/nltk_data/corpora/reuters/training/2801
inflating: /root/nltk_data/corpora/reuters/training/2802
inflating: /root/nltk_data/corpora/reuters/training/2803
inflating: /root/nltk_data/corpora/reuters/training/2806
inflating: /root/nltk_data/corpora/reuters/training/2807
inflating: /root/nltk_data/corpora/reuters/training/2809
inflating: /root/nltk_data/corpora/reuters/training/281
inflating: /root/nltk_data/corpora/reuters/training/2810
inflating: /root/nltk_data/corpora/reuters/training/2811
inflating: /root/nltk_data/corpora/reuters/training/2813
inflating: /root/nltk_data/corpora/reuters/training/2815
inflating: /root/nltk_data/corpora/reuters/training/2818
inflating: /root/nltk_data/corpora/reuters/training/2819
inflating: /root/nltk_data/corpora/reuters/training/2820
inflating: /root/nltk_data/corpora/reuters/training/2825
inflating: /root/nltk_data/corpora/reuters/training/2826
inflating: /root/nltk_data/corpora/reuters/training/2827
inflating: /root/nltk_data/corpora/reuters/training/2828
inflating: /root/nltk_data/corpora/reuters/training/283
inflating: /root/nltk_data/corpora/reuters/training/2833
inflating: /root/nltk_data/corpora/reuters/training/2834
inflating: /root/nltk_data/corpora/reuters/training/2836
inflating: /root/nltk_data/corpora/reuters/training/2838
inflating: /root/nltk_data/corpora/reuters/training/2839
inflating: /root/nltk_data/corpora/reuters/training/284
inflating: /root/nltk_data/corpora/reuters/training/2843
inflating: /root/nltk_data/corpora/reuters/training/2847
inflating: /root/nltk_data/corpora/reuters/training/2848
```

```
inflating: /root/nltk_data/corpora/reuters/training/2849
inflating: /root/nltk_data/corpora/reuters/training/2850
inflating: /root/nltk_data/corpora/reuters/training/2851
inflating: /root/nltk_data/corpora/reuters/training/2853
inflating: /root/nltk_data/corpora/reuters/training/2856
inflating: /root/nltk_data/corpora/reuters/training/2857
inflating: /root/nltk_data/corpora/reuters/training/2858
inflating: /root/nltk_data/corpora/reuters/training/2862
inflating: /root/nltk_data/corpora/reuters/training/2864
inflating: /root/nltk_data/corpora/reuters/training/2866
inflating: /root/nltk_data/corpora/reuters/training/2867
inflating: /root/nltk_data/corpora/reuters/training/2870
inflating: /root/nltk_data/corpora/reuters/training/2872
inflating: /root/nltk_data/corpora/reuters/training/2873
inflating: /root/nltk_data/corpora/reuters/training/2877
inflating: /root/nltk_data/corpora/reuters/training/2880
inflating: /root/nltk_data/corpora/reuters/training/2881
inflating: /root/nltk_data/corpora/reuters/training/2884
inflating: /root/nltk_data/corpora/reuters/training/2886
inflating: /root/nltk_data/corpora/reuters/training/2887
inflating: /root/nltk_data/corpora/reuters/training/2890
inflating: /root/nltk_data/corpora/reuters/training/2896
inflating: /root/nltk_data/corpora/reuters/training/2897
inflating: /root/nltk_data/corpora/reuters/training/2898
inflating: /root/nltk_data/corpora/reuters/training/29
inflating: /root/nltk_data/corpora/reuters/training/290
inflating: /root/nltk_data/corpora/reuters/training/2902
inflating: /root/nltk_data/corpora/reuters/training/2907
inflating: /root/nltk_data/corpora/reuters/training/2908
inflating: /root/nltk_data/corpora/reuters/training/2909
inflating: /root/nltk_data/corpora/reuters/training/2912
inflating: /root/nltk_data/corpora/reuters/training/2913
inflating: /root/nltk_data/corpora/reuters/training/2915
inflating: /root/nltk_data/corpora/reuters/training/2916
inflating: /root/nltk_data/corpora/reuters/training/2917
inflating: /root/nltk_data/corpora/reuters/training/2918
inflating: /root/nltk_data/corpora/reuters/training/2919
inflating: /root/nltk_data/corpora/reuters/training/2920
inflating: /root/nltk_data/corpora/reuters/training/2922
inflating: /root/nltk_data/corpora/reuters/training/2923
inflating: /root/nltk_data/corpora/reuters/training/2924
inflating: /root/nltk_data/corpora/reuters/training/2925
inflating: /root/nltk_data/corpora/reuters/training/2926
inflating: /root/nltk_data/corpora/reuters/training/2927
inflating: /root/nltk_data/corpora/reuters/training/2928
inflating: /root/nltk_data/corpora/reuters/training/2929
inflating: /root/nltk_data/corpora/reuters/training/293
inflating: /root/nltk_data/corpora/reuters/training/2930
inflating: /root/nltk_data/corpora/reuters/training/2933
inflating: /root/nltk_data/corpora/reuters/training/2936
inflating: /root/nltk_data/corpora/reuters/training/2938
inflating: /root/nltk_data/corpora/reuters/training/2939
inflating: /root/nltk_data/corpora/reuters/training/2940
inflating: /root/nltk_data/corpora/reuters/training/2942
inflating: /root/nltk_data/corpora/reuters/training/2944
inflating: /root/nltk_data/corpora/reuters/training/2945
inflating: /root/nltk_data/corpora/reuters/training/2946
inflating: /root/nltk_data/corpora/reuters/training/2947
inflating: /root/nltk_data/corpora/reuters/training/2948
inflating: /root/nltk_data/corpora/reuters/training/2950
inflating: /root/nltk_data/corpora/reuters/training/2951
inflating: /root/nltk_data/corpora/reuters/training/2952
inflating: /root/nltk_data/corpora/reuters/training/2953
inflating: /root/nltk_data/corpora/reuters/training/2954
inflating: /root/nltk_data/corpora/reuters/training/2955
inflating: /root/nltk_data/corpora/reuters/training/2956
inflating: /root/nltk_data/corpora/reuters/training/2957
inflating: /root/nltk_data/corpora/reuters/training/2958
inflating: /root/nltk_data/corpora/reuters/training/2959
inflating: /root/nltk_data/corpora/reuters/training/296
inflating: /root/nltk_data/corpora/reuters/training/2960
inflating: /root/nltk_data/corpora/reuters/training/2961
inflating: /root/nltk_data/corpora/reuters/training/2962
inflating: /root/nltk_data/corpora/reuters/training/2963
inflating: /root/nltk_data/corpora/reuters/training/2964
inflating: /root/nltk_data/corpora/reuters/training/2965
inflating: /root/nltk_data/corpora/reuters/training/2966
inflating: /root/nltk_data/corpora/reuters/training/2968
inflating: /root/nltk_data/corpora/reuters/training/297
inflating: /root/nltk_data/corpora/reuters/training/2970
inflating: /root/nltk_data/corpora/reuters/training/2971
inflating: /root/nltk_data/corpora/reuters/training/2973
inflating: /root/nltk_data/corpora/reuters/training/2974
inflating: /root/nltk_data/corpora/reuters/training/2975
inflating: /root/nltk_data/corpora/reuters/training/2976
inflating: /root/nltk_data/corpora/reuters/training/2977
inflating: /root/nltk_data/corpora/reuters/training/2978
inflating: /root/nltk_data/corpora/reuters/training/2979
inflating: /root/nltk_data/corpora/reuters/training/298
inflating: /root/nltk_data/corpora/reuters/training/2982
inflating: /root/nltk_data/corpora/reuters/training/2983
inflating: /root/nltk_data/corpora/reuters/training/2985
inflating: /root/nltk_data/corpora/reuters/training/299
inflating: /root/nltk_data/corpora/reuters/training/2990
inflating: /root/nltk_data/corpora/reuters/training/2993
inflating: /root/nltk_data/corpora/reuters/training/2994
inflating: /root/nltk_data/corpora/reuters/training/2995
inflating: /root/nltk_data/corpora/reuters/training/2996
inflating: /root/nltk_data/corpora/reuters/training/2998
inflating: /root/nltk_data/corpora/reuters/training/2999
inflating: /root/nltk_data/corpora/reuters/training/30
inflating: /root/nltk_data/corpora/reuters/training/3001
inflating: /root/nltk_data/corpora/reuters/training/3002
inflating: /root/nltk_data/corpora/reuters/training/3003
inflating: /root/nltk_data/corpora/reuters/training/3006
inflating: /root/nltk_data/corpora/reuters/training/3007
inflating: /root/nltk_data/corpora/reuters/training/3008
inflating: /root/nltk_data/corpora/reuters/training/3009
inflating: /root/nltk_data/corpora/reuters/training/3010
inflating: /root/nltk_data/corpora/reuters/training/3013
inflating: /root/nltk_data/corpora/reuters/training/3015
inflating: /root/nltk_data/corpora/reuters/training/3016
inflating: /root/nltk_data/corpora/reuters/training/3017
inflating: /root/nltk_data/corpora/reuters/training/3019
inflating: /root/nltk_data/corpora/reuters/training/302
inflating: /root/nltk_data/corpora/reuters/training/3020
inflating: /root/nltk_data/corpora/reuters/training/3023
inflating: /root/nltk_data/corpora/reuters/training/3024
inflating: /root/nltk_data/corpora/reuters/training/3025
inflating: /root/nltk_data/corpora/reuters/training/3026
inflating: /root/nltk_data/corpora/reuters/training/3028
inflating: /root/nltk_data/corpora/reuters/training/303
inflating: /root/nltk_data/corpora/reuters/training/3031
inflating: /root/nltk_data/corpora/reuters/training/3034
inflating: /root/nltk_data/corpora/reuters/training/3035
inflating: /root/nltk_data/corpora/reuters/training/3036
inflating: /root/nltk_data/corpora/reuters/training/3037
inflating: /root/nltk_data/corpora/reuters/training/3038
inflating: /root/nltk_data/corpora/reuters/training/3039
inflating: /root/nltk_data/corpora/reuters/training/304
inflating: /root/nltk_data/corpora/reuters/training/3040
inflating: /root/nltk_data/corpora/reuters/training/3041
inflating: /root/nltk_data/corpora/reuters/training/3042
inflating: /root/nltk_data/corpora/reuters/training/3043
inflating: /root/nltk_data/corpora/reuters/training/3044
```

```
inflating: /root/nltk_data/corpora/reuters/training/3046
inflating: /root/nltk_data/corpora/reuters/training/3047
inflating: /root/nltk_data/corpora/reuters/training/3048
inflating: /root/nltk_data/corpora/reuters/training/3049
inflating: /root/nltk_data/corpora/reuters/training/3051
inflating: /root/nltk_data/corpora/reuters/training/3053
inflating: /root/nltk_data/corpora/reuters/training/3056
inflating: /root/nltk_data/corpora/reuters/training/3057
inflating: /root/nltk_data/corpora/reuters/training/3060
inflating: /root/nltk_data/corpora/reuters/training/3061
inflating: /root/nltk_data/corpora/reuters/training/3062
inflating: /root/nltk_data/corpora/reuters/training/3063
inflating: /root/nltk_data/corpora/reuters/training/3065
inflating: /root/nltk_data/corpora/reuters/training/3067
inflating: /root/nltk_data/corpora/reuters/training/3068
inflating: /root/nltk_data/corpora/reuters/training/3069
inflating: /root/nltk_data/corpora/reuters/training/307
inflating: /root/nltk_data/corpora/reuters/training/3071
inflating: /root/nltk_data/corpora/reuters/training/3072
inflating: /root/nltk_data/corpora/reuters/training/3073
inflating: /root/nltk_data/corpora/reuters/training/3074
inflating: /root/nltk_data/corpora/reuters/training/3076
inflating: /root/nltk_data/corpora/reuters/training/3077
inflating: /root/nltk_data/corpora/reuters/training/3078
inflating: /root/nltk_data/corpora/reuters/training/3080
inflating: /root/nltk_data/corpora/reuters/training/3082
inflating: /root/nltk_data/corpora/reuters/training/3083
inflating: /root/nltk_data/corpora/reuters/training/3084
inflating: /root/nltk_data/corpora/reuters/training/3085
inflating: /root/nltk_data/corpora/reuters/training/3086
inflating: /root/nltk_data/corpora/reuters/training/3087
inflating: /root/nltk_data/corpora/reuters/training/3088
inflating: /root/nltk_data/corpora/reuters/training/3089
inflating: /root/nltk_data/corpora/reuters/training/309
inflating: /root/nltk_data/corpora/reuters/training/3091
inflating: /root/nltk_data/corpora/reuters/training/3092
inflating: /root/nltk_data/corpora/reuters/training/3093
inflating: /root/nltk_data/corpora/reuters/training/3094
inflating: /root/nltk_data/corpora/reuters/training/3095
inflating: /root/nltk_data/corpora/reuters/training/3096
inflating: /root/nltk_data/corpora/reuters/training/3098
inflating: /root/nltk_data/corpora/reuters/training/310
inflating: /root/nltk_data/corpora/reuters/training/3100
inflating: /root/nltk_data/corpora/reuters/training/3103
inflating: /root/nltk_data/corpora/reuters/training/3104
inflating: /root/nltk_data/corpora/reuters/training/3105
inflating: /root/nltk_data/corpora/reuters/training/3107
inflating: /root/nltk_data/corpora/reuters/training/3109
inflating: /root/nltk_data/corpora/reuters/training/311
inflating: /root/nltk_data/corpora/reuters/training/3110
inflating: /root/nltk_data/corpora/reuters/training/3111
inflating: /root/nltk_data/corpora/reuters/training/3114
inflating: /root/nltk_data/corpora/reuters/training/3115
inflating: /root/nltk_data/corpora/reuters/training/3117
inflating: /root/nltk_data/corpora/reuters/training/3118
inflating: /root/nltk_data/corpora/reuters/training/312
inflating: /root/nltk_data/corpora/reuters/training/3120
inflating: /root/nltk_data/corpora/reuters/training/3121
inflating: /root/nltk_data/corpora/reuters/training/3122
inflating: /root/nltk_data/corpora/reuters/training/3125
inflating: /root/nltk_data/corpora/reuters/training/3128
inflating: /root/nltk_data/corpora/reuters/training/3129
inflating: /root/nltk_data/corpora/reuters/training/313
inflating: /root/nltk_data/corpora/reuters/training/3131
inflating: /root/nltk_data/corpora/reuters/training/3132
inflating: /root/nltk_data/corpora/reuters/training/3133
inflating: /root/nltk_data/corpora/reuters/training/3134
inflating: /root/nltk_data/corpora/reuters/training/3135
inflating: /root/nltk_data/corpora/reuters/training/3137
inflating: /root/nltk_data/corpora/reuters/training/3138
inflating: /root/nltk_data/corpora/reuters/training/3139
inflating: /root/nltk_data/corpora/reuters/training/314
inflating: /root/nltk_data/corpora/reuters/training/3142
inflating: /root/nltk_data/corpora/reuters/training/3145
inflating: /root/nltk_data/corpora/reuters/training/3146
inflating: /root/nltk_data/corpora/reuters/training/3149
inflating: /root/nltk_data/corpora/reuters/training/315
inflating: /root/nltk_data/corpora/reuters/training/3152
inflating: /root/nltk_data/corpora/reuters/training/3153
inflating: /root/nltk_data/corpora/reuters/training/3154
inflating: /root/nltk_data/corpora/reuters/training/3155
inflating: /root/nltk_data/corpora/reuters/training/3157
inflating: /root/nltk_data/corpora/reuters/training/3159
inflating: /root/nltk_data/corpora/reuters/training/3161
inflating: /root/nltk_data/corpora/reuters/training/3162
inflating: /root/nltk_data/corpora/reuters/training/3163
inflating: /root/nltk_data/corpora/reuters/training/3164
inflating: /root/nltk_data/corpora/reuters/training/3166
inflating: /root/nltk_data/corpora/reuters/training/3168
inflating: /root/nltk_data/corpora/reuters/training/3169
inflating: /root/nltk_data/corpora/reuters/training/317
inflating: /root/nltk_data/corpora/reuters/training/3174
inflating: /root/nltk_data/corpora/reuters/training/3176
inflating: /root/nltk_data/corpora/reuters/training/3180
inflating: /root/nltk_data/corpora/reuters/training/3181
inflating: /root/nltk_data/corpora/reuters/training/3183
inflating: /root/nltk_data/corpora/reuters/training/3184
inflating: /root/nltk_data/corpora/reuters/training/3186
inflating: /root/nltk_data/corpora/reuters/training/3187
inflating: /root/nltk_data/corpora/reuters/training/3189
inflating: /root/nltk_data/corpora/reuters/training/3190
inflating: /root/nltk_data/corpora/reuters/training/3191
inflating: /root/nltk_data/corpora/reuters/training/3192
inflating: /root/nltk_data/corpora/reuters/training/3194
inflating: /root/nltk_data/corpora/reuters/training/3198
inflating: /root/nltk_data/corpora/reuters/training/3199
inflating: /root/nltk_data/corpora/reuters/training/320
inflating: /root/nltk_data/corpora/reuters/training/3204
inflating: /root/nltk_data/corpora/reuters/training/3205
inflating: /root/nltk_data/corpora/reuters/training/3206
inflating: /root/nltk_data/corpora/reuters/training/3207
inflating: /root/nltk_data/corpora/reuters/training/3210
inflating: /root/nltk_data/corpora/reuters/training/3211
inflating: /root/nltk_data/corpora/reuters/training/3216
inflating: /root/nltk_data/corpora/reuters/training/3217
inflating: /root/nltk_data/corpora/reuters/training/3219
inflating: /root/nltk_data/corpora/reuters/training/322
inflating: /root/nltk_data/corpora/reuters/training/3222
inflating: /root/nltk_data/corpora/reuters/training/3223
inflating: /root/nltk_data/corpora/reuters/training/3225
inflating: /root/nltk_data/corpora/reuters/training/3228
inflating: /root/nltk_data/corpora/reuters/training/3229
inflating: /root/nltk_data/corpora/reuters/training/323
inflating: /root/nltk_data/corpora/reuters/training/3235
inflating: /root/nltk_data/corpora/reuters/training/3239
inflating: /root/nltk_data/corpora/reuters/training/3241
inflating: /root/nltk_data/corpora/reuters/training/3242
inflating: /root/nltk_data/corpora/reuters/training/3243
inflating: /root/nltk_data/corpora/reuters/training/3246
inflating: /root/nltk_data/corpora/reuters/training/3247
inflating: /root/nltk_data/corpora/reuters/training/3249
inflating: /root/nltk_data/corpora/reuters/training/3251
inflating: /root/nltk_data/corpora/reuters/training/3252
inflating: /root/nltk_data/corpora/reuters/training/3253
```

```
inflating: /root/nltk_data/corpora/reuters/training/3254
inflating: /root/nltk_data/corpora/reuters/training/3256
inflating: /root/nltk_data/corpora/reuters/training/3258
inflating: /root/nltk_data/corpora/reuters/training/3259
inflating: /root/nltk_data/corpora/reuters/training/3261
inflating: /root/nltk_data/corpora/reuters/training/3263
inflating: /root/nltk_data/corpora/reuters/training/3264
inflating: /root/nltk_data/corpora/reuters/training/3266
inflating: /root/nltk_data/corpora/reuters/training/3267
inflating: /root/nltk_data/corpora/reuters/training/3269
inflating: /root/nltk_data/corpora/reuters/training/327
inflating: /root/nltk_data/corpora/reuters/training/3271
inflating: /root/nltk_data/corpora/reuters/training/3272
inflating: /root/nltk_data/corpora/reuters/training/3273
inflating: /root/nltk_data/corpora/reuters/training/3275
inflating: /root/nltk_data/corpora/reuters/training/3276
inflating: /root/nltk_data/corpora/reuters/training/3277
inflating: /root/nltk_data/corpora/reuters/training/3278
inflating: /root/nltk_data/corpora/reuters/training/3282
inflating: /root/nltk_data/corpora/reuters/training/3283
inflating: /root/nltk_data/corpora/reuters/training/3285
inflating: /root/nltk_data/corpora/reuters/training/3287
inflating: /root/nltk_data/corpora/reuters/training/3290
inflating: /root/nltk_data/corpora/reuters/training/3292
inflating: /root/nltk_data/corpora/reuters/training/3295
inflating: /root/nltk_data/corpora/reuters/training/3298
inflating: /root/nltk_data/corpora/reuters/training/3299
inflating: /root/nltk_data/corpora/reuters/training/330
inflating: /root/nltk_data/corpora/reuters/training/3302
inflating: /root/nltk_data/corpora/reuters/training/3303
inflating: /root/nltk_data/corpora/reuters/training/3305
inflating: /root/nltk_data/corpora/reuters/training/3306
inflating: /root/nltk_data/corpora/reuters/training/3309
inflating: /root/nltk_data/corpora/reuters/training/331
inflating: /root/nltk_data/corpora/reuters/training/3310
inflating: /root/nltk_data/corpora/reuters/training/3311
inflating: /root/nltk_data/corpora/reuters/training/3312
inflating: /root/nltk_data/corpora/reuters/training/3313
inflating: /root/nltk_data/corpora/reuters/training/3314
inflating: /root/nltk_data/corpora/reuters/training/3315
inflating: /root/nltk_data/corpora/reuters/training/3317
inflating: /root/nltk_data/corpora/reuters/training/3318
inflating: /root/nltk_data/corpora/reuters/training/332
inflating: /root/nltk_data/corpora/reuters/training/3320
inflating: /root/nltk_data/corpora/reuters/training/3322
inflating: /root/nltk_data/corpora/reuters/training/3323
inflating: /root/nltk_data/corpora/reuters/training/3324
inflating: /root/nltk_data/corpora/reuters/training/3327
inflating: /root/nltk_data/corpora/reuters/training/3329
inflating: /root/nltk_data/corpora/reuters/training/3330
inflating: /root/nltk_data/corpora/reuters/training/3332
inflating: /root/nltk_data/corpora/reuters/training/3333
inflating: /root/nltk_data/corpora/reuters/training/3334
inflating: /root/nltk_data/corpora/reuters/training/3335
inflating: /root/nltk_data/corpora/reuters/training/3336
inflating: /root/nltk_data/corpora/reuters/training/3337
inflating: /root/nltk_data/corpora/reuters/training/3338
inflating: /root/nltk_data/corpora/reuters/training/3339
inflating: /root/nltk_data/corpora/reuters/training/334
inflating: /root/nltk_data/corpora/reuters/training/3340
inflating: /root/nltk_data/corpora/reuters/training/3341
inflating: /root/nltk_data/corpora/reuters/training/3342
inflating: /root/nltk_data/corpora/reuters/training/3343
inflating: /root/nltk_data/corpora/reuters/training/3344
inflating: /root/nltk_data/corpora/reuters/training/3345
inflating: /root/nltk_data/corpora/reuters/training/3347
inflating: /root/nltk_data/corpora/reuters/training/3348
inflating: /root/nltk_data/corpora/reuters/training/3349
inflating: /root/nltk_data/corpora/reuters/training/3351
inflating: /root/nltk_data/corpora/reuters/training/3352
inflating: /root/nltk_data/corpora/reuters/training/3353
inflating: /root/nltk_data/corpora/reuters/training/3354
inflating: /root/nltk_data/corpora/reuters/training/3356
inflating: /root/nltk_data/corpora/reuters/training/3358
inflating: /root/nltk_data/corpora/reuters/training/3359
inflating: /root/nltk_data/corpora/reuters/training/336
inflating: /root/nltk_data/corpora/reuters/training/3360
inflating: /root/nltk_data/corpora/reuters/training/3363
inflating: /root/nltk_data/corpora/reuters/training/3364
inflating: /root/nltk_data/corpora/reuters/training/3365
inflating: /root/nltk_data/corpora/reuters/training/3367
inflating: /root/nltk_data/corpora/reuters/training/3370
inflating: /root/nltk_data/corpora/reuters/training/3371
inflating: /root/nltk_data/corpora/reuters/training/3372
inflating: /root/nltk_data/corpora/reuters/training/3377
inflating: /root/nltk_data/corpora/reuters/training/3378
inflating: /root/nltk_data/corpora/reuters/training/338
inflating: /root/nltk_data/corpora/reuters/training/3380
inflating: /root/nltk_data/corpora/reuters/training/3381
inflating: /root/nltk_data/corpora/reuters/training/3382
inflating: /root/nltk_data/corpora/reuters/training/3384
inflating: /root/nltk_data/corpora/reuters/training/3386
inflating: /root/nltk_data/corpora/reuters/training/3387
inflating: /root/nltk_data/corpora/reuters/training/3389
inflating: /root/nltk_data/corpora/reuters/training/3390
inflating: /root/nltk_data/corpora/reuters/training/3394
inflating: /root/nltk_data/corpora/reuters/training/3395
inflating: /root/nltk_data/corpora/reuters/training/3396
inflating: /root/nltk_data/corpora/reuters/training/3399
inflating: /root/nltk_data/corpora/reuters/training/3400
inflating: /root/nltk_data/corpora/reuters/training/3401
inflating: /root/nltk_data/corpora/reuters/training/3405
inflating: /root/nltk_data/corpora/reuters/training/3406
inflating: /root/nltk_data/corpora/reuters/training/3409
inflating: /root/nltk_data/corpora/reuters/training/341
inflating: /root/nltk_data/corpora/reuters/training/3411
inflating: /root/nltk_data/corpora/reuters/training/3413
inflating: /root/nltk_data/corpora/reuters/training/3415
inflating: /root/nltk_data/corpora/reuters/training/3416
inflating: /root/nltk_data/corpora/reuters/training/3417
inflating: /root/nltk_data/corpora/reuters/training/3418
inflating: /root/nltk_data/corpora/reuters/training/3419
inflating: /root/nltk_data/corpora/reuters/training/342
inflating: /root/nltk_data/corpora/reuters/training/3420
inflating: /root/nltk_data/corpora/reuters/training/3421
inflating: /root/nltk_data/corpora/reuters/training/3422
inflating: /root/nltk_data/corpora/reuters/training/3425
inflating: /root/nltk_data/corpora/reuters/training/3429
inflating: /root/nltk_data/corpora/reuters/training/343
inflating: /root/nltk_data/corpora/reuters/training/3430
inflating: /root/nltk_data/corpora/reuters/training/3433
inflating: /root/nltk_data/corpora/reuters/training/3435
inflating: /root/nltk_data/corpora/reuters/training/3438
inflating: /root/nltk_data/corpora/reuters/training/3440
inflating: /root/nltk_data/corpora/reuters/training/3441
inflating: /root/nltk_data/corpora/reuters/training/3442
inflating: /root/nltk_data/corpora/reuters/training/3445
inflating: /root/nltk_data/corpora/reuters/training/3446
inflating: /root/nltk_data/corpora/reuters/training/3448
inflating: /root/nltk_data/corpora/reuters/training/3449
inflating: /root/nltk_data/corpora/reuters/training/345
inflating: /root/nltk_data/corpora/reuters/training/3452
inflating: /root/nltk_data/corpora/reuters/training/3453
inflating: /root/nltk_data/corpora/reuters/training/3454
inflating: /root/nltk_data/corpora/reuters/training/3455
```

```
inflating: /root/nltk_data/corpora/reuters/training/3455
inflating: /root/nltk_data/corpora/reuters/training/3456
inflating: /root/nltk_data/corpora/reuters/training/3458
inflating: /root/nltk_data/corpora/reuters/training/3459
inflating: /root/nltk_data/corpora/reuters/training/346
inflating: /root/nltk_data/corpora/reuters/training/3460
inflating: /root/nltk_data/corpora/reuters/training/3461
inflating: /root/nltk_data/corpora/reuters/training/3463
inflating: /root/nltk_data/corpora/reuters/training/3466
inflating: /root/nltk_data/corpora/reuters/training/3467
inflating: /root/nltk_data/corpora/reuters/training/3468
inflating: /root/nltk_data/corpora/reuters/training/3469
inflating: /root/nltk_data/corpora/reuters/training/3472
inflating: /root/nltk_data/corpora/reuters/training/3473
inflating: /root/nltk_data/corpora/reuters/training/3474
inflating: /root/nltk_data/corpora/reuters/training/3476
inflating: /root/nltk_data/corpora/reuters/training/3477
inflating: /root/nltk_data/corpora/reuters/training/3480
inflating: /root/nltk_data/corpora/reuters/training/3481
inflating: /root/nltk_data/corpora/reuters/training/3482
inflating: /root/nltk_data/corpora/reuters/training/3483
inflating: /root/nltk_data/corpora/reuters/training/3486
inflating: /root/nltk_data/corpora/reuters/training/3488
inflating: /root/nltk_data/corpora/reuters/training/349
inflating: /root/nltk_data/corpora/reuters/training/3490
inflating: /root/nltk_data/corpora/reuters/training/3491
inflating: /root/nltk_data/corpora/reuters/training/3492
inflating: /root/nltk_data/corpora/reuters/training/3493
inflating: /root/nltk_data/corpora/reuters/training/3497
inflating: /root/nltk_data/corpora/reuters/training/3499
inflating: /root/nltk_data/corpora/reuters/training/3504
inflating: /root/nltk_data/corpora/reuters/training/3505
inflating: /root/nltk_data/corpora/reuters/training/3507
inflating: /root/nltk_data/corpora/reuters/training/3509
inflating: /root/nltk_data/corpora/reuters/training/3511
inflating: /root/nltk_data/corpora/reuters/training/3512
inflating: /root/nltk_data/corpora/reuters/training/3514
inflating: /root/nltk_data/corpora/reuters/training/3517
inflating: /root/nltk_data/corpora/reuters/training/352
inflating: /root/nltk_data/corpora/reuters/training/3520
inflating: /root/nltk_data/corpora/reuters/training/3524
inflating: /root/nltk_data/corpora/reuters/training/3526
inflating: /root/nltk_data/corpora/reuters/training/3528
inflating: /root/nltk_data/corpora/reuters/training/353
inflating: /root/nltk_data/corpora/reuters/training/3531
inflating: /root/nltk_data/corpora/reuters/training/3532
inflating: /root/nltk_data/corpora/reuters/training/3533
inflating: /root/nltk_data/corpora/reuters/training/3534
inflating: /root/nltk_data/corpora/reuters/training/3535
inflating: /root/nltk_data/corpora/reuters/training/3538
inflating: /root/nltk_data/corpora/reuters/training/3539
inflating: /root/nltk_data/corpora/reuters/training/354
inflating: /root/nltk_data/corpora/reuters/training/3540
inflating: /root/nltk_data/corpora/reuters/training/3541
inflating: /root/nltk_data/corpora/reuters/training/3542
inflating: /root/nltk_data/corpora/reuters/training/3543
inflating: /root/nltk_data/corpora/reuters/training/3545
inflating: /root/nltk_data/corpora/reuters/training/3547
inflating: /root/nltk_data/corpora/reuters/training/3553
inflating: /root/nltk_data/corpora/reuters/training/3554
inflating: /root/nltk_data/corpora/reuters/training/3556
inflating: /root/nltk_data/corpora/reuters/training/3557
inflating: /root/nltk_data/corpora/reuters/training/3559
inflating: /root/nltk_data/corpora/reuters/training/356
inflating: /root/nltk_data/corpora/reuters/training/3563
inflating: /root/nltk_data/corpora/reuters/training/3564
inflating: /root/nltk_data/corpora/reuters/training/3565
inflating: /root/nltk_data/corpora/reuters/training/3569
inflating: /root/nltk_data/corpora/reuters/training/3570
inflating: /root/nltk_data/corpora/reuters/training/3571
inflating: /root/nltk_data/corpora/reuters/training/3572
inflating: /root/nltk_data/corpora/reuters/training/3573
inflating: /root/nltk_data/corpora/reuters/training/3574
inflating: /root/nltk_data/corpora/reuters/training/3575
inflating: /root/nltk_data/corpora/reuters/training/3577
inflating: /root/nltk_data/corpora/reuters/training/3578
inflating: /root/nltk_data/corpora/reuters/training/3579
inflating: /root/nltk_data/corpora/reuters/training/3580
inflating: /root/nltk_data/corpora/reuters/training/3582
inflating: /root/nltk_data/corpora/reuters/training/3584
inflating: /root/nltk_data/corpora/reuters/training/3585
inflating: /root/nltk_data/corpora/reuters/training/3586
inflating: /root/nltk_data/corpora/reuters/training/3587
inflating: /root/nltk_data/corpora/reuters/training/3589
inflating: /root/nltk_data/corpora/reuters/training/3591
inflating: /root/nltk_data/corpora/reuters/training/3592
inflating: /root/nltk_data/corpora/reuters/training/3593
inflating: /root/nltk_data/corpora/reuters/training/3594
inflating: /root/nltk_data/corpora/reuters/training/3595
inflating: /root/nltk_data/corpora/reuters/training/3596
inflating: /root/nltk_data/corpora/reuters/training/3597
inflating: /root/nltk_data/corpora/reuters/training/3598
inflating: /root/nltk_data/corpora/reuters/training/3599
inflating: /root/nltk_data/corpora/reuters/training/36
inflating: /root/nltk_data/corpora/reuters/training/3601
inflating: /root/nltk_data/corpora/reuters/training/3602
inflating: /root/nltk_data/corpora/reuters/training/3603
inflating: /root/nltk_data/corpora/reuters/training/3604
inflating: /root/nltk_data/corpora/reuters/training/3605
inflating: /root/nltk_data/corpora/reuters/training/3607
inflating: /root/nltk_data/corpora/reuters/training/3609
inflating: /root/nltk_data/corpora/reuters/training/361
inflating: /root/nltk_data/corpora/reuters/training/3610
inflating: /root/nltk_data/corpora/reuters/training/3611
inflating: /root/nltk_data/corpora/reuters/training/3612
inflating: /root/nltk_data/corpora/reuters/training/3613
inflating: /root/nltk_data/corpora/reuters/training/3615
inflating: /root/nltk_data/corpora/reuters/training/3616
inflating: /root/nltk_data/corpora/reuters/training/3619
inflating: /root/nltk_data/corpora/reuters/training/362
inflating: /root/nltk_data/corpora/reuters/training/3620
inflating: /root/nltk_data/corpora/reuters/training/3621
inflating: /root/nltk_data/corpora/reuters/training/3622
inflating: /root/nltk_data/corpora/reuters/training/3623
inflating: /root/nltk_data/corpora/reuters/training/3624
inflating: /root/nltk_data/corpora/reuters/training/3625
inflating: /root/nltk_data/corpora/reuters/training/3626
inflating: /root/nltk_data/corpora/reuters/training/3629
inflating: /root/nltk_data/corpora/reuters/training/3630
inflating: /root/nltk_data/corpora/reuters/training/3631
inflating: /root/nltk_data/corpora/reuters/training/3632
inflating: /root/nltk_data/corpora/reuters/training/3633
inflating: /root/nltk_data/corpora/reuters/training/3634
inflating: /root/nltk_data/corpora/reuters/training/3637
inflating: /root/nltk_data/corpora/reuters/training/3639
inflating: /root/nltk_data/corpora/reuters/training/3640
inflating: /root/nltk_data/corpora/reuters/training/3643
inflating: /root/nltk_data/corpora/reuters/training/3645
inflating: /root/nltk_data/corpora/reuters/training/3646
inflating: /root/nltk_data/corpora/reuters/training/3647
inflating: /root/nltk_data/corpora/reuters/training/3649
inflating: /root/nltk_data/corpora/reuters/training/3651
inflating: /root/nltk_data/corpora/reuters/training/3652
inflating: /root/nltk_data/corpora/reuters/training/366
inflating: /root/nltk_data/corpora/reuters/training/3660
```

```
inflating: /root/nltk_data/corpora/reuters/training/3661
inflating: /root/nltk_data/corpora/reuters/training/3666
inflating: /root/nltk_data/corpora/reuters/training/3667
inflating: /root/nltk_data/corpora/reuters/training/3668
inflating: /root/nltk_data/corpora/reuters/training/3669
inflating: /root/nltk_data/corpora/reuters/training/3672
inflating: /root/nltk_data/corpora/reuters/training/3673
inflating: /root/nltk_data/corpora/reuters/training/3674
inflating: /root/nltk_data/corpora/reuters/training/3676
inflating: /root/nltk_data/corpora/reuters/training/368
inflating: /root/nltk_data/corpora/reuters/training/3680
inflating: /root/nltk_data/corpora/reuters/training/3681
inflating: /root/nltk_data/corpora/reuters/training/3684
inflating: /root/nltk_data/corpora/reuters/training/3688
inflating: /root/nltk_data/corpora/reuters/training/369
inflating: /root/nltk_data/corpora/reuters/training/3690
inflating: /root/nltk_data/corpora/reuters/training/3694
inflating: /root/nltk_data/corpora/reuters/training/3695
inflating: /root/nltk_data/corpora/reuters/training/3697
inflating: /root/nltk_data/corpora/reuters/training/3699
inflating: /root/nltk_data/corpora/reuters/training/37
inflating: /root/nltk_data/corpora/reuters/training/3701
inflating: /root/nltk_data/corpora/reuters/training/3702
inflating: /root/nltk_data/corpora/reuters/training/3704
inflating: /root/nltk_data/corpora/reuters/training/3707
inflating: /root/nltk_data/corpora/reuters/training/3708
inflating: /root/nltk_data/corpora/reuters/training/371
inflating: /root/nltk_data/corpora/reuters/training/3711
inflating: /root/nltk_data/corpora/reuters/training/3713
inflating: /root/nltk_data/corpora/reuters/training/3717
inflating: /root/nltk_data/corpora/reuters/training/372
inflating: /root/nltk_data/corpora/reuters/training/3720
inflating: /root/nltk_data/corpora/reuters/training/3721
inflating: /root/nltk_data/corpora/reuters/training/3723
inflating: /root/nltk_data/corpora/reuters/training/3729
inflating: /root/nltk_data/corpora/reuters/training/3730
inflating: /root/nltk_data/corpora/reuters/training/3734
inflating: /root/nltk_data/corpora/reuters/training/3735
inflating: /root/nltk_data/corpora/reuters/training/3737
inflating: /root/nltk_data/corpora/reuters/training/3738
inflating: /root/nltk_data/corpora/reuters/training/374
inflating: /root/nltk_data/corpora/reuters/training/3740
inflating: /root/nltk_data/corpora/reuters/training/3741
inflating: /root/nltk_data/corpora/reuters/training/3742
inflating: /root/nltk_data/corpora/reuters/training/3743
inflating: /root/nltk_data/corpora/reuters/training/3744
inflating: /root/nltk_data/corpora/reuters/training/3747
inflating: /root/nltk_data/corpora/reuters/training/3748
inflating: /root/nltk_data/corpora/reuters/training/3750
inflating: /root/nltk_data/corpora/reuters/training/3751
inflating: /root/nltk_data/corpora/reuters/training/3755
inflating: /root/nltk_data/corpora/reuters/training/3757
inflating: /root/nltk_data/corpora/reuters/training/3758
inflating: /root/nltk_data/corpora/reuters/training/3759
inflating: /root/nltk_data/corpora/reuters/training/376
inflating: /root/nltk_data/corpora/reuters/training/3760
inflating: /root/nltk_data/corpora/reuters/training/3762
inflating: /root/nltk_data/corpora/reuters/training/3763
inflating: /root/nltk_data/corpora/reuters/training/3764
inflating: /root/nltk_data/corpora/reuters/training/3765
inflating: /root/nltk_data/corpora/reuters/training/3766
inflating: /root/nltk_data/corpora/reuters/training/3768
inflating: /root/nltk_data/corpora/reuters/training/377
inflating: /root/nltk_data/corpora/reuters/training/3770
inflating: /root/nltk_data/corpora/reuters/training/3771
inflating: /root/nltk_data/corpora/reuters/training/3772
inflating: /root/nltk_data/corpora/reuters/training/3774
inflating: /root/nltk_data/corpora/reuters/training/3775
inflating: /root/nltk_data/corpora/reuters/training/3777
inflating: /root/nltk_data/corpora/reuters/training/3779
inflating: /root/nltk_data/corpora/reuters/training/378
inflating: /root/nltk_data/corpora/reuters/training/3782
inflating: /root/nltk_data/corpora/reuters/training/3784
inflating: /root/nltk_data/corpora/reuters/training/3785
inflating: /root/nltk_data/corpora/reuters/training/3786
inflating: /root/nltk_data/corpora/reuters/training/3787
inflating: /root/nltk_data/corpora/reuters/training/3788
inflating: /root/nltk_data/corpora/reuters/training/379
inflating: /root/nltk_data/corpora/reuters/training/3790
inflating: /root/nltk_data/corpora/reuters/training/3792
inflating: /root/nltk_data/corpora/reuters/training/3793
inflating: /root/nltk_data/corpora/reuters/training/3795
inflating: /root/nltk_data/corpora/reuters/training/3796
inflating: /root/nltk_data/corpora/reuters/training/3797
inflating: /root/nltk_data/corpora/reuters/training/3798
inflating: /root/nltk_data/corpora/reuters/training/38
inflating: /root/nltk_data/corpora/reuters/training/3800
inflating: /root/nltk_data/corpora/reuters/training/3801
inflating: /root/nltk_data/corpora/reuters/training/3806
inflating: /root/nltk_data/corpora/reuters/training/3807
inflating: /root/nltk_data/corpora/reuters/training/3808
inflating: /root/nltk_data/corpora/reuters/training/381
inflating: /root/nltk_data/corpora/reuters/training/3817
inflating: /root/nltk_data/corpora/reuters/training/3821
inflating: /root/nltk_data/corpora/reuters/training/3825
inflating: /root/nltk_data/corpora/reuters/training/3828
inflating: /root/nltk_data/corpora/reuters/training/383
inflating: /root/nltk_data/corpora/reuters/training/3830
inflating: /root/nltk_data/corpora/reuters/training/3831
inflating: /root/nltk_data/corpora/reuters/training/3833
inflating: /root/nltk_data/corpora/reuters/training/3838
inflating: /root/nltk_data/corpora/reuters/training/3839
inflating: /root/nltk_data/corpora/reuters/training/3840
inflating: /root/nltk_data/corpora/reuters/training/3841
inflating: /root/nltk_data/corpora/reuters/training/3843
inflating: /root/nltk_data/corpora/reuters/training/3846
inflating: /root/nltk_data/corpora/reuters/training/3847
inflating: /root/nltk_data/corpora/reuters/training/3850
inflating: /root/nltk_data/corpora/reuters/training/3851
inflating: /root/nltk_data/corpora/reuters/training/3853
inflating: /root/nltk_data/corpora/reuters/training/3855
inflating: /root/nltk_data/corpora/reuters/training/3856
inflating: /root/nltk_data/corpora/reuters/training/3857
inflating: /root/nltk_data/corpora/reuters/training/3858
inflating: /root/nltk_data/corpora/reuters/training/3860
inflating: /root/nltk_data/corpora/reuters/training/3862
inflating: /root/nltk_data/corpora/reuters/training/3864
inflating: /root/nltk_data/corpora/reuters/training/3865
inflating: /root/nltk_data/corpora/reuters/training/3867
inflating: /root/nltk_data/corpora/reuters/training/3868
inflating: /root/nltk_data/corpora/reuters/training/3869
inflating: /root/nltk_data/corpora/reuters/training/387
inflating: /root/nltk_data/corpora/reuters/training/3870
inflating: /root/nltk_data/corpora/reuters/training/3872
inflating: /root/nltk_data/corpora/reuters/training/3873
inflating: /root/nltk_data/corpora/reuters/training/3874
inflating: /root/nltk_data/corpora/reuters/training/3875
inflating: /root/nltk_data/corpora/reuters/training/3879
inflating: /root/nltk_data/corpora/reuters/training/388
inflating: /root/nltk_data/corpora/reuters/training/3881
inflating: /root/nltk_data/corpora/reuters/training/3884
inflating: /root/nltk_data/corpora/reuters/training/3886
inflating: /root/nltk_data/corpora/reuters/training/3888
inflating: /root/nltk_data/corpora/reuters/training/389
```

```
inflating: /root/nltk_data/corpora/reuters/training/3890
inflating: /root/nltk_data/corpora/reuters/training/3892
inflating: /root/nltk_data/corpora/reuters/training/3893
inflating: /root/nltk_data/corpora/reuters/training/3894
inflating: /root/nltk_data/corpora/reuters/training/3895
inflating: /root/nltk_data/corpora/reuters/training/3897
inflating: /root/nltk_data/corpora/reuters/training/3898
inflating: /root/nltk_data/corpora/reuters/training/3899
inflating: /root/nltk_data/corpora/reuters/training/390
inflating: /root/nltk_data/corpora/reuters/training/3900
inflating: /root/nltk_data/corpora/reuters/training/3901
inflating: /root/nltk_data/corpora/reuters/training/3902
inflating: /root/nltk_data/corpora/reuters/training/3903
inflating: /root/nltk_data/corpora/reuters/training/3904
inflating: /root/nltk_data/corpora/reuters/training/3905
inflating: /root/nltk_data/corpora/reuters/training/3906
inflating: /root/nltk_data/corpora/reuters/training/3908
inflating: /root/nltk_data/corpora/reuters/training/3909
inflating: /root/nltk_data/corpora/reuters/training/3910
inflating: /root/nltk_data/corpora/reuters/training/3912
inflating: /root/nltk_data/corpora/reuters/training/3913
inflating: /root/nltk_data/corpora/reuters/training/3916
inflating: /root/nltk_data/corpora/reuters/training/3917
inflating: /root/nltk_data/corpora/reuters/training/3918
inflating: /root/nltk_data/corpora/reuters/training/392
inflating: /root/nltk_data/corpora/reuters/training/3920
inflating: /root/nltk_data/corpora/reuters/training/3922
inflating: /root/nltk_data/corpora/reuters/training/3926
inflating: /root/nltk_data/corpora/reuters/training/3928
inflating: /root/nltk_data/corpora/reuters/training/393
inflating: /root/nltk_data/corpora/reuters/training/3930
inflating: /root/nltk_data/corpora/reuters/training/3931
inflating: /root/nltk_data/corpora/reuters/training/3933
inflating: /root/nltk_data/corpora/reuters/training/3934
inflating: /root/nltk_data/corpora/reuters/training/3937
inflating: /root/nltk_data/corpora/reuters/training/3939
inflating: /root/nltk_data/corpora/reuters/training/3942
inflating: /root/nltk_data/corpora/reuters/training/3944
inflating: /root/nltk_data/corpora/reuters/training/3946
inflating: /root/nltk_data/corpora/reuters/training/3948
inflating: /root/nltk_data/corpora/reuters/training/3949
inflating: /root/nltk_data/corpora/reuters/training/395
inflating: /root/nltk_data/corpora/reuters/training/3950
inflating: /root/nltk_data/corpora/reuters/training/3951
inflating: /root/nltk_data/corpora/reuters/training/3952
inflating: /root/nltk_data/corpora/reuters/training/3954
inflating: /root/nltk_data/corpora/reuters/training/3955
inflating: /root/nltk_data/corpora/reuters/training/3957
inflating: /root/nltk_data/corpora/reuters/training/3958
inflating: /root/nltk_data/corpora/reuters/training/3959
inflating: /root/nltk_data/corpora/reuters/training/396
inflating: /root/nltk_data/corpora/reuters/training/3960
inflating: /root/nltk_data/corpora/reuters/training/3961
inflating: /root/nltk_data/corpora/reuters/training/3964
inflating: /root/nltk_data/corpora/reuters/training/3968
inflating: /root/nltk_data/corpora/reuters/training/3970
inflating: /root/nltk_data/corpora/reuters/training/3971
inflating: /root/nltk_data/corpora/reuters/training/3973
inflating: /root/nltk_data/corpora/reuters/training/3976
inflating: /root/nltk_data/corpora/reuters/training/3977
inflating: /root/nltk_data/corpora/reuters/training/3979
inflating: /root/nltk_data/corpora/reuters/training/398
inflating: /root/nltk_data/corpora/reuters/training/3980
inflating: /root/nltk_data/corpora/reuters/training/3981
inflating: /root/nltk_data/corpora/reuters/training/3982
inflating: /root/nltk_data/corpora/reuters/training/3983
inflating: /root/nltk_data/corpora/reuters/training/3984
inflating: /root/nltk_data/corpora/reuters/training/3985
inflating: /root/nltk_data/corpora/reuters/training/3989
inflating: /root/nltk_data/corpora/reuters/training/399
inflating: /root/nltk_data/corpora/reuters/training/3990
inflating: /root/nltk_data/corpora/reuters/training/3991
inflating: /root/nltk_data/corpora/reuters/training/3993
inflating: /root/nltk_data/corpora/reuters/training/3994
inflating: /root/nltk_data/corpora/reuters/training/3995
inflating: /root/nltk_data/corpora/reuters/training/3996
inflating: /root/nltk_data/corpora/reuters/training/3997
inflating: /root/nltk_data/corpora/reuters/training/3999
inflating: /root/nltk_data/corpora/reuters/training/40
inflating: /root/nltk_data/corpora/reuters/training/400
inflating: /root/nltk_data/corpora/reuters/training/4000
inflating: /root/nltk_data/corpora/reuters/training/4005
inflating: /root/nltk_data/corpora/reuters/training/401
inflating: /root/nltk_data/corpora/reuters/training/4012
inflating: /root/nltk_data/corpora/reuters/training/4014
inflating: /root/nltk_data/corpora/reuters/training/4015
inflating: /root/nltk_data/corpora/reuters/training/4016
inflating: /root/nltk_data/corpora/reuters/training/4019
inflating: /root/nltk_data/corpora/reuters/training/402
inflating: /root/nltk_data/corpora/reuters/training/4021
inflating: /root/nltk_data/corpora/reuters/training/4022
inflating: /root/nltk_data/corpora/reuters/training/4023
inflating: /root/nltk_data/corpora/reuters/training/4025
inflating: /root/nltk_data/corpora/reuters/training/4026
inflating: /root/nltk_data/corpora/reuters/training/4027
inflating: /root/nltk_data/corpora/reuters/training/4028
inflating: /root/nltk_data/corpora/reuters/training/4029
inflating: /root/nltk_data/corpora/reuters/training/4031
inflating: /root/nltk_data/corpora/reuters/training/4032
inflating: /root/nltk_data/corpora/reuters/training/4033
inflating: /root/nltk_data/corpora/reuters/training/4035
inflating: /root/nltk_data/corpora/reuters/training/4036
inflating: /root/nltk_data/corpora/reuters/training/4038
inflating: /root/nltk_data/corpora/reuters/training/4039
inflating: /root/nltk_data/corpora/reuters/training/404
inflating: /root/nltk_data/corpora/reuters/training/4040
inflating: /root/nltk_data/corpora/reuters/training/4041
inflating: /root/nltk_data/corpora/reuters/training/4043
inflating: /root/nltk_data/corpora/reuters/training/4045
inflating: /root/nltk_data/corpora/reuters/training/4047
inflating: /root/nltk_data/corpora/reuters/training/4048
inflating: /root/nltk_data/corpora/reuters/training/4049
inflating: /root/nltk_data/corpora/reuters/training/405
inflating: /root/nltk_data/corpora/reuters/training/4051
inflating: /root/nltk_data/corpora/reuters/training/4052
inflating: /root/nltk_data/corpora/reuters/training/4054
inflating: /root/nltk_data/corpora/reuters/training/4055
inflating: /root/nltk_data/corpora/reuters/training/4056
inflating: /root/nltk_data/corpora/reuters/training/4057
inflating: /root/nltk_data/corpora/reuters/training/4058
inflating: /root/nltk_data/corpora/reuters/training/406
inflating: /root/nltk_data/corpora/reuters/training/4061
inflating: /root/nltk_data/corpora/reuters/training/4062
inflating: /root/nltk_data/corpora/reuters/training/4063
inflating: /root/nltk_data/corpora/reuters/training/4064
inflating: /root/nltk_data/corpora/reuters/training/4065
inflating: /root/nltk_data/corpora/reuters/training/4066
inflating: /root/nltk_data/corpora/reuters/training/4067
inflating: /root/nltk_data/corpora/reuters/training/4069
inflating: /root/nltk_data/corpora/reuters/training/407
inflating: /root/nltk_data/corpora/reuters/training/4071
inflating: /root/nltk_data/corpora/reuters/training/4074
inflating: /root/nltk_data/corpora/reuters/training/4075
inflating: /root/nltk_data/corpora/reuters/training/4076
inflating: /root/nltk_data/corpora/reuters/training/4077
```

```
inflating: /root/nltk_data/corpora/reuters/training/4078
inflating: /root/nltk_data/corpora/reuters/training/408
inflating: /root/nltk_data/corpora/reuters/training/4080
inflating: /root/nltk_data/corpora/reuters/training/4081
inflating: /root/nltk_data/corpora/reuters/training/4083
inflating: /root/nltk_data/corpora/reuters/training/4084
inflating: /root/nltk_data/corpora/reuters/training/4085
inflating: /root/nltk_data/corpora/reuters/training/4087
inflating: /root/nltk_data/corpora/reuters/training/4088
inflating: /root/nltk_data/corpora/reuters/training/4089
inflating: /root/nltk_data/corpora/reuters/training/409
inflating: /root/nltk_data/corpora/reuters/training/4090
inflating: /root/nltk_data/corpora/reuters/training/4091
inflating: /root/nltk_data/corpora/reuters/training/4092
inflating: /root/nltk_data/corpora/reuters/training/4093
inflating: /root/nltk_data/corpora/reuters/training/4094
inflating: /root/nltk_data/corpora/reuters/training/4096
inflating: /root/nltk_data/corpora/reuters/training/4097
inflating: /root/nltk_data/corpora/reuters/training/4098
inflating: /root/nltk_data/corpora/reuters/training/4099
inflating: /root/nltk_data/corpora/reuters/training/41
inflating: /root/nltk_data/corpora/reuters/training/410
inflating: /root/nltk_data/corpora/reuters/training/4100
inflating: /root/nltk_data/corpora/reuters/training/4101
inflating: /root/nltk_data/corpora/reuters/training/4103
inflating: /root/nltk_data/corpora/reuters/training/4105
inflating: /root/nltk_data/corpora/reuters/training/4106
inflating: /root/nltk_data/corpora/reuters/training/4107
inflating: /root/nltk_data/corpora/reuters/training/411
inflating: /root/nltk_data/corpora/reuters/training/4111
inflating: /root/nltk_data/corpora/reuters/training/4113
inflating: /root/nltk_data/corpora/reuters/training/4114
inflating: /root/nltk_data/corpora/reuters/training/4115
inflating: /root/nltk_data/corpora/reuters/training/4117
inflating: /root/nltk_data/corpora/reuters/training/4120
inflating: /root/nltk_data/corpora/reuters/training/4122
inflating: /root/nltk_data/corpora/reuters/training/4123
inflating: /root/nltk_data/corpora/reuters/training/4125
inflating: /root/nltk_data/corpora/reuters/training/4127
inflating: /root/nltk_data/corpora/reuters/training/4129
inflating: /root/nltk_data/corpora/reuters/training/4131
inflating: /root/nltk_data/corpora/reuters/training/4132
inflating: /root/nltk_data/corpora/reuters/training/4133
inflating: /root/nltk_data/corpora/reuters/training/4135
inflating: /root/nltk_data/corpora/reuters/training/4138
inflating: /root/nltk_data/corpora/reuters/training/4139
inflating: /root/nltk_data/corpora/reuters/training/4141
inflating: /root/nltk_data/corpora/reuters/training/4142
inflating: /root/nltk_data/corpora/reuters/training/4143
inflating: /root/nltk_data/corpora/reuters/training/4144
inflating: /root/nltk_data/corpora/reuters/training/4145
inflating: /root/nltk_data/corpora/reuters/training/4146
inflating: /root/nltk_data/corpora/reuters/training/4147
inflating: /root/nltk_data/corpora/reuters/training/4148
inflating: /root/nltk_data/corpora/reuters/training/4149
inflating: /root/nltk_data/corpora/reuters/training/4150
inflating: /root/nltk_data/corpora/reuters/training/4152
inflating: /root/nltk_data/corpora/reuters/training/4153
inflating: /root/nltk_data/corpora/reuters/training/4155
inflating: /root/nltk_data/corpora/reuters/training/4156
inflating: /root/nltk_data/corpora/reuters/training/4157
inflating: /root/nltk_data/corpora/reuters/training/4158
inflating: /root/nltk_data/corpora/reuters/training/4159
inflating: /root/nltk_data/corpora/reuters/training/416
inflating: /root/nltk_data/corpora/reuters/training/4160
inflating: /root/nltk_data/corpora/reuters/training/4161
inflating: /root/nltk_data/corpora/reuters/training/4162
inflating: /root/nltk_data/corpora/reuters/training/4164
inflating: /root/nltk_data/corpora/reuters/training/4165
inflating: /root/nltk_data/corpora/reuters/training/4167
inflating: /root/nltk_data/corpora/reuters/training/4169
inflating: /root/nltk_data/corpora/reuters/training/417
inflating: /root/nltk_data/corpora/reuters/training/4170
inflating: /root/nltk_data/corpora/reuters/training/4171
inflating: /root/nltk_data/corpora/reuters/training/4172
inflating: /root/nltk_data/corpora/reuters/training/4173
inflating: /root/nltk_data/corpora/reuters/training/4174
inflating: /root/nltk_data/corpora/reuters/training/4175
inflating: /root/nltk_data/corpora/reuters/training/4176
inflating: /root/nltk_data/corpora/reuters/training/4177
inflating: /root/nltk_data/corpora/reuters/training/4178
inflating: /root/nltk_data/corpora/reuters/training/418
inflating: /root/nltk_data/corpora/reuters/training/4182
inflating: /root/nltk_data/corpora/reuters/training/4184
inflating: /root/nltk_data/corpora/reuters/training/4185
inflating: /root/nltk_data/corpora/reuters/training/4188
inflating: /root/nltk_data/corpora/reuters/training/4189
inflating: /root/nltk_data/corpora/reuters/training/4192
inflating: /root/nltk_data/corpora/reuters/training/4193
inflating: /root/nltk_data/corpora/reuters/training/4194
inflating: /root/nltk_data/corpora/reuters/training/4199
inflating: /root/nltk_data/corpora/reuters/training/42
inflating: /root/nltk_data/corpora/reuters/training/420
inflating: /root/nltk_data/corpora/reuters/training/4201
inflating: /root/nltk_data/corpora/reuters/training/4203
inflating: /root/nltk_data/corpora/reuters/training/4204
inflating: /root/nltk_data/corpora/reuters/training/4205
inflating: /root/nltk_data/corpora/reuters/training/4206
inflating: /root/nltk_data/corpora/reuters/training/4209
inflating: /root/nltk_data/corpora/reuters/training/4210
inflating: /root/nltk_data/corpora/reuters/training/4212
inflating: /root/nltk_data/corpora/reuters/training/4215
inflating: /root/nltk_data/corpora/reuters/training/4216
inflating: /root/nltk_data/corpora/reuters/training/4218
inflating: /root/nltk_data/corpora/reuters/training/4219
inflating: /root/nltk_data/corpora/reuters/training/422
inflating: /root/nltk_data/corpora/reuters/training/4220
inflating: /root/nltk_data/corpora/reuters/training/4223
inflating: /root/nltk_data/corpora/reuters/training/4227
inflating: /root/nltk_data/corpora/reuters/training/4228
inflating: /root/nltk_data/corpora/reuters/training/4229
inflating: /root/nltk_data/corpora/reuters/training/423
inflating: /root/nltk_data/corpora/reuters/training/4231
inflating: /root/nltk_data/corpora/reuters/training/4232
inflating: /root/nltk_data/corpora/reuters/training/4233
inflating: /root/nltk_data/corpora/reuters/training/4235
inflating: /root/nltk_data/corpora/reuters/training/4237
inflating: /root/nltk_data/corpora/reuters/training/4238
inflating: /root/nltk_data/corpora/reuters/training/4239
inflating: /root/nltk_data/corpora/reuters/training/424
inflating: /root/nltk_data/corpora/reuters/training/4240
inflating: /root/nltk_data/corpora/reuters/training/4241
inflating: /root/nltk_data/corpora/reuters/training/4245
inflating: /root/nltk_data/corpora/reuters/training/4246
inflating: /root/nltk_data/corpora/reuters/training/4249
inflating: /root/nltk_data/corpora/reuters/training/425
inflating: /root/nltk_data/corpora/reuters/training/4253
inflating: /root/nltk_data/corpora/reuters/training/4255
inflating: /root/nltk_data/corpora/reuters/training/4257
inflating: /root/nltk_data/corpora/reuters/training/4258
inflating: /root/nltk_data/corpora/reuters/training/4259
inflating: /root/nltk_data/corpora/reuters/training/4261
inflating: /root/nltk_data/corpora/reuters/training/4266
inflating: /root/nltk_data/corpora/reuters/training/4267
```

```
inflating: /root/nltk_data/corpora/reuters/training/4273
inflating: /root/nltk_data/corpora/reuters/training/4275
inflating: /root/nltk_data/corpora/reuters/training/4276
inflating: /root/nltk_data/corpora/reuters/training/4277
inflating: /root/nltk_data/corpora/reuters/training/4278
inflating: /root/nltk_data/corpora/reuters/training/4279
inflating: /root/nltk_data/corpora/reuters/training/4280
inflating: /root/nltk_data/corpora/reuters/training/4281
inflating: /root/nltk_data/corpora/reuters/training/4282
inflating: /root/nltk_data/corpora/reuters/training/4284
inflating: /root/nltk_data/corpora/reuters/training/4285
inflating: /root/nltk_data/corpora/reuters/training/4286
inflating: /root/nltk_data/corpora/reuters/training/4288
inflating: /root/nltk_data/corpora/reuters/training/4289
inflating: /root/nltk_data/corpora/reuters/training/4290
inflating: /root/nltk_data/corpora/reuters/training/4291
inflating: /root/nltk_data/corpora/reuters/training/4292
inflating: /root/nltk_data/corpora/reuters/training/4293
inflating: /root/nltk_data/corpora/reuters/training/4296
inflating: /root/nltk_data/corpora/reuters/training/4297
inflating: /root/nltk_data/corpora/reuters/training/4298
inflating: /root/nltk_data/corpora/reuters/training/4301
inflating: /root/nltk_data/corpora/reuters/training/4302
inflating: /root/nltk_data/corpora/reuters/training/4303
inflating: /root/nltk_data/corpora/reuters/training/4304
inflating: /root/nltk_data/corpora/reuters/training/4305
inflating: /root/nltk_data/corpora/reuters/training/4306
inflating: /root/nltk_data/corpora/reuters/training/4307
inflating: /root/nltk_data/corpora/reuters/training/4310
inflating: /root/nltk_data/corpora/reuters/training/4313
inflating: /root/nltk_data/corpora/reuters/training/4314
inflating: /root/nltk_data/corpora/reuters/training/4315
inflating: /root/nltk_data/corpora/reuters/training/4318
inflating: /root/nltk_data/corpora/reuters/training/4323
inflating: /root/nltk_data/corpora/reuters/training/4327
inflating: /root/nltk_data/corpora/reuters/training/4328
inflating: /root/nltk_data/corpora/reuters/training/4330
inflating: /root/nltk_data/corpora/reuters/training/4331
inflating: /root/nltk_data/corpora/reuters/training/4333
inflating: /root/nltk_data/corpora/reuters/training/4337
inflating: /root/nltk_data/corpora/reuters/training/4338
inflating: /root/nltk_data/corpora/reuters/training/4339
inflating: /root/nltk_data/corpora/reuters/training/4340
inflating: /root/nltk_data/corpora/reuters/training/4341
inflating: /root/nltk_data/corpora/reuters/training/4343
inflating: /root/nltk_data/corpora/reuters/training/4344
inflating: /root/nltk_data/corpora/reuters/training/4345
inflating: /root/nltk_data/corpora/reuters/training/4346
inflating: /root/nltk_data/corpora/reuters/training/4347
inflating: /root/nltk_data/corpora/reuters/training/4349
inflating: /root/nltk_data/corpora/reuters/training/435
inflating: /root/nltk_data/corpora/reuters/training/4351
inflating: /root/nltk_data/corpora/reuters/training/4353
inflating: /root/nltk_data/corpora/reuters/training/4356
inflating: /root/nltk_data/corpora/reuters/training/4357
inflating: /root/nltk_data/corpora/reuters/training/4358
inflating: /root/nltk_data/corpora/reuters/training/436
inflating: /root/nltk_data/corpora/reuters/training/4360
inflating: /root/nltk_data/corpora/reuters/training/4361
inflating: /root/nltk_data/corpora/reuters/training/4363
inflating: /root/nltk_data/corpora/reuters/training/4364
inflating: /root/nltk_data/corpora/reuters/training/4365
inflating: /root/nltk_data/corpora/reuters/training/4366
inflating: /root/nltk_data/corpora/reuters/training/4367
inflating: /root/nltk_data/corpora/reuters/training/4369
inflating: /root/nltk_data/corpora/reuters/training/4370
inflating: /root/nltk_data/corpora/reuters/training/4371
inflating: /root/nltk_data/corpora/reuters/training/4372
inflating: /root/nltk_data/corpora/reuters/training/4374
inflating: /root/nltk_data/corpora/reuters/training/4376
inflating: /root/nltk_data/corpora/reuters/training/4378
inflating: /root/nltk_data/corpora/reuters/training/438
inflating: /root/nltk_data/corpora/reuters/training/4380
inflating: /root/nltk_data/corpora/reuters/training/4382
inflating: /root/nltk_data/corpora/reuters/training/4384
inflating: /root/nltk_data/corpora/reuters/training/4385
inflating: /root/nltk_data/corpora/reuters/training/4386
inflating: /root/nltk_data/corpora/reuters/training/4387
inflating: /root/nltk_data/corpora/reuters/training/4388
inflating: /root/nltk_data/corpora/reuters/training/4389
inflating: /root/nltk_data/corpora/reuters/training/439
inflating: /root/nltk_data/corpora/reuters/training/4392
inflating: /root/nltk_data/corpora/reuters/training/4397
inflating: /root/nltk_data/corpora/reuters/training/4398
inflating: /root/nltk_data/corpora/reuters/training/4399
inflating: /root/nltk_data/corpora/reuters/training/44
inflating: /root/nltk_data/corpora/reuters/training/440
inflating: /root/nltk_data/corpora/reuters/training/4401
inflating: /root/nltk_data/corpora/reuters/training/4404
inflating: /root/nltk_data/corpora/reuters/training/4405
inflating: /root/nltk_data/corpora/reuters/training/4407
inflating: /root/nltk_data/corpora/reuters/training/4408
inflating: /root/nltk_data/corpora/reuters/training/4409
inflating: /root/nltk_data/corpora/reuters/training/441
inflating: /root/nltk_data/corpora/reuters/training/4410
inflating: /root/nltk_data/corpora/reuters/training/4411
inflating: /root/nltk_data/corpora/reuters/training/4412
inflating: /root/nltk_data/corpora/reuters/training/4419
inflating: /root/nltk_data/corpora/reuters/training/442
inflating: /root/nltk_data/corpora/reuters/training/4420
inflating: /root/nltk_data/corpora/reuters/training/4425
inflating: /root/nltk_data/corpora/reuters/training/4429
inflating: /root/nltk_data/corpora/reuters/training/443
inflating: /root/nltk_data/corpora/reuters/training/4431
inflating: /root/nltk_data/corpora/reuters/training/4433
inflating: /root/nltk_data/corpora/reuters/training/4434
inflating: /root/nltk_data/corpora/reuters/training/4436
inflating: /root/nltk_data/corpora/reuters/training/4437
inflating: /root/nltk_data/corpora/reuters/training/4438
inflating: /root/nltk_data/corpora/reuters/training/4439
inflating: /root/nltk_data/corpora/reuters/training/444
inflating: /root/nltk_data/corpora/reuters/training/4440
inflating: /root/nltk_data/corpora/reuters/training/4445
inflating: /root/nltk_data/corpora/reuters/training/4446
inflating: /root/nltk_data/corpora/reuters/training/4447
inflating: /root/nltk_data/corpora/reuters/training/4449
inflating: /root/nltk_data/corpora/reuters/training/4451
inflating: /root/nltk_data/corpora/reuters/training/4452
inflating: /root/nltk_data/corpora/reuters/training/4453
inflating: /root/nltk_data/corpora/reuters/training/4454
inflating: /root/nltk_data/corpora/reuters/training/4455
inflating: /root/nltk_data/corpora/reuters/training/4457
inflating: /root/nltk_data/corpora/reuters/training/4458
inflating: /root/nltk_data/corpora/reuters/training/4459
inflating: /root/nltk_data/corpora/reuters/training/4460
inflating: /root/nltk_data/corpora/reuters/training/4462
inflating: /root/nltk_data/corpora/reuters/training/4465
inflating: /root/nltk_data/corpora/reuters/training/4466
inflating: /root/nltk_data/corpora/reuters/training/4467
inflating: /root/nltk_data/corpora/reuters/training/4468
inflating: /root/nltk_data/corpora/reuters/training/447
inflating: /root/nltk_data/corpora/reuters/training/4470
inflating: /root/nltk_data/corpora/reuters/training/4472
inflating: /root/nltk_data/corpora/reuters/training/4473
inflating: /root/nltk_data/corpora/reuters/training/4474
```

```
inflating: /root/nltk_data/corpora/reuters/training/4474
inflating: /root/nltk_data/corpora/reuters/training/4475
inflating: /root/nltk_data/corpora/reuters/training/4476
inflating: /root/nltk_data/corpora/reuters/training/4477
inflating: /root/nltk_data/corpora/reuters/training/4478
inflating: /root/nltk_data/corpora/reuters/training/448
inflating: /root/nltk_data/corpora/reuters/training/4480
inflating: /root/nltk_data/corpora/reuters/training/4481
inflating: /root/nltk_data/corpora/reuters/training/4483
inflating: /root/nltk_data/corpora/reuters/training/4484
inflating: /root/nltk_data/corpora/reuters/training/4486
inflating: /root/nltk_data/corpora/reuters/training/4487
inflating: /root/nltk_data/corpora/reuters/training/4488
inflating: /root/nltk_data/corpora/reuters/training/449
inflating: /root/nltk_data/corpora/reuters/training/4490
inflating: /root/nltk_data/corpora/reuters/training/4493
inflating: /root/nltk_data/corpora/reuters/training/4494
inflating: /root/nltk_data/corpora/reuters/training/4498
inflating: /root/nltk_data/corpora/reuters/training/4499
inflating: /root/nltk_data/corpora/reuters/training/45
inflating: /root/nltk_data/corpora/reuters/training/450
inflating: /root/nltk_data/corpora/reuters/training/4503
inflating: /root/nltk_data/corpora/reuters/training/4504
inflating: /root/nltk_data/corpora/reuters/training/4505
inflating: /root/nltk_data/corpora/reuters/training/4507
inflating: /root/nltk_data/corpora/reuters/training/4509
inflating: /root/nltk_data/corpora/reuters/training/4510
inflating: /root/nltk_data/corpora/reuters/training/4512
inflating: /root/nltk_data/corpora/reuters/training/4513
inflating: /root/nltk_data/corpora/reuters/training/4514
inflating: /root/nltk_data/corpora/reuters/training/4515
inflating: /root/nltk_data/corpora/reuters/training/4516
inflating: /root/nltk_data/corpora/reuters/training/4517
inflating: /root/nltk_data/corpora/reuters/training/4518
inflating: /root/nltk_data/corpora/reuters/training/4519
inflating: /root/nltk_data/corpora/reuters/training/4521
inflating: /root/nltk_data/corpora/reuters/training/4523
inflating: /root/nltk_data/corpora/reuters/training/4524
inflating: /root/nltk_data/corpora/reuters/training/4525
inflating: /root/nltk_data/corpora/reuters/training/4532
inflating: /root/nltk_data/corpora/reuters/training/4533
inflating: /root/nltk_data/corpora/reuters/training/4534
inflating: /root/nltk_data/corpora/reuters/training/4536
inflating: /root/nltk_data/corpora/reuters/training/4537
inflating: /root/nltk_data/corpora/reuters/training/4541
inflating: /root/nltk_data/corpora/reuters/training/4544
inflating: /root/nltk_data/corpora/reuters/training/4548
inflating: /root/nltk_data/corpora/reuters/training/4549
inflating: /root/nltk_data/corpora/reuters/training/455
inflating: /root/nltk_data/corpora/reuters/training/4551
inflating: /root/nltk_data/corpora/reuters/training/4552
inflating: /root/nltk_data/corpora/reuters/training/4554
inflating: /root/nltk_data/corpora/reuters/training/4555
inflating: /root/nltk_data/corpora/reuters/training/4558
inflating: /root/nltk_data/corpora/reuters/training/456
inflating: /root/nltk_data/corpora/reuters/training/4560
inflating: /root/nltk_data/corpora/reuters/training/4562
inflating: /root/nltk_data/corpora/reuters/training/4564
inflating: /root/nltk_data/corpora/reuters/training/4569
inflating: /root/nltk_data/corpora/reuters/training/4573
inflating: /root/nltk_data/corpora/reuters/training/4574
inflating: /root/nltk_data/corpora/reuters/training/4576
inflating: /root/nltk_data/corpora/reuters/training/4577
inflating: /root/nltk_data/corpora/reuters/training/4578
inflating: /root/nltk_data/corpora/reuters/training/4579
inflating: /root/nltk_data/corpora/reuters/training/4581
inflating: /root/nltk_data/corpora/reuters/training/4583
inflating: /root/nltk_data/corpora/reuters/training/4584
inflating: /root/nltk_data/corpora/reuters/training/4585
inflating: /root/nltk_data/corpora/reuters/training/4586
inflating: /root/nltk_data/corpora/reuters/training/459
inflating: /root/nltk_data/corpora/reuters/training/4590
inflating: /root/nltk_data/corpora/reuters/training/4591
inflating: /root/nltk_data/corpora/reuters/training/4592
inflating: /root/nltk_data/corpora/reuters/training/4593
inflating: /root/nltk_data/corpora/reuters/training/4595
inflating: /root/nltk_data/corpora/reuters/training/4599
inflating: /root/nltk_data/corpora/reuters/training/46
inflating: /root/nltk_data/corpora/reuters/training/4600
inflating: /root/nltk_data/corpora/reuters/training/4602
inflating: /root/nltk_data/corpora/reuters/training/4603
inflating: /root/nltk_data/corpora/reuters/training/4604
inflating: /root/nltk_data/corpora/reuters/training/4605
inflating: /root/nltk_data/corpora/reuters/training/4606
inflating: /root/nltk_data/corpora/reuters/training/4609
inflating: /root/nltk_data/corpora/reuters/training/461
inflating: /root/nltk_data/corpora/reuters/training/4610
inflating: /root/nltk_data/corpora/reuters/training/4613
inflating: /root/nltk_data/corpora/reuters/training/4615
inflating: /root/nltk_data/corpora/reuters/training/4616
inflating: /root/nltk_data/corpora/reuters/training/4617
inflating: /root/nltk_data/corpora/reuters/training/4620
inflating: /root/nltk_data/corpora/reuters/training/4621
inflating: /root/nltk_data/corpora/reuters/training/4622
inflating: /root/nltk_data/corpora/reuters/training/4623
inflating: /root/nltk_data/corpora/reuters/training/4625
inflating: /root/nltk_data/corpora/reuters/training/4629
inflating: /root/nltk_data/corpora/reuters/training/4630
inflating: /root/nltk_data/corpora/reuters/training/4631
inflating: /root/nltk_data/corpora/reuters/training/4632
inflating: /root/nltk_data/corpora/reuters/training/4633
inflating: /root/nltk_data/corpora/reuters/training/4634
inflating: /root/nltk_data/corpora/reuters/training/4635
inflating: /root/nltk_data/corpora/reuters/training/4636
inflating: /root/nltk_data/corpora/reuters/training/4637
inflating: /root/nltk_data/corpora/reuters/training/4638
inflating: /root/nltk_data/corpora/reuters/training/464
inflating: /root/nltk_data/corpora/reuters/training/4640
inflating: /root/nltk_data/corpora/reuters/training/4641
inflating: /root/nltk_data/corpora/reuters/training/4644
inflating: /root/nltk_data/corpora/reuters/training/4648
inflating: /root/nltk_data/corpora/reuters/training/4649
inflating: /root/nltk_data/corpora/reuters/training/4650
inflating: /root/nltk_data/corpora/reuters/training/4652
inflating: /root/nltk_data/corpora/reuters/training/4654
inflating: /root/nltk_data/corpora/reuters/training/4656
inflating: /root/nltk_data/corpora/reuters/training/4657
inflating: /root/nltk_data/corpora/reuters/training/4658
inflating: /root/nltk_data/corpora/reuters/training/4659
inflating: /root/nltk_data/corpora/reuters/training/466
inflating: /root/nltk_data/corpora/reuters/training/4660
inflating: /root/nltk_data/corpora/reuters/training/4662
inflating: /root/nltk_data/corpora/reuters/training/4664
inflating: /root/nltk_data/corpora/reuters/training/4665
inflating: /root/nltk_data/corpora/reuters/training/4666
inflating: /root/nltk_data/corpora/reuters/training/467
inflating: /root/nltk_data/corpora/reuters/training/4671
inflating: /root/nltk_data/corpora/reuters/training/4675
inflating: /root/nltk_data/corpora/reuters/training/4678
inflating: /root/nltk_data/corpora/reuters/training/4679
inflating: /root/nltk_data/corpora/reuters/training/4680
inflating: /root/nltk_data/corpora/reuters/training/4681
inflating: /root/nltk_data/corpora/reuters/training/4682
inflating: /root/nltk_data/corpora/reuters/training/4686
inflating: /root/nltk_data/corpora/reuters/training/4687
```

```
inflating: /root/nltk_data/corpora/reuters/training/4689
inflating: /root/nltk_data/corpora/reuters/training/4690
inflating: /root/nltk_data/corpora/reuters/training/4691
inflating: /root/nltk_data/corpora/reuters/training/4692
inflating: /root/nltk_data/corpora/reuters/training/4695
inflating: /root/nltk_data/corpora/reuters/training/4696
inflating: /root/nltk_data/corpora/reuters/training/4697
inflating: /root/nltk_data/corpora/reuters/training/4698
inflating: /root/nltk_data/corpora/reuters/training/47
inflating: /root/nltk_data/corpora/reuters/training/470
inflating: /root/nltk_data/corpora/reuters/training/4700
inflating: /root/nltk_data/corpora/reuters/training/4703
inflating: /root/nltk_data/corpora/reuters/training/4708
inflating: /root/nltk_data/corpora/reuters/training/4709
inflating: /root/nltk_data/corpora/reuters/training/4711
inflating: /root/nltk_data/corpora/reuters/training/4712
inflating: /root/nltk_data/corpora/reuters/training/4713
inflating: /root/nltk_data/corpora/reuters/training/4714
inflating: /root/nltk_data/corpora/reuters/training/4717
inflating: /root/nltk_data/corpora/reuters/training/4719
inflating: /root/nltk_data/corpora/reuters/training/4720
inflating: /root/nltk_data/corpora/reuters/training/4729
inflating: /root/nltk_data/corpora/reuters/training/473
inflating: /root/nltk_data/corpora/reuters/training/4732
inflating: /root/nltk_data/corpora/reuters/training/4733
inflating: /root/nltk_data/corpora/reuters/training/4735
inflating: /root/nltk_data/corpora/reuters/training/4739
inflating: /root/nltk_data/corpora/reuters/training/474
inflating: /root/nltk_data/corpora/reuters/training/4740
inflating: /root/nltk_data/corpora/reuters/training/4741
inflating: /root/nltk_data/corpora/reuters/training/4742
inflating: /root/nltk_data/corpora/reuters/training/4743
inflating: /root/nltk_data/corpora/reuters/training/4744
inflating: /root/nltk_data/corpora/reuters/training/4745
inflating: /root/nltk_data/corpora/reuters/training/4747
inflating: /root/nltk_data/corpora/reuters/training/4748
inflating: /root/nltk_data/corpora/reuters/training/475
inflating: /root/nltk_data/corpora/reuters/training/4751
inflating: /root/nltk_data/corpora/reuters/training/4755
inflating: /root/nltk_data/corpora/reuters/training/4756
inflating: /root/nltk_data/corpora/reuters/training/4758
inflating: /root/nltk_data/corpora/reuters/training/4759
inflating: /root/nltk_data/corpora/reuters/training/4761
inflating: /root/nltk_data/corpora/reuters/training/4763
inflating: /root/nltk_data/corpora/reuters/training/4765
inflating: /root/nltk_data/corpora/reuters/training/4768
inflating: /root/nltk_data/corpora/reuters/training/4769
inflating: /root/nltk_data/corpora/reuters/training/4770
inflating: /root/nltk_data/corpora/reuters/training/4771
inflating: /root/nltk_data/corpora/reuters/training/4772
inflating: /root/nltk_data/corpora/reuters/training/4773
inflating: /root/nltk_data/corpora/reuters/training/4774
inflating: /root/nltk_data/corpora/reuters/training/4775
inflating: /root/nltk_data/corpora/reuters/training/4776
inflating: /root/nltk_data/corpora/reuters/training/4777
inflating: /root/nltk_data/corpora/reuters/training/4778
inflating: /root/nltk_data/corpora/reuters/training/4779
inflating: /root/nltk_data/corpora/reuters/training/478
inflating: /root/nltk_data/corpora/reuters/training/4780
inflating: /root/nltk_data/corpora/reuters/training/4781
inflating: /root/nltk_data/corpora/reuters/training/4782
inflating: /root/nltk_data/corpora/reuters/training/4783
inflating: /root/nltk_data/corpora/reuters/training/4785
inflating: /root/nltk_data/corpora/reuters/training/4787
inflating: /root/nltk_data/corpora/reuters/training/4788
inflating: /root/nltk_data/corpora/reuters/training/4791
inflating: /root/nltk_data/corpora/reuters/training/4792
inflating: /root/nltk_data/corpora/reuters/training/4794
inflating: /root/nltk_data/corpora/reuters/training/4795
inflating: /root/nltk_data/corpora/reuters/training/4796
inflating: /root/nltk_data/corpora/reuters/training/4797
inflating: /root/nltk_data/corpora/reuters/training/4798
inflating: /root/nltk_data/corpora/reuters/training/48
inflating: /root/nltk_data/corpora/reuters/training/480
inflating: /root/nltk_data/corpora/reuters/training/4801
inflating: /root/nltk_data/corpora/reuters/training/4803
inflating: /root/nltk_data/corpora/reuters/training/4804
inflating: /root/nltk_data/corpora/reuters/training/4805
inflating: /root/nltk_data/corpora/reuters/training/4806
inflating: /root/nltk_data/corpora/reuters/training/4807
inflating: /root/nltk_data/corpora/reuters/training/4809
inflating: /root/nltk_data/corpora/reuters/training/4810
inflating: /root/nltk_data/corpora/reuters/training/4812
inflating: /root/nltk_data/corpora/reuters/training/4814
inflating: /root/nltk_data/corpora/reuters/training/4816
inflating: /root/nltk_data/corpora/reuters/training/4818
inflating: /root/nltk_data/corpora/reuters/training/4819
inflating: /root/nltk_data/corpora/reuters/training/4820
inflating: /root/nltk_data/corpora/reuters/training/4821
inflating: /root/nltk_data/corpora/reuters/training/4822
inflating: /root/nltk_data/corpora/reuters/training/4824
inflating: /root/nltk_data/corpora/reuters/training/4825
inflating: /root/nltk_data/corpora/reuters/training/4827
inflating: /root/nltk_data/corpora/reuters/training/4829
inflating: /root/nltk_data/corpora/reuters/training/483
inflating: /root/nltk_data/corpora/reuters/training/4831
inflating: /root/nltk_data/corpora/reuters/training/4832
inflating: /root/nltk_data/corpora/reuters/training/4833
inflating: /root/nltk_data/corpora/reuters/training/4835
inflating: /root/nltk_data/corpora/reuters/training/484
inflating: /root/nltk_data/corpora/reuters/training/4841
inflating: /root/nltk_data/corpora/reuters/training/4843
inflating: /root/nltk_data/corpora/reuters/training/4845
inflating: /root/nltk_data/corpora/reuters/training/4848
inflating: /root/nltk_data/corpora/reuters/training/4850
inflating: /root/nltk_data/corpora/reuters/training/486
inflating: /root/nltk_data/corpora/reuters/training/4867
inflating: /root/nltk_data/corpora/reuters/training/488
inflating: /root/nltk_data/corpora/reuters/training/4881
inflating: /root/nltk_data/corpora/reuters/training/4882
inflating: /root/nltk_data/corpora/reuters/training/4884
inflating: /root/nltk_data/corpora/reuters/training/4887
inflating: /root/nltk_data/corpora/reuters/training/4889
inflating: /root/nltk_data/corpora/reuters/training/489
inflating: /root/nltk_data/corpora/reuters/training/4890
inflating: /root/nltk_data/corpora/reuters/training/4892
inflating: /root/nltk_data/corpora/reuters/training/4893
inflating: /root/nltk_data/corpora/reuters/training/4894
inflating: /root/nltk_data/corpora/reuters/training/4898
inflating: /root/nltk_data/corpora/reuters/training/49
inflating: /root/nltk_data/corpora/reuters/training/490
inflating: /root/nltk_data/corpora/reuters/training/4900
inflating: /root/nltk_data/corpora/reuters/training/4901
inflating: /root/nltk_data/corpora/reuters/training/4903
inflating: /root/nltk_data/corpora/reuters/training/4904
inflating: /root/nltk_data/corpora/reuters/training/4905
inflating: /root/nltk_data/corpora/reuters/training/4908
inflating: /root/nltk_data/corpora/reuters/training/4910
inflating: /root/nltk_data/corpora/reuters/training/4911
inflating: /root/nltk_data/corpora/reuters/training/4914
inflating: /root/nltk_data/corpora/reuters/training/4918
inflating: /root/nltk_data/corpora/reuters/training/4920
inflating: /root/nltk_data/corpora/reuters/training/4922
inflating: /root/nltk_data/corpora/reuters/training/4925
```

```
inflating: /root/nltk_data/corpora/reuters/training/4929
inflating: /root/nltk_data/corpora/reuters/training/4930
inflating: /root/nltk_data/corpora/reuters/training/4933
inflating: /root/nltk_data/corpora/reuters/training/4934
inflating: /root/nltk_data/corpora/reuters/training/4937
inflating: /root/nltk_data/corpora/reuters/training/4938
inflating: /root/nltk_data/corpora/reuters/training/4939
inflating: /root/nltk_data/corpora/reuters/training/494
inflating: /root/nltk_data/corpora/reuters/training/4941
inflating: /root/nltk_data/corpora/reuters/training/4944
inflating: /root/nltk_data/corpora/reuters/training/4951
inflating: /root/nltk_data/corpora/reuters/training/4952
inflating: /root/nltk_data/corpora/reuters/training/4953
inflating: /root/nltk_data/corpora/reuters/training/4954
inflating: /root/nltk_data/corpora/reuters/training/4955
inflating: /root/nltk_data/corpora/reuters/training/4956
inflating: /root/nltk_data/corpora/reuters/training/4959
inflating: /root/nltk_data/corpora/reuters/training/496
inflating: /root/nltk_data/corpora/reuters/training/4960
inflating: /root/nltk_data/corpora/reuters/training/4962
inflating: /root/nltk_data/corpora/reuters/training/4963
inflating: /root/nltk_data/corpora/reuters/training/4964
inflating: /root/nltk_data/corpora/reuters/training/4965
inflating: /root/nltk_data/corpora/reuters/training/497
inflating: /root/nltk_data/corpora/reuters/training/4970
inflating: /root/nltk_data/corpora/reuters/training/4972
inflating: /root/nltk_data/corpora/reuters/training/4974
inflating: /root/nltk_data/corpora/reuters/training/4976
inflating: /root/nltk_data/corpora/reuters/training/498
inflating: /root/nltk_data/corpora/reuters/training/4981
inflating: /root/nltk_data/corpora/reuters/training/4983
inflating: /root/nltk_data/corpora/reuters/training/4987
inflating: /root/nltk_data/corpora/reuters/training/4988
inflating: /root/nltk_data/corpora/reuters/training/4990
inflating: /root/nltk_data/corpora/reuters/training/4991
inflating: /root/nltk_data/corpora/reuters/training/4993
inflating: /root/nltk_data/corpora/reuters/training/4998
inflating: /root/nltk_data/corpora/reuters/training/4999
inflating: /root/nltk_data/corpora/reuters/training/5
inflating: /root/nltk_data/corpora/reuters/training/50
inflating: /root/nltk_data/corpora/reuters/training/5001
inflating: /root/nltk_data/corpora/reuters/training/5002
inflating: /root/nltk_data/corpora/reuters/training/5003
inflating: /root/nltk_data/corpora/reuters/training/5009
inflating: /root/nltk_data/corpora/reuters/training/501
inflating: /root/nltk_data/corpora/reuters/training/5012
inflating: /root/nltk_data/corpora/reuters/training/5013
inflating: /root/nltk_data/corpora/reuters/training/5015
inflating: /root/nltk_data/corpora/reuters/training/5017
inflating: /root/nltk_data/corpora/reuters/training/5018
inflating: /root/nltk_data/corpora/reuters/training/502
inflating: /root/nltk_data/corpora/reuters/training/5020
inflating: /root/nltk_data/corpora/reuters/training/5021
inflating: /root/nltk_data/corpora/reuters/training/5029
inflating: /root/nltk_data/corpora/reuters/training/5031
inflating: /root/nltk_data/corpora/reuters/training/5033
inflating: /root/nltk_data/corpora/reuters/training/5034
inflating: /root/nltk_data/corpora/reuters/training/5037
inflating: /root/nltk_data/corpora/reuters/training/504
inflating: /root/nltk_data/corpora/reuters/training/5040
inflating: /root/nltk_data/corpora/reuters/training/5042
inflating: /root/nltk_data/corpora/reuters/training/5043
inflating: /root/nltk_data/corpora/reuters/training/5044
inflating: /root/nltk_data/corpora/reuters/training/5046
inflating: /root/nltk_data/corpora/reuters/training/5048
inflating: /root/nltk_data/corpora/reuters/training/505
inflating: /root/nltk_data/corpora/reuters/training/5050
inflating: /root/nltk_data/corpora/reuters/training/5051
inflating: /root/nltk_data/corpora/reuters/training/5052
inflating: /root/nltk_data/corpora/reuters/training/5053
inflating: /root/nltk_data/corpora/reuters/training/5054
inflating: /root/nltk_data/corpora/reuters/training/5055
inflating: /root/nltk_data/corpora/reuters/training/5056
inflating: /root/nltk_data/corpora/reuters/training/5057
inflating: /root/nltk_data/corpora/reuters/training/5058
inflating: /root/nltk_data/corpora/reuters/training/5059
inflating: /root/nltk_data/corpora/reuters/training/506
inflating: /root/nltk_data/corpora/reuters/training/5060
inflating: /root/nltk_data/corpora/reuters/training/5061
inflating: /root/nltk_data/corpora/reuters/training/5062
inflating: /root/nltk_data/corpora/reuters/training/5068
inflating: /root/nltk_data/corpora/reuters/training/5070
inflating: /root/nltk_data/corpora/reuters/training/5071
inflating: /root/nltk_data/corpora/reuters/training/5072
inflating: /root/nltk_data/corpora/reuters/training/5076
inflating: /root/nltk_data/corpora/reuters/training/5077
inflating: /root/nltk_data/corpora/reuters/training/5078
inflating: /root/nltk_data/corpora/reuters/training/508
inflating: /root/nltk_data/corpora/reuters/training/5080
inflating: /root/nltk_data/corpora/reuters/training/5081
inflating: /root/nltk_data/corpora/reuters/training/5082
inflating: /root/nltk_data/corpora/reuters/training/5083
inflating: /root/nltk_data/corpora/reuters/training/5084
inflating: /root/nltk_data/corpora/reuters/training/5085
inflating: /root/nltk_data/corpora/reuters/training/5086
inflating: /root/nltk_data/corpora/reuters/training/5089
inflating: /root/nltk_data/corpora/reuters/training/509
inflating: /root/nltk_data/corpora/reuters/training/5090
inflating: /root/nltk_data/corpora/reuters/training/5091
inflating: /root/nltk_data/corpora/reuters/training/5092
inflating: /root/nltk_data/corpora/reuters/training/5094
inflating: /root/nltk_data/corpora/reuters/training/5095
inflating: /root/nltk_data/corpora/reuters/training/5097
inflating: /root/nltk_data/corpora/reuters/training/5100
inflating: /root/nltk_data/corpora/reuters/training/5102
inflating: /root/nltk_data/corpora/reuters/training/5105
inflating: /root/nltk_data/corpora/reuters/training/5109
inflating: /root/nltk_data/corpora/reuters/training/511
inflating: /root/nltk_data/corpora/reuters/training/5112
inflating: /root/nltk_data/corpora/reuters/training/5113
inflating: /root/nltk_data/corpora/reuters/training/5114
inflating: /root/nltk_data/corpora/reuters/training/5115
inflating: /root/nltk_data/corpora/reuters/training/5116
inflating: /root/nltk_data/corpora/reuters/training/5117
inflating: /root/nltk_data/corpora/reuters/training/5118
inflating: /root/nltk_data/corpora/reuters/training/5119
inflating: /root/nltk_data/corpora/reuters/training/512
inflating: /root/nltk_data/corpora/reuters/training/5121
inflating: /root/nltk_data/corpora/reuters/training/5123
inflating: /root/nltk_data/corpora/reuters/training/5125
inflating: /root/nltk_data/corpora/reuters/training/5126
inflating: /root/nltk_data/corpora/reuters/training/5127
inflating: /root/nltk_data/corpora/reuters/training/5128
inflating: /root/nltk_data/corpora/reuters/training/513
inflating: /root/nltk_data/corpora/reuters/training/5132
inflating: /root/nltk_data/corpora/reuters/training/5134
inflating: /root/nltk_data/corpora/reuters/training/5138
inflating: /root/nltk_data/corpora/reuters/training/5139
inflating: /root/nltk_data/corpora/reuters/training/514
inflating: /root/nltk_data/corpora/reuters/training/5141
inflating: /root/nltk_data/corpora/reuters/training/5142
inflating: /root/nltk_data/corpora/reuters/training/5145
inflating: /root/nltk_data/corpora/reuters/training/5146
inflating: /root/nltk_data/corpora/reuters/training/5148
inflating: /root/nltk_data/corpora/reuters/training/5149
```

```
inflating: /root/nltk_data/corpora/reuters/training/5149
inflating: /root/nltk_data/corpora/reuters/training/5150
inflating: /root/nltk_data/corpora/reuters/training/5152
inflating: /root/nltk_data/corpora/reuters/training/5153
inflating: /root/nltk_data/corpora/reuters/training/5154
inflating: /root/nltk_data/corpora/reuters/training/5156
inflating: /root/nltk_data/corpora/reuters/training/516
inflating: /root/nltk_data/corpora/reuters/training/5160
inflating: /root/nltk_data/corpora/reuters/training/5162
inflating: /root/nltk_data/corpora/reuters/training/5166
inflating: /root/nltk_data/corpora/reuters/training/5167
inflating: /root/nltk_data/corpora/reuters/training/5168
inflating: /root/nltk_data/corpora/reuters/training/5169
inflating: /root/nltk_data/corpora/reuters/training/517
inflating: /root/nltk_data/corpora/reuters/training/5171
inflating: /root/nltk_data/corpora/reuters/training/5172
inflating: /root/nltk_data/corpora/reuters/training/5175
inflating: /root/nltk_data/corpora/reuters/training/5176
inflating: /root/nltk_data/corpora/reuters/training/5177
inflating: /root/nltk_data/corpora/reuters/training/5178
inflating: /root/nltk_data/corpora/reuters/training/518
inflating: /root/nltk_data/corpora/reuters/training/5181
inflating: /root/nltk_data/corpora/reuters/training/5183
inflating: /root/nltk_data/corpora/reuters/training/5185
inflating: /root/nltk_data/corpora/reuters/training/5188
inflating: /root/nltk_data/corpora/reuters/training/5189
inflating: /root/nltk_data/corpora/reuters/training/519
inflating: /root/nltk_data/corpora/reuters/training/5190
inflating: /root/nltk_data/corpora/reuters/training/5191
inflating: /root/nltk_data/corpora/reuters/training/5192
inflating: /root/nltk_data/corpora/reuters/training/5193
inflating: /root/nltk_data/corpora/reuters/training/5194
inflating: /root/nltk_data/corpora/reuters/training/5195
inflating: /root/nltk_data/corpora/reuters/training/5196
inflating: /root/nltk_data/corpora/reuters/training/5201
inflating: /root/nltk_data/corpora/reuters/training/5203
inflating: /root/nltk_data/corpora/reuters/training/5204
inflating: /root/nltk_data/corpora/reuters/training/5205
inflating: /root/nltk_data/corpora/reuters/training/5206
inflating: /root/nltk_data/corpora/reuters/training/5207
inflating: /root/nltk_data/corpora/reuters/training/5209
inflating: /root/nltk_data/corpora/reuters/training/521
inflating: /root/nltk_data/corpora/reuters/training/5210
inflating: /root/nltk_data/corpora/reuters/training/5212
inflating: /root/nltk_data/corpora/reuters/training/5214
inflating: /root/nltk_data/corpora/reuters/training/5215
inflating: /root/nltk_data/corpora/reuters/training/5216
inflating: /root/nltk_data/corpora/reuters/training/5218
inflating: /root/nltk_data/corpora/reuters/training/5219
inflating: /root/nltk_data/corpora/reuters/training/522
inflating: /root/nltk_data/corpora/reuters/training/5220
inflating: /root/nltk_data/corpora/reuters/training/5222
inflating: /root/nltk_data/corpora/reuters/training/5223
inflating: /root/nltk_data/corpora/reuters/training/5226
inflating: /root/nltk_data/corpora/reuters/training/5228
inflating: /root/nltk_data/corpora/reuters/training/5231
inflating: /root/nltk_data/corpora/reuters/training/5232
inflating: /root/nltk_data/corpora/reuters/training/5234
inflating: /root/nltk_data/corpora/reuters/training/5235
inflating: /root/nltk_data/corpora/reuters/training/5236
inflating: /root/nltk_data/corpora/reuters/training/5237
inflating: /root/nltk_data/corpora/reuters/training/5238
inflating: /root/nltk_data/corpora/reuters/training/524
inflating: /root/nltk_data/corpora/reuters/training/5240
inflating: /root/nltk_data/corpora/reuters/training/5241
inflating: /root/nltk_data/corpora/reuters/training/5243
inflating: /root/nltk_data/corpora/reuters/training/5244
inflating: /root/nltk_data/corpora/reuters/training/5245
inflating: /root/nltk_data/corpora/reuters/training/525
inflating: /root/nltk_data/corpora/reuters/training/5250
inflating: /root/nltk_data/corpora/reuters/training/5251
inflating: /root/nltk_data/corpora/reuters/training/5253
inflating: /root/nltk_data/corpora/reuters/training/5254
inflating: /root/nltk_data/corpora/reuters/training/5255
inflating: /root/nltk_data/corpora/reuters/training/5256
inflating: /root/nltk_data/corpora/reuters/training/5257
inflating: /root/nltk_data/corpora/reuters/training/5258
inflating: /root/nltk_data/corpora/reuters/training/5260
inflating: /root/nltk_data/corpora/reuters/training/5264
inflating: /root/nltk_data/corpora/reuters/training/5266
inflating: /root/nltk_data/corpora/reuters/training/5268
inflating: /root/nltk_data/corpora/reuters/training/5269
inflating: /root/nltk_data/corpora/reuters/training/5270
inflating: /root/nltk_data/corpora/reuters/training/5271
inflating: /root/nltk_data/corpora/reuters/training/5272
inflating: /root/nltk_data/corpora/reuters/training/5273
inflating: /root/nltk_data/corpora/reuters/training/5274
inflating: /root/nltk_data/corpora/reuters/training/5276
inflating: /root/nltk_data/corpora/reuters/training/5277
inflating: /root/nltk_data/corpora/reuters/training/5278
inflating: /root/nltk_data/corpora/reuters/training/5279
inflating: /root/nltk_data/corpora/reuters/training/528
inflating: /root/nltk_data/corpora/reuters/training/5280
inflating: /root/nltk_data/corpora/reuters/training/5281
inflating: /root/nltk_data/corpora/reuters/training/5283
inflating: /root/nltk_data/corpora/reuters/training/5284
inflating: /root/nltk_data/corpora/reuters/training/5285
inflating: /root/nltk_data/corpora/reuters/training/5286
inflating: /root/nltk_data/corpora/reuters/training/5287
inflating: /root/nltk_data/corpora/reuters/training/5288
inflating: /root/nltk_data/corpora/reuters/training/529
inflating: /root/nltk_data/corpora/reuters/training/5290
inflating: /root/nltk_data/corpora/reuters/training/5291
inflating: /root/nltk_data/corpora/reuters/training/5292
inflating: /root/nltk_data/corpora/reuters/training/5293
inflating: /root/nltk_data/corpora/reuters/training/5295
inflating: /root/nltk_data/corpora/reuters/training/5296
inflating: /root/nltk_data/corpora/reuters/training/5297
inflating: /root/nltk_data/corpora/reuters/training/5299
inflating: /root/nltk_data/corpora/reuters/training/53
inflating: /root/nltk_data/corpora/reuters/training/530
inflating: /root/nltk_data/corpora/reuters/training/5300
inflating: /root/nltk_data/corpora/reuters/training/5303
inflating: /root/nltk_data/corpora/reuters/training/5304
inflating: /root/nltk_data/corpora/reuters/training/5308
inflating: /root/nltk_data/corpora/reuters/training/5309
inflating: /root/nltk_data/corpora/reuters/training/531
inflating: /root/nltk_data/corpora/reuters/training/5312
inflating: /root/nltk_data/corpora/reuters/training/5313
inflating: /root/nltk_data/corpora/reuters/training/5314
inflating: /root/nltk_data/corpora/reuters/training/5315
inflating: /root/nltk_data/corpora/reuters/training/5316
inflating: /root/nltk_data/corpora/reuters/training/5317
inflating: /root/nltk_data/corpora/reuters/training/5318
inflating: /root/nltk_data/corpora/reuters/training/532
inflating: /root/nltk_data/corpora/reuters/training/5320
inflating: /root/nltk_data/corpora/reuters/training/5321
inflating: /root/nltk_data/corpora/reuters/training/5323
inflating: /root/nltk_data/corpora/reuters/training/5326
inflating: /root/nltk_data/corpora/reuters/training/5327
inflating: /root/nltk_data/corpora/reuters/training/5328
inflating: /root/nltk_data/corpora/reuters/training/5329
inflating: /root/nltk_data/corpora/reuters/training/533
inflating: /root/nltk_data/corpora/reuters/training/5330
inflating: /root/nltk_data/corpora/reuters/training/5331
```

```
inflating: /root/nltk_data/corpora/reuters/training/5332
inflating: /root/nltk_data/corpora/reuters/training/5333
inflating: /root/nltk_data/corpora/reuters/training/5334
inflating: /root/nltk_data/corpora/reuters/training/5335
inflating: /root/nltk_data/corpora/reuters/training/5336
inflating: /root/nltk_data/corpora/reuters/training/5338
inflating: /root/nltk_data/corpora/reuters/training/5339
inflating: /root/nltk_data/corpora/reuters/training/534
inflating: /root/nltk_data/corpora/reuters/training/5340
inflating: /root/nltk_data/corpora/reuters/training/5342
inflating: /root/nltk_data/corpora/reuters/training/5344
inflating: /root/nltk_data/corpora/reuters/training/5345
inflating: /root/nltk_data/corpora/reuters/training/5346
inflating: /root/nltk_data/corpora/reuters/training/5349
inflating: /root/nltk_data/corpora/reuters/training/535
inflating: /root/nltk_data/corpora/reuters/training/5350
inflating: /root/nltk_data/corpora/reuters/training/5351
inflating: /root/nltk_data/corpora/reuters/training/5352
inflating: /root/nltk_data/corpora/reuters/training/5354
inflating: /root/nltk_data/corpora/reuters/training/5355
inflating: /root/nltk_data/corpora/reuters/training/5356
inflating: /root/nltk_data/corpora/reuters/training/5357
inflating: /root/nltk_data/corpora/reuters/training/5358
inflating: /root/nltk_data/corpora/reuters/training/5359
inflating: /root/nltk_data/corpora/reuters/training/536
inflating: /root/nltk_data/corpora/reuters/training/5360
inflating: /root/nltk_data/corpora/reuters/training/5362
inflating: /root/nltk_data/corpora/reuters/training/5363
inflating: /root/nltk_data/corpora/reuters/training/5366
inflating: /root/nltk_data/corpora/reuters/training/5367
inflating: /root/nltk_data/corpora/reuters/training/5371
inflating: /root/nltk_data/corpora/reuters/training/5373
inflating: /root/nltk_data/corpora/reuters/training/5375
inflating: /root/nltk_data/corpora/reuters/training/5376
inflating: /root/nltk_data/corpora/reuters/training/5377
inflating: /root/nltk_data/corpora/reuters/training/5378
inflating: /root/nltk_data/corpora/reuters/training/5379
inflating: /root/nltk_data/corpora/reuters/training/538
inflating: /root/nltk_data/corpora/reuters/training/5382
inflating: /root/nltk_data/corpora/reuters/training/5383
inflating: /root/nltk_data/corpora/reuters/training/5385
inflating: /root/nltk_data/corpora/reuters/training/5388
inflating: /root/nltk_data/corpora/reuters/training/5389
inflating: /root/nltk_data/corpora/reuters/training/539
inflating: /root/nltk_data/corpora/reuters/training/5390
inflating: /root/nltk_data/corpora/reuters/training/5391
inflating: /root/nltk_data/corpora/reuters/training/5392
inflating: /root/nltk_data/corpora/reuters/training/5394
inflating: /root/nltk_data/corpora/reuters/training/5396
inflating: /root/nltk_data/corpora/reuters/training/5398
inflating: /root/nltk_data/corpora/reuters/training/5400
inflating: /root/nltk_data/corpora/reuters/training/5402
inflating: /root/nltk_data/corpora/reuters/training/5404
inflating: /root/nltk_data/corpora/reuters/training/5407
inflating: /root/nltk_data/corpora/reuters/training/5408
inflating: /root/nltk_data/corpora/reuters/training/5409
inflating: /root/nltk_data/corpora/reuters/training/541
inflating: /root/nltk_data/corpora/reuters/training/5410
inflating: /root/nltk_data/corpora/reuters/training/5411
inflating: /root/nltk_data/corpora/reuters/training/5412
inflating: /root/nltk_data/corpora/reuters/training/5414
inflating: /root/nltk_data/corpora/reuters/training/5415
inflating: /root/nltk_data/corpora/reuters/training/5416
inflating: /root/nltk_data/corpora/reuters/training/5417
inflating: /root/nltk_data/corpora/reuters/training/5421
inflating: /root/nltk_data/corpora/reuters/training/5422
inflating: /root/nltk_data/corpora/reuters/training/5423
inflating: /root/nltk_data/corpora/reuters/training/5426
inflating: /root/nltk_data/corpora/reuters/training/5429
inflating: /root/nltk_data/corpora/reuters/training/543
inflating: /root/nltk_data/corpora/reuters/training/5431
inflating: /root/nltk_data/corpora/reuters/training/5434
inflating: /root/nltk_data/corpora/reuters/training/5435
inflating: /root/nltk_data/corpora/reuters/training/5437
inflating: /root/nltk_data/corpora/reuters/training/5438
inflating: /root/nltk_data/corpora/reuters/training/5439
inflating: /root/nltk_data/corpora/reuters/training/544
inflating: /root/nltk_data/corpora/reuters/training/5440
inflating: /root/nltk_data/corpora/reuters/training/5441
inflating: /root/nltk_data/corpora/reuters/training/5442
inflating: /root/nltk_data/corpora/reuters/training/5443
inflating: /root/nltk_data/corpora/reuters/training/5444
inflating: /root/nltk_data/corpora/reuters/training/5445
inflating: /root/nltk_data/corpora/reuters/training/5447
inflating: /root/nltk_data/corpora/reuters/training/545
inflating: /root/nltk_data/corpora/reuters/training/5451
inflating: /root/nltk_data/corpora/reuters/training/5452
inflating: /root/nltk_data/corpora/reuters/training/5453
inflating: /root/nltk_data/corpora/reuters/training/5454
inflating: /root/nltk_data/corpora/reuters/training/5455
inflating: /root/nltk_data/corpora/reuters/training/5456
inflating: /root/nltk_data/corpora/reuters/training/5458
inflating: /root/nltk_data/corpora/reuters/training/546
inflating: /root/nltk_data/corpora/reuters/training/5460
inflating: /root/nltk_data/corpora/reuters/training/5461
inflating: /root/nltk_data/corpora/reuters/training/5464
inflating: /root/nltk_data/corpora/reuters/training/5465
inflating: /root/nltk_data/corpora/reuters/training/5467
inflating: /root/nltk_data/corpora/reuters/training/5469
inflating: /root/nltk_data/corpora/reuters/training/547
inflating: /root/nltk_data/corpora/reuters/training/5470
inflating: /root/nltk_data/corpora/reuters/training/5471
inflating: /root/nltk_data/corpora/reuters/training/5472
inflating: /root/nltk_data/corpora/reuters/training/5473
inflating: /root/nltk_data/corpora/reuters/training/5474
inflating: /root/nltk_data/corpora/reuters/training/5475
inflating: /root/nltk_data/corpora/reuters/training/5476
inflating: /root/nltk_data/corpora/reuters/training/5477
inflating: /root/nltk_data/corpora/reuters/training/548
inflating: /root/nltk_data/corpora/reuters/training/5481
inflating: /root/nltk_data/corpora/reuters/training/5483
inflating: /root/nltk_data/corpora/reuters/training/5485
inflating: /root/nltk_data/corpora/reuters/training/5487
inflating: /root/nltk_data/corpora/reuters/training/5490
inflating: /root/nltk_data/corpora/reuters/training/5491
inflating: /root/nltk_data/corpora/reuters/training/5492
inflating: /root/nltk_data/corpora/reuters/training/5494
inflating: /root/nltk_data/corpora/reuters/training/5495
inflating: /root/nltk_data/corpora/reuters/training/5496
inflating: /root/nltk_data/corpora/reuters/training/5498
inflating: /root/nltk_data/corpora/reuters/training/5499
inflating: /root/nltk_data/corpora/reuters/training/550
inflating: /root/nltk_data/corpora/reuters/training/5500
inflating: /root/nltk_data/corpora/reuters/training/5501
inflating: /root/nltk_data/corpora/reuters/training/5502
inflating: /root/nltk_data/corpora/reuters/training/5505
inflating: /root/nltk_data/corpora/reuters/training/5506
inflating: /root/nltk_data/corpora/reuters/training/5507
inflating: /root/nltk_data/corpora/reuters/training/5508
inflating: /root/nltk_data/corpora/reuters/training/551
inflating: /root/nltk_data/corpora/reuters/training/5511
inflating: /root/nltk_data/corpora/reuters/training/5516
inflating: /root/nltk_data/corpora/reuters/training/5517
inflating: /root/nltk_data/corpora/reuters/training/5518
```

```
inflating: /root/nltk_data/corpora/reuters/training/552
inflating: /root/nltk_data/corpora/reuters/training/5520
inflating: /root/nltk_data/corpora/reuters/training/5522
inflating: /root/nltk_data/corpora/reuters/training/5523
inflating: /root/nltk_data/corpora/reuters/training/5525
inflating: /root/nltk_data/corpora/reuters/training/5526
inflating: /root/nltk_data/corpora/reuters/training/5528
inflating: /root/nltk_data/corpora/reuters/training/553
inflating: /root/nltk_data/corpora/reuters/training/5530
inflating: /root/nltk_data/corpora/reuters/training/5531
inflating: /root/nltk_data/corpora/reuters/training/5532
inflating: /root/nltk_data/corpora/reuters/training/5533
inflating: /root/nltk_data/corpora/reuters/training/5535
inflating: /root/nltk_data/corpora/reuters/training/5537
inflating: /root/nltk_data/corpora/reuters/training/5538
inflating: /root/nltk_data/corpora/reuters/training/554
inflating: /root/nltk_data/corpora/reuters/training/5540
inflating: /root/nltk_data/corpora/reuters/training/5541
inflating: /root/nltk_data/corpora/reuters/training/5542
inflating: /root/nltk_data/corpora/reuters/training/5543
inflating: /root/nltk_data/corpora/reuters/training/5544
inflating: /root/nltk_data/corpora/reuters/training/5545
inflating: /root/nltk_data/corpora/reuters/training/5546
inflating: /root/nltk_data/corpora/reuters/training/5549
inflating: /root/nltk_data/corpora/reuters/training/555
inflating: /root/nltk_data/corpora/reuters/training/5553
inflating: /root/nltk_data/corpora/reuters/training/5554
inflating: /root/nltk_data/corpora/reuters/training/5555
inflating: /root/nltk_data/corpora/reuters/training/5556
inflating: /root/nltk_data/corpora/reuters/training/5558
inflating: /root/nltk_data/corpora/reuters/training/5559
inflating: /root/nltk_data/corpora/reuters/training/5561
inflating: /root/nltk_data/corpora/reuters/training/5563
inflating: /root/nltk_data/corpora/reuters/training/5564
inflating: /root/nltk_data/corpora/reuters/training/5565
inflating: /root/nltk_data/corpora/reuters/training/5566
inflating: /root/nltk_data/corpora/reuters/training/5567
inflating: /root/nltk_data/corpora/reuters/training/5568
inflating: /root/nltk_data/corpora/reuters/training/557
inflating: /root/nltk_data/corpora/reuters/training/5570
inflating: /root/nltk_data/corpora/reuters/training/5571
inflating: /root/nltk_data/corpora/reuters/training/5572
inflating: /root/nltk_data/corpora/reuters/training/5573
inflating: /root/nltk_data/corpora/reuters/training/5575
inflating: /root/nltk_data/corpora/reuters/training/5576
inflating: /root/nltk_data/corpora/reuters/training/5579
inflating: /root/nltk_data/corpora/reuters/training/558
inflating: /root/nltk_data/corpora/reuters/training/5582
inflating: /root/nltk_data/corpora/reuters/training/5583
inflating: /root/nltk_data/corpora/reuters/training/5585
inflating: /root/nltk_data/corpora/reuters/training/5586
inflating: /root/nltk_data/corpora/reuters/training/559
inflating: /root/nltk_data/corpora/reuters/training/5592
inflating: /root/nltk_data/corpora/reuters/training/5593
inflating: /root/nltk_data/corpora/reuters/training/5594
inflating: /root/nltk_data/corpora/reuters/training/5595
inflating: /root/nltk_data/corpora/reuters/training/5596
inflating: /root/nltk_data/corpora/reuters/training/5597
inflating: /root/nltk_data/corpora/reuters/training/5598
inflating: /root/nltk_data/corpora/reuters/training/56
inflating: /root/nltk_data/corpora/reuters/training/560
inflating: /root/nltk_data/corpora/reuters/training/5600
inflating: /root/nltk_data/corpora/reuters/training/5602
inflating: /root/nltk_data/corpora/reuters/training/5603
inflating: /root/nltk_data/corpora/reuters/training/5604
inflating: /root/nltk_data/corpora/reuters/training/5606
inflating: /root/nltk_data/corpora/reuters/training/5608
inflating: /root/nltk_data/corpora/reuters/training/5609
inflating: /root/nltk_data/corpora/reuters/training/5610
inflating: /root/nltk_data/corpora/reuters/training/5611
inflating: /root/nltk_data/corpora/reuters/training/5615
inflating: /root/nltk_data/corpora/reuters/training/5617
inflating: /root/nltk_data/corpora/reuters/training/562
inflating: /root/nltk_data/corpora/reuters/training/5620
inflating: /root/nltk_data/corpora/reuters/training/5624
inflating: /root/nltk_data/corpora/reuters/training/5627
inflating: /root/nltk_data/corpora/reuters/training/5629
inflating: /root/nltk_data/corpora/reuters/training/5630
inflating: /root/nltk_data/corpora/reuters/training/5632
inflating: /root/nltk_data/corpora/reuters/training/5634
inflating: /root/nltk_data/corpora/reuters/training/5635
inflating: /root/nltk_data/corpora/reuters/training/5636
inflating: /root/nltk_data/corpora/reuters/training/5637
inflating: /root/nltk_data/corpora/reuters/training/5639
inflating: /root/nltk_data/corpora/reuters/training/5640
inflating: /root/nltk_data/corpora/reuters/training/5641
inflating: /root/nltk_data/corpora/reuters/training/5643
inflating: /root/nltk_data/corpora/reuters/training/5644
inflating: /root/nltk_data/corpora/reuters/training/5647
inflating: /root/nltk_data/corpora/reuters/training/5648
inflating: /root/nltk_data/corpora/reuters/training/5650
inflating: /root/nltk_data/corpora/reuters/training/5652
inflating: /root/nltk_data/corpora/reuters/training/5654
inflating: /root/nltk_data/corpora/reuters/training/5655
inflating: /root/nltk_data/corpora/reuters/training/5657
inflating: /root/nltk_data/corpora/reuters/training/5658
inflating: /root/nltk_data/corpora/reuters/training/5659
inflating: /root/nltk_data/corpora/reuters/training/5660
inflating: /root/nltk_data/corpora/reuters/training/5661
inflating: /root/nltk_data/corpora/reuters/training/5662
inflating: /root/nltk_data/corpora/reuters/training/5663
inflating: /root/nltk_data/corpora/reuters/training/5664
inflating: /root/nltk_data/corpora/reuters/training/5665
inflating: /root/nltk_data/corpora/reuters/training/5667
inflating: /root/nltk_data/corpora/reuters/training/5669
inflating: /root/nltk_data/corpora/reuters/training/5670
inflating: /root/nltk_data/corpora/reuters/training/5671
inflating: /root/nltk_data/corpora/reuters/training/5672
inflating: /root/nltk_data/corpora/reuters/training/5675
inflating: /root/nltk_data/corpora/reuters/training/5677
inflating: /root/nltk_data/corpora/reuters/training/5678
inflating: /root/nltk_data/corpora/reuters/training/5679
inflating: /root/nltk_data/corpora/reuters/training/5680
inflating: /root/nltk_data/corpora/reuters/training/5682
inflating: /root/nltk_data/corpora/reuters/training/5683
inflating: /root/nltk_data/corpora/reuters/training/5684
inflating: /root/nltk_data/corpora/reuters/training/5686
inflating: /root/nltk_data/corpora/reuters/training/5687
inflating: /root/nltk_data/corpora/reuters/training/5688
inflating: /root/nltk_data/corpora/reuters/training/5689
inflating: /root/nltk_data/corpora/reuters/training/5690
inflating: /root/nltk_data/corpora/reuters/training/5692
inflating: /root/nltk_data/corpora/reuters/training/5693
inflating: /root/nltk_data/corpora/reuters/training/5694
inflating: /root/nltk_data/corpora/reuters/training/57
inflating: /root/nltk_data/corpora/reuters/training/5701
inflating: /root/nltk_data/corpora/reuters/training/5702
inflating: /root/nltk_data/corpora/reuters/training/5706
inflating: /root/nltk_data/corpora/reuters/training/5707
inflating: /root/nltk_data/corpora/reuters/training/5708
inflating: /root/nltk_data/corpora/reuters/training/5710
inflating: /root/nltk_data/corpora/reuters/training/5711
inflating: /root/nltk_data/corpora/reuters/training/5712
inflating: /root/nltk_data/corpora/reuters/training/5713
inflating: /root/nltk_data/corpora/reuters/training/5715
```

inflating: /root/nltk_data/corpora/reuters/training/5716
inflating: /root/nltk_data/corpora/reuters/training/5717
inflating: /root/nltk_data/corpora/reuters/training/5719
inflating: /root/nltk_data/corpora/reuters/training/5720
inflating: /root/nltk_data/corpora/reuters/training/5721
inflating: /root/nltk_data/corpora/reuters/training/5722
inflating: /root/nltk_data/corpora/reuters/training/5723
inflating: /root/nltk_data/corpora/reuters/training/5724
inflating: /root/nltk_data/corpora/reuters/training/5725
inflating: /root/nltk_data/corpora/reuters/training/5727
inflating: /root/nltk_data/corpora/reuters/training/5729
inflating: /root/nltk_data/corpora/reuters/training/5731
inflating: /root/nltk_data/corpora/reuters/training/5732
inflating: /root/nltk_data/corpora/reuters/training/5734
inflating: /root/nltk_data/corpora/reuters/training/5737
inflating: /root/nltk_data/corpora/reuters/training/5738
inflating: /root/nltk_data/corpora/reuters/training/5739
inflating: /root/nltk_data/corpora/reuters/training/5740
inflating: /root/nltk_data/corpora/reuters/training/5741
inflating: /root/nltk_data/corpora/reuters/training/5742
inflating: /root/nltk_data/corpora/reuters/training/5743
inflating: /root/nltk_data/corpora/reuters/training/5744
inflating: /root/nltk_data/corpora/reuters/training/5746
inflating: /root/nltk_data/corpora/reuters/training/5749
inflating: /root/nltk_data/corpora/reuters/training/5750
inflating: /root/nltk_data/corpora/reuters/training/5751
inflating: /root/nltk_data/corpora/reuters/training/5752
inflating: /root/nltk_data/corpora/reuters/training/5754
inflating: /root/nltk_data/corpora/reuters/training/5755
inflating: /root/nltk_data/corpora/reuters/training/5756
inflating: /root/nltk_data/corpora/reuters/training/5757
inflating: /root/nltk_data/corpora/reuters/training/5758
inflating: /root/nltk_data/corpora/reuters/training/5761
inflating: /root/nltk_data/corpora/reuters/training/5762
inflating: /root/nltk_data/corpora/reuters/training/5764
inflating: /root/nltk_data/corpora/reuters/training/5765
inflating: /root/nltk_data/corpora/reuters/training/5767
inflating: /root/nltk_data/corpora/reuters/training/5768
inflating: /root/nltk_data/corpora/reuters/training/5769
inflating: /root/nltk_data/corpora/reuters/training/5778
inflating: /root/nltk_data/corpora/reuters/training/5779
inflating: /root/nltk_data/corpora/reuters/training/5780
inflating: /root/nltk_data/corpora/reuters/training/5782
inflating: /root/nltk_data/corpora/reuters/training/5783
inflating: /root/nltk_data/corpora/reuters/training/5785
inflating: /root/nltk_data/corpora/reuters/training/5786
inflating: /root/nltk_data/corpora/reuters/training/5787
inflating: /root/nltk_data/corpora/reuters/training/5788
inflating: /root/nltk_data/corpora/reuters/training/5789
inflating: /root/nltk_data/corpora/reuters/training/5791
inflating: /root/nltk_data/corpora/reuters/training/5792
inflating: /root/nltk_data/corpora/reuters/training/5793
inflating: /root/nltk_data/corpora/reuters/training/5796
inflating: /root/nltk_data/corpora/reuters/training/5797
inflating: /root/nltk_data/corpora/reuters/training/5798
inflating: /root/nltk_data/corpora/reuters/training/58
inflating: /root/nltk_data/corpora/reuters/training/5800
inflating: /root/nltk_data/corpora/reuters/training/5803
inflating: /root/nltk_data/corpora/reuters/training/5804
inflating: /root/nltk_data/corpora/reuters/training/5805
inflating: /root/nltk_data/corpora/reuters/training/5807
inflating: /root/nltk_data/corpora/reuters/training/5808
inflating: /root/nltk_data/corpora/reuters/training/5810
inflating: /root/nltk_data/corpora/reuters/training/5811
inflating: /root/nltk_data/corpora/reuters/training/5812
inflating: /root/nltk_data/corpora/reuters/training/5814
inflating: /root/nltk_data/corpora/reuters/training/5815
inflating: /root/nltk_data/corpora/reuters/training/5818
inflating: /root/nltk_data/corpora/reuters/training/5819
inflating: /root/nltk_data/corpora/reuters/training/5822
inflating: /root/nltk_data/corpora/reuters/training/5826
inflating: /root/nltk_data/corpora/reuters/training/5827
inflating: /root/nltk_data/corpora/reuters/training/5830
inflating: /root/nltk_data/corpora/reuters/training/5833
inflating: /root/nltk_data/corpora/reuters/training/5835
inflating: /root/nltk_data/corpora/reuters/training/5836
inflating: /root/nltk_data/corpora/reuters/training/5838
inflating: /root/nltk_data/corpora/reuters/training/5839
inflating: /root/nltk_data/corpora/reuters/training/5841
inflating: /root/nltk_data/corpora/reuters/training/5842
inflating: /root/nltk_data/corpora/reuters/training/5845
inflating: /root/nltk_data/corpora/reuters/training/5847
inflating: /root/nltk_data/corpora/reuters/training/5848
inflating: /root/nltk_data/corpora/reuters/training/5849
inflating: /root/nltk_data/corpora/reuters/training/5850
inflating: /root/nltk_data/corpora/reuters/training/5852
inflating: /root/nltk_data/corpora/reuters/training/5854
inflating: /root/nltk_data/corpora/reuters/training/5858
inflating: /root/nltk_data/corpora/reuters/training/5861
inflating: /root/nltk_data/corpora/reuters/training/5866
inflating: /root/nltk_data/corpora/reuters/training/5868
inflating: /root/nltk_data/corpora/reuters/training/5870
inflating: /root/nltk_data/corpora/reuters/training/5873
inflating: /root/nltk_data/corpora/reuters/training/5875
inflating: /root/nltk_data/corpora/reuters/training/5876
inflating: /root/nltk_data/corpora/reuters/training/5879
inflating: /root/nltk_data/corpora/reuters/training/5880
inflating: /root/nltk_data/corpora/reuters/training/5883
inflating: /root/nltk_data/corpora/reuters/training/5886
inflating: /root/nltk_data/corpora/reuters/training/5887
inflating: /root/nltk_data/corpora/reuters/training/5888
inflating: /root/nltk_data/corpora/reuters/training/5889
inflating: /root/nltk_data/corpora/reuters/training/5894
inflating: /root/nltk_data/corpora/reuters/training/5895
inflating: /root/nltk_data/corpora/reuters/training/5896
inflating: /root/nltk_data/corpora/reuters/training/5898
inflating: /root/nltk_data/corpora/reuters/training/5899
inflating: /root/nltk_data/corpora/reuters/training/59
inflating: /root/nltk_data/corpora/reuters/training/5900
inflating: /root/nltk_data/corpora/reuters/training/5901
inflating: /root/nltk_data/corpora/reuters/training/5902
inflating: /root/nltk_data/corpora/reuters/training/5905
inflating: /root/nltk_data/corpora/reuters/training/5906
inflating: /root/nltk_data/corpora/reuters/training/5908
inflating: /root/nltk_data/corpora/reuters/training/5909
inflating: /root/nltk_data/corpora/reuters/training/5910
inflating: /root/nltk_data/corpora/reuters/training/5911
inflating: /root/nltk_data/corpora/reuters/training/5912
inflating: /root/nltk_data/corpora/reuters/training/5913
inflating: /root/nltk_data/corpora/reuters/training/5914
inflating: /root/nltk_data/corpora/reuters/training/5915
inflating: /root/nltk_data/corpora/reuters/training/5919
inflating: /root/nltk_data/corpora/reuters/training/5922
inflating: /root/nltk_data/corpora/reuters/training/5925
inflating: /root/nltk_data/corpora/reuters/training/5926
inflating: /root/nltk_data/corpora/reuters/training/5930
inflating: /root/nltk_data/corpora/reuters/training/5931
inflating: /root/nltk_data/corpora/reuters/training/5932
inflating: /root/nltk_data/corpora/reuters/training/5933
inflating: /root/nltk_data/corpora/reuters/training/5934
inflating: /root/nltk_data/corpora/reuters/training/5935
inflating: /root/nltk_data/corpora/reuters/training/5937
inflating: /root/nltk_data/corpora/reuters/training/5938
inflating: /root/nltk_data/corpora/reuters/training/5939

```
inflating: /root/nltk_data/corpora/reuters/training/5941
inflating: /root/nltk_data/corpora/reuters/training/5945
inflating: /root/nltk_data/corpora/reuters/training/5946
inflating: /root/nltk_data/corpora/reuters/training/5947
inflating: /root/nltk_data/corpora/reuters/training/5949
inflating: /root/nltk_data/corpora/reuters/training/5950
inflating: /root/nltk_data/corpora/reuters/training/5951
inflating: /root/nltk_data/corpora/reuters/training/5954
inflating: /root/nltk_data/corpora/reuters/training/5955
inflating: /root/nltk_data/corpora/reuters/training/5957
inflating: /root/nltk_data/corpora/reuters/training/5958
inflating: /root/nltk_data/corpora/reuters/training/5960
inflating: /root/nltk_data/corpora/reuters/training/5964
inflating: /root/nltk_data/corpora/reuters/training/5965
inflating: /root/nltk_data/corpora/reuters/training/5967
inflating: /root/nltk_data/corpora/reuters/training/5968
inflating: /root/nltk_data/corpora/reuters/training/5969
inflating: /root/nltk_data/corpora/reuters/training/5972
inflating: /root/nltk_data/corpora/reuters/training/5973
inflating: /root/nltk_data/corpora/reuters/training/5974
inflating: /root/nltk_data/corpora/reuters/training/5975
inflating: /root/nltk_data/corpora/reuters/training/5976
inflating: /root/nltk_data/corpora/reuters/training/5977
inflating: /root/nltk_data/corpora/reuters/training/5980
inflating: /root/nltk_data/corpora/reuters/training/5981
inflating: /root/nltk_data/corpora/reuters/training/5982
inflating: /root/nltk_data/corpora/reuters/training/5984
inflating: /root/nltk_data/corpora/reuters/training/5985
inflating: /root/nltk_data/corpora/reuters/training/5986
inflating: /root/nltk_data/corpora/reuters/training/5991
inflating: /root/nltk_data/corpora/reuters/training/5993
inflating: /root/nltk_data/corpora/reuters/training/5994
inflating: /root/nltk_data/corpora/reuters/training/5998
inflating: /root/nltk_data/corpora/reuters/training/6
inflating: /root/nltk_data/corpora/reuters/training/6000
inflating: /root/nltk_data/corpora/reuters/training/6003
inflating: /root/nltk_data/corpora/reuters/training/6005
inflating: /root/nltk_data/corpora/reuters/training/6006
inflating: /root/nltk_data/corpora/reuters/training/6010
inflating: /root/nltk_data/corpora/reuters/training/6013
inflating: /root/nltk_data/corpora/reuters/training/6014
inflating: /root/nltk_data/corpora/reuters/training/6017
inflating: /root/nltk_data/corpora/reuters/training/6018
inflating: /root/nltk_data/corpora/reuters/training/6020
inflating: /root/nltk_data/corpora/reuters/training/6023
inflating: /root/nltk_data/corpora/reuters/training/6025
inflating: /root/nltk_data/corpora/reuters/training/6026
inflating: /root/nltk_data/corpora/reuters/training/6027
inflating: /root/nltk_data/corpora/reuters/training/6030
inflating: /root/nltk_data/corpora/reuters/training/6033
inflating: /root/nltk_data/corpora/reuters/training/6034
inflating: /root/nltk_data/corpora/reuters/training/6035
inflating: /root/nltk_data/corpora/reuters/training/6037
inflating: /root/nltk_data/corpora/reuters/training/6041
inflating: /root/nltk_data/corpora/reuters/training/6042
inflating: /root/nltk_data/corpora/reuters/training/6043
inflating: /root/nltk_data/corpora/reuters/training/6044
inflating: /root/nltk_data/corpora/reuters/training/6046
inflating: /root/nltk_data/corpora/reuters/training/6047
inflating: /root/nltk_data/corpora/reuters/training/6048
inflating: /root/nltk_data/corpora/reuters/training/6051
inflating: /root/nltk_data/corpora/reuters/training/6054
inflating: /root/nltk_data/corpora/reuters/training/6055
inflating: /root/nltk_data/corpora/reuters/training/6056
inflating: /root/nltk_data/corpora/reuters/training/6058
inflating: /root/nltk_data/corpora/reuters/training/6059
inflating: /root/nltk_data/corpora/reuters/training/6060
inflating: /root/nltk_data/corpora/reuters/training/6062
inflating: /root/nltk_data/corpora/reuters/training/6063
inflating: /root/nltk_data/corpora/reuters/training/6065
inflating: /root/nltk_data/corpora/reuters/training/6067
inflating: /root/nltk_data/corpora/reuters/training/6068
inflating: /root/nltk_data/corpora/reuters/training/6069
inflating: /root/nltk_data/corpora/reuters/training/6070
inflating: /root/nltk_data/corpora/reuters/training/6073
inflating: /root/nltk_data/corpora/reuters/training/6075
inflating: /root/nltk_data/corpora/reuters/training/6078
inflating: /root/nltk_data/corpora/reuters/training/6079
inflating: /root/nltk_data/corpora/reuters/training/6081
inflating: /root/nltk_data/corpora/reuters/training/6082
inflating: /root/nltk_data/corpora/reuters/training/6083
inflating: /root/nltk_data/corpora/reuters/training/6084
inflating: /root/nltk_data/corpora/reuters/training/6085
inflating: /root/nltk_data/corpora/reuters/training/6086
inflating: /root/nltk_data/corpora/reuters/training/6088
inflating: /root/nltk_data/corpora/reuters/training/6090
inflating: /root/nltk_data/corpora/reuters/training/6091
inflating: /root/nltk_data/corpora/reuters/training/6092
inflating: /root/nltk_data/corpora/reuters/training/6093
inflating: /root/nltk_data/corpora/reuters/training/6094
inflating: /root/nltk_data/corpora/reuters/training/6096
inflating: /root/nltk_data/corpora/reuters/training/6101
inflating: /root/nltk_data/corpora/reuters/training/6102
inflating: /root/nltk_data/corpora/reuters/training/6103
inflating: /root/nltk_data/corpora/reuters/training/6104
inflating: /root/nltk_data/corpora/reuters/training/6105
inflating: /root/nltk_data/corpora/reuters/training/6106
inflating: /root/nltk_data/corpora/reuters/training/6107
inflating: /root/nltk_data/corpora/reuters/training/6108
inflating: /root/nltk_data/corpora/reuters/training/6111
inflating: /root/nltk_data/corpora/reuters/training/6113
inflating: /root/nltk_data/corpora/reuters/training/6114
inflating: /root/nltk_data/corpora/reuters/training/6115
inflating: /root/nltk_data/corpora/reuters/training/6116
inflating: /root/nltk_data/corpora/reuters/training/6119
inflating: /root/nltk_data/corpora/reuters/training/6120
inflating: /root/nltk_data/corpora/reuters/training/6121
inflating: /root/nltk_data/corpora/reuters/training/6123
inflating: /root/nltk_data/corpora/reuters/training/6124
inflating: /root/nltk_data/corpora/reuters/training/6125
inflating: /root/nltk_data/corpora/reuters/training/6127
inflating: /root/nltk_data/corpora/reuters/training/6128
inflating: /root/nltk_data/corpora/reuters/training/6129
inflating: /root/nltk_data/corpora/reuters/training/6130
inflating: /root/nltk_data/corpora/reuters/training/6134
inflating: /root/nltk_data/corpora/reuters/training/6135
inflating: /root/nltk_data/corpora/reuters/training/6137
inflating: /root/nltk_data/corpora/reuters/training/6138
inflating: /root/nltk_data/corpora/reuters/training/6142
inflating: /root/nltk_data/corpora/reuters/training/6146
inflating: /root/nltk_data/corpora/reuters/training/6147
inflating: /root/nltk_data/corpora/reuters/training/6149
inflating: /root/nltk_data/corpora/reuters/training/6153
inflating: /root/nltk_data/corpora/reuters/training/6154
inflating: /root/nltk_data/corpora/reuters/training/6155
inflating: /root/nltk_data/corpora/reuters/training/6156
inflating: /root/nltk_data/corpora/reuters/training/6157
inflating: /root/nltk_data/corpora/reuters/training/6158
inflating: /root/nltk_data/corpora/reuters/training/6159
inflating: /root/nltk_data/corpora/reuters/training/6163
inflating: /root/nltk_data/corpora/reuters/training/6165
inflating: /root/nltk_data/corpora/reuters/training/6166
inflating: /root/nltk_data/corpora/reuters/training/6169
inflating: /root/nltk_data/corpora/reuters/training/6172
```

```
inflating: /root/nltk_data/corpora/reuters/training/6173
inflating: /root/nltk_data/corpora/reuters/training/6175
inflating: /root/nltk_data/corpora/reuters/training/6176
inflating: /root/nltk_data/corpora/reuters/training/6177
inflating: /root/nltk_data/corpora/reuters/training/6178
inflating: /root/nltk_data/corpora/reuters/training/6179
inflating: /root/nltk_data/corpora/reuters/training/6180
inflating: /root/nltk_data/corpora/reuters/training/6181
inflating: /root/nltk_data/corpora/reuters/training/6184
inflating: /root/nltk_data/corpora/reuters/training/6185
inflating: /root/nltk_data/corpora/reuters/training/6186
inflating: /root/nltk_data/corpora/reuters/training/6187
inflating: /root/nltk_data/corpora/reuters/training/6188
inflating: /root/nltk_data/corpora/reuters/training/6189
inflating: /root/nltk_data/corpora/reuters/training/6197
inflating: /root/nltk_data/corpora/reuters/training/6198
inflating: /root/nltk_data/corpora/reuters/training/6199
inflating: /root/nltk_data/corpora/reuters/training/6201
inflating: /root/nltk_data/corpora/reuters/training/6202
inflating: /root/nltk_data/corpora/reuters/training/6204
inflating: /root/nltk_data/corpora/reuters/training/6205
inflating: /root/nltk_data/corpora/reuters/training/6208
inflating: /root/nltk_data/corpora/reuters/training/6209
inflating: /root/nltk_data/corpora/reuters/training/6211
inflating: /root/nltk_data/corpora/reuters/training/6212
inflating: /root/nltk_data/corpora/reuters/training/6213
inflating: /root/nltk_data/corpora/reuters/training/6214
inflating: /root/nltk_data/corpora/reuters/training/6215
inflating: /root/nltk_data/corpora/reuters/training/6217
inflating: /root/nltk_data/corpora/reuters/training/6218
inflating: /root/nltk_data/corpora/reuters/training/6220
inflating: /root/nltk_data/corpora/reuters/training/6221
inflating: /root/nltk_data/corpora/reuters/training/6222
inflating: /root/nltk_data/corpora/reuters/training/6223
inflating: /root/nltk_data/corpora/reuters/training/6225
inflating: /root/nltk_data/corpora/reuters/training/6227
inflating: /root/nltk_data/corpora/reuters/training/6231
inflating: /root/nltk_data/corpora/reuters/training/6232
inflating: /root/nltk_data/corpora/reuters/training/6233
inflating: /root/nltk_data/corpora/reuters/training/6234
inflating: /root/nltk_data/corpora/reuters/training/6235
inflating: /root/nltk_data/corpora/reuters/training/6236
inflating: /root/nltk_data/corpora/reuters/training/6239
inflating: /root/nltk_data/corpora/reuters/training/6240
inflating: /root/nltk_data/corpora/reuters/training/6242
inflating: /root/nltk_data/corpora/reuters/training/6243
inflating: /root/nltk_data/corpora/reuters/training/6244
inflating: /root/nltk_data/corpora/reuters/training/6246
inflating: /root/nltk_data/corpora/reuters/training/6248
inflating: /root/nltk_data/corpora/reuters/training/6252
inflating: /root/nltk_data/corpora/reuters/training/6253
inflating: /root/nltk_data/corpora/reuters/training/6254
inflating: /root/nltk_data/corpora/reuters/training/6257
inflating: /root/nltk_data/corpora/reuters/training/6259
inflating: /root/nltk_data/corpora/reuters/training/6264
inflating: /root/nltk_data/corpora/reuters/training/6265
inflating: /root/nltk_data/corpora/reuters/training/6267
inflating: /root/nltk_data/corpora/reuters/training/6269
inflating: /root/nltk_data/corpora/reuters/training/6271
inflating: /root/nltk_data/corpora/reuters/training/6273
inflating: /root/nltk_data/corpora/reuters/training/6274
inflating: /root/nltk_data/corpora/reuters/training/6275
inflating: /root/nltk_data/corpora/reuters/training/6276
inflating: /root/nltk_data/corpora/reuters/training/6277
inflating: /root/nltk_data/corpora/reuters/training/6278
inflating: /root/nltk_data/corpora/reuters/training/6279
inflating: /root/nltk_data/corpora/reuters/training/6280
inflating: /root/nltk_data/corpora/reuters/training/6282
inflating: /root/nltk_data/corpora/reuters/training/6283
inflating: /root/nltk_data/corpora/reuters/training/6287
inflating: /root/nltk_data/corpora/reuters/training/6288
inflating: /root/nltk_data/corpora/reuters/training/6293
inflating: /root/nltk_data/corpora/reuters/training/6294
inflating: /root/nltk_data/corpora/reuters/training/6295
inflating: /root/nltk_data/corpora/reuters/training/6296
inflating: /root/nltk_data/corpora/reuters/training/6297
inflating: /root/nltk_data/corpora/reuters/training/6298
inflating: /root/nltk_data/corpora/reuters/training/6299
inflating: /root/nltk_data/corpora/reuters/training/6301
inflating: /root/nltk_data/corpora/reuters/training/6305
inflating: /root/nltk_data/corpora/reuters/training/6306
inflating: /root/nltk_data/corpora/reuters/training/6309
inflating: /root/nltk_data/corpora/reuters/training/6310
inflating: /root/nltk_data/corpora/reuters/training/6311
inflating: /root/nltk_data/corpora/reuters/training/6312
inflating: /root/nltk_data/corpora/reuters/training/6313
inflating: /root/nltk_data/corpora/reuters/training/6315
inflating: /root/nltk_data/corpora/reuters/training/6317
inflating: /root/nltk_data/corpora/reuters/training/6323
inflating: /root/nltk_data/corpora/reuters/training/6324
inflating: /root/nltk_data/corpora/reuters/training/6325
inflating: /root/nltk_data/corpora/reuters/training/6326
inflating: /root/nltk_data/corpora/reuters/training/6331
inflating: /root/nltk_data/corpora/reuters/training/6333
inflating: /root/nltk_data/corpora/reuters/training/6335
inflating: /root/nltk_data/corpora/reuters/training/6337
inflating: /root/nltk_data/corpora/reuters/training/6338
inflating: /root/nltk_data/corpora/reuters/training/6339
inflating: /root/nltk_data/corpora/reuters/training/6340
inflating: /root/nltk_data/corpora/reuters/training/6341
inflating: /root/nltk_data/corpora/reuters/training/6342
inflating: /root/nltk_data/corpora/reuters/training/6344
inflating: /root/nltk_data/corpora/reuters/training/6346
inflating: /root/nltk_data/corpora/reuters/training/6347
inflating: /root/nltk_data/corpora/reuters/training/6348
inflating: /root/nltk_data/corpora/reuters/training/6349
inflating: /root/nltk_data/corpora/reuters/training/6350
inflating: /root/nltk_data/corpora/reuters/training/6352
inflating: /root/nltk_data/corpora/reuters/training/6353
inflating: /root/nltk_data/corpora/reuters/training/6354
inflating: /root/nltk_data/corpora/reuters/training/6355
inflating: /root/nltk_data/corpora/reuters/training/6356
inflating: /root/nltk_data/corpora/reuters/training/6357
inflating: /root/nltk_data/corpora/reuters/training/6358
inflating: /root/nltk_data/corpora/reuters/training/6359
inflating: /root/nltk_data/corpora/reuters/training/6361
inflating: /root/nltk_data/corpora/reuters/training/6362
inflating: /root/nltk_data/corpora/reuters/training/6364
inflating: /root/nltk_data/corpora/reuters/training/6365
inflating: /root/nltk_data/corpora/reuters/training/6366
inflating: /root/nltk_data/corpora/reuters/training/6368
inflating: /root/nltk_data/corpora/reuters/training/6369
inflating: /root/nltk_data/corpora/reuters/training/6371
inflating: /root/nltk_data/corpora/reuters/training/6372
inflating: /root/nltk_data/corpora/reuters/training/6373
inflating: /root/nltk_data/corpora/reuters/training/6375
inflating: /root/nltk_data/corpora/reuters/training/6376
inflating: /root/nltk_data/corpora/reuters/training/6382
inflating: /root/nltk_data/corpora/reuters/training/6384
inflating: /root/nltk_data/corpora/reuters/training/6385
inflating: /root/nltk_data/corpora/reuters/training/6386
inflating: /root/nltk_data/corpora/reuters/training/6394
inflating: /root/nltk_data/corpora/reuters/training/6398
inflating: /root/nltk_data/corpora/reuters/training/6399
inflating: /root/nltk_data/corpora/reuters/training/64
```

```
inflating: /root/nltk_data/corpora/reuters/training/6400
inflating: /root/nltk_data/corpora/reuters/training/6402
inflating: /root/nltk_data/corpora/reuters/training/6404
inflating: /root/nltk_data/corpora/reuters/training/6405
inflating: /root/nltk_data/corpora/reuters/training/6406
inflating: /root/nltk_data/corpora/reuters/training/6407
inflating: /root/nltk_data/corpora/reuters/training/6410
inflating: /root/nltk_data/corpora/reuters/training/6412
inflating: /root/nltk_data/corpora/reuters/training/6413
inflating: /root/nltk_data/corpora/reuters/training/6414
inflating: /root/nltk_data/corpora/reuters/training/6416
inflating: /root/nltk_data/corpora/reuters/training/6417
inflating: /root/nltk_data/corpora/reuters/training/6419
inflating: /root/nltk_data/corpora/reuters/training/6420
inflating: /root/nltk_data/corpora/reuters/training/6421
inflating: /root/nltk_data/corpora/reuters/training/6422
inflating: /root/nltk_data/corpora/reuters/training/6425
inflating: /root/nltk_data/corpora/reuters/training/6426
inflating: /root/nltk_data/corpora/reuters/training/6427
inflating: /root/nltk_data/corpora/reuters/training/6428
inflating: /root/nltk_data/corpora/reuters/training/6430
inflating: /root/nltk_data/corpora/reuters/training/6432
inflating: /root/nltk_data/corpora/reuters/training/6433
inflating: /root/nltk_data/corpora/reuters/training/6434
inflating: /root/nltk_data/corpora/reuters/training/6435
inflating: /root/nltk_data/corpora/reuters/training/6436
inflating: /root/nltk_data/corpora/reuters/training/6437
inflating: /root/nltk_data/corpora/reuters/training/6438
inflating: /root/nltk_data/corpora/reuters/training/6440
inflating: /root/nltk_data/corpora/reuters/training/6441
inflating: /root/nltk_data/corpora/reuters/training/6442
inflating: /root/nltk_data/corpora/reuters/training/6443
inflating: /root/nltk_data/corpora/reuters/training/6445
inflating: /root/nltk_data/corpora/reuters/training/6447
inflating: /root/nltk_data/corpora/reuters/training/6449
inflating: /root/nltk_data/corpora/reuters/training/6450
inflating: /root/nltk_data/corpora/reuters/training/6451
inflating: /root/nltk_data/corpora/reuters/training/6452
inflating: /root/nltk_data/corpora/reuters/training/6453
inflating: /root/nltk_data/corpora/reuters/training/6454
inflating: /root/nltk_data/corpora/reuters/training/6455
inflating: /root/nltk_data/corpora/reuters/training/6456
inflating: /root/nltk_data/corpora/reuters/training/6457
inflating: /root/nltk_data/corpora/reuters/training/6458
inflating: /root/nltk_data/corpora/reuters/training/6459
inflating: /root/nltk_data/corpora/reuters/training/6460
inflating: /root/nltk_data/corpora/reuters/training/6461
inflating: /root/nltk_data/corpora/reuters/training/6463
inflating: /root/nltk_data/corpora/reuters/training/6464
inflating: /root/nltk_data/corpora/reuters/training/6465
inflating: /root/nltk_data/corpora/reuters/training/6466
inflating: /root/nltk_data/corpora/reuters/training/6471
inflating: /root/nltk_data/corpora/reuters/training/6473
inflating: /root/nltk_data/corpora/reuters/training/6474
inflating: /root/nltk_data/corpora/reuters/training/6475
inflating: /root/nltk_data/corpora/reuters/training/6476
inflating: /root/nltk_data/corpora/reuters/training/6477
inflating: /root/nltk_data/corpora/reuters/training/6478
inflating: /root/nltk_data/corpora/reuters/training/6481
inflating: /root/nltk_data/corpora/reuters/training/6483
inflating: /root/nltk_data/corpora/reuters/training/6484
inflating: /root/nltk_data/corpora/reuters/training/6485
inflating: /root/nltk_data/corpora/reuters/training/6486
inflating: /root/nltk_data/corpora/reuters/training/6488
inflating: /root/nltk_data/corpora/reuters/training/6491
inflating: /root/nltk_data/corpora/reuters/training/6492
inflating: /root/nltk_data/corpora/reuters/training/6493
inflating: /root/nltk_data/corpora/reuters/training/6494
inflating: /root/nltk_data/corpora/reuters/training/6495
inflating: /root/nltk_data/corpora/reuters/training/6496
inflating: /root/nltk_data/corpora/reuters/training/6497
inflating: /root/nltk_data/corpora/reuters/training/6498
inflating: /root/nltk_data/corpora/reuters/training/6499
inflating: /root/nltk_data/corpora/reuters/training/65
inflating: /root/nltk_data/corpora/reuters/training/6500
inflating: /root/nltk_data/corpora/reuters/training/6502
inflating: /root/nltk_data/corpora/reuters/training/6503
inflating: /root/nltk_data/corpora/reuters/training/6504
inflating: /root/nltk_data/corpora/reuters/training/6505
inflating: /root/nltk_data/corpora/reuters/training/6506
inflating: /root/nltk_data/corpora/reuters/training/6507
inflating: /root/nltk_data/corpora/reuters/training/6509
inflating: /root/nltk_data/corpora/reuters/training/6511
inflating: /root/nltk_data/corpora/reuters/training/6514
inflating: /root/nltk_data/corpora/reuters/training/6515
inflating: /root/nltk_data/corpora/reuters/training/6520
inflating: /root/nltk_data/corpora/reuters/training/6521
inflating: /root/nltk_data/corpora/reuters/training/6522
inflating: /root/nltk_data/corpora/reuters/training/6523
inflating: /root/nltk_data/corpora/reuters/training/6524
inflating: /root/nltk_data/corpora/reuters/training/6525
inflating: /root/nltk_data/corpora/reuters/training/6527
inflating: /root/nltk_data/corpora/reuters/training/6528
inflating: /root/nltk_data/corpora/reuters/training/6529
inflating: /root/nltk_data/corpora/reuters/training/6531
inflating: /root/nltk_data/corpora/reuters/training/6532
inflating: /root/nltk_data/corpora/reuters/training/6533
inflating: /root/nltk_data/corpora/reuters/training/6534
inflating: /root/nltk_data/corpora/reuters/training/6535
inflating: /root/nltk_data/corpora/reuters/training/6537
inflating: /root/nltk_data/corpora/reuters/training/6538
inflating: /root/nltk_data/corpora/reuters/training/6539
inflating: /root/nltk_data/corpora/reuters/training/6540
inflating: /root/nltk_data/corpora/reuters/training/6541
inflating: /root/nltk_data/corpora/reuters/training/6542
inflating: /root/nltk_data/corpora/reuters/training/6543
inflating: /root/nltk_data/corpora/reuters/training/6546
inflating: /root/nltk_data/corpora/reuters/training/6547
inflating: /root/nltk_data/corpora/reuters/training/6549
inflating: /root/nltk_data/corpora/reuters/training/6551
inflating: /root/nltk_data/corpora/reuters/training/6552
inflating: /root/nltk_data/corpora/reuters/training/6556
inflating: /root/nltk_data/corpora/reuters/training/6557
inflating: /root/nltk_data/corpora/reuters/training/6559
inflating: /root/nltk_data/corpora/reuters/training/6561
inflating: /root/nltk_data/corpora/reuters/training/6562
inflating: /root/nltk_data/corpora/reuters/training/6568
inflating: /root/nltk_data/corpora/reuters/training/6572
inflating: /root/nltk_data/corpora/reuters/training/6573
inflating: /root/nltk_data/corpora/reuters/training/6575
inflating: /root/nltk_data/corpora/reuters/training/6576
inflating: /root/nltk_data/corpora/reuters/training/6577
inflating: /root/nltk_data/corpora/reuters/training/6578
inflating: /root/nltk_data/corpora/reuters/training/6581
inflating: /root/nltk_data/corpora/reuters/training/6583
inflating: /root/nltk_data/corpora/reuters/training/6584
inflating: /root/nltk_data/corpora/reuters/training/6585
inflating: /root/nltk_data/corpora/reuters/training/6586
inflating: /root/nltk_data/corpora/reuters/training/6588
inflating: /root/nltk_data/corpora/reuters/training/6590
inflating: /root/nltk_data/corpora/reuters/training/6591
inflating: /root/nltk_data/corpora/reuters/training/6592
inflating: /root/nltk_data/corpora/reuters/training/6593
inflating: /root/nltk_data/corpora/reuters/training/6594
```

```
inflating: /root/nltk_data/corpora/reuters/training/6595
inflating: /root/nltk_data/corpora/reuters/training/6596
inflating: /root/nltk_data/corpora/reuters/training/6597
inflating: /root/nltk_data/corpora/reuters/training/6598
inflating: /root/nltk_data/corpora/reuters/training/66
inflating: /root/nltk_data/corpora/reuters/training/6601
inflating: /root/nltk_data/corpora/reuters/training/6603
inflating: /root/nltk_data/corpora/reuters/training/6604
inflating: /root/nltk_data/corpora/reuters/training/6606
inflating: /root/nltk_data/corpora/reuters/training/6607
inflating: /root/nltk_data/corpora/reuters/training/6609
inflating: /root/nltk_data/corpora/reuters/training/6611
inflating: /root/nltk_data/corpora/reuters/training/6612
inflating: /root/nltk_data/corpora/reuters/training/6616
inflating: /root/nltk_data/corpora/reuters/training/6619
inflating: /root/nltk_data/corpora/reuters/training/6620
inflating: /root/nltk_data/corpora/reuters/training/6621
inflating: /root/nltk_data/corpora/reuters/training/6622
inflating: /root/nltk_data/corpora/reuters/training/6623
inflating: /root/nltk_data/corpora/reuters/training/6626
inflating: /root/nltk_data/corpora/reuters/training/6629
inflating: /root/nltk_data/corpora/reuters/training/6630
inflating: /root/nltk_data/corpora/reuters/training/6632
inflating: /root/nltk_data/corpora/reuters/training/6635
inflating: /root/nltk_data/corpora/reuters/training/6636
inflating: /root/nltk_data/corpora/reuters/training/6637
inflating: /root/nltk_data/corpora/reuters/training/6638
inflating: /root/nltk_data/corpora/reuters/training/6639
inflating: /root/nltk_data/corpora/reuters/training/6641
inflating: /root/nltk_data/corpora/reuters/training/6642
inflating: /root/nltk_data/corpora/reuters/training/6647
inflating: /root/nltk_data/corpora/reuters/training/6650
inflating: /root/nltk_data/corpora/reuters/training/6652
inflating: /root/nltk_data/corpora/reuters/training/6653
inflating: /root/nltk_data/corpora/reuters/training/6654
inflating: /root/nltk_data/corpora/reuters/training/6655
inflating: /root/nltk_data/corpora/reuters/training/6656
inflating: /root/nltk_data/corpora/reuters/training/6657
inflating: /root/nltk_data/corpora/reuters/training/6658
inflating: /root/nltk_data/corpora/reuters/training/6659
inflating: /root/nltk_data/corpora/reuters/training/6660
inflating: /root/nltk_data/corpora/reuters/training/6661
inflating: /root/nltk_data/corpora/reuters/training/6665
inflating: /root/nltk_data/corpora/reuters/training/6666
inflating: /root/nltk_data/corpora/reuters/training/6668
inflating: /root/nltk_data/corpora/reuters/training/6669
inflating: /root/nltk_data/corpora/reuters/training/6670
inflating: /root/nltk_data/corpora/reuters/training/6672
inflating: /root/nltk_data/corpora/reuters/training/6673
inflating: /root/nltk_data/corpora/reuters/training/6674
inflating: /root/nltk_data/corpora/reuters/training/6675
inflating: /root/nltk_data/corpora/reuters/training/6677
inflating: /root/nltk_data/corpora/reuters/training/6679
inflating: /root/nltk_data/corpora/reuters/training/6680
inflating: /root/nltk_data/corpora/reuters/training/6681
inflating: /root/nltk_data/corpora/reuters/training/6682
inflating: /root/nltk_data/corpora/reuters/training/6683
inflating: /root/nltk_data/corpora/reuters/training/6684
inflating: /root/nltk_data/corpora/reuters/training/6685
inflating: /root/nltk_data/corpora/reuters/training/6686
inflating: /root/nltk_data/corpora/reuters/training/6687
inflating: /root/nltk_data/corpora/reuters/training/6689
inflating: /root/nltk_data/corpora/reuters/training/6694
inflating: /root/nltk_data/corpora/reuters/training/6695
inflating: /root/nltk_data/corpora/reuters/training/6696
inflating: /root/nltk_data/corpora/reuters/training/6697
inflating: /root/nltk_data/corpora/reuters/training/6699
inflating: /root/nltk_data/corpora/reuters/training/6701
inflating: /root/nltk_data/corpora/reuters/training/6702
inflating: /root/nltk_data/corpora/reuters/training/6705
inflating: /root/nltk_data/corpora/reuters/training/6707
inflating: /root/nltk_data/corpora/reuters/training/6708
inflating: /root/nltk_data/corpora/reuters/training/6709
inflating: /root/nltk_data/corpora/reuters/training/6710
inflating: /root/nltk_data/corpora/reuters/training/6711
inflating: /root/nltk_data/corpora/reuters/training/6712
inflating: /root/nltk_data/corpora/reuters/training/6713
inflating: /root/nltk_data/corpora/reuters/training/6714
inflating: /root/nltk_data/corpora/reuters/training/6716
inflating: /root/nltk_data/corpora/reuters/training/6719
inflating: /root/nltk_data/corpora/reuters/training/6720
inflating: /root/nltk_data/corpora/reuters/training/6721
inflating: /root/nltk_data/corpora/reuters/training/6722
inflating: /root/nltk_data/corpora/reuters/training/6723
inflating: /root/nltk_data/corpora/reuters/training/6725
inflating: /root/nltk_data/corpora/reuters/training/6726
inflating: /root/nltk_data/corpora/reuters/training/6729
inflating: /root/nltk_data/corpora/reuters/training/6730
inflating: /root/nltk_data/corpora/reuters/training/6731
inflating: /root/nltk_data/corpora/reuters/training/6732
inflating: /root/nltk_data/corpora/reuters/training/6733
inflating: /root/nltk_data/corpora/reuters/training/6734
inflating: /root/nltk_data/corpora/reuters/training/6735
inflating: /root/nltk_data/corpora/reuters/training/6737
inflating: /root/nltk_data/corpora/reuters/training/6738
inflating: /root/nltk_data/corpora/reuters/training/6739
inflating: /root/nltk_data/corpora/reuters/training/6740
inflating: /root/nltk_data/corpora/reuters/training/6741
inflating: /root/nltk_data/corpora/reuters/training/6742
inflating: /root/nltk_data/corpora/reuters/training/6743
inflating: /root/nltk_data/corpora/reuters/training/6744
inflating: /root/nltk_data/corpora/reuters/training/6746
inflating: /root/nltk_data/corpora/reuters/training/6747
inflating: /root/nltk_data/corpora/reuters/training/6748
inflating: /root/nltk_data/corpora/reuters/training/6751
inflating: /root/nltk_data/corpora/reuters/training/6752
inflating: /root/nltk_data/corpora/reuters/training/6757
inflating: /root/nltk_data/corpora/reuters/training/6758
inflating: /root/nltk_data/corpora/reuters/training/6759
inflating: /root/nltk_data/corpora/reuters/training/6760
inflating: /root/nltk_data/corpora/reuters/training/6761
inflating: /root/nltk_data/corpora/reuters/training/6762
inflating: /root/nltk_data/corpora/reuters/training/6763
inflating: /root/nltk_data/corpora/reuters/training/6764
inflating: /root/nltk_data/corpora/reuters/training/6765
inflating: /root/nltk_data/corpora/reuters/training/6767
inflating: /root/nltk_data/corpora/reuters/training/677
inflating: /root/nltk_data/corpora/reuters/training/6770
inflating: /root/nltk_data/corpora/reuters/training/6771
inflating: /root/nltk_data/corpora/reuters/training/6773
inflating: /root/nltk_data/corpora/reuters/training/6774
inflating: /root/nltk_data/corpora/reuters/training/6776
inflating: /root/nltk_data/corpora/reuters/training/6777
inflating: /root/nltk_data/corpora/reuters/training/6778
inflating: /root/nltk_data/corpora/reuters/training/6779
inflating: /root/nltk_data/corpora/reuters/training/6783
inflating: /root/nltk_data/corpora/reuters/training/6784
inflating: /root/nltk_data/corpora/reuters/training/6785
inflating: /root/nltk_data/corpora/reuters/training/6787
inflating: /root/nltk_data/corpora/reuters/training/6788
inflating: /root/nltk_data/corpora/reuters/training/6789
inflating: /root/nltk_data/corpora/reuters/training/679
inflating: /root/nltk_data/corpora/reuters/training/6790
inflating: /root/nltk_data/corpora/reuters/training/6791
inflating: /root/nltk_data/corpora/reuters/training/6794
```

```
inflating: /root/nltk_data/corpora/reuters/training/6794
inflating: /root/nltk_data/corpora/reuters/training/6797
inflating: /root/nltk_data/corpora/reuters/training/6798
inflating: /root/nltk_data/corpora/reuters/training/6799
inflating: /root/nltk_data/corpora/reuters/training/68
inflating: /root/nltk_data/corpora/reuters/training/680
inflating: /root/nltk_data/corpora/reuters/training/6801
inflating: /root/nltk_data/corpora/reuters/training/6802
inflating: /root/nltk_data/corpora/reuters/training/6804
inflating: /root/nltk_data/corpora/reuters/training/6805
inflating: /root/nltk_data/corpora/reuters/training/6807
inflating: /root/nltk_data/corpora/reuters/training/6809
inflating: /root/nltk_data/corpora/reuters/training/6810
inflating: /root/nltk_data/corpora/reuters/training/6811
inflating: /root/nltk_data/corpora/reuters/training/6813
inflating: /root/nltk_data/corpora/reuters/training/6814
inflating: /root/nltk_data/corpora/reuters/training/6815
inflating: /root/nltk_data/corpora/reuters/training/6816
inflating: /root/nltk_data/corpora/reuters/training/6817
inflating: /root/nltk_data/corpora/reuters/training/6818
inflating: /root/nltk_data/corpora/reuters/training/6819
inflating: /root/nltk_data/corpora/reuters/training/682
inflating: /root/nltk_data/corpora/reuters/training/6825
inflating: /root/nltk_data/corpora/reuters/training/6827
inflating: /root/nltk_data/corpora/reuters/training/6828
inflating: /root/nltk_data/corpora/reuters/training/6830
inflating: /root/nltk_data/corpora/reuters/training/6831
inflating: /root/nltk_data/corpora/reuters/training/6836
inflating: /root/nltk_data/corpora/reuters/training/6837
inflating: /root/nltk_data/corpora/reuters/training/684
inflating: /root/nltk_data/corpora/reuters/training/6840
inflating: /root/nltk_data/corpora/reuters/training/6841
inflating: /root/nltk_data/corpora/reuters/training/6844
inflating: /root/nltk_data/corpora/reuters/training/6846
inflating: /root/nltk_data/corpora/reuters/training/6848
inflating: /root/nltk_data/corpora/reuters/training/6849
inflating: /root/nltk_data/corpora/reuters/training/685
inflating: /root/nltk_data/corpora/reuters/training/6850
inflating: /root/nltk_data/corpora/reuters/training/6852
inflating: /root/nltk_data/corpora/reuters/training/6853
inflating: /root/nltk_data/corpora/reuters/training/6855
inflating: /root/nltk_data/corpora/reuters/training/6856
inflating: /root/nltk_data/corpora/reuters/training/6859
inflating: /root/nltk_data/corpora/reuters/training/686
inflating: /root/nltk_data/corpora/reuters/training/6860
inflating: /root/nltk_data/corpora/reuters/training/6861
inflating: /root/nltk_data/corpora/reuters/training/6862
inflating: /root/nltk_data/corpora/reuters/training/6863
inflating: /root/nltk_data/corpora/reuters/training/6864
inflating: /root/nltk_data/corpora/reuters/training/6865
inflating: /root/nltk_data/corpora/reuters/training/6866
inflating: /root/nltk_data/corpora/reuters/training/6869
inflating: /root/nltk_data/corpora/reuters/training/687
inflating: /root/nltk_data/corpora/reuters/training/6871
inflating: /root/nltk_data/corpora/reuters/training/6872
inflating: /root/nltk_data/corpora/reuters/training/6873
inflating: /root/nltk_data/corpora/reuters/training/6874
inflating: /root/nltk_data/corpora/reuters/training/6875
inflating: /root/nltk_data/corpora/reuters/training/6876
inflating: /root/nltk_data/corpora/reuters/training/6877
inflating: /root/nltk_data/corpora/reuters/training/688
inflating: /root/nltk_data/corpora/reuters/training/6880
inflating: /root/nltk_data/corpora/reuters/training/6881
inflating: /root/nltk_data/corpora/reuters/training/6882
inflating: /root/nltk_data/corpora/reuters/training/6884
inflating: /root/nltk_data/corpora/reuters/training/6886
inflating: /root/nltk_data/corpora/reuters/training/6887
inflating: /root/nltk_data/corpora/reuters/training/6888
inflating: /root/nltk_data/corpora/reuters/training/689
inflating: /root/nltk_data/corpora/reuters/training/6890
inflating: /root/nltk_data/corpora/reuters/training/6893
inflating: /root/nltk_data/corpora/reuters/training/6894
inflating: /root/nltk_data/corpora/reuters/training/6897
inflating: /root/nltk_data/corpora/reuters/training/6898
inflating: /root/nltk_data/corpora/reuters/training/69
inflating: /root/nltk_data/corpora/reuters/training/690
inflating: /root/nltk_data/corpora/reuters/training/6903
inflating: /root/nltk_data/corpora/reuters/training/6905
inflating: /root/nltk_data/corpora/reuters/training/6906
inflating: /root/nltk_data/corpora/reuters/training/6907
inflating: /root/nltk_data/corpora/reuters/training/6908
inflating: /root/nltk_data/corpora/reuters/training/6909
inflating: /root/nltk_data/corpora/reuters/training/691
inflating: /root/nltk_data/corpora/reuters/training/6912
inflating: /root/nltk_data/corpora/reuters/training/6913
inflating: /root/nltk_data/corpora/reuters/training/6914
inflating: /root/nltk_data/corpora/reuters/training/6917
inflating: /root/nltk_data/corpora/reuters/training/692
inflating: /root/nltk_data/corpora/reuters/training/6920
inflating: /root/nltk_data/corpora/reuters/training/6922
inflating: /root/nltk_data/corpora/reuters/training/6926
inflating: /root/nltk_data/corpora/reuters/training/6927
inflating: /root/nltk_data/corpora/reuters/training/6928
inflating: /root/nltk_data/corpora/reuters/training/6929
inflating: /root/nltk_data/corpora/reuters/training/693
inflating: /root/nltk_data/corpora/reuters/training/6930
inflating: /root/nltk_data/corpora/reuters/training/6931
inflating: /root/nltk_data/corpora/reuters/training/6934
inflating: /root/nltk_data/corpora/reuters/training/6935
inflating: /root/nltk_data/corpora/reuters/training/6939
inflating: /root/nltk_data/corpora/reuters/training/694
inflating: /root/nltk_data/corpora/reuters/training/6941
inflating: /root/nltk_data/corpora/reuters/training/6942
inflating: /root/nltk_data/corpora/reuters/training/6943
inflating: /root/nltk_data/corpora/reuters/training/6944
inflating: /root/nltk_data/corpora/reuters/training/6945
inflating: /root/nltk_data/corpora/reuters/training/6946
inflating: /root/nltk_data/corpora/reuters/training/695
inflating: /root/nltk_data/corpora/reuters/training/6951
inflating: /root/nltk_data/corpora/reuters/training/6952
inflating: /root/nltk_data/corpora/reuters/training/6953
inflating: /root/nltk_data/corpora/reuters/training/6954
inflating: /root/nltk_data/corpora/reuters/training/6955
inflating: /root/nltk_data/corpora/reuters/training/6957
inflating: /root/nltk_data/corpora/reuters/training/6959
inflating: /root/nltk_data/corpora/reuters/training/696
inflating: /root/nltk_data/corpora/reuters/training/6960
inflating: /root/nltk_data/corpora/reuters/training/6961
inflating: /root/nltk_data/corpora/reuters/training/6964
inflating: /root/nltk_data/corpora/reuters/training/6965
inflating: /root/nltk_data/corpora/reuters/training/6968
inflating: /root/nltk_data/corpora/reuters/training/6969
inflating: /root/nltk_data/corpora/reuters/training/697
inflating: /root/nltk_data/corpora/reuters/training/6970
inflating: /root/nltk_data/corpora/reuters/training/6972
inflating: /root/nltk_data/corpora/reuters/training/6975
inflating: /root/nltk_data/corpora/reuters/training/6976
inflating: /root/nltk_data/corpora/reuters/training/698
inflating: /root/nltk_data/corpora/reuters/training/6981
inflating: /root/nltk_data/corpora/reuters/training/6982
inflating: /root/nltk_data/corpora/reuters/training/6983
inflating: /root/nltk_data/corpora/reuters/training/6985
inflating: /root/nltk_data/corpora/reuters/training/6986
inflating: /root/nltk_data/corpora/reuters/training/699
inflating: /root/nltk_data/corpora/reuters/training/6993
```

```
inflating: /root/nltk_data/corpora/reuters/training/6994
inflating: /root/nltk_data/corpora/reuters/training/6995
inflating: /root/nltk_data/corpora/reuters/training/6996
inflating: /root/nltk_data/corpora/reuters/training/6998
inflating: /root/nltk_data/corpora/reuters/training/6999
inflating: /root/nltk_data/corpora/reuters/training/700
inflating: /root/nltk_data/corpora/reuters/training/7000
inflating: /root/nltk_data/corpora/reuters/training/7001
inflating: /root/nltk_data/corpora/reuters/training/7003
inflating: /root/nltk_data/corpora/reuters/training/7004
inflating: /root/nltk_data/corpora/reuters/training/7005
inflating: /root/nltk_data/corpora/reuters/training/7006
inflating: /root/nltk_data/corpora/reuters/training/7007
inflating: /root/nltk_data/corpora/reuters/training/7008
inflating: /root/nltk_data/corpora/reuters/training/7009
inflating: /root/nltk_data/corpora/reuters/training/701
inflating: /root/nltk_data/corpora/reuters/training/7010
inflating: /root/nltk_data/corpora/reuters/training/7011
inflating: /root/nltk_data/corpora/reuters/training/7012
inflating: /root/nltk_data/corpora/reuters/training/7013
inflating: /root/nltk_data/corpora/reuters/training/7014
inflating: /root/nltk_data/corpora/reuters/training/7015
inflating: /root/nltk_data/corpora/reuters/training/702
inflating: /root/nltk_data/corpora/reuters/training/7020
inflating: /root/nltk_data/corpora/reuters/training/7023
inflating: /root/nltk_data/corpora/reuters/training/7024
inflating: /root/nltk_data/corpora/reuters/training/7027
inflating: /root/nltk_data/corpora/reuters/training/7029
inflating: /root/nltk_data/corpora/reuters/training/703
inflating: /root/nltk_data/corpora/reuters/training/7030
inflating: /root/nltk_data/corpora/reuters/training/7031
inflating: /root/nltk_data/corpora/reuters/training/7032
inflating: /root/nltk_data/corpora/reuters/training/7036
inflating: /root/nltk_data/corpora/reuters/training/7037
inflating: /root/nltk_data/corpora/reuters/training/7039
inflating: /root/nltk_data/corpora/reuters/training/704
inflating: /root/nltk_data/corpora/reuters/training/7041
inflating: /root/nltk_data/corpora/reuters/training/7043
inflating: /root/nltk_data/corpora/reuters/training/7045
inflating: /root/nltk_data/corpora/reuters/training/7046
inflating: /root/nltk_data/corpora/reuters/training/7047
inflating: /root/nltk_data/corpora/reuters/training/7048
inflating: /root/nltk_data/corpora/reuters/training/7049
inflating: /root/nltk_data/corpora/reuters/training/7052
inflating: /root/nltk_data/corpora/reuters/training/7057
inflating: /root/nltk_data/corpora/reuters/training/7058
inflating: /root/nltk_data/corpora/reuters/training/706
inflating: /root/nltk_data/corpora/reuters/training/7060
inflating: /root/nltk_data/corpora/reuters/training/7061
inflating: /root/nltk_data/corpora/reuters/training/7062
inflating: /root/nltk_data/corpora/reuters/training/7063
inflating: /root/nltk_data/corpora/reuters/training/7064
inflating: /root/nltk_data/corpora/reuters/training/7065
inflating: /root/nltk_data/corpora/reuters/training/7066
inflating: /root/nltk_data/corpora/reuters/training/7067
inflating: /root/nltk_data/corpora/reuters/training/7068
inflating: /root/nltk_data/corpora/reuters/training/707
inflating: /root/nltk_data/corpora/reuters/training/7070
inflating: /root/nltk_data/corpora/reuters/training/7071
inflating: /root/nltk_data/corpora/reuters/training/7072
inflating: /root/nltk_data/corpora/reuters/training/7073
inflating: /root/nltk_data/corpora/reuters/training/7074
inflating: /root/nltk_data/corpora/reuters/training/7075
inflating: /root/nltk_data/corpora/reuters/training/7076
inflating: /root/nltk_data/corpora/reuters/training/7077
inflating: /root/nltk_data/corpora/reuters/training/7078
inflating: /root/nltk_data/corpora/reuters/training/7079
inflating: /root/nltk_data/corpora/reuters/training/708
inflating: /root/nltk_data/corpora/reuters/training/7080
inflating: /root/nltk_data/corpora/reuters/training/7083
inflating: /root/nltk_data/corpora/reuters/training/7085
inflating: /root/nltk_data/corpora/reuters/training/7087
inflating: /root/nltk_data/corpora/reuters/training/7088
inflating: /root/nltk_data/corpora/reuters/training/709
inflating: /root/nltk_data/corpora/reuters/training/7090
inflating: /root/nltk_data/corpora/reuters/training/7092
inflating: /root/nltk_data/corpora/reuters/training/7093
inflating: /root/nltk_data/corpora/reuters/training/7096
inflating: /root/nltk_data/corpora/reuters/training/7097
inflating: /root/nltk_data/corpora/reuters/training/7098
inflating: /root/nltk_data/corpora/reuters/training/7099
inflating: /root/nltk_data/corpora/reuters/training/71
inflating: /root/nltk_data/corpora/reuters/training/7100
inflating: /root/nltk_data/corpora/reuters/training/7101
inflating: /root/nltk_data/corpora/reuters/training/7102
inflating: /root/nltk_data/corpora/reuters/training/7103
inflating: /root/nltk_data/corpora/reuters/training/7104
inflating: /root/nltk_data/corpora/reuters/training/7106
inflating: /root/nltk_data/corpora/reuters/training/7107
inflating: /root/nltk_data/corpora/reuters/training/7109
inflating: /root/nltk_data/corpora/reuters/training/7110
inflating: /root/nltk_data/corpora/reuters/training/7111
inflating: /root/nltk_data/corpora/reuters/training/7113
inflating: /root/nltk_data/corpora/reuters/training/7115
inflating: /root/nltk_data/corpora/reuters/training/7117
inflating: /root/nltk_data/corpora/reuters/training/7119
inflating: /root/nltk_data/corpora/reuters/training/712
inflating: /root/nltk_data/corpora/reuters/training/7120
inflating: /root/nltk_data/corpora/reuters/training/7123
inflating: /root/nltk_data/corpora/reuters/training/7124
inflating: /root/nltk_data/corpora/reuters/training/7126
inflating: /root/nltk_data/corpora/reuters/training/7128
inflating: /root/nltk_data/corpora/reuters/training/7129
inflating: /root/nltk_data/corpora/reuters/training/7130
inflating: /root/nltk_data/corpora/reuters/training/7132
inflating: /root/nltk_data/corpora/reuters/training/7133
inflating: /root/nltk_data/corpora/reuters/training/7135
inflating: /root/nltk_data/corpora/reuters/training/7136
inflating: /root/nltk_data/corpora/reuters/training/7139
inflating: /root/nltk_data/corpora/reuters/training/714
inflating: /root/nltk_data/corpora/reuters/training/7141
inflating: /root/nltk_data/corpora/reuters/training/7143
inflating: /root/nltk_data/corpora/reuters/training/7146
inflating: /root/nltk_data/corpora/reuters/training/7148
inflating: /root/nltk_data/corpora/reuters/training/7149
inflating: /root/nltk_data/corpora/reuters/training/715
inflating: /root/nltk_data/corpora/reuters/training/7150
inflating: /root/nltk_data/corpora/reuters/training/7152
inflating: /root/nltk_data/corpora/reuters/training/7154
inflating: /root/nltk_data/corpora/reuters/training/7155
inflating: /root/nltk_data/corpora/reuters/training/7157
inflating: /root/nltk_data/corpora/reuters/training/7158
inflating: /root/nltk_data/corpora/reuters/training/7159
inflating: /root/nltk_data/corpora/reuters/training/7160
inflating: /root/nltk_data/corpora/reuters/training/7161
inflating: /root/nltk_data/corpora/reuters/training/7166
inflating: /root/nltk_data/corpora/reuters/training/7167
inflating: /root/nltk_data/corpora/reuters/training/7168
inflating: /root/nltk_data/corpora/reuters/training/7174
inflating: /root/nltk_data/corpora/reuters/training/7177
inflating: /root/nltk_data/corpora/reuters/training/718
inflating: /root/nltk_data/corpora/reuters/training/7184
inflating: /root/nltk_data/corpora/reuters/training/7185
inflating: /root/nltk_data/corpora/reuters/training/7188
```

```
inflating: /root/nltk_data/corpora/reuters/training/7189
inflating: /root/nltk_data/corpora/reuters/training/7190
inflating: /root/nltk_data/corpora/reuters/training/7191
inflating: /root/nltk_data/corpora/reuters/training/7192
inflating: /root/nltk_data/corpora/reuters/training/7193
inflating: /root/nltk_data/corpora/reuters/training/7196
inflating: /root/nltk_data/corpora/reuters/training/7199
inflating: /root/nltk_data/corpora/reuters/training/7204
inflating: /root/nltk_data/corpora/reuters/training/7205
inflating: /root/nltk_data/corpora/reuters/training/7207
inflating: /root/nltk_data/corpora/reuters/training/7208
inflating: /root/nltk_data/corpora/reuters/training/7209
inflating: /root/nltk_data/corpora/reuters/training/7212
inflating: /root/nltk_data/corpora/reuters/training/7215
inflating: /root/nltk_data/corpora/reuters/training/7216
inflating: /root/nltk_data/corpora/reuters/training/7217
inflating: /root/nltk_data/corpora/reuters/training/7218
inflating: /root/nltk_data/corpora/reuters/training/722
inflating: /root/nltk_data/corpora/reuters/training/7220
inflating: /root/nltk_data/corpora/reuters/training/7221
inflating: /root/nltk_data/corpora/reuters/training/7222
inflating: /root/nltk_data/corpora/reuters/training/7225
inflating: /root/nltk_data/corpora/reuters/training/7226
inflating: /root/nltk_data/corpora/reuters/training/7229
inflating: /root/nltk_data/corpora/reuters/training/7234
inflating: /root/nltk_data/corpora/reuters/training/7235
inflating: /root/nltk_data/corpora/reuters/training/724
inflating: /root/nltk_data/corpora/reuters/training/7240
inflating: /root/nltk_data/corpora/reuters/training/7244
inflating: /root/nltk_data/corpora/reuters/training/7245
inflating: /root/nltk_data/corpora/reuters/training/7247
inflating: /root/nltk_data/corpora/reuters/training/725
inflating: /root/nltk_data/corpora/reuters/training/7250
inflating: /root/nltk_data/corpora/reuters/training/7253
inflating: /root/nltk_data/corpora/reuters/training/7256
inflating: /root/nltk_data/corpora/reuters/training/7259
inflating: /root/nltk_data/corpora/reuters/training/7260
inflating: /root/nltk_data/corpora/reuters/training/7261
inflating: /root/nltk_data/corpora/reuters/training/7263
inflating: /root/nltk_data/corpora/reuters/training/7264
inflating: /root/nltk_data/corpora/reuters/training/7266
inflating: /root/nltk_data/corpora/reuters/training/7267
inflating: /root/nltk_data/corpora/reuters/training/7272
inflating: /root/nltk_data/corpora/reuters/training/7275
inflating: /root/nltk_data/corpora/reuters/training/7276
inflating: /root/nltk_data/corpora/reuters/training/7277
inflating: /root/nltk_data/corpora/reuters/training/7278
inflating: /root/nltk_data/corpora/reuters/training/7279
inflating: /root/nltk_data/corpora/reuters/training/728
inflating: /root/nltk_data/corpora/reuters/training/7286
inflating: /root/nltk_data/corpora/reuters/training/7287
inflating: /root/nltk_data/corpora/reuters/training/7296
inflating: /root/nltk_data/corpora/reuters/training/7297
inflating: /root/nltk_data/corpora/reuters/training/730
inflating: /root/nltk_data/corpora/reuters/training/7302
inflating: /root/nltk_data/corpora/reuters/training/7304
inflating: /root/nltk_data/corpora/reuters/training/7305
inflating: /root/nltk_data/corpora/reuters/training/7308
inflating: /root/nltk_data/corpora/reuters/training/731
inflating: /root/nltk_data/corpora/reuters/training/7310
inflating: /root/nltk_data/corpora/reuters/training/7311
inflating: /root/nltk_data/corpora/reuters/training/7312
inflating: /root/nltk_data/corpora/reuters/training/7313
inflating: /root/nltk_data/corpora/reuters/training/7314
inflating: /root/nltk_data/corpora/reuters/training/7315
inflating: /root/nltk_data/corpora/reuters/training/7317
inflating: /root/nltk_data/corpora/reuters/training/7318
inflating: /root/nltk_data/corpora/reuters/training/7319
inflating: /root/nltk_data/corpora/reuters/training/732
inflating: /root/nltk_data/corpora/reuters/training/7320
inflating: /root/nltk_data/corpora/reuters/training/7321
inflating: /root/nltk_data/corpora/reuters/training/7322
inflating: /root/nltk_data/corpora/reuters/training/7323
inflating: /root/nltk_data/corpora/reuters/training/7326
inflating: /root/nltk_data/corpora/reuters/training/7329
inflating: /root/nltk_data/corpora/reuters/training/733
inflating: /root/nltk_data/corpora/reuters/training/7330
inflating: /root/nltk_data/corpora/reuters/training/7331
inflating: /root/nltk_data/corpora/reuters/training/7332
inflating: /root/nltk_data/corpora/reuters/training/7336
inflating: /root/nltk_data/corpora/reuters/training/7337
inflating: /root/nltk_data/corpora/reuters/training/7338
inflating: /root/nltk_data/corpora/reuters/training/7339
inflating: /root/nltk_data/corpora/reuters/training/734
inflating: /root/nltk_data/corpora/reuters/training/7341
inflating: /root/nltk_data/corpora/reuters/training/7342
inflating: /root/nltk_data/corpora/reuters/training/7344
inflating: /root/nltk_data/corpora/reuters/training/7346
inflating: /root/nltk_data/corpora/reuters/training/7349
inflating: /root/nltk_data/corpora/reuters/training/735
inflating: /root/nltk_data/corpora/reuters/training/7352
inflating: /root/nltk_data/corpora/reuters/training/7353
inflating: /root/nltk_data/corpora/reuters/training/7354
inflating: /root/nltk_data/corpora/reuters/training/7355
inflating: /root/nltk_data/corpora/reuters/training/7356
inflating: /root/nltk_data/corpora/reuters/training/7357
inflating: /root/nltk_data/corpora/reuters/training/7358
inflating: /root/nltk_data/corpora/reuters/training/7359
inflating: /root/nltk_data/corpora/reuters/training/736
inflating: /root/nltk_data/corpora/reuters/training/7362
inflating: /root/nltk_data/corpora/reuters/training/7363
inflating: /root/nltk_data/corpora/reuters/training/7366
inflating: /root/nltk_data/corpora/reuters/training/7367
inflating: /root/nltk_data/corpora/reuters/training/7370
inflating: /root/nltk_data/corpora/reuters/training/7375
inflating: /root/nltk_data/corpora/reuters/training/7379
inflating: /root/nltk_data/corpora/reuters/training/7381
inflating: /root/nltk_data/corpora/reuters/training/7382
inflating: /root/nltk_data/corpora/reuters/training/7385
inflating: /root/nltk_data/corpora/reuters/training/7386
inflating: /root/nltk_data/corpora/reuters/training/7387
inflating: /root/nltk_data/corpora/reuters/training/7388
inflating: /root/nltk_data/corpora/reuters/training/7389
inflating: /root/nltk_data/corpora/reuters/training/739
inflating: /root/nltk_data/corpora/reuters/training/7390
inflating: /root/nltk_data/corpora/reuters/training/7391
inflating: /root/nltk_data/corpora/reuters/training/7394
inflating: /root/nltk_data/corpora/reuters/training/7395
inflating: /root/nltk_data/corpora/reuters/training/7397
inflating: /root/nltk_data/corpora/reuters/training/7398
inflating: /root/nltk_data/corpora/reuters/training/74
inflating: /root/nltk_data/corpora/reuters/training/7404
inflating: /root/nltk_data/corpora/reuters/training/7406
inflating: /root/nltk_data/corpora/reuters/training/7407
inflating: /root/nltk_data/corpora/reuters/training/7408
inflating: /root/nltk_data/corpora/reuters/training/741
inflating: /root/nltk_data/corpora/reuters/training/7410
inflating: /root/nltk_data/corpora/reuters/training/7413
inflating: /root/nltk_data/corpora/reuters/training/7419
inflating: /root/nltk_data/corpora/reuters/training/742
inflating: /root/nltk_data/corpora/reuters/training/7420
inflating: /root/nltk_data/corpora/reuters/training/7421
inflating: /root/nltk_data/corpora/reuters/training/7423
inflating: /root/nltk_data/corpora/reuters/training/7424
inflating: /root/nltk_data/corpora/reuters/training/7425
```

```
  inflating: /root/nltk_data/corpora/reuters/training/7426
  inflating: /root/nltk_data/corpora/reuters/training/7427
  inflating: /root/nltk_data/corpora/reuters/training/7428
  inflating: /root/nltk_data/corpora/reuters/training/7429
  inflating: /root/nltk_data/corpora/reuters/training/743
  inflating: /root/nltk_data/corpora/reuters/training/7432
  inflating: /root/nltk_data/corpora/reuters/training/7433
  inflating: /root/nltk_data/corpora/reuters/training/7437
  inflating: /root/nltk_data/corpora/reuters/training/7438
  inflating: /root/nltk_data/corpora/reuters/training/7439
  inflating: /root/nltk_data/corpora/reuters/training/7441
  inflating: /root/nltk_data/corpora/reuters/training/7442
  inflating: /root/nltk_data/corpora/reuters/training/7443
  inflating: /root/nltk_data/corpora/reuters/training/7444
  inflating: /root/nltk_data/corpora/reuters/training/7447
  inflating: /root/nltk_data/corpora/reuters/training/7449
  inflating: /root/nltk_data/corpora/reuters/training/745
  inflating: /root/nltk_data/corpora/reuters/training/7452
  inflating: /root/nltk_data/corpora/reuters/training/7453
  inflating: /root/nltk_data/corpora/reuters/training/7455
  inflating: /root/nltk_data/corpora/reuters/training/7459
  inflating: /root/nltk_data/corpora/reuters/training/7460
  inflating: /root/nltk_data/corpora/reuters/training/7462
  inflating: /root/nltk_data/corpora/reuters/training/7464
  inflating: /root/nltk_data/corpora/reuters/training/7465
  inflating: /root/nltk_data/corpora/reuters/training/7466
  inflating: /root/nltk_data/corpora/reuters/training/7468
  inflating: /root/nltk_data/corpora/reuters/training/7469
  inflating: /root/nltk_data/corpora/reuters/training/7470
  inflating: /root/nltk_data/corpora/reuters/training/7471
  inflating: /root/nltk_data/corpora/reuters/training/7472
  inflating: /root/nltk_data/corpora/reuters/training/7474
  inflating: /root/nltk_data/corpora/reuters/training/7475
  inflating: /root/nltk_data/corpora/reuters/training/7476
  inflating: /root/nltk_data/corpora/reuters/training/7477
  inflating: /root/nltk_data/corpora/reuters/training/748
  inflating: /root/nltk_data/corpora/reuters/training/7481
  inflating: /root/nltk_data/corpora/reuters/training/7482
  inflating: /root/nltk_data/corpora/reuters/training/7483
  inflating: /root/nltk_data/corpora/reuters/training/7484
  inflating: /root/nltk_data/corpora/reuters/training/7487
  inflating: /root/nltk_data/corpora/reuters/training/7488
  inflating: /root/nltk_data/corpora/reuters/training/7490
  inflating: /root/nltk_data/corpora/reuters/training/7491
  inflating: /root/nltk_data/corpora/reuters/training/7494
  inflating: /root/nltk_data/corpora/reuters/training/7495
  inflating: /root/nltk_data/corpora/reuters/training/7496
  inflating: /root/nltk_data/corpora/reuters/training/7497
  inflating: /root/nltk_data/corpora/reuters/training/7498
  inflating: /root/nltk_data/corpora/reuters/training/7499
  inflating: /root/nltk_data/corpora/reuters/training/75
  inflating: /root/nltk_data/corpora/reuters/training/7500
  inflating: /root/nltk_data/corpora/reuters/training/7501
  inflating: /root/nltk_data/corpora/reuters/training/7504
  inflating: /root/nltk_data/corpora/reuters/training/7505
  inflating: /root/nltk_data/corpora/reuters/training/7508
  inflating: /root/nltk_data/corpora/reuters/training/7510
  inflating: /root/nltk_data/corpora/reuters/training/7512
  inflating: /root/nltk_data/corpora/reuters/training/7514
  inflating: /root/nltk_data/corpora/reuters/training/7515
  inflating: /root/nltk_data/corpora/reuters/training/7516
  inflating: /root/nltk_data/corpora/reuters/training/7517
  inflating: /root/nltk_data/corpora/reuters/training/7518
  inflating: /root/nltk_data/corpora/reuters/training/7519
  inflating: /root/nltk_data/corpora/reuters/training/7520
  inflating: /root/nltk_data/corpora/reuters/training/7521
  inflating: /root/nltk_data/corpora/reuters/training/7528
  inflating: /root/nltk_data/corpora/reuters/training/7529
  inflating: /root/nltk_data/corpora/reuters/training/7531
  inflating: /root/nltk_data/corpora/reuters/training/7533
  inflating: /root/nltk_data/corpora/reuters/training/7534
  inflating: /root/nltk_data/corpora/reuters/training/7537
  inflating: /root/nltk_data/corpora/reuters/training/7538
  inflating: /root/nltk_data/corpora/reuters/training/7539
  inflating: /root/nltk_data/corpora/reuters/training/754
  inflating: /root/nltk_data/corpora/reuters/training/7541
  inflating: /root/nltk_data/corpora/reuters/training/7543
  inflating: /root/nltk_data/corpora/reuters/training/7544
  inflating: /root/nltk_data/corpora/reuters/training/7545
  inflating: /root/nltk_data/corpora/reuters/training/7548
  inflating: /root/nltk_data/corpora/reuters/training/7552
  inflating: /root/nltk_data/corpora/reuters/training/7554
  inflating: /root/nltk_data/corpora/reuters/training/7555
  inflating: /root/nltk_data/corpora/reuters/training/7557
  inflating: /root/nltk_data/corpora/reuters/training/7558
  inflating: /root/nltk_data/corpora/reuters/training/7559
  inflating: /root/nltk_data/corpora/reuters/training/756
  inflating: /root/nltk_data/corpora/reuters/training/7562
  inflating: /root/nltk_data/corpora/reuters/training/7564
  inflating: /root/nltk_data/corpora/reuters/training/7565
  inflating: /root/nltk_data/corpora/reuters/training/7566
  inflating: /root/nltk_data/corpora/reuters/training/7568
  inflating: /root/nltk_data/corpora/reuters/training/7569
  inflating: /root/nltk_data/corpora/reuters/training/757
  inflating: /root/nltk_data/corpora/reuters/training/7570
  inflating: /root/nltk_data/corpora/reuters/training/7571
  inflating: /root/nltk_data/corpora/reuters/training/7574
  inflating: /root/nltk_data/corpora/reuters/training/7576
  inflating: /root/nltk_data/corpora/reuters/training/7577
  inflating: /root/nltk_data/corpora/reuters/training/7579
  inflating: /root/nltk_data/corpora/reuters/training/758
  inflating: /root/nltk_data/corpora/reuters/training/7580
  inflating: /root/nltk_data/corpora/reuters/training/7581
  inflating: /root/nltk_data/corpora/reuters/training/7583
  inflating: /root/nltk_data/corpora/reuters/training/7584
  inflating: /root/nltk_data/corpora/reuters/training/7585
  inflating: /root/nltk_data/corpora/reuters/training/7587
  inflating: /root/nltk_data/corpora/reuters/training/7590
  inflating: /root/nltk_data/corpora/reuters/training/7591
  inflating: /root/nltk_data/corpora/reuters/training/7592
  inflating: /root/nltk_data/corpora/reuters/training/7593
  inflating: /root/nltk_data/corpora/reuters/training/7596
  inflating: /root/nltk_data/corpora/reuters/training/7599
  inflating: /root/nltk_data/corpora/reuters/training/76
  inflating: /root/nltk_data/corpora/reuters/training/7600
  inflating: /root/nltk_data/corpora/reuters/training/7602
  inflating: /root/nltk_data/corpora/reuters/training/7603
  inflating: /root/nltk_data/corpora/reuters/training/7604
  inflating: /root/nltk_data/corpora/reuters/training/7605
  inflating: /root/nltk_data/corpora/reuters/training/7606
  inflating: /root/nltk_data/corpora/reuters/training/7609
  inflating: /root/nltk_data/corpora/reuters/training/7611
  inflating: /root/nltk_data/corpora/reuters/training/7614
  inflating: /root/nltk_data/corpora/reuters/training/7618
  inflating: /root/nltk_data/corpora/reuters/training/762
  inflating: /root/nltk_data/corpora/reuters/training/7620
  inflating: /root/nltk_data/corpora/reuters/training/7621
  inflating: /root/nltk_data/corpora/reuters/training/7622
  inflating: /root/nltk_data/corpora/reuters/training/7623
  inflating: /root/nltk_data/corpora/reuters/training/7625
  inflating: /root/nltk_data/corpora/reuters/training/7626
  inflating: /root/nltk_data/corpora/reuters/training/7627
  inflating: /root/nltk_data/corpora/reuters/training/7628
  inflating: /root/nltk_data/corpora/reuters/training/7629
```

```
inflating: /root/nltk_data/corpora/reuters/training/7631
inflating: /root/nltk_data/corpora/reuters/training/7632
inflating: /root/nltk_data/corpora/reuters/training/7633
inflating: /root/nltk_data/corpora/reuters/training/7634
inflating: /root/nltk_data/corpora/reuters/training/7635
inflating: /root/nltk_data/corpora/reuters/training/7636
inflating: /root/nltk_data/corpora/reuters/training/7637
inflating: /root/nltk_data/corpora/reuters/training/7638
inflating: /root/nltk_data/corpora/reuters/training/7639
inflating: /root/nltk_data/corpora/reuters/training/764
inflating: /root/nltk_data/corpora/reuters/training/7640
inflating: /root/nltk_data/corpora/reuters/training/7641
inflating: /root/nltk_data/corpora/reuters/training/7642
inflating: /root/nltk_data/corpora/reuters/training/7643
inflating: /root/nltk_data/corpora/reuters/training/7644
inflating: /root/nltk_data/corpora/reuters/training/7645
inflating: /root/nltk_data/corpora/reuters/training/7647
inflating: /root/nltk_data/corpora/reuters/training/7650
inflating: /root/nltk_data/corpora/reuters/training/7652
inflating: /root/nltk_data/corpora/reuters/training/7657
inflating: /root/nltk_data/corpora/reuters/training/7659
inflating: /root/nltk_data/corpora/reuters/training/7662
inflating: /root/nltk_data/corpora/reuters/training/7669
inflating: /root/nltk_data/corpora/reuters/training/767
inflating: /root/nltk_data/corpora/reuters/training/7670
inflating: /root/nltk_data/corpora/reuters/training/7671
inflating: /root/nltk_data/corpora/reuters/training/7672
inflating: /root/nltk_data/corpora/reuters/training/7673
inflating: /root/nltk_data/corpora/reuters/training/7674
inflating: /root/nltk_data/corpora/reuters/training/7676
inflating: /root/nltk_data/corpora/reuters/training/7678
inflating: /root/nltk_data/corpora/reuters/training/7679
inflating: /root/nltk_data/corpora/reuters/training/768
inflating: /root/nltk_data/corpora/reuters/training/7681
inflating: /root/nltk_data/corpora/reuters/training/7682
inflating: /root/nltk_data/corpora/reuters/training/7684
inflating: /root/nltk_data/corpora/reuters/training/7686
inflating: /root/nltk_data/corpora/reuters/training/7687
inflating: /root/nltk_data/corpora/reuters/training/7690
inflating: /root/nltk_data/corpora/reuters/training/7693
inflating: /root/nltk_data/corpora/reuters/training/7694
inflating: /root/nltk_data/corpora/reuters/training/7696
inflating: /root/nltk_data/corpora/reuters/training/7698
inflating: /root/nltk_data/corpora/reuters/training/77
inflating: /root/nltk_data/corpora/reuters/training/7700
inflating: /root/nltk_data/corpora/reuters/training/7702
inflating: /root/nltk_data/corpora/reuters/training/7704
inflating: /root/nltk_data/corpora/reuters/training/7707
inflating: /root/nltk_data/corpora/reuters/training/7709
inflating: /root/nltk_data/corpora/reuters/training/7710
inflating: /root/nltk_data/corpora/reuters/training/7715
inflating: /root/nltk_data/corpora/reuters/training/7716
inflating: /root/nltk_data/corpora/reuters/training/7722
inflating: /root/nltk_data/corpora/reuters/training/7723
inflating: /root/nltk_data/corpora/reuters/training/7725
inflating: /root/nltk_data/corpora/reuters/training/7726
inflating: /root/nltk_data/corpora/reuters/training/7728
inflating: /root/nltk_data/corpora/reuters/training/7729
inflating: /root/nltk_data/corpora/reuters/training/773
inflating: /root/nltk_data/corpora/reuters/training/7730
inflating: /root/nltk_data/corpora/reuters/training/7732
inflating: /root/nltk_data/corpora/reuters/training/7733
inflating: /root/nltk_data/corpora/reuters/training/7734
inflating: /root/nltk_data/corpora/reuters/training/7735
inflating: /root/nltk_data/corpora/reuters/training/7736
inflating: /root/nltk_data/corpora/reuters/training/7737
inflating: /root/nltk_data/corpora/reuters/training/7738
inflating: /root/nltk_data/corpora/reuters/training/7739
inflating: /root/nltk_data/corpora/reuters/training/774
inflating: /root/nltk_data/corpora/reuters/training/7742
inflating: /root/nltk_data/corpora/reuters/training/7743
inflating: /root/nltk_data/corpora/reuters/training/7744
inflating: /root/nltk_data/corpora/reuters/training/7746
inflating: /root/nltk_data/corpora/reuters/training/7747
inflating: /root/nltk_data/corpora/reuters/training/7749
inflating: /root/nltk_data/corpora/reuters/training/7752
inflating: /root/nltk_data/corpora/reuters/training/7755
inflating: /root/nltk_data/corpora/reuters/training/7757
inflating: /root/nltk_data/corpora/reuters/training/776
inflating: /root/nltk_data/corpora/reuters/training/7761
inflating: /root/nltk_data/corpora/reuters/training/7762
inflating: /root/nltk_data/corpora/reuters/training/7763
inflating: /root/nltk_data/corpora/reuters/training/7764
inflating: /root/nltk_data/corpora/reuters/training/7765
inflating: /root/nltk_data/corpora/reuters/training/7766
inflating: /root/nltk_data/corpora/reuters/training/7767
inflating: /root/nltk_data/corpora/reuters/training/7768
inflating: /root/nltk_data/corpora/reuters/training/7769
inflating: /root/nltk_data/corpora/reuters/training/7771
inflating: /root/nltk_data/corpora/reuters/training/7773
inflating: /root/nltk_data/corpora/reuters/training/7774
inflating: /root/nltk_data/corpora/reuters/training/7775
inflating: /root/nltk_data/corpora/reuters/training/7776
inflating: /root/nltk_data/corpora/reuters/training/778
inflating: /root/nltk_data/corpora/reuters/training/7782
inflating: /root/nltk_data/corpora/reuters/training/7784
inflating: /root/nltk_data/corpora/reuters/training/7785
inflating: /root/nltk_data/corpora/reuters/training/7789
inflating: /root/nltk_data/corpora/reuters/training/7790
inflating: /root/nltk_data/corpora/reuters/training/7791
inflating: /root/nltk_data/corpora/reuters/training/7792
inflating: /root/nltk_data/corpora/reuters/training/7795
inflating: /root/nltk_data/corpora/reuters/training/7796
inflating: /root/nltk_data/corpora/reuters/training/7797
inflating: /root/nltk_data/corpora/reuters/training/78
inflating: /root/nltk_data/corpora/reuters/training/780
inflating: /root/nltk_data/corpora/reuters/training/7802
inflating: /root/nltk_data/corpora/reuters/training/7804
inflating: /root/nltk_data/corpora/reuters/training/7805
inflating: /root/nltk_data/corpora/reuters/training/7807
inflating: /root/nltk_data/corpora/reuters/training/7809
inflating: /root/nltk_data/corpora/reuters/training/781
inflating: /root/nltk_data/corpora/reuters/training/7810
inflating: /root/nltk_data/corpora/reuters/training/7813
inflating: /root/nltk_data/corpora/reuters/training/7814
inflating: /root/nltk_data/corpora/reuters/training/7816
inflating: /root/nltk_data/corpora/reuters/training/7817
inflating: /root/nltk_data/corpora/reuters/training/7818
inflating: /root/nltk_data/corpora/reuters/training/7819
inflating: /root/nltk_data/corpora/reuters/training/7820
inflating: /root/nltk_data/corpora/reuters/training/7821
inflating: /root/nltk_data/corpora/reuters/training/7823
inflating: /root/nltk_data/corpora/reuters/training/7824
inflating: /root/nltk_data/corpora/reuters/training/7825
inflating: /root/nltk_data/corpora/reuters/training/7826
inflating: /root/nltk_data/corpora/reuters/training/7828
inflating: /root/nltk_data/corpora/reuters/training/783
inflating: /root/nltk_data/corpora/reuters/training/7830
inflating: /root/nltk_data/corpora/reuters/training/7831
inflating: /root/nltk_data/corpora/reuters/training/7835
inflating: /root/nltk_data/corpora/reuters/training/7836
inflating: /root/nltk_data/corpora/reuters/training/7837
inflating: /root/nltk_data/corpora/reuters/training/7838
inflating: /root/nltk_data/corpora/reuters/training/7839
```

```
inflating: /root/nltk_data/corpora/reuters/training/7840
inflating: /root/nltk_data/corpora/reuters/training/7841
inflating: /root/nltk_data/corpora/reuters/training/7842
inflating: /root/nltk_data/corpora/reuters/training/7843
inflating: /root/nltk_data/corpora/reuters/training/7845
inflating: /root/nltk_data/corpora/reuters/training/7847
inflating: /root/nltk_data/corpora/reuters/training/7848
inflating: /root/nltk_data/corpora/reuters/training/7850
inflating: /root/nltk_data/corpora/reuters/training/7851
inflating: /root/nltk_data/corpora/reuters/training/7854
inflating: /root/nltk_data/corpora/reuters/training/7855
inflating: /root/nltk_data/corpora/reuters/training/7857
inflating: /root/nltk_data/corpora/reuters/training/7858
inflating: /root/nltk_data/corpora/reuters/training/7859
inflating: /root/nltk_data/corpora/reuters/training/7860
inflating: /root/nltk_data/corpora/reuters/training/7861
inflating: /root/nltk_data/corpora/reuters/training/7862
inflating: /root/nltk_data/corpora/reuters/training/7864
inflating: /root/nltk_data/corpora/reuters/training/7865
inflating: /root/nltk_data/corpora/reuters/training/7866
inflating: /root/nltk_data/corpora/reuters/training/7867
inflating: /root/nltk_data/corpora/reuters/training/7869
inflating: /root/nltk_data/corpora/reuters/training/787
inflating: /root/nltk_data/corpora/reuters/training/7870
inflating: /root/nltk_data/corpora/reuters/training/7871
inflating: /root/nltk_data/corpora/reuters/training/7872
inflating: /root/nltk_data/corpora/reuters/training/7873
inflating: /root/nltk_data/corpora/reuters/training/7874
inflating: /root/nltk_data/corpora/reuters/training/7875
inflating: /root/nltk_data/corpora/reuters/training/7876
inflating: /root/nltk_data/corpora/reuters/training/7877
inflating: /root/nltk_data/corpora/reuters/training/7878
inflating: /root/nltk_data/corpora/reuters/training/7879
inflating: /root/nltk_data/corpora/reuters/training/7880
inflating: /root/nltk_data/corpora/reuters/training/7881
inflating: /root/nltk_data/corpora/reuters/training/7883
inflating: /root/nltk_data/corpora/reuters/training/7884
inflating: /root/nltk_data/corpora/reuters/training/7887
inflating: /root/nltk_data/corpora/reuters/training/7888
inflating: /root/nltk_data/corpora/reuters/training/7891
inflating: /root/nltk_data/corpora/reuters/training/7893
inflating: /root/nltk_data/corpora/reuters/training/7894
inflating: /root/nltk_data/corpora/reuters/training/7895
inflating: /root/nltk_data/corpora/reuters/training/7896
inflating: /root/nltk_data/corpora/reuters/training/7897
inflating: /root/nltk_data/corpora/reuters/training/7899
inflating: /root/nltk_data/corpora/reuters/training/7901
inflating: /root/nltk_data/corpora/reuters/training/7902
inflating: /root/nltk_data/corpora/reuters/training/7903
inflating: /root/nltk_data/corpora/reuters/training/7904
inflating: /root/nltk_data/corpora/reuters/training/7907
inflating: /root/nltk_data/corpora/reuters/training/791
inflating: /root/nltk_data/corpora/reuters/training/7912
inflating: /root/nltk_data/corpora/reuters/training/7913
inflating: /root/nltk_data/corpora/reuters/training/7914
inflating: /root/nltk_data/corpora/reuters/training/7917
inflating: /root/nltk_data/corpora/reuters/training/7918
inflating: /root/nltk_data/corpora/reuters/training/7920
inflating: /root/nltk_data/corpora/reuters/training/7921
inflating: /root/nltk_data/corpora/reuters/training/7922
inflating: /root/nltk_data/corpora/reuters/training/7924
inflating: /root/nltk_data/corpora/reuters/training/7927
inflating: /root/nltk_data/corpora/reuters/training/7928
inflating: /root/nltk_data/corpora/reuters/training/793
inflating: /root/nltk_data/corpora/reuters/training/7930
inflating: /root/nltk_data/corpora/reuters/training/7934
inflating: /root/nltk_data/corpora/reuters/training/7935
inflating: /root/nltk_data/corpora/reuters/training/7937
inflating: /root/nltk_data/corpora/reuters/training/7939
inflating: /root/nltk_data/corpora/reuters/training/7940
inflating: /root/nltk_data/corpora/reuters/training/7942
inflating: /root/nltk_data/corpora/reuters/training/7943
inflating: /root/nltk_data/corpora/reuters/training/7948
inflating: /root/nltk_data/corpora/reuters/training/7949
inflating: /root/nltk_data/corpora/reuters/training/7950
inflating: /root/nltk_data/corpora/reuters/training/7951
inflating: /root/nltk_data/corpora/reuters/training/7953
inflating: /root/nltk_data/corpora/reuters/training/7954
inflating: /root/nltk_data/corpora/reuters/training/7955
inflating: /root/nltk_data/corpora/reuters/training/7956
inflating: /root/nltk_data/corpora/reuters/training/7957
inflating: /root/nltk_data/corpora/reuters/training/7958
inflating: /root/nltk_data/corpora/reuters/training/7959
inflating: /root/nltk_data/corpora/reuters/training/7961
inflating: /root/nltk_data/corpora/reuters/training/7962
inflating: /root/nltk_data/corpora/reuters/training/7963
inflating: /root/nltk_data/corpora/reuters/training/7964
inflating: /root/nltk_data/corpora/reuters/training/7966
inflating: /root/nltk_data/corpora/reuters/training/7967
inflating: /root/nltk_data/corpora/reuters/training/7968
inflating: /root/nltk_data/corpora/reuters/training/7969
inflating: /root/nltk_data/corpora/reuters/training/797
inflating: /root/nltk_data/corpora/reuters/training/7970
inflating: /root/nltk_data/corpora/reuters/training/7971
inflating: /root/nltk_data/corpora/reuters/training/7972
inflating: /root/nltk_data/corpora/reuters/training/7974
inflating: /root/nltk_data/corpora/reuters/training/7975
inflating: /root/nltk_data/corpora/reuters/training/7977
inflating: /root/nltk_data/corpora/reuters/training/7979
inflating: /root/nltk_data/corpora/reuters/training/798
inflating: /root/nltk_data/corpora/reuters/training/7980
inflating: /root/nltk_data/corpora/reuters/training/7983
inflating: /root/nltk_data/corpora/reuters/training/7984
inflating: /root/nltk_data/corpora/reuters/training/7985
inflating: /root/nltk_data/corpora/reuters/training/7987
inflating: /root/nltk_data/corpora/reuters/training/7990
inflating: /root/nltk_data/corpora/reuters/training/7991
inflating: /root/nltk_data/corpora/reuters/training/7992
inflating: /root/nltk_data/corpora/reuters/training/7993
inflating: /root/nltk_data/corpora/reuters/training/7994
inflating: /root/nltk_data/corpora/reuters/training/7997
inflating: /root/nltk_data/corpora/reuters/training/80
inflating: /root/nltk_data/corpora/reuters/training/800
inflating: /root/nltk_data/corpora/reuters/training/8000
inflating: /root/nltk_data/corpora/reuters/training/8001
inflating: /root/nltk_data/corpora/reuters/training/8002
inflating: /root/nltk_data/corpora/reuters/training/8003
inflating: /root/nltk_data/corpora/reuters/training/8004
inflating: /root/nltk_data/corpora/reuters/training/8005
inflating: /root/nltk_data/corpora/reuters/training/8007
inflating: /root/nltk_data/corpora/reuters/training/8008
inflating: /root/nltk_data/corpora/reuters/training/8009
inflating: /root/nltk_data/corpora/reuters/training/8010
inflating: /root/nltk_data/corpora/reuters/training/8011
inflating: /root/nltk_data/corpora/reuters/training/8012
inflating: /root/nltk_data/corpora/reuters/training/8015
inflating: /root/nltk_data/corpora/reuters/training/8016
inflating: /root/nltk_data/corpora/reuters/training/8017
inflating: /root/nltk_data/corpora/reuters/training/8020
inflating: /root/nltk_data/corpora/reuters/training/8021
inflating: /root/nltk_data/corpora/reuters/training/8022
inflating: /root/nltk_data/corpora/reuters/training/8023
inflating: /root/nltk_data/corpora/reuters/training/8025
inflating: /root/nltk_data/corpora/reuters/training/8026
inflating: /root/nltk_data/corpora/reuters/training/8027
```

```
inflating: /root/nltk_data/corpora/reuters/training/8028
inflating: /root/nltk_data/corpora/reuters/training/8029
inflating: /root/nltk_data/corpora/reuters/training/8030
inflating: /root/nltk_data/corpora/reuters/training/8032
inflating: /root/nltk_data/corpora/reuters/training/8036
inflating: /root/nltk_data/corpora/reuters/training/8039
inflating: /root/nltk_data/corpora/reuters/training/8040
inflating: /root/nltk_data/corpora/reuters/training/8041
inflating: /root/nltk_data/corpora/reuters/training/8044
inflating: /root/nltk_data/corpora/reuters/training/8045
inflating: /root/nltk_data/corpora/reuters/training/8047
inflating: /root/nltk_data/corpora/reuters/training/8048
inflating: /root/nltk_data/corpora/reuters/training/8050
inflating: /root/nltk_data/corpora/reuters/training/8051
inflating: /root/nltk_data/corpora/reuters/training/8055
inflating: /root/nltk_data/corpora/reuters/training/8056
inflating: /root/nltk_data/corpora/reuters/training/8060
inflating: /root/nltk_data/corpora/reuters/training/8063
inflating: /root/nltk_data/corpora/reuters/training/8064
inflating: /root/nltk_data/corpora/reuters/training/8068
inflating: /root/nltk_data/corpora/reuters/training/8069
inflating: /root/nltk_data/corpora/reuters/training/8074
inflating: /root/nltk_data/corpora/reuters/training/8080
inflating: /root/nltk_data/corpora/reuters/training/8085
inflating: /root/nltk_data/corpora/reuters/training/8086
inflating: /root/nltk_data/corpora/reuters/training/8087
inflating: /root/nltk_data/corpora/reuters/training/8088
inflating: /root/nltk_data/corpora/reuters/training/8089
inflating: /root/nltk_data/corpora/reuters/training/8091
inflating: /root/nltk_data/corpora/reuters/training/8094
inflating: /root/nltk_data/corpora/reuters/training/8096
inflating: /root/nltk_data/corpora/reuters/training/8097
inflating: /root/nltk_data/corpora/reuters/training/8098
inflating: /root/nltk_data/corpora/reuters/training/81
inflating: /root/nltk_data/corpora/reuters/training/8100
inflating: /root/nltk_data/corpora/reuters/training/8102
inflating: /root/nltk_data/corpora/reuters/training/8103
inflating: /root/nltk_data/corpora/reuters/training/8105
inflating: /root/nltk_data/corpora/reuters/training/8106
inflating: /root/nltk_data/corpora/reuters/training/8108
inflating: /root/nltk_data/corpora/reuters/training/8109
inflating: /root/nltk_data/corpora/reuters/training/811
inflating: /root/nltk_data/corpora/reuters/training/8111
inflating: /root/nltk_data/corpora/reuters/training/8112
inflating: /root/nltk_data/corpora/reuters/training/8113
inflating: /root/nltk_data/corpora/reuters/training/8115
inflating: /root/nltk_data/corpora/reuters/training/8117
inflating: /root/nltk_data/corpora/reuters/training/8119
inflating: /root/nltk_data/corpora/reuters/training/8120
inflating: /root/nltk_data/corpora/reuters/training/8123
inflating: /root/nltk_data/corpora/reuters/training/8125
inflating: /root/nltk_data/corpora/reuters/training/8126
inflating: /root/nltk_data/corpora/reuters/training/8130
inflating: /root/nltk_data/corpora/reuters/training/8131
inflating: /root/nltk_data/corpora/reuters/training/8132
inflating: /root/nltk_data/corpora/reuters/training/8134
inflating: /root/nltk_data/corpora/reuters/training/8135
inflating: /root/nltk_data/corpora/reuters/training/8137
inflating: /root/nltk_data/corpora/reuters/training/8138
inflating: /root/nltk_data/corpora/reuters/training/8140
inflating: /root/nltk_data/corpora/reuters/training/8141
inflating: /root/nltk_data/corpora/reuters/training/8144
inflating: /root/nltk_data/corpora/reuters/training/8145
inflating: /root/nltk_data/corpora/reuters/training/8147
inflating: /root/nltk_data/corpora/reuters/training/8149
inflating: /root/nltk_data/corpora/reuters/training/8151
inflating: /root/nltk_data/corpora/reuters/training/8153
inflating: /root/nltk_data/corpora/reuters/training/8155
inflating: /root/nltk_data/corpora/reuters/training/8156
inflating: /root/nltk_data/corpora/reuters/training/8158
inflating: /root/nltk_data/corpora/reuters/training/8159
inflating: /root/nltk_data/corpora/reuters/training/816
inflating: /root/nltk_data/corpora/reuters/training/8160
inflating: /root/nltk_data/corpora/reuters/training/8161
inflating: /root/nltk_data/corpora/reuters/training/8162
inflating: /root/nltk_data/corpora/reuters/training/8163
inflating: /root/nltk_data/corpora/reuters/training/8164
inflating: /root/nltk_data/corpora/reuters/training/8166
inflating: /root/nltk_data/corpora/reuters/training/8167
inflating: /root/nltk_data/corpora/reuters/training/8168
inflating: /root/nltk_data/corpora/reuters/training/8169
inflating: /root/nltk_data/corpora/reuters/training/817
inflating: /root/nltk_data/corpora/reuters/training/8172
inflating: /root/nltk_data/corpora/reuters/training/8173
inflating: /root/nltk_data/corpora/reuters/training/8174
inflating: /root/nltk_data/corpora/reuters/training/8176
inflating: /root/nltk_data/corpora/reuters/training/8177
inflating: /root/nltk_data/corpora/reuters/training/8178
inflating: /root/nltk_data/corpora/reuters/training/8179
inflating: /root/nltk_data/corpora/reuters/training/818
inflating: /root/nltk_data/corpora/reuters/training/8180
inflating: /root/nltk_data/corpora/reuters/training/8182
inflating: /root/nltk_data/corpora/reuters/training/8184
inflating: /root/nltk_data/corpora/reuters/training/8186
inflating: /root/nltk_data/corpora/reuters/training/8188
inflating: /root/nltk_data/corpora/reuters/training/8189
inflating: /root/nltk_data/corpora/reuters/training/8190
inflating: /root/nltk_data/corpora/reuters/training/8193
inflating: /root/nltk_data/corpora/reuters/training/8194
inflating: /root/nltk_data/corpora/reuters/training/8197
inflating: /root/nltk_data/corpora/reuters/training/8198
inflating: /root/nltk_data/corpora/reuters/training/8199
inflating: /root/nltk_data/corpora/reuters/training/82
inflating: /root/nltk_data/corpora/reuters/training/820
inflating: /root/nltk_data/corpora/reuters/training/8200
inflating: /root/nltk_data/corpora/reuters/training/8204
inflating: /root/nltk_data/corpora/reuters/training/8205
inflating: /root/nltk_data/corpora/reuters/training/8206
inflating: /root/nltk_data/corpora/reuters/training/8209
inflating: /root/nltk_data/corpora/reuters/training/8210
inflating: /root/nltk_data/corpora/reuters/training/8211
inflating: /root/nltk_data/corpora/reuters/training/8212
inflating: /root/nltk_data/corpora/reuters/training/8213
inflating: /root/nltk_data/corpora/reuters/training/8214
inflating: /root/nltk_data/corpora/reuters/training/8217
inflating: /root/nltk_data/corpora/reuters/training/8218
inflating: /root/nltk_data/corpora/reuters/training/8219
inflating: /root/nltk_data/corpora/reuters/training/8220
inflating: /root/nltk_data/corpora/reuters/training/8221
inflating: /root/nltk_data/corpora/reuters/training/8222
inflating: /root/nltk_data/corpora/reuters/training/8223
inflating: /root/nltk_data/corpora/reuters/training/8224
inflating: /root/nltk_data/corpora/reuters/training/8225
inflating: /root/nltk_data/corpora/reuters/training/8226
inflating: /root/nltk_data/corpora/reuters/training/8227
inflating: /root/nltk_data/corpora/reuters/training/8229
inflating: /root/nltk_data/corpora/reuters/training/8234
inflating: /root/nltk_data/corpora/reuters/training/8236
inflating: /root/nltk_data/corpora/reuters/training/8240
inflating: /root/nltk_data/corpora/reuters/training/8244
inflating: /root/nltk_data/corpora/reuters/training/8245
inflating: /root/nltk_data/corpora/reuters/training/8246
inflating: /root/nltk_data/corpora/reuters/training/8247
inflating: /root/nltk_data/corpora/reuters/training/8252
```

```
inflating: /root/nltk_data/corpora/reuters/training/8253
inflating: /root/nltk_data/corpora/reuters/training/8254
inflating: /root/nltk_data/corpora/reuters/training/8255
inflating: /root/nltk_data/corpora/reuters/training/8257
inflating: /root/nltk_data/corpora/reuters/training/8258
inflating: /root/nltk_data/corpora/reuters/training/8259
inflating: /root/nltk_data/corpora/reuters/training/8262
inflating: /root/nltk_data/corpora/reuters/training/8264
inflating: /root/nltk_data/corpora/reuters/training/8269
inflating: /root/nltk_data/corpora/reuters/training/8270
inflating: /root/nltk_data/corpora/reuters/training/8273
inflating: /root/nltk_data/corpora/reuters/training/8275
inflating: /root/nltk_data/corpora/reuters/training/8276
inflating: /root/nltk_data/corpora/reuters/training/8279
inflating: /root/nltk_data/corpora/reuters/training/828
inflating: /root/nltk_data/corpora/reuters/training/8286
inflating: /root/nltk_data/corpora/reuters/training/8287
inflating: /root/nltk_data/corpora/reuters/training/8288
inflating: /root/nltk_data/corpora/reuters/training/8289
inflating: /root/nltk_data/corpora/reuters/training/829
inflating: /root/nltk_data/corpora/reuters/training/8291
inflating: /root/nltk_data/corpora/reuters/training/8292
inflating: /root/nltk_data/corpora/reuters/training/8293
inflating: /root/nltk_data/corpora/reuters/training/8294
inflating: /root/nltk_data/corpora/reuters/training/8296
inflating: /root/nltk_data/corpora/reuters/training/8297
inflating: /root/nltk_data/corpora/reuters/training/8298
inflating: /root/nltk_data/corpora/reuters/training/8299
inflating: /root/nltk_data/corpora/reuters/training/83
inflating: /root/nltk_data/corpora/reuters/training/830
inflating: /root/nltk_data/corpora/reuters/training/8300
inflating: /root/nltk_data/corpora/reuters/training/8302
inflating: /root/nltk_data/corpora/reuters/training/8306
inflating: /root/nltk_data/corpora/reuters/training/8308
inflating: /root/nltk_data/corpora/reuters/training/8309
inflating: /root/nltk_data/corpora/reuters/training/8311
inflating: /root/nltk_data/corpora/reuters/training/8313
inflating: /root/nltk_data/corpora/reuters/training/8314
inflating: /root/nltk_data/corpora/reuters/training/8317
inflating: /root/nltk_data/corpora/reuters/training/8318
inflating: /root/nltk_data/corpora/reuters/training/8319
inflating: /root/nltk_data/corpora/reuters/training/8320
inflating: /root/nltk_data/corpora/reuters/training/8322
inflating: /root/nltk_data/corpora/reuters/training/8326
inflating: /root/nltk_data/corpora/reuters/training/8327
inflating: /root/nltk_data/corpora/reuters/training/833
inflating: /root/nltk_data/corpora/reuters/training/8331
inflating: /root/nltk_data/corpora/reuters/training/8334
inflating: /root/nltk_data/corpora/reuters/training/8335
inflating: /root/nltk_data/corpora/reuters/training/8337
inflating: /root/nltk_data/corpora/reuters/training/8338
inflating: /root/nltk_data/corpora/reuters/training/8339
inflating: /root/nltk_data/corpora/reuters/training/834
inflating: /root/nltk_data/corpora/reuters/training/8340
inflating: /root/nltk_data/corpora/reuters/training/8342
inflating: /root/nltk_data/corpora/reuters/training/8343
inflating: /root/nltk_data/corpora/reuters/training/8344
inflating: /root/nltk_data/corpora/reuters/training/8347
inflating: /root/nltk_data/corpora/reuters/training/8349
inflating: /root/nltk_data/corpora/reuters/training/835
inflating: /root/nltk_data/corpora/reuters/training/8350
inflating: /root/nltk_data/corpora/reuters/training/8351
inflating: /root/nltk_data/corpora/reuters/training/8353
inflating: /root/nltk_data/corpora/reuters/training/8355
inflating: /root/nltk_data/corpora/reuters/training/8359
inflating: /root/nltk_data/corpora/reuters/training/836
inflating: /root/nltk_data/corpora/reuters/training/8361
inflating: /root/nltk_data/corpora/reuters/training/8363
inflating: /root/nltk_data/corpora/reuters/training/8364
inflating: /root/nltk_data/corpora/reuters/training/8366
inflating: /root/nltk_data/corpora/reuters/training/8367
inflating: /root/nltk_data/corpora/reuters/training/8368
inflating: /root/nltk_data/corpora/reuters/training/837
inflating: /root/nltk_data/corpora/reuters/training/8370
inflating: /root/nltk_data/corpora/reuters/training/8373
inflating: /root/nltk_data/corpora/reuters/training/8374
inflating: /root/nltk_data/corpora/reuters/training/8378
inflating: /root/nltk_data/corpora/reuters/training/8387
inflating: /root/nltk_data/corpora/reuters/training/8388
inflating: /root/nltk_data/corpora/reuters/training/839
inflating: /root/nltk_data/corpora/reuters/training/8391
inflating: /root/nltk_data/corpora/reuters/training/8394
inflating: /root/nltk_data/corpora/reuters/training/8396
inflating: /root/nltk_data/corpora/reuters/training/8399
inflating: /root/nltk_data/corpora/reuters/training/840
inflating: /root/nltk_data/corpora/reuters/training/8400
inflating: /root/nltk_data/corpora/reuters/training/8401
inflating: /root/nltk_data/corpora/reuters/training/8402
inflating: /root/nltk_data/corpora/reuters/training/8405
inflating: /root/nltk_data/corpora/reuters/training/8406
inflating: /root/nltk_data/corpora/reuters/training/8407
inflating: /root/nltk_data/corpora/reuters/training/8409
inflating: /root/nltk_data/corpora/reuters/training/8413
inflating: /root/nltk_data/corpora/reuters/training/8415
inflating: /root/nltk_data/corpora/reuters/training/8416
inflating: /root/nltk_data/corpora/reuters/training/842
inflating: /root/nltk_data/corpora/reuters/training/8421
inflating: /root/nltk_data/corpora/reuters/training/8422
inflating: /root/nltk_data/corpora/reuters/training/8424
inflating: /root/nltk_data/corpora/reuters/training/8425
inflating: /root/nltk_data/corpora/reuters/training/8426
inflating: /root/nltk_data/corpora/reuters/training/8427
inflating: /root/nltk_data/corpora/reuters/training/8428
inflating: /root/nltk_data/corpora/reuters/training/8429
inflating: /root/nltk_data/corpora/reuters/training/843
inflating: /root/nltk_data/corpora/reuters/training/8430
inflating: /root/nltk_data/corpora/reuters/training/8432
inflating: /root/nltk_data/corpora/reuters/training/8433
inflating: /root/nltk_data/corpora/reuters/training/8435
inflating: /root/nltk_data/corpora/reuters/training/8438
inflating: /root/nltk_data/corpora/reuters/training/844
inflating: /root/nltk_data/corpora/reuters/training/8440
inflating: /root/nltk_data/corpora/reuters/training/8441
inflating: /root/nltk_data/corpora/reuters/training/8443
inflating: /root/nltk_data/corpora/reuters/training/8446
inflating: /root/nltk_data/corpora/reuters/training/8447
inflating: /root/nltk_data/corpora/reuters/training/8448
inflating: /root/nltk_data/corpora/reuters/training/8449
inflating: /root/nltk_data/corpora/reuters/training/8451
inflating: /root/nltk_data/corpora/reuters/training/8453
inflating: /root/nltk_data/corpora/reuters/training/8454
inflating: /root/nltk_data/corpora/reuters/training/8457
inflating: /root/nltk_data/corpora/reuters/training/8458
inflating: /root/nltk_data/corpora/reuters/training/846
inflating: /root/nltk_data/corpora/reuters/training/8467
inflating: /root/nltk_data/corpora/reuters/training/8472
inflating: /root/nltk_data/corpora/reuters/training/8475
inflating: /root/nltk_data/corpora/reuters/training/8478
inflating: /root/nltk_data/corpora/reuters/training/8479
inflating: /root/nltk_data/corpora/reuters/training/848
inflating: /root/nltk_data/corpora/reuters/training/8480
inflating: /root/nltk_data/corpora/reuters/training/8481
inflating: /root/nltk_data/corpora/reuters/training/8487
inflating: /root/nltk_data/corpora/reuters/training/8488
inflating: /root/nltk_data/corpora/reuters/training/8489
```

```
inflating: /root/nltk_data/corpora/reuters/training/8489
inflating: /root/nltk_data/corpora/reuters/training/849
inflating: /root/nltk_data/corpora/reuters/training/8490
inflating: /root/nltk_data/corpora/reuters/training/8491
inflating: /root/nltk_data/corpora/reuters/training/8493
inflating: /root/nltk_data/corpora/reuters/training/8496
inflating: /root/nltk_data/corpora/reuters/training/8497
inflating: /root/nltk_data/corpora/reuters/training/85
inflating: /root/nltk_data/corpora/reuters/training/850
inflating: /root/nltk_data/corpora/reuters/training/8501
inflating: /root/nltk_data/corpora/reuters/training/8502
inflating: /root/nltk_data/corpora/reuters/training/8506
inflating: /root/nltk_data/corpora/reuters/training/8509
inflating: /root/nltk_data/corpora/reuters/training/851
inflating: /root/nltk_data/corpora/reuters/training/8510
inflating: /root/nltk_data/corpora/reuters/training/8512
inflating: /root/nltk_data/corpora/reuters/training/8513
inflating: /root/nltk_data/corpora/reuters/training/8514
inflating: /root/nltk_data/corpora/reuters/training/8516
inflating: /root/nltk_data/corpora/reuters/training/8518
inflating: /root/nltk_data/corpora/reuters/training/852
inflating: /root/nltk_data/corpora/reuters/training/8520
inflating: /root/nltk_data/corpora/reuters/training/8521
inflating: /root/nltk_data/corpora/reuters/training/8522
inflating: /root/nltk_data/corpora/reuters/training/8526
inflating: /root/nltk_data/corpora/reuters/training/8528
inflating: /root/nltk_data/corpora/reuters/training/853
inflating: /root/nltk_data/corpora/reuters/training/8530
inflating: /root/nltk_data/corpora/reuters/training/8533
inflating: /root/nltk_data/corpora/reuters/training/8534
inflating: /root/nltk_data/corpora/reuters/training/8535
inflating: /root/nltk_data/corpora/reuters/training/8536
inflating: /root/nltk_data/corpora/reuters/training/8540
inflating: /root/nltk_data/corpora/reuters/training/8541
inflating: /root/nltk_data/corpora/reuters/training/8543
inflating: /root/nltk_data/corpora/reuters/training/8544
inflating: /root/nltk_data/corpora/reuters/training/8545
inflating: /root/nltk_data/corpora/reuters/training/8546
inflating: /root/nltk_data/corpora/reuters/training/8547
inflating: /root/nltk_data/corpora/reuters/training/855
inflating: /root/nltk_data/corpora/reuters/training/8553
inflating: /root/nltk_data/corpora/reuters/training/8554
inflating: /root/nltk_data/corpora/reuters/training/8555
inflating: /root/nltk_data/corpora/reuters/training/8556
inflating: /root/nltk_data/corpora/reuters/training/8558
inflating: /root/nltk_data/corpora/reuters/training/8559
inflating: /root/nltk_data/corpora/reuters/training/856
inflating: /root/nltk_data/corpora/reuters/training/8561
inflating: /root/nltk_data/corpora/reuters/training/8562
inflating: /root/nltk_data/corpora/reuters/training/8563
inflating: /root/nltk_data/corpora/reuters/training/8565
inflating: /root/nltk_data/corpora/reuters/training/8567
inflating: /root/nltk_data/corpora/reuters/training/8569
inflating: /root/nltk_data/corpora/reuters/training/857
inflating: /root/nltk_data/corpora/reuters/training/8571
inflating: /root/nltk_data/corpora/reuters/training/8572
inflating: /root/nltk_data/corpora/reuters/training/8574
inflating: /root/nltk_data/corpora/reuters/training/8577
inflating: /root/nltk_data/corpora/reuters/training/8578
inflating: /root/nltk_data/corpora/reuters/training/8579
inflating: /root/nltk_data/corpora/reuters/training/8580
inflating: /root/nltk_data/corpora/reuters/training/8581
inflating: /root/nltk_data/corpora/reuters/training/8582
inflating: /root/nltk_data/corpora/reuters/training/8584
inflating: /root/nltk_data/corpora/reuters/training/8585
inflating: /root/nltk_data/corpora/reuters/training/8586
inflating: /root/nltk_data/corpora/reuters/training/8587
inflating: /root/nltk_data/corpora/reuters/training/8588
inflating: /root/nltk_data/corpora/reuters/training/859
inflating: /root/nltk_data/corpora/reuters/training/8590
inflating: /root/nltk_data/corpora/reuters/training/8591
inflating: /root/nltk_data/corpora/reuters/training/8592
inflating: /root/nltk_data/corpora/reuters/training/8595
inflating: /root/nltk_data/corpora/reuters/training/8596
inflating: /root/nltk_data/corpora/reuters/training/8597
inflating: /root/nltk_data/corpora/reuters/training/8598
inflating: /root/nltk_data/corpora/reuters/training/8599
inflating: /root/nltk_data/corpora/reuters/training/86
inflating: /root/nltk_data/corpora/reuters/training/8600
inflating: /root/nltk_data/corpora/reuters/training/8602
inflating: /root/nltk_data/corpora/reuters/training/8603
inflating: /root/nltk_data/corpora/reuters/training/8604
inflating: /root/nltk_data/corpora/reuters/training/8605
inflating: /root/nltk_data/corpora/reuters/training/8606
inflating: /root/nltk_data/corpora/reuters/training/8607
inflating: /root/nltk_data/corpora/reuters/training/8608
inflating: /root/nltk_data/corpora/reuters/training/8609
inflating: /root/nltk_data/corpora/reuters/training/861
inflating: /root/nltk_data/corpora/reuters/training/8610
inflating: /root/nltk_data/corpora/reuters/training/8611
inflating: /root/nltk_data/corpora/reuters/training/8612
inflating: /root/nltk_data/corpora/reuters/training/8613
inflating: /root/nltk_data/corpora/reuters/training/8614
inflating: /root/nltk_data/corpora/reuters/training/8615
inflating: /root/nltk_data/corpora/reuters/training/8616
inflating: /root/nltk_data/corpora/reuters/training/8617
inflating: /root/nltk_data/corpora/reuters/training/8618
inflating: /root/nltk_data/corpora/reuters/training/8619
inflating: /root/nltk_data/corpora/reuters/training/8621
inflating: /root/nltk_data/corpora/reuters/training/8623
inflating: /root/nltk_data/corpora/reuters/training/8624
inflating: /root/nltk_data/corpora/reuters/training/8628
inflating: /root/nltk_data/corpora/reuters/training/8629
inflating: /root/nltk_data/corpora/reuters/training/8630
inflating: /root/nltk_data/corpora/reuters/training/8632
inflating: /root/nltk_data/corpora/reuters/training/8635
inflating: /root/nltk_data/corpora/reuters/training/8636
inflating: /root/nltk_data/corpora/reuters/training/8637
inflating: /root/nltk_data/corpora/reuters/training/8638
inflating: /root/nltk_data/corpora/reuters/training/8641
inflating: /root/nltk_data/corpora/reuters/training/8643
inflating: /root/nltk_data/corpora/reuters/training/8644
inflating: /root/nltk_data/corpora/reuters/training/865
inflating: /root/nltk_data/corpora/reuters/training/8654
inflating: /root/nltk_data/corpora/reuters/training/8656
inflating: /root/nltk_data/corpora/reuters/training/8657
inflating: /root/nltk_data/corpora/reuters/training/8658
inflating: /root/nltk_data/corpora/reuters/training/866
inflating: /root/nltk_data/corpora/reuters/training/8661
inflating: /root/nltk_data/corpora/reuters/training/8662
inflating: /root/nltk_data/corpora/reuters/training/8663
inflating: /root/nltk_data/corpora/reuters/training/8664
inflating: /root/nltk_data/corpora/reuters/training/8665
inflating: /root/nltk_data/corpora/reuters/training/8666
inflating: /root/nltk_data/corpora/reuters/training/8667
inflating: /root/nltk_data/corpora/reuters/training/8668
inflating: /root/nltk_data/corpora/reuters/training/867
inflating: /root/nltk_data/corpora/reuters/training/8670
inflating: /root/nltk_data/corpora/reuters/training/8671
inflating: /root/nltk_data/corpora/reuters/training/8672
inflating: /root/nltk_data/corpora/reuters/training/8673
inflating: /root/nltk_data/corpora/reuters/training/8674
inflating: /root/nltk_data/corpora/reuters/training/8675
inflating: /root/nltk_data/corpora/reuters/training/8676
inflating: /root/nltk_data/corpora/reuters/training/8677
```

```
inflating: /root/nltk_data/corpora/reuters/training/8678
inflating: /root/nltk_data/corpora/reuters/training/868
inflating: /root/nltk_data/corpora/reuters/training/8681
inflating: /root/nltk_data/corpora/reuters/training/8683
inflating: /root/nltk_data/corpora/reuters/training/8684
inflating: /root/nltk_data/corpora/reuters/training/8685
inflating: /root/nltk_data/corpora/reuters/training/8686
inflating: /root/nltk_data/corpora/reuters/training/8688
inflating: /root/nltk_data/corpora/reuters/training/8689
inflating: /root/nltk_data/corpora/reuters/training/8690
inflating: /root/nltk_data/corpora/reuters/training/8691
inflating: /root/nltk_data/corpora/reuters/training/8692
inflating: /root/nltk_data/corpora/reuters/training/8694
inflating: /root/nltk_data/corpora/reuters/training/8696
inflating: /root/nltk_data/corpora/reuters/training/8699
inflating: /root/nltk_data/corpora/reuters/training/87
inflating: /root/nltk_data/corpora/reuters/training/870
inflating: /root/nltk_data/corpora/reuters/training/8701
inflating: /root/nltk_data/corpora/reuters/training/8702
inflating: /root/nltk_data/corpora/reuters/training/8704
inflating: /root/nltk_data/corpora/reuters/training/8705
inflating: /root/nltk_data/corpora/reuters/training/8706
inflating: /root/nltk_data/corpora/reuters/training/8707
inflating: /root/nltk_data/corpora/reuters/training/8708
inflating: /root/nltk_data/corpora/reuters/training/8709
inflating: /root/nltk_data/corpora/reuters/training/8710
inflating: /root/nltk_data/corpora/reuters/training/8711
inflating: /root/nltk_data/corpora/reuters/training/8713
inflating: /root/nltk_data/corpora/reuters/training/8714
inflating: /root/nltk_data/corpora/reuters/training/8716
inflating: /root/nltk_data/corpora/reuters/training/8717
inflating: /root/nltk_data/corpora/reuters/training/872
inflating: /root/nltk_data/corpora/reuters/training/8722
inflating: /root/nltk_data/corpora/reuters/training/8723
inflating: /root/nltk_data/corpora/reuters/training/8724
inflating: /root/nltk_data/corpora/reuters/training/8725
inflating: /root/nltk_data/corpora/reuters/training/8726
inflating: /root/nltk_data/corpora/reuters/training/8728
inflating: /root/nltk_data/corpora/reuters/training/873
inflating: /root/nltk_data/corpora/reuters/training/8730
inflating: /root/nltk_data/corpora/reuters/training/8731
inflating: /root/nltk_data/corpora/reuters/training/8732
inflating: /root/nltk_data/corpora/reuters/training/8733
inflating: /root/nltk_data/corpora/reuters/training/8734
inflating: /root/nltk_data/corpora/reuters/training/8735
inflating: /root/nltk_data/corpora/reuters/training/8736
inflating: /root/nltk_data/corpora/reuters/training/8738
inflating: /root/nltk_data/corpora/reuters/training/874
inflating: /root/nltk_data/corpora/reuters/training/8740
inflating: /root/nltk_data/corpora/reuters/training/8741
inflating: /root/nltk_data/corpora/reuters/training/8742
inflating: /root/nltk_data/corpora/reuters/training/8744
inflating: /root/nltk_data/corpora/reuters/training/8745
inflating: /root/nltk_data/corpora/reuters/training/8746
inflating: /root/nltk_data/corpora/reuters/training/8747
inflating: /root/nltk_data/corpora/reuters/training/8748
inflating: /root/nltk_data/corpora/reuters/training/8749
inflating: /root/nltk_data/corpora/reuters/training/875
inflating: /root/nltk_data/corpora/reuters/training/8750
inflating: /root/nltk_data/corpora/reuters/training/8755
inflating: /root/nltk_data/corpora/reuters/training/8756
inflating: /root/nltk_data/corpora/reuters/training/8757
inflating: /root/nltk_data/corpora/reuters/training/8758
inflating: /root/nltk_data/corpora/reuters/training/8759
inflating: /root/nltk_data/corpora/reuters/training/8760
inflating: /root/nltk_data/corpora/reuters/training/8761
inflating: /root/nltk_data/corpora/reuters/training/8765
inflating: /root/nltk_data/corpora/reuters/training/8767
inflating: /root/nltk_data/corpora/reuters/training/8768
inflating: /root/nltk_data/corpora/reuters/training/8771
inflating: /root/nltk_data/corpora/reuters/training/8777
inflating: /root/nltk_data/corpora/reuters/training/8778
inflating: /root/nltk_data/corpora/reuters/training/8780
inflating: /root/nltk_data/corpora/reuters/training/8781
inflating: /root/nltk_data/corpora/reuters/training/8784
inflating: /root/nltk_data/corpora/reuters/training/8785
inflating: /root/nltk_data/corpora/reuters/training/8788
inflating: /root/nltk_data/corpora/reuters/training/879
inflating: /root/nltk_data/corpora/reuters/training/8791
inflating: /root/nltk_data/corpora/reuters/training/8793
inflating: /root/nltk_data/corpora/reuters/training/8794
inflating: /root/nltk_data/corpora/reuters/training/8795
inflating: /root/nltk_data/corpora/reuters/training/8796
inflating: /root/nltk_data/corpora/reuters/training/8797
inflating: /root/nltk_data/corpora/reuters/training/8799
inflating: /root/nltk_data/corpora/reuters/training/880
inflating: /root/nltk_data/corpora/reuters/training/8800
inflating: /root/nltk_data/corpora/reuters/training/8802
inflating: /root/nltk_data/corpora/reuters/training/8806
inflating: /root/nltk_data/corpora/reuters/training/8808
inflating: /root/nltk_data/corpora/reuters/training/8810
inflating: /root/nltk_data/corpora/reuters/training/8812
inflating: /root/nltk_data/corpora/reuters/training/8814
inflating: /root/nltk_data/corpora/reuters/training/8815
inflating: /root/nltk_data/corpora/reuters/training/882
inflating: /root/nltk_data/corpora/reuters/training/8820
inflating: /root/nltk_data/corpora/reuters/training/8822
inflating: /root/nltk_data/corpora/reuters/training/8823
inflating: /root/nltk_data/corpora/reuters/training/8826
inflating: /root/nltk_data/corpora/reuters/training/8829
inflating: /root/nltk_data/corpora/reuters/training/883
inflating: /root/nltk_data/corpora/reuters/training/8830
inflating: /root/nltk_data/corpora/reuters/training/8835
inflating: /root/nltk_data/corpora/reuters/training/8836
inflating: /root/nltk_data/corpora/reuters/training/8837
inflating: /root/nltk_data/corpora/reuters/training/884
inflating: /root/nltk_data/corpora/reuters/training/8845
inflating: /root/nltk_data/corpora/reuters/training/8846
inflating: /root/nltk_data/corpora/reuters/training/8849
inflating: /root/nltk_data/corpora/reuters/training/885
inflating: /root/nltk_data/corpora/reuters/training/8850
inflating: /root/nltk_data/corpora/reuters/training/8853
inflating: /root/nltk_data/corpora/reuters/training/8854
inflating: /root/nltk_data/corpora/reuters/training/8855
inflating: /root/nltk_data/corpora/reuters/training/8856
inflating: /root/nltk_data/corpora/reuters/training/8857
inflating: /root/nltk_data/corpora/reuters/training/8858
inflating: /root/nltk_data/corpora/reuters/training/8859
inflating: /root/nltk_data/corpora/reuters/training/886
inflating: /root/nltk_data/corpora/reuters/training/8860
inflating: /root/nltk_data/corpora/reuters/training/8861
inflating: /root/nltk_data/corpora/reuters/training/8863
inflating: /root/nltk_data/corpora/reuters/training/8865
inflating: /root/nltk_data/corpora/reuters/training/8866
inflating: /root/nltk_data/corpora/reuters/training/8868
inflating: /root/nltk_data/corpora/reuters/training/887
inflating: /root/nltk_data/corpora/reuters/training/8873
inflating: /root/nltk_data/corpora/reuters/training/8877
inflating: /root/nltk_data/corpora/reuters/training/8879
inflating: /root/nltk_data/corpora/reuters/training/8880
inflating: /root/nltk_data/corpora/reuters/training/8882
inflating: /root/nltk_data/corpora/reuters/training/8884
inflating: /root/nltk_data/corpora/reuters/training/889
inflating: /root/nltk_data/corpora/reuters/training/8890
```

```
inflating: /root/nltk_data/corpora/reuters/training/8892
inflating: /root/nltk_data/corpora/reuters/training/8893
inflating: /root/nltk_data/corpora/reuters/training/8895
inflating: /root/nltk_data/corpora/reuters/training/8898
inflating: /root/nltk_data/corpora/reuters/training/8899
inflating: /root/nltk_data/corpora/reuters/training/89
inflating: /root/nltk_data/corpora/reuters/training/890
inflating: /root/nltk_data/corpora/reuters/training/8902
inflating: /root/nltk_data/corpora/reuters/training/8903
inflating: /root/nltk_data/corpora/reuters/training/8905
inflating: /root/nltk_data/corpora/reuters/training/8906
inflating: /root/nltk_data/corpora/reuters/training/8908
inflating: /root/nltk_data/corpora/reuters/training/8909
inflating: /root/nltk_data/corpora/reuters/training/8911
inflating: /root/nltk_data/corpora/reuters/training/8912
inflating: /root/nltk_data/corpora/reuters/training/8914
inflating: /root/nltk_data/corpora/reuters/training/8915
inflating: /root/nltk_data/corpora/reuters/training/8916
inflating: /root/nltk_data/corpora/reuters/training/8918
inflating: /root/nltk_data/corpora/reuters/training/8919
inflating: /root/nltk_data/corpora/reuters/training/8922
inflating: /root/nltk_data/corpora/reuters/training/8925
inflating: /root/nltk_data/corpora/reuters/training/8928
inflating: /root/nltk_data/corpora/reuters/training/8930
inflating: /root/nltk_data/corpora/reuters/training/8932
inflating: /root/nltk_data/corpora/reuters/training/8933
inflating: /root/nltk_data/corpora/reuters/training/8935
inflating: /root/nltk_data/corpora/reuters/training/894
inflating: /root/nltk_data/corpora/reuters/training/8940
inflating: /root/nltk_data/corpora/reuters/training/8941
inflating: /root/nltk_data/corpora/reuters/training/8943
inflating: /root/nltk_data/corpora/reuters/training/8944
inflating: /root/nltk_data/corpora/reuters/training/8945
inflating: /root/nltk_data/corpora/reuters/training/8946
inflating: /root/nltk_data/corpora/reuters/training/8947
inflating: /root/nltk_data/corpora/reuters/training/8948
inflating: /root/nltk_data/corpora/reuters/training/895
inflating: /root/nltk_data/corpora/reuters/training/8950
inflating: /root/nltk_data/corpora/reuters/training/8951
inflating: /root/nltk_data/corpora/reuters/training/8956
inflating: /root/nltk_data/corpora/reuters/training/8959
inflating: /root/nltk_data/corpora/reuters/training/896
inflating: /root/nltk_data/corpora/reuters/training/8961
inflating: /root/nltk_data/corpora/reuters/training/8962
inflating: /root/nltk_data/corpora/reuters/training/8964
inflating: /root/nltk_data/corpora/reuters/training/8965
inflating: /root/nltk_data/corpora/reuters/training/8968
inflating: /root/nltk_data/corpora/reuters/training/8969
inflating: /root/nltk_data/corpora/reuters/training/897
inflating: /root/nltk_data/corpora/reuters/training/8970
inflating: /root/nltk_data/corpora/reuters/training/8971
inflating: /root/nltk_data/corpora/reuters/training/8972
inflating: /root/nltk_data/corpora/reuters/training/8974
inflating: /root/nltk_data/corpora/reuters/training/8975
inflating: /root/nltk_data/corpora/reuters/training/8976
inflating: /root/nltk_data/corpora/reuters/training/8977
inflating: /root/nltk_data/corpora/reuters/training/8978
inflating: /root/nltk_data/corpora/reuters/training/8979
inflating: /root/nltk_data/corpora/reuters/training/8981
inflating: /root/nltk_data/corpora/reuters/training/8982
inflating: /root/nltk_data/corpora/reuters/training/8983
inflating: /root/nltk_data/corpora/reuters/training/8984
inflating: /root/nltk_data/corpora/reuters/training/8985
inflating: /root/nltk_data/corpora/reuters/training/8986
inflating: /root/nltk_data/corpora/reuters/training/8987
inflating: /root/nltk_data/corpora/reuters/training/899
inflating: /root/nltk_data/corpora/reuters/training/8991
inflating: /root/nltk_data/corpora/reuters/training/8992
inflating: /root/nltk_data/corpora/reuters/training/8993
inflating: /root/nltk_data/corpora/reuters/training/8995
inflating: /root/nltk_data/corpora/reuters/training/8996
inflating: /root/nltk_data/corpora/reuters/training/8997
inflating: /root/nltk_data/corpora/reuters/training/8998
inflating: /root/nltk_data/corpora/reuters/training/8999
inflating: /root/nltk_data/corpora/reuters/training/9
inflating: /root/nltk_data/corpora/reuters/training/900
inflating: /root/nltk_data/corpora/reuters/training/9003
inflating: /root/nltk_data/corpora/reuters/training/9007
inflating: /root/nltk_data/corpora/reuters/training/901
inflating: /root/nltk_data/corpora/reuters/training/9010
inflating: /root/nltk_data/corpora/reuters/training/9012
inflating: /root/nltk_data/corpora/reuters/training/9014
inflating: /root/nltk_data/corpora/reuters/training/9015
inflating: /root/nltk_data/corpora/reuters/training/9018
inflating: /root/nltk_data/corpora/reuters/training/9020
inflating: /root/nltk_data/corpora/reuters/training/9021
inflating: /root/nltk_data/corpora/reuters/training/9022
inflating: /root/nltk_data/corpora/reuters/training/9025
inflating: /root/nltk_data/corpora/reuters/training/9027
inflating: /root/nltk_data/corpora/reuters/training/9029
inflating: /root/nltk_data/corpora/reuters/training/903
inflating: /root/nltk_data/corpora/reuters/training/9030
inflating: /root/nltk_data/corpora/reuters/training/9031
inflating: /root/nltk_data/corpora/reuters/training/9032
inflating: /root/nltk_data/corpora/reuters/training/9033
inflating: /root/nltk_data/corpora/reuters/training/9034
inflating: /root/nltk_data/corpora/reuters/training/9036
inflating: /root/nltk_data/corpora/reuters/training/9039
inflating: /root/nltk_data/corpora/reuters/training/904
inflating: /root/nltk_data/corpora/reuters/training/9040
inflating: /root/nltk_data/corpora/reuters/training/9041
inflating: /root/nltk_data/corpora/reuters/training/9044
inflating: /root/nltk_data/corpora/reuters/training/9045
inflating: /root/nltk_data/corpora/reuters/training/9047
inflating: /root/nltk_data/corpora/reuters/training/9048
inflating: /root/nltk_data/corpora/reuters/training/9049
inflating: /root/nltk_data/corpora/reuters/training/9051
inflating: /root/nltk_data/corpora/reuters/training/9053
inflating: /root/nltk_data/corpora/reuters/training/9054
inflating: /root/nltk_data/corpora/reuters/training/9055
inflating: /root/nltk_data/corpora/reuters/training/9056
inflating: /root/nltk_data/corpora/reuters/training/9058
inflating: /root/nltk_data/corpora/reuters/training/9059
inflating: /root/nltk_data/corpora/reuters/training/9060
inflating: /root/nltk_data/corpora/reuters/training/9061
inflating: /root/nltk_data/corpora/reuters/training/9064
inflating: /root/nltk_data/corpora/reuters/training/9065
inflating: /root/nltk_data/corpora/reuters/training/9066
inflating: /root/nltk_data/corpora/reuters/training/9067
inflating: /root/nltk_data/corpora/reuters/training/9068
inflating: /root/nltk_data/corpora/reuters/training/9069
inflating: /root/nltk_data/corpora/reuters/training/9072
inflating: /root/nltk_data/corpora/reuters/training/9073
inflating: /root/nltk_data/corpora/reuters/training/9074
inflating: /root/nltk_data/corpora/reuters/training/9075
inflating: /root/nltk_data/corpora/reuters/training/9076
inflating: /root/nltk_data/corpora/reuters/training/9077
inflating: /root/nltk_data/corpora/reuters/training/9079
inflating: /root/nltk_data/corpora/reuters/training/908
inflating: /root/nltk_data/corpora/reuters/training/9081
inflating: /root/nltk_data/corpora/reuters/training/9087
inflating: /root/nltk_data/corpora/reuters/training/9088
inflating: /root/nltk_data/corpora/reuters/training/9090
inflating: /root/nltk_data/corpora/reuters/training/9091
inflating: /root/nltk_data/corpora/reuters/training/9093
```

```
inflating: /root/nltk_data/corpora/reuters/training/9094
inflating: /root/nltk_data/corpora/reuters/training/9095
inflating: /root/nltk_data/corpora/reuters/training/9097
inflating: /root/nltk_data/corpora/reuters/training/9098
inflating: /root/nltk_data/corpora/reuters/training/9100
inflating: /root/nltk_data/corpora/reuters/training/9101
inflating: /root/nltk_data/corpora/reuters/training/9102
inflating: /root/nltk_data/corpora/reuters/training/9103
inflating: /root/nltk_data/corpora/reuters/training/9104
inflating: /root/nltk_data/corpora/reuters/training/9107
inflating: /root/nltk_data/corpora/reuters/training/9108
inflating: /root/nltk_data/corpora/reuters/training/9110
inflating: /root/nltk_data/corpora/reuters/training/9111
inflating: /root/nltk_data/corpora/reuters/training/9112
inflating: /root/nltk_data/corpora/reuters/training/9113
inflating: /root/nltk_data/corpora/reuters/training/9115
inflating: /root/nltk_data/corpora/reuters/training/9116
inflating: /root/nltk_data/corpora/reuters/training/9117
inflating: /root/nltk_data/corpora/reuters/training/9118
inflating: /root/nltk_data/corpora/reuters/training/9120
inflating: /root/nltk_data/corpora/reuters/training/9124
inflating: /root/nltk_data/corpora/reuters/training/9127
inflating: /root/nltk_data/corpora/reuters/training/9128
inflating: /root/nltk_data/corpora/reuters/training/9129
inflating: /root/nltk_data/corpora/reuters/training/913
inflating: /root/nltk_data/corpora/reuters/training/9131
inflating: /root/nltk_data/corpora/reuters/training/9132
inflating: /root/nltk_data/corpora/reuters/training/9133
inflating: /root/nltk_data/corpora/reuters/training/9134
inflating: /root/nltk_data/corpora/reuters/training/9135
inflating: /root/nltk_data/corpora/reuters/training/9137
inflating: /root/nltk_data/corpora/reuters/training/9138
inflating: /root/nltk_data/corpora/reuters/training/9139
inflating: /root/nltk_data/corpora/reuters/training/9142
inflating: /root/nltk_data/corpora/reuters/training/9143
inflating: /root/nltk_data/corpora/reuters/training/9146
inflating: /root/nltk_data/corpora/reuters/training/9147
inflating: /root/nltk_data/corpora/reuters/training/9149
inflating: /root/nltk_data/corpora/reuters/training/915
inflating: /root/nltk_data/corpora/reuters/training/9152
inflating: /root/nltk_data/corpora/reuters/training/9153
inflating: /root/nltk_data/corpora/reuters/training/9155
inflating: /root/nltk_data/corpora/reuters/training/9156
inflating: /root/nltk_data/corpora/reuters/training/9158
inflating: /root/nltk_data/corpora/reuters/training/9159
inflating: /root/nltk_data/corpora/reuters/training/9160
inflating: /root/nltk_data/corpora/reuters/training/9161
inflating: /root/nltk_data/corpora/reuters/training/9162
inflating: /root/nltk_data/corpora/reuters/training/9163
inflating: /root/nltk_data/corpora/reuters/training/9164
inflating: /root/nltk_data/corpora/reuters/training/9165
inflating: /root/nltk_data/corpora/reuters/training/9166
inflating: /root/nltk_data/corpora/reuters/training/9170
inflating: /root/nltk_data/corpora/reuters/training/9171
inflating: /root/nltk_data/corpora/reuters/training/9172
inflating: /root/nltk_data/corpora/reuters/training/9175
inflating: /root/nltk_data/corpora/reuters/training/9177
inflating: /root/nltk_data/corpora/reuters/training/9178
inflating: /root/nltk_data/corpora/reuters/training/918
inflating: /root/nltk_data/corpora/reuters/training/9181
inflating: /root/nltk_data/corpora/reuters/training/9182
inflating: /root/nltk_data/corpora/reuters/training/9184
inflating: /root/nltk_data/corpora/reuters/training/9187
inflating: /root/nltk_data/corpora/reuters/training/9189
inflating: /root/nltk_data/corpora/reuters/training/9190
inflating: /root/nltk_data/corpora/reuters/training/9192
inflating: /root/nltk_data/corpora/reuters/training/9193
inflating: /root/nltk_data/corpora/reuters/training/9194
inflating: /root/nltk_data/corpora/reuters/training/9195
inflating: /root/nltk_data/corpora/reuters/training/9196
inflating: /root/nltk_data/corpora/reuters/training/9197
inflating: /root/nltk_data/corpora/reuters/training/92
inflating: /root/nltk_data/corpora/reuters/training/920
inflating: /root/nltk_data/corpora/reuters/training/9201
inflating: /root/nltk_data/corpora/reuters/training/9202
inflating: /root/nltk_data/corpora/reuters/training/9203
inflating: /root/nltk_data/corpora/reuters/training/9204
inflating: /root/nltk_data/corpora/reuters/training/9205
inflating: /root/nltk_data/corpora/reuters/training/9206
inflating: /root/nltk_data/corpora/reuters/training/9207
inflating: /root/nltk_data/corpora/reuters/training/9208
inflating: /root/nltk_data/corpora/reuters/training/921
inflating: /root/nltk_data/corpora/reuters/training/9210
inflating: /root/nltk_data/corpora/reuters/training/9213
inflating: /root/nltk_data/corpora/reuters/training/9214
inflating: /root/nltk_data/corpora/reuters/training/9216
inflating: /root/nltk_data/corpora/reuters/training/9217
inflating: /root/nltk_data/corpora/reuters/training/9218
inflating: /root/nltk_data/corpora/reuters/training/922
inflating: /root/nltk_data/corpora/reuters/training/9220
inflating: /root/nltk_data/corpora/reuters/training/9222
inflating: /root/nltk_data/corpora/reuters/training/9224
inflating: /root/nltk_data/corpora/reuters/training/9228
inflating: /root/nltk_data/corpora/reuters/training/9229
inflating: /root/nltk_data/corpora/reuters/training/9230
inflating: /root/nltk_data/corpora/reuters/training/9231
inflating: /root/nltk_data/corpora/reuters/training/9232
inflating: /root/nltk_data/corpora/reuters/training/9234
inflating: /root/nltk_data/corpora/reuters/training/9235
inflating: /root/nltk_data/corpora/reuters/training/9237
inflating: /root/nltk_data/corpora/reuters/training/9238
inflating: /root/nltk_data/corpora/reuters/training/9239
inflating: /root/nltk_data/corpora/reuters/training/924
inflating: /root/nltk_data/corpora/reuters/training/9240
inflating: /root/nltk_data/corpora/reuters/training/9242
inflating: /root/nltk_data/corpora/reuters/training/9243
inflating: /root/nltk_data/corpora/reuters/training/9246
inflating: /root/nltk_data/corpora/reuters/training/925
inflating: /root/nltk_data/corpora/reuters/training/9250
inflating: /root/nltk_data/corpora/reuters/training/9252
inflating: /root/nltk_data/corpora/reuters/training/9253
inflating: /root/nltk_data/corpora/reuters/training/9255
inflating: /root/nltk_data/corpora/reuters/training/9258
inflating: /root/nltk_data/corpora/reuters/training/9259
inflating: /root/nltk_data/corpora/reuters/training/926
inflating: /root/nltk_data/corpora/reuters/training/9260
inflating: /root/nltk_data/corpora/reuters/training/9261
inflating: /root/nltk_data/corpora/reuters/training/9263
inflating: /root/nltk_data/corpora/reuters/training/9264
inflating: /root/nltk_data/corpora/reuters/training/9265
inflating: /root/nltk_data/corpora/reuters/training/9266
inflating: /root/nltk_data/corpora/reuters/training/9267
inflating: /root/nltk_data/corpora/reuters/training/9268
inflating: /root/nltk_data/corpora/reuters/training/9270
inflating: /root/nltk_data/corpora/reuters/training/9271
inflating: /root/nltk_data/corpora/reuters/training/9272
inflating: /root/nltk_data/corpora/reuters/training/9273
inflating: /root/nltk_data/corpora/reuters/training/9274
inflating: /root/nltk_data/corpora/reuters/training/9277
inflating: /root/nltk_data/corpora/reuters/training/9278
inflating: /root/nltk_data/corpora/reuters/training/9279
inflating: /root/nltk_data/corpora/reuters/training/9280
inflating: /root/nltk_data/corpora/reuters/training/9282
inflating: /root/nltk_data/corpora/reuters/training/9283
```

```
inflating: /root/nltk_data/corpora/reuters/training/9284
inflating: /root/nltk_data/corpora/reuters/training/9285
inflating: /root/nltk_data/corpora/reuters/training/9287
inflating: /root/nltk_data/corpora/reuters/training/9288
inflating: /root/nltk_data/corpora/reuters/training/929
inflating: /root/nltk_data/corpora/reuters/training/9290
inflating: /root/nltk_data/corpora/reuters/training/9293
inflating: /root/nltk_data/corpora/reuters/training/9294
inflating: /root/nltk_data/corpora/reuters/training/9295
inflating: /root/nltk_data/corpora/reuters/training/9296
inflating: /root/nltk_data/corpora/reuters/training/9299
inflating: /root/nltk_data/corpora/reuters/training/93
inflating: /root/nltk_data/corpora/reuters/training/930
inflating: /root/nltk_data/corpora/reuters/training/9301
inflating: /root/nltk_data/corpora/reuters/training/9302
inflating: /root/nltk_data/corpora/reuters/training/9304
inflating: /root/nltk_data/corpora/reuters/training/9305
inflating: /root/nltk_data/corpora/reuters/training/9306
inflating: /root/nltk_data/corpora/reuters/training/9308
inflating: /root/nltk_data/corpora/reuters/training/9311
inflating: /root/nltk_data/corpora/reuters/training/9312
inflating: /root/nltk_data/corpora/reuters/training/9313
inflating: /root/nltk_data/corpora/reuters/training/9314
inflating: /root/nltk_data/corpora/reuters/training/9315
inflating: /root/nltk_data/corpora/reuters/training/9316
inflating: /root/nltk_data/corpora/reuters/training/932
inflating: /root/nltk_data/corpora/reuters/training/9324
inflating: /root/nltk_data/corpora/reuters/training/9326
inflating: /root/nltk_data/corpora/reuters/training/9327
inflating: /root/nltk_data/corpora/reuters/training/9328
inflating: /root/nltk_data/corpora/reuters/training/9330
inflating: /root/nltk_data/corpora/reuters/training/9332
inflating: /root/nltk_data/corpora/reuters/training/9333
inflating: /root/nltk_data/corpora/reuters/training/9334
inflating: /root/nltk_data/corpora/reuters/training/9335
inflating: /root/nltk_data/corpora/reuters/training/9336
inflating: /root/nltk_data/corpora/reuters/training/9337
inflating: /root/nltk_data/corpora/reuters/training/9338
inflating: /root/nltk_data/corpora/reuters/training/9339
inflating: /root/nltk_data/corpora/reuters/training/934
inflating: /root/nltk_data/corpora/reuters/training/9343
inflating: /root/nltk_data/corpora/reuters/training/9345
inflating: /root/nltk_data/corpora/reuters/training/9346
inflating: /root/nltk_data/corpora/reuters/training/9348
inflating: /root/nltk_data/corpora/reuters/training/9349
inflating: /root/nltk_data/corpora/reuters/training/9352
inflating: /root/nltk_data/corpora/reuters/training/9354
inflating: /root/nltk_data/corpora/reuters/training/9356
inflating: /root/nltk_data/corpora/reuters/training/9357
inflating: /root/nltk_data/corpora/reuters/training/9361
inflating: /root/nltk_data/corpora/reuters/training/9362
inflating: /root/nltk_data/corpora/reuters/training/9363
inflating: /root/nltk_data/corpora/reuters/training/9364
inflating: /root/nltk_data/corpora/reuters/training/9365
inflating: /root/nltk_data/corpora/reuters/training/9366
inflating: /root/nltk_data/corpora/reuters/training/9369
inflating: /root/nltk_data/corpora/reuters/training/937
inflating: /root/nltk_data/corpora/reuters/training/9371
inflating: /root/nltk_data/corpora/reuters/training/9372
inflating: /root/nltk_data/corpora/reuters/training/9373
inflating: /root/nltk_data/corpora/reuters/training/9374
inflating: /root/nltk_data/corpora/reuters/training/9377
inflating: /root/nltk_data/corpora/reuters/training/9378
inflating: /root/nltk_data/corpora/reuters/training/938
inflating: /root/nltk_data/corpora/reuters/training/9381
inflating: /root/nltk_data/corpora/reuters/training/9382
inflating: /root/nltk_data/corpora/reuters/training/9383
inflating: /root/nltk_data/corpora/reuters/training/9385
inflating: /root/nltk_data/corpora/reuters/training/9388
inflating: /root/nltk_data/corpora/reuters/training/9389
inflating: /root/nltk_data/corpora/reuters/training/939
inflating: /root/nltk_data/corpora/reuters/training/9392
inflating: /root/nltk_data/corpora/reuters/training/9393
inflating: /root/nltk_data/corpora/reuters/training/9397
inflating: /root/nltk_data/corpora/reuters/training/9398
inflating: /root/nltk_data/corpora/reuters/training/9399
inflating: /root/nltk_data/corpora/reuters/training/94
inflating: /root/nltk_data/corpora/reuters/training/940
inflating: /root/nltk_data/corpora/reuters/training/9403
inflating: /root/nltk_data/corpora/reuters/training/9405
inflating: /root/nltk_data/corpora/reuters/training/9406
inflating: /root/nltk_data/corpora/reuters/training/9408
inflating: /root/nltk_data/corpora/reuters/training/941
inflating: /root/nltk_data/corpora/reuters/training/9410
inflating: /root/nltk_data/corpora/reuters/training/9412
inflating: /root/nltk_data/corpora/reuters/training/9413
inflating: /root/nltk_data/corpora/reuters/training/9414
inflating: /root/nltk_data/corpora/reuters/training/9415
inflating: /root/nltk_data/corpora/reuters/training/9417
inflating: /root/nltk_data/corpora/reuters/training/942
inflating: /root/nltk_data/corpora/reuters/training/9422
inflating: /root/nltk_data/corpora/reuters/training/9423
inflating: /root/nltk_data/corpora/reuters/training/9425
inflating: /root/nltk_data/corpora/reuters/training/9426
inflating: /root/nltk_data/corpora/reuters/training/9427
inflating: /root/nltk_data/corpora/reuters/training/9428
inflating: /root/nltk_data/corpora/reuters/training/9429
inflating: /root/nltk_data/corpora/reuters/training/943
inflating: /root/nltk_data/corpora/reuters/training/9431
inflating: /root/nltk_data/corpora/reuters/training/9432
inflating: /root/nltk_data/corpora/reuters/training/9433
inflating: /root/nltk_data/corpora/reuters/training/9434
inflating: /root/nltk_data/corpora/reuters/training/9435
inflating: /root/nltk_data/corpora/reuters/training/9436
inflating: /root/nltk_data/corpora/reuters/training/9437
inflating: /root/nltk_data/corpora/reuters/training/9438
inflating: /root/nltk_data/corpora/reuters/training/944
inflating: /root/nltk_data/corpora/reuters/training/9441
inflating: /root/nltk_data/corpora/reuters/training/9443
inflating: /root/nltk_data/corpora/reuters/training/9444
inflating: /root/nltk_data/corpora/reuters/training/9445
inflating: /root/nltk_data/corpora/reuters/training/945
inflating: /root/nltk_data/corpora/reuters/training/9450
inflating: /root/nltk_data/corpora/reuters/training/9453
inflating: /root/nltk_data/corpora/reuters/training/9454
inflating: /root/nltk_data/corpora/reuters/training/9463
inflating: /root/nltk_data/corpora/reuters/training/9465
inflating: /root/nltk_data/corpora/reuters/training/9469
inflating: /root/nltk_data/corpora/reuters/training/9470
inflating: /root/nltk_data/corpora/reuters/training/9472
inflating: /root/nltk_data/corpora/reuters/training/9473
inflating: /root/nltk_data/corpora/reuters/training/9475
inflating: /root/nltk_data/corpora/reuters/training/9477
inflating: /root/nltk_data/corpora/reuters/training/9478
inflating: /root/nltk_data/corpora/reuters/training/9479
inflating: /root/nltk_data/corpora/reuters/training/9484
inflating: /root/nltk_data/corpora/reuters/training/9486
inflating: /root/nltk_data/corpora/reuters/training/9487
inflating: /root/nltk_data/corpora/reuters/training/9488
inflating: /root/nltk_data/corpora/reuters/training/9489
inflating: /root/nltk_data/corpora/reuters/training/949
inflating: /root/nltk_data/corpora/reuters/training/9491
inflating: /root/nltk_data/corpora/reuters/training/9493
inflating: /root/nltk_data/corpora/reuters/training/9497
```

```
inflating: /root/nltk_data/corpora/reuters/training/9498
inflating: /root/nltk_data/corpora/reuters/training/9499
inflating: /root/nltk_data/corpora/reuters/training/95
inflating: /root/nltk_data/corpora/reuters/training/9502
inflating: /root/nltk_data/corpora/reuters/training/9503
inflating: /root/nltk_data/corpora/reuters/training/9505
inflating: /root/nltk_data/corpora/reuters/training/9509
inflating: /root/nltk_data/corpora/reuters/training/9510
inflating: /root/nltk_data/corpora/reuters/training/9511
inflating: /root/nltk_data/corpora/reuters/training/9515
inflating: /root/nltk_data/corpora/reuters/training/9519
inflating: /root/nltk_data/corpora/reuters/training/952
inflating: /root/nltk_data/corpora/reuters/training/9521
inflating: /root/nltk_data/corpora/reuters/training/9525
inflating: /root/nltk_data/corpora/reuters/training/9526
inflating: /root/nltk_data/corpora/reuters/training/9527
inflating: /root/nltk_data/corpora/reuters/training/953
inflating: /root/nltk_data/corpora/reuters/training/9530
inflating: /root/nltk_data/corpora/reuters/training/9531
inflating: /root/nltk_data/corpora/reuters/training/9532
inflating: /root/nltk_data/corpora/reuters/training/9534
inflating: /root/nltk_data/corpora/reuters/training/9535
inflating: /root/nltk_data/corpora/reuters/training/9538
inflating: /root/nltk_data/corpora/reuters/training/9539
inflating: /root/nltk_data/corpora/reuters/training/9543
inflating: /root/nltk_data/corpora/reuters/training/9544
inflating: /root/nltk_data/corpora/reuters/training/9547
inflating: /root/nltk_data/corpora/reuters/training/9549
inflating: /root/nltk_data/corpora/reuters/training/955
inflating: /root/nltk_data/corpora/reuters/training/9550
inflating: /root/nltk_data/corpora/reuters/training/9551
inflating: /root/nltk_data/corpora/reuters/training/9554
inflating: /root/nltk_data/corpora/reuters/training/9556
inflating: /root/nltk_data/corpora/reuters/training/9557
inflating: /root/nltk_data/corpora/reuters/training/9558
inflating: /root/nltk_data/corpora/reuters/training/9559
inflating: /root/nltk_data/corpora/reuters/training/9560
inflating: /root/nltk_data/corpora/reuters/training/9561
inflating: /root/nltk_data/corpora/reuters/training/9562
inflating: /root/nltk_data/corpora/reuters/training/9563
inflating: /root/nltk_data/corpora/reuters/training/9565
inflating: /root/nltk_data/corpora/reuters/training/9566
inflating: /root/nltk_data/corpora/reuters/training/9570
inflating: /root/nltk_data/corpora/reuters/training/9572
inflating: /root/nltk_data/corpora/reuters/training/9573
inflating: /root/nltk_data/corpora/reuters/training/9576
inflating: /root/nltk_data/corpora/reuters/training/9577
inflating: /root/nltk_data/corpora/reuters/training/9578
inflating: /root/nltk_data/corpora/reuters/training/958
inflating: /root/nltk_data/corpora/reuters/training/9580
inflating: /root/nltk_data/corpora/reuters/training/9582
inflating: /root/nltk_data/corpora/reuters/training/9583
inflating: /root/nltk_data/corpora/reuters/training/9585
inflating: /root/nltk_data/corpora/reuters/training/9587
inflating: /root/nltk_data/corpora/reuters/training/9588
inflating: /root/nltk_data/corpora/reuters/training/9590
inflating: /root/nltk_data/corpora/reuters/training/9592
inflating: /root/nltk_data/corpora/reuters/training/9597
inflating: /root/nltk_data/corpora/reuters/training/9598
inflating: /root/nltk_data/corpora/reuters/training/9599
inflating: /root/nltk_data/corpora/reuters/training/96
inflating: /root/nltk_data/corpora/reuters/training/960
inflating: /root/nltk_data/corpora/reuters/training/9600
inflating: /root/nltk_data/corpora/reuters/training/9603
inflating: /root/nltk_data/corpora/reuters/training/9604
inflating: /root/nltk_data/corpora/reuters/training/9605
inflating: /root/nltk_data/corpora/reuters/training/9606
inflating: /root/nltk_data/corpora/reuters/training/9607
inflating: /root/nltk_data/corpora/reuters/training/9608
inflating: /root/nltk_data/corpora/reuters/training/9610
inflating: /root/nltk_data/corpora/reuters/training/9611
inflating: /root/nltk_data/corpora/reuters/training/9612
inflating: /root/nltk_data/corpora/reuters/training/9613
inflating: /root/nltk_data/corpora/reuters/training/9614
inflating: /root/nltk_data/corpora/reuters/training/9615
inflating: /root/nltk_data/corpora/reuters/training/9617
inflating: /root/nltk_data/corpora/reuters/training/9618
inflating: /root/nltk_data/corpora/reuters/training/9619
inflating: /root/nltk_data/corpora/reuters/training/9625
inflating: /root/nltk_data/corpora/reuters/training/9626
inflating: /root/nltk_data/corpora/reuters/training/9628
inflating: /root/nltk_data/corpora/reuters/training/963
inflating: /root/nltk_data/corpora/reuters/training/9634
inflating: /root/nltk_data/corpora/reuters/training/9635
inflating: /root/nltk_data/corpora/reuters/training/9637
inflating: /root/nltk_data/corpora/reuters/training/9638
inflating: /root/nltk_data/corpora/reuters/training/9639
inflating: /root/nltk_data/corpora/reuters/training/9640
inflating: /root/nltk_data/corpora/reuters/training/9641
inflating: /root/nltk_data/corpora/reuters/training/9642
inflating: /root/nltk_data/corpora/reuters/training/9644
inflating: /root/nltk_data/corpora/reuters/training/9647
inflating: /root/nltk_data/corpora/reuters/training/9650
inflating: /root/nltk_data/corpora/reuters/training/9652
inflating: /root/nltk_data/corpora/reuters/training/9654
inflating: /root/nltk_data/corpora/reuters/training/9656
inflating: /root/nltk_data/corpora/reuters/training/9657
inflating: /root/nltk_data/corpora/reuters/training/9661
inflating: /root/nltk_data/corpora/reuters/training/9667
inflating: /root/nltk_data/corpora/reuters/training/9668
inflating: /root/nltk_data/corpora/reuters/training/9671
inflating: /root/nltk_data/corpora/reuters/training/9674
inflating: /root/nltk_data/corpora/reuters/training/9675
inflating: /root/nltk_data/corpora/reuters/training/9677
inflating: /root/nltk_data/corpora/reuters/training/9678
inflating: /root/nltk_data/corpora/reuters/training/968
inflating: /root/nltk_data/corpora/reuters/training/9680
inflating: /root/nltk_data/corpora/reuters/training/9681
inflating: /root/nltk_data/corpora/reuters/training/9682
inflating: /root/nltk_data/corpora/reuters/training/9686
inflating: /root/nltk_data/corpora/reuters/training/9689
inflating: /root/nltk_data/corpora/reuters/training/969
inflating: /root/nltk_data/corpora/reuters/training/9697
inflating: /root/nltk_data/corpora/reuters/training/9698
inflating: /root/nltk_data/corpora/reuters/training/9699
inflating: /root/nltk_data/corpora/reuters/training/97
inflating: /root/nltk_data/corpora/reuters/training/9700
inflating: /root/nltk_data/corpora/reuters/training/9701
inflating: /root/nltk_data/corpora/reuters/training/9704
inflating: /root/nltk_data/corpora/reuters/training/9705
inflating: /root/nltk_data/corpora/reuters/training/9706
inflating: /root/nltk_data/corpora/reuters/training/9707
inflating: /root/nltk_data/corpora/reuters/training/971
inflating: /root/nltk_data/corpora/reuters/training/9712
inflating: /root/nltk_data/corpora/reuters/training/9718
inflating: /root/nltk_data/corpora/reuters/training/972
inflating: /root/nltk_data/corpora/reuters/training/9720
inflating: /root/nltk_data/corpora/reuters/training/9723
inflating: /root/nltk_data/corpora/reuters/training/9727
inflating: /root/nltk_data/corpora/reuters/training/9728
inflating: /root/nltk_data/corpora/reuters/training/9729
inflating: /root/nltk_data/corpora/reuters/training/9730
inflating: /root/nltk_data/corpora/reuters/training/9731
inflating: /root/nltk_data/corpora/reuters/training/9733
inflating: /root/nltk_data/corpora/reuters/training/9734
```

```
  inflating: /root/nltk_data/corpora/reuters/training/9735
  inflating: /root/nltk_data/corpora/reuters/training/9736
  inflating: /root/nltk_data/corpora/reuters/training/9737
  inflating: /root/nltk_data/corpora/reuters/training/9738
  inflating: /root/nltk_data/corpora/reuters/training/974
  inflating: /root/nltk_data/corpora/reuters/training/9741
  inflating: /root/nltk_data/corpora/reuters/training/9742
  inflating: /root/nltk_data/corpora/reuters/training/9744
  inflating: /root/nltk_data/corpora/reuters/training/9745
  inflating: /root/nltk_data/corpora/reuters/training/9746
  inflating: /root/nltk_data/corpora/reuters/training/9747
  inflating: /root/nltk_data/corpora/reuters/training/9748
  inflating: /root/nltk_data/corpora/reuters/training/9749
  inflating: /root/nltk_data/corpora/reuters/training/9751
  inflating: /root/nltk_data/corpora/reuters/training/9752
  inflating: /root/nltk_data/corpora/reuters/training/9753
  inflating: /root/nltk_data/corpora/reuters/training/9754
  inflating: /root/nltk_data/corpora/reuters/training/9755
  inflating: /root/nltk_data/corpora/reuters/training/9756
  inflating: /root/nltk_data/corpora/reuters/training/976
  inflating: /root/nltk_data/corpora/reuters/training/9760
  inflating: /root/nltk_data/corpora/reuters/training/9761
  inflating: /root/nltk_data/corpora/reuters/training/9763
  inflating: /root/nltk_data/corpora/reuters/training/9764
  inflating: /root/nltk_data/corpora/reuters/training/9767
  inflating: /root/nltk_data/corpora/reuters/training/9768
  inflating: /root/nltk_data/corpora/reuters/training/9769
  inflating: /root/nltk_data/corpora/reuters/training/977
  inflating: /root/nltk_data/corpora/reuters/training/9770
  inflating: /root/nltk_data/corpora/reuters/training/9771
  inflating: /root/nltk_data/corpora/reuters/training/9772
  inflating: /root/nltk_data/corpora/reuters/training/9773
  inflating: /root/nltk_data/corpora/reuters/training/9777
  inflating: /root/nltk_data/corpora/reuters/training/9779
  inflating: /root/nltk_data/corpora/reuters/training/978
  inflating: /root/nltk_data/corpora/reuters/training/9781
  inflating: /root/nltk_data/corpora/reuters/training/9782
  inflating: /root/nltk_data/corpora/reuters/training/9783
  inflating: /root/nltk_data/corpora/reuters/training/9784
  inflating: /root/nltk_data/corpora/reuters/training/9787
  inflating: /root/nltk_data/corpora/reuters/training/9792
  inflating: /root/nltk_data/corpora/reuters/training/9793
  inflating: /root/nltk_data/corpora/reuters/training/9795
  inflating: /root/nltk_data/corpora/reuters/training/9797
  inflating: /root/nltk_data/corpora/reuters/training/9799
  inflating: /root/nltk_data/corpora/reuters/training/98
  inflating: /root/nltk_data/corpora/reuters/training/9801
  inflating: /root/nltk_data/corpora/reuters/training/9804
  inflating: /root/nltk_data/corpora/reuters/training/9805
  inflating: /root/nltk_data/corpora/reuters/training/9807
  inflating: /root/nltk_data/corpora/reuters/training/9809
  inflating: /root/nltk_data/corpora/reuters/training/981
  inflating: /root/nltk_data/corpora/reuters/training/9810
  inflating: /root/nltk_data/corpora/reuters/training/9812
  inflating: /root/nltk_data/corpora/reuters/training/9814
  inflating: /root/nltk_data/corpora/reuters/training/9815
  inflating: /root/nltk_data/corpora/reuters/training/9816
  inflating: /root/nltk_data/corpora/reuters/training/9818
  inflating: /root/nltk_data/corpora/reuters/training/982
  inflating: /root/nltk_data/corpora/reuters/training/9821
  inflating: /root/nltk_data/corpora/reuters/training/9822
  inflating: /root/nltk_data/corpora/reuters/training/9823
  inflating: /root/nltk_data/corpora/reuters/training/9824
  inflating: /root/nltk_data/corpora/reuters/training/9825
  inflating: /root/nltk_data/corpora/reuters/training/9827
  inflating: /root/nltk_data/corpora/reuters/training/9829
  inflating: /root/nltk_data/corpora/reuters/training/983
  inflating: /root/nltk_data/corpora/reuters/training/9833
  inflating: /root/nltk_data/corpora/reuters/training/9834
  inflating: /root/nltk_data/corpora/reuters/training/9836
  inflating: /root/nltk_data/corpora/reuters/training/9837
  inflating: /root/nltk_data/corpora/reuters/training/9839
  inflating: /root/nltk_data/corpora/reuters/training/984
  inflating: /root/nltk_data/corpora/reuters/training/9841
  inflating: /root/nltk_data/corpora/reuters/training/9844
  inflating: /root/nltk_data/corpora/reuters/training/9847
  inflating: /root/nltk_data/corpora/reuters/training/9848
  inflating: /root/nltk_data/corpora/reuters/training/9849
  inflating: /root/nltk_data/corpora/reuters/training/985
  inflating: /root/nltk_data/corpora/reuters/training/9850
  inflating: /root/nltk_data/corpora/reuters/training/9851
  inflating: /root/nltk_data/corpora/reuters/training/9852
  inflating: /root/nltk_data/corpora/reuters/training/9853
  inflating: /root/nltk_data/corpora/reuters/training/9855
  inflating: /root/nltk_data/corpora/reuters/training/9857
  inflating: /root/nltk_data/corpora/reuters/training/9858
  inflating: /root/nltk_data/corpora/reuters/training/986
  inflating: /root/nltk_data/corpora/reuters/training/9860
  inflating: /root/nltk_data/corpora/reuters/training/9861
  inflating: /root/nltk_data/corpora/reuters/training/9862
  inflating: /root/nltk_data/corpora/reuters/training/9864
  inflating: /root/nltk_data/corpora/reuters/training/9865
  inflating: /root/nltk_data/corpora/reuters/training/9866
  inflating: /root/nltk_data/corpora/reuters/training/9867
  inflating: /root/nltk_data/corpora/reuters/training/9868
  inflating: /root/nltk_data/corpora/reuters/training/9869
  inflating: /root/nltk_data/corpora/reuters/training/9871
  inflating: /root/nltk_data/corpora/reuters/training/9872
  inflating: /root/nltk_data/corpora/reuters/training/9874
  inflating: /root/nltk_data/corpora/reuters/training/9875
  inflating: /root/nltk_data/corpora/reuters/training/9878
  inflating: /root/nltk_data/corpora/reuters/training/988
  inflating: /root/nltk_data/corpora/reuters/training/9880
  inflating: /root/nltk_data/corpora/reuters/training/9884
  inflating: /root/nltk_data/corpora/reuters/training/9891
  inflating: /root/nltk_data/corpora/reuters/training/9892
  inflating: /root/nltk_data/corpora/reuters/training/9893
  inflating: /root/nltk_data/corpora/reuters/training/9894
  inflating: /root/nltk_data/corpora/reuters/training/9896
  inflating: /root/nltk_data/corpora/reuters/training/9897
  inflating: /root/nltk_data/corpora/reuters/training/99
  inflating: /root/nltk_data/corpora/reuters/training/9901
  inflating: /root/nltk_data/corpora/reuters/training/9902
  inflating: /root/nltk_data/corpora/reuters/training/9903
  inflating: /root/nltk_data/corpora/reuters/training/9904
  inflating: /root/nltk_data/corpora/reuters/training/9905
  inflating: /root/nltk_data/corpora/reuters/training/9906
  inflating: /root/nltk_data/corpora/reuters/training/9907
  inflating: /root/nltk_data/corpora/reuters/training/9909
  inflating: /root/nltk_data/corpora/reuters/training/991
  inflating: /root/nltk_data/corpora/reuters/training/9912
  inflating: /root/nltk_data/corpora/reuters/training/9913
  inflating: /root/nltk_data/corpora/reuters/training/9914
  inflating: /root/nltk_data/corpora/reuters/training/9915
  inflating: /root/nltk_data/corpora/reuters/training/9918
  inflating: /root/nltk_data/corpora/reuters/training/9919
  inflating: /root/nltk_data/corpora/reuters/training/9920
  inflating: /root/nltk_data/corpora/reuters/training/9923
  inflating: /root/nltk_data/corpora/reuters/training/9925
  inflating: /root/nltk_data/corpora/reuters/training/9926
  inflating: /root/nltk_data/corpora/reuters/training/9927
  inflating: /root/nltk_data/corpora/reuters/training/9928
  inflating: /root/nltk_data/corpora/reuters/training/9929
  inflating: /root/nltk_data/corpora/reuters/training/9930
```

```
inflating: /root/nltk_data/corpora/reuters/training/9933
inflating: /root/nltk_data/corpora/reuters/training/9934
inflating: /root/nltk_data/corpora/reuters/training/9936
inflating: /root/nltk_data/corpora/reuters/training/9937
inflating: /root/nltk_data/corpora/reuters/training/9939
inflating: /root/nltk_data/corpora/reuters/training/9940
inflating: /root/nltk_data/corpora/reuters/training/9941
inflating: /root/nltk_data/corpora/reuters/training/9942
inflating: /root/nltk_data/corpora/reuters/training/9943
inflating: /root/nltk_data/corpora/reuters/training/9946
inflating: /root/nltk_data/corpora/reuters/training/9947
inflating: /root/nltk_data/corpora/reuters/training/995
inflating: /root/nltk_data/corpora/reuters/training/9952
inflating: /root/nltk_data/corpora/reuters/training/9953
inflating: /root/nltk_data/corpora/reuters/training/9954
inflating: /root/nltk_data/corpora/reuters/training/9955
inflating: /root/nltk_data/corpora/reuters/training/9956
inflating: /root/nltk_data/corpora/reuters/training/9957
inflating: /root/nltk_data/corpora/reuters/training/9958
inflating: /root/nltk_data/corpora/reuters/training/9959
inflating: /root/nltk_data/corpora/reuters/training/9961
inflating: /root/nltk_data/corpora/reuters/training/9963
inflating: /root/nltk_data/corpora/reuters/training/9964
inflating: /root/nltk_data/corpora/reuters/training/9965
inflating: /root/nltk_data/corpora/reuters/training/9967
inflating: /root/nltk_data/corpora/reuters/training/9970
inflating: /root/nltk_data/corpora/reuters/training/9971
inflating: /root/nltk_data/corpora/reuters/training/9972
inflating: /root/nltk_data/corpora/reuters/training/9973
inflating: /root/nltk_data/corpora/reuters/training/9974
inflating: /root/nltk_data/corpora/reuters/training/9975
inflating: /root/nltk_data/corpora/reuters/training/9976
inflating: /root/nltk_data/corpora/reuters/training/9977
inflating: /root/nltk_data/corpora/reuters/training/9978
inflating: /root/nltk_data/corpora/reuters/training/998
inflating: /root/nltk_data/corpora/reuters/training/9981
inflating: /root/nltk_data/corpora/reuters/training/9982
inflating: /root/nltk_data/corpora/reuters/training/9984
inflating: /root/nltk_data/corpora/reuters/training/9985
inflating: /root/nltk_data/corpora/reuters/training/9988
inflating: /root/nltk_data/corpora/reuters/training/9989
inflating: /root/nltk_data/corpora/reuters/training/999
inflating: /root/nltk_data/corpora/reuters/training/9992
inflating: /root/nltk_data/corpora/reuters/training/9993
inflating: /root/nltk_data/corpora/reuters/training/9994
inflating: /root/nltk_data/corpora/reuters/training/9995
```

**Let's have a look what these documents are like....**

In [6]:

```python
reuters_corpus = read_corpus()
pprint.pprint(reuters_corpus[:3], compact=True, width=100)
```

```
[['<START>', 'japan', 'to', 'revise', 'long', '-', 'term', 'energy', 'demand', 'downwards', 'the',
  'ministry', 'of', 'international', 'trade', 'and', 'industry', '(', 'miti', ')', 'will', 'revise',
  'its', 'long', '-', 'term', 'energy', 'supply', '/', 'demand', 'outlook', 'by', 'august', 'to',
  'meet', 'a', 'forecast', 'downtrend', 'in', 'japanese', 'energy', 'demand', ',', 'ministry',
  'officials', 'said', '.', 'miti', 'is', 'expected', 'to', 'lower', 'the', 'projection', 'for',
  'primary', 'energy', 'supplies', 'in', 'the', 'year', '2000', 'to', '550', 'mln', 'kilolitres',
  '(', 'kl', ')', 'from', '600', 'mln', ',', 'they', 'said', '.', 'the', 'decision', 'follows',
  'the', 'emergence', 'of', 'structural', 'changes', 'in', 'japanese', 'industry', 'following',
  'the', 'rise', 'in', 'the', 'value', 'of', 'the', 'yen', 'and', 'a', 'decline', 'in', 'domestic',
  'electric', 'power', 'demand', '.', 'miti', 'is', 'planning', 'to', 'work', 'out', 'a', 'revised',
  'energy', 'supply', '/', 'demand', 'outlook', 'through', 'deliberations', 'of', 'committee',
  'meetings', 'of', 'the', 'agency', 'of', 'natural', 'resources', 'and', 'energy', ',', 'the',
  'officials', 'said', '.', 'they', 'said', 'miti', 'will', 'also', 'review', 'the', 'breakdown',
  'of', 'energy', 'supply', 'sources', ',', 'including', 'oil', ',', 'nuclear', ',', 'coal', 'and',
  'natural', 'gas', '.', 'nuclear', 'energy', 'provided', 'the', 'bulk', 'of', 'japan', "'", 's',
  'electric', 'power', 'in', 'the', 'fiscal', 'year', 'ended', 'march', '31', ',', 'supplying',
  'an', 'estimated', '27', 'pct', 'on', 'a', 'kilowatt', '/', 'hour', 'basis', ',', 'followed',
  'by', 'oil', '(', '23', 'pct', ')', 'and', 'liquefied', 'natural', 'gas', '(', '21', 'pct', '),',
  'they', 'noted', '.', '<END>'],
 ['<START>', 'energy', '/', 'u', '.', 's', '.', 'petrochemical', 'industry', 'cheap', 'oil',
  'feedstocks', ',', 'the', 'weakened', 'u', '.', 's', '.', 'dollar', 'and', 'a', 'plant',
  'utilization', 'rate', 'approaching', '90', 'pct', 'will', 'propel', 'the', 'streamlined', 'u',
  '.', 's', '.', 'petrochemical', 'industry', 'to', 'record', 'profits', 'this', 'year', ',',
  'with', 'growth', 'expected', 'through', 'at', 'least', '1990', ',', 'major', 'company',
  'executives', 'predicted', '.', 'this', 'bullish', 'outlook', 'for', 'chemical', 'manufacturing',
  'and', 'an', 'industrywide', 'move', 'to', 'shed', 'unrelated', 'businesses', 'has', 'prompted',
  'gaf', 'corp', '&', 'lt', ';', 'gaf', '>,', 'privately', '-', 'held', 'cain', 'chemical', 'inc',
  ',', 'and', 'other', 'firms', 'to', 'aggressively', 'seek', 'acquisitions', 'of', 'petrochemical',
  'plants', '.', 'oil', 'companies', 'such', 'as', 'ashland', 'oil', 'inc', '&', 'lt', ';', 'ash',
  '>,', 'the', 'kentucky', '-', 'based', 'oil', 'refiner', 'and', 'marketer', ',', 'are', 'also',
  'shopping', 'for', 'money', '-', 'making', 'petrochemical', 'businesses', 'to', 'buy', '.', '""',
  'i', 'see', 'us', 'poised', 'at', 'the', 'threshold', 'of', 'a', 'golden', 'period', ',"', 'said',
  'paul', 'oreffice', ',', 'chairman', 'of', 'giant', 'dow', 'chemical', 'co', '&', 'lt', ';',
  'dow', '>,', 'adding', ',', '""', 'there', '"\'', 's', 'no', 'major', 'plant', 'capacity', 'being',
  'added', 'around', 'the', 'world', 'now', '.', 'the', 'whole', 'game', 'is', 'bringing', 'out',
  'new', 'products', 'and', 'improving', 'the', 'old', 'ones', '."', 'analysts', 'say', 'the',
  'chemical', 'industry', '"\'', 's', 'biggest', 'customers', ',', 'automobile', 'manufacturers',
  'and', 'home', 'builders', 'that', 'use', 'a', 'lot', 'of', 'paints', 'and', 'plastics', ',',
  'are', 'expected', 'to', 'buy', 'quantities', 'this', 'year', '.', 'u', '.', 's', '.',
  'petrochemical', 'plants', 'are', 'currently', 'operating', 'at', 'about', '90', 'pct',
  'capacity', ',', 'reflecting', 'tighter', 'supply', 'that', 'could', 'hike', 'product', 'prices',
  'by', '30', 'to', '40', 'pct', 'this', 'year', ',', 'said', 'john', 'dosher', ',', 'managing',
  'director', 'of', 'pace', 'consultants', 'inc', 'of', 'houston', '.', 'demand', 'for', 'some',
  'products', 'such', 'as', 'styrene', 'could', 'push', 'profit', 'margins', 'up', 'by', 'as',
  'much', 'as', '300', 'pct', ',', 'he', 'said', '.', 'oreffice', ',', 'speaking', 'at', 'a',
  'meeting', 'of', 'chemical', 'engineers', 'in', 'houston', ',', 'said', 'dow', 'would', 'easily',
  'top', 'the', '741', 'mln', 'dlrs', 'it', 'earned', 'last', 'year', 'and', 'predicted', 'it',
  'would', 'have', 'the', 'best', 'year', 'in', 'its', 'history', '.', 'in', '1985', ',', 'when',
  'oil', 'prices', 'were', 'still', 'above', '25', 'dlrs', 'a', 'barrel', 'and', 'chemical',
  'exports', 'were', 'adversely', 'affected', 'by', 'the', 'strong', 'u', '.', 's', '.', 'dollar',
  ',', 'dow', 'had', 'profits', 'of', '58', 'mln', 'dlrs', '.', '""', 'i', 'believe', 'the',
  'entire', 'chemical', 'industry', 'is', 'headed', 'for', 'a', 'record', 'year', 'or', 'close',
  'to', 'it', ',"', 'oreffice', 'said', '.', 'gaf', 'chairman', 'samuel', 'heyman', 'estimated',
  'that', 'the', 'u', '.', 's', '.', 'chemical', 'industry', 'would', 'report', 'a', '20', 'pct',
  'gain', 'in', 'profits', 'during', '1987', '.', 'last', 'year', ',', 'the', 'domestic',
  'industry', 'earned', 'a', 'total', 'of', '13', 'billion', 'dlrs', ',', 'a', '54', 'pct', 'leap',
  'from', '1985', '.', 'the', 'turn', 'in', 'the', 'fortunes', 'of', 'the', 'once', '-', 'sickly',
  'chemical', 'industry', 'has', 'been', 'brought', 'about', 'by', 'a', 'combination', 'of', 'luck',
  'and', 'planning', ',', 'said', 'pace', "'", 's', 'john', 'dosher', '.', 'dosher', 'said', 'last',
  'year', "'", 's', 'fall', 'in', 'oil', 'prices', 'made', 'feedstocks', 'dramatically', 'cheaper',
  'and', 'at', 'the', 'same', 'time', 'the', 'american', 'dollar', 'was', 'weakening', 'against',
  'foreign', 'currencies', '.', 'that', 'helped', 'boost', 'u', '.', 's', '.', 'chemical',
  'exports', '.', 'also', 'helping', 'to', 'bring', 'supply', 'and', 'demand', 'into', 'balance',
  'has', 'been', 'the', 'gradual', 'market', 'absorption', 'of', 'the', 'extra', 'chemical',
  'manufacturing', 'capacity', 'created', 'by', 'middle', 'eastern', 'oil', 'producers', 'in',
  'the', 'early', '1980s', '.', 'finally', ',', 'virtually', 'all', 'major', 'u', '.', 's', '.',
  'chemical', 'manufacturers', 'have', 'embarked', 'on', 'an', 'extensive', 'corporate',
  'restructuring', 'program', 'to', 'mothball', 'inefficient', 'plants', ',', 'trim', 'the',
  'payroll', 'and', 'eliminate', 'unrelated', 'businesses', '.', 'the', 'restructuring', 'touched',
  'off', 'a', 'flurry', 'of', 'friendly', 'and', 'hostile', 'takeover', 'attempts', '.', 'gaf', ',',
  'which', 'made', 'an', 'unsuccessful', 'attempt', 'in', '1985', 'to', 'acquire', 'union',
  'carbide', 'corp', '&', 'lt', ';', 'uk', '>,', 'recently', 'offered', 'three', 'billion', 'dlrs',
  'for', 'borg', 'warner', 'corp', '&', 'lt', ';', 'bor', '>,', 'a', 'chicago', 'manufacturer',
  'of', 'plastics', 'and', 'chemicals', '.', 'another', 'industry', 'powerhouse', ',', 'w', '.',
  'r', '.', 'grace', '&', 'lt', ';', 'gra', '>', 'has', 'divested', 'its', 'retailing', ',',
  'restaurant', 'and', 'fertilizer', 'businesses', 'to', 'raise', 'cash', 'for', 'chemical',
  'acquisitions', '.', 'but', 'some', 'experts', 'worry', 'that', 'the', 'chemical', 'industry',
  'may', 'be', 'headed', 'for', 'trouble', 'if', 'companies', 'continue', 'turning', 'their',
  'back', 'on', 'the', 'manufacturing', 'of', 'staple', 'petrochemical', 'commodities', ',', 'such',
  'as', 'ethylene', ',', 'in', 'favor', 'of', 'more', 'profitable', 'specialty', 'chemicals',
```

```
'that', 'are', 'custom', '-', 'designed', 'for', 'a', 'small', 'group', 'of', 'buyers', '.', '"',
'companies', 'like', 'dupont', '&', 'lt', ';', 'dd', '>', 'and', 'monsanto', 'co', '&', 'lt', ';',
'mtc', '>', 'spent', 'the', 'past', 'two', 'or', 'three', 'years', 'trying', 'to', 'get', 'out',
'of', 'the', 'commodity', 'chemical', 'business', 'in', 'reaction', 'to', 'how', 'badly', 'the',
'market', 'had', 'deteriorated', ',"', 'dosher', 'said', '.', '"', 'but', 'i', 'think', 'they',
'will', 'eventually', 'kill', 'the', 'margins', 'on', 'the', 'profitable', 'chemicals', 'in',
'the', 'niche', 'market', '."', 'some', 'top', 'chemical', 'executives', 'share', 'the',
'concern', '.', '"', 'the', 'challenge', 'for', 'our', 'industry', 'is', 'to', 'keep', 'from',
'getting', 'carried', 'away', 'and', 'repeating', 'past', 'mistakes', ',"', 'gaf', '"', 's',
'heyman', 'cautioned', '.', '"', 'the', 'shift', 'from', 'commodity', 'chemicals', 'may', 'be',
'ill', '-', 'advised', '.', 'specialty', 'businesses', 'do', 'not', 'stay', 'special', 'long',
'."', 'houston', '-', 'based', 'cain', 'chemical', ',', 'created', 'this', 'month', 'by', 'the',
'sterling', 'investment', 'banking', 'group', ',', 'believes', 'it', 'can', 'generate', '700',
'mln', 'dlrs', 'in', 'annual', 'sales', 'by', 'bucking', 'the', 'industry', 'trend', '.',
'chairman', 'gordon', 'cain', ',', 'who', 'previously', 'led', 'a', 'leveraged', 'buyout', 'of',
'dupont', '"', 's', 'conoco', 'inc', '"', 's', 'chemical', 'business', ',', 'has', 'spent', '1',
'.', '1', 'billion', 'dlrs', 'since', 'january', 'to', 'buy', 'seven', 'petrochemical', 'plants',
'along', 'the', 'texas', 'gulf', 'coast', '.', 'the', 'plants', 'produce', 'only', 'basic',
'commodity', 'petrochemicals', 'that', 'are', 'the', 'building', 'blocks', 'of', 'specialty',
'products', '.', '"', 'this', 'kind', 'of', 'commodity', 'chemical', 'business', 'will', 'never',
'be', 'a', 'glamorous', ',', 'high', '-', 'margin', 'business', ',"', 'cain', 'said', ',',
'adding', 'that', 'demand', 'is', 'expected', 'to', 'grow', 'by', 'about', 'three', 'pct',
'annually', '.', 'garo', 'armen', ',', 'an', 'analyst', 'with', 'dean', 'witter', 'reynolds', ',',
'said', 'chemical', 'makers', 'have', 'also', 'benefitted', 'by', 'increasing', 'demand', 'for',
'plastics', 'as', 'prices', 'become', 'more', 'competitive', 'with', 'aluminum', ',', 'wood',
'and', 'steel', 'products', '.', 'armen', 'estimated', 'the', 'upturn', 'in', 'the', 'chemical',
'business', 'could', 'last', 'as', 'long', 'as', 'four', 'or', 'five', 'years', ',', 'provided',
'the', 'u', 's', '.', 'economy', 'continues', 'its', 'modest', 'rate', 'of', 'growth', '.',
'<END>'],
['<START>', 'turkey', 'calls', 'for', 'dialogue', 'to', 'solve', 'dispute', 'turkey', 'said',
'today', 'its', 'disputes', 'with', 'greece', ',', 'including', 'rights', 'on', 'the',
'continental', 'shelf', 'in', 'the', 'aegean', 'sea', ',', 'should', 'be', 'solved', 'through',
'negotiations', '.', 'a', 'foreign', 'ministry', 'statement', 'said', 'the', 'latest', 'crisis',
'between', 'the', 'two', 'nato', 'members', 'stemmed', 'from', 'the', 'continental', 'shelf',
'dispute', 'and', 'an', 'agreement', 'on', 'this', 'issue', 'would', 'effect', 'the', 'security',
',', 'economy', 'and', 'other', 'rights', 'of', 'both', 'countries', '.', '"', 'as', 'the',
'issue', 'is', 'basicly', 'political', ',', 'a', 'solution', 'can', 'only', 'be', 'found', 'by',
'bilateral', 'negotiations', ',"', 'the', 'statement', 'said', '.', 'greece', 'has', 'repeatedly',
'said', 'the', 'issue', 'was', 'legal', 'and', 'could', 'be', 'solved', 'at', 'the',
'international', 'court', 'of', 'justice', '.', 'the', 'two', 'countries', 'approached', 'armed',
'confrontation', 'last', 'month', 'after', 'greece', 'announced', 'it', 'planned', 'oil',
'exploration', 'work', 'in', 'the', 'aegean', 'and', 'turkey', 'said', 'it', 'would', 'also',
'search', 'for', 'oil', '.', 'a', 'face', '-', 'off', 'was', 'averted', 'when', 'turkey',
'confined', 'its', 'research', 'to', 'territorrial', 'waters', '.', '"', 'the', 'latest',
'crises', 'created', 'an', 'historic', 'opportunity', 'to', 'solve', 'the', 'disputes', 'between',
'the', 'two', 'countries', ',"', 'the', 'foreign', 'ministry', 'statement', 'said', '.', 'turkey',
'"', 's', 'ambassador', 'in', 'athens', ',', 'nazmi', 'akiman', ',', 'was', 'due', 'to', 'meet',
'prime', 'minister', 'andreas', 'papandreou', 'today', 'for', 'the', 'greek', 'reply', 'to', 'a',
'message', 'sent', 'last', 'week', 'by', 'turkish', 'prime', 'minister', 'turgut', 'ozal', '.',
'the', 'contents', 'of', 'the', 'message', 'were', 'not', 'disclosed', '.', '<END>']]
```

In [7]:
```python
type(reuters_corpus)
```

Out[7]:

list

In [8]:
```python
len(reuters_corpus)
```

Out[8]:

578

In [9]:
```python
len(reuters_corpus[0])
```

Out[9]:

209

In [10]:
```python
for i in range(5):
    print(len(reuters_corpus[i]))
```

209
1001
225
240
350

## Question 1.1: Implement `distinct_words` [code] (2 points)

Write a method to work out the distinct words (word types) that occur in the corpus. You can do this with `for` loops, but it's more efficient to do it with Python list comprehensions. In particular, this may be useful to flatten a list of lists. If you're not familiar with Python list comprehensions in general, here's more information.

Your returned `corpus_words` should be sorted. You can use python's `sorted` function for this.

You may find it useful to use Python sets to remove duplicate words.

# List comprehension way

```python
#The list of lists
list_of_lists = [range(4), range(7)]


#flatten the lists
flattened_list = [y for x in list_of_lists for y in x]
```

In [11]:
```python
list1 = [range(5), range(16,20)]
```

In [12]:
```python
[y for i in list1 for y in i]
```

Out[12]:

```
[0, 1, 2, 3, 4, 16, 17, 18, 19]
```

In [13]:

```python
len(set([word for wordlist in reuters_corpus for word in wordlist]))
```

Out[13]:

```
8185
```

In [14]:

```python
def distinct_words(corpus):
    """ Determine a list of distinct words for the corpus.
        Params:
            corpus (list of list of strings): corpus of documents
        Return:
            corpus_words (list of strings): sorted list of distinct words across the corpus
            num_corpus_words (integer): number of distinct words across the corpus
    """
    corpus_words = []
    num_corpus_words = -1

    # ------------------
    # Write your implementation here..

    corpus_words = sorted(list(set([word for wordlist in corpus for word in wordlist])))
    num_corpus_words = len(corpus_words)


    # ------------------

    return corpus_words, num_corpus_words
```

In [15]:

```python
["{} All that glitters isn't gold {}".format(START_TOKEN, END_TOKEN).split(" "), "{} All's well that ends well {}".format(START_TOKEN, END_TOKEN).split(" ")]
```

Out[15]:

```
[['<START>', 'All', 'that', 'glitters', "isn't", 'gold', '<END>'],
 ['<START>', "All's", 'well', 'that', 'ends', 'well', '<END>']]
```

In [16]:

```python
distinct_words([['i', 'don','know'], ['hey', 'how']])
```

Out[16]:

```
(['don', 'hey', 'how', 'i', 'know'], 5)
```

In [17]:

```python
# --------------------
# Run this sanity check
# Note that this not an exhaustive check for correctness.
# --------------------

# Define toy corpus
test_corpus = ["{} All that glitters isn't gold {}".format(START_TOKEN, END_TOKEN).split(" "), "{} All's well that ends well {}".format(START_TOKEN, END_TOKEN).split(" ")]
test_corpus_words, num_corpus_words = distinct_words(test_corpus)

# Correct answers
ans_test_corpus_words = sorted([START_TOKEN, "All", "ends", "that", "gold", "All's", "glitters", "isn't", "well", END_TOKEN])
ans_num_corpus_words = len(ans_test_corpus_words)

# Test correct number of words
assert(num_corpus_words == ans_num_corpus_words), "Incorrect number of distinct words. Correct: {}. Yours: {}".format(ans_num_corpus_words, num_corpus_words)

# Test correct words
assert (test_corpus_words == ans_test_corpus_words), "Incorrect corpus_words.\nCorrect: {}\nYours:   {}".format(str(ans_test_corpus_words), str(test_corpus_words))

# Print Success
print ("-" * 80)
print ("Passed All Tests!")
print ("-" * 80)
```

```
--------------------------------------------------------------------------------
Passed All Tests!
--------------------------------------------------------------------------------
```

## Question 1.2: Implement `compute_co_occurrence_matrix` [code] (3 points)

Write a method that constructs a co-occurrence matrix for a certain window-size $n$ (with a default of 4), considering words $n$ before and $n$ after the word in the center of the window. Here, we start to use `numpy (np)` to represent vectors, matrices, and tensors.

In [18]:

```python
words, num_words = distinct_words(test_corpus)
word2ind = {}
print(words[:5])
print(num_words)

for num in range(num_words):
    word2ind[words[num]] = num

print(word2ind)
```

```
['<END>', '<START>', 'All', "All's", 'ends']
10
{'<END>': 0, '<START>': 1, 'All': 2, "All's": 3, 'ends': 4, 'glitters': 5, 'gold': 6, "isn't": 7, 'that': 8, 'well': 9}
```

In [19]:

```python
def compute_co_occurrence_matrix(corpus, window_size=4):
    """ Compute co-occurrence matrix for the given corpus and window_size (default of 4).

        Note: Each word in a document should be at the center of a window. Words near edges will have a smaller
              number of co-occurring words.

              For example, if we take the document "<START> All that glitters is not gold <END>" with window size of 4,
              "All" will co-occur with "<START>", "that", "glitters", "is", and "not".

        Params:
            corpus (list of list of strings): corpus of documents
            window_size (int): size of context window
        Return:
            M (a symmetric numpy matrix of shape (number of unique words in the corpus , number of unique words in the corpus)):
                Co-occurence matrix of word counts.
                The ordering of the words in the rows/columns should be the same as the ordering of the words given by the distinct_words function.
            word2ind (dict): dictionary that maps word to index (i.e. row/column number) for matrix M.
    """
    words, num_words = distinct_words(corpus)
    M = None
    word2ind = {}

    # ------------------
    # Write your implementation here.

    for i in range(num_words):
        word2ind[words[i]] = i
    M = np.zeros((num_words, num_words))
```

```
        for line in corpus:
            for i in range(len(line)):
                target = line[i]
                target_index = word2ind[target]

                left = max(i - window_size, 0)
                right = min(i + window_size, len(line) - 1)
                for j in range(left, i):
                    window_word = line[j]
                    M[target_index][word2ind[window_word]] += 1
                    M[word2ind[window_word]][target_index] += 1

    # ------------------

    return M, word2ind
```

In [20]:

```
# --------------------
# Run this sanity check
# Note that this is not an exhaustive check for correctness.
# --------------------

# Define toy corpus and get student's co-occurrence matrix
test_corpus = ["{} All that glitters isn't gold {}".format(START_TOKEN, END_TOKEN).split(" "), "{} All's well that ends well {}".format(START_TOKEN, END_TOKEN).split(" ")]
M_test, word2ind_test = compute_co_occurrence_matrix(test_corpus, window_size=1)

# Correct M and word2ind
M_test_ans = np.array(
    [[0., 0., 0., 0., 0., 0., 1., 0., 0., 1.,],
     [0., 0., 1., 1., 0., 0., 0., 0., 0., 0.,],
     [0., 1., 0., 0., 0., 0., 0., 0., 1., 0.,],
     [0., 1., 0., 0., 0., 0., 0., 0., 0., 1.,],
     [0., 0., 0., 0., 0., 0., 0., 0., 1., 1.,],
     [0., 0., 0., 0., 0., 0., 0., 1., 1., 0.,],
     [1., 0., 0., 0., 0., 0., 0., 1., 0., 0.,],
     [0., 0., 0., 0., 0., 1., 1., 0., 0., 0.,],
     [0., 0., 1., 0., 1., 1., 0., 0., 0., 1.,],
     [1., 0., 0., 1., 1., 0., 0., 0., 1., 0.,]]
)
ans_test_corpus_words = sorted([START_TOKEN, "All", "ends", "that", "gold", "All's", "glitters", "isn't", "well", END_TOKEN])
word2ind_ans = dict(zip(ans_test_corpus_words, range(len(ans_test_corpus_words))))

# Test correct word2ind
assert (word2ind_ans == word2ind_test), "Your word2ind is incorrect:\nCorrect: {}\nYours: {}".format(word2ind_ans, word2ind_test)

# Test correct M shape
assert (M_test.shape == M_test_ans.shape), "M matrix has incorrect shape.\nCorrect: {}\nYours: {}".format(M_test.shape, M_test_ans.shape)

# Test correct M values
for w1 in word2ind_ans.keys():
    idx1 = word2ind_ans[w1]
    for w2 in word2ind_ans.keys():
        idx2 = word2ind_ans[w2]
        student = M_test[idx1, idx2]
        correct = M_test_ans[idx1, idx2]
        if student != correct:
            print("Correct M:")
            print(M_test_ans)
            print("Your M: ")
            print(M_test)
            raise AssertionError("Incorrect count at index ({}, {})=({}, {}) in matrix M. Yours has {} but should have {}.".format(idx1, idx2, w1, w2, student, correct))

# Print Success
print ("-" * 80)
print("Passed All Tests!")
print ("-" * 80)
```

```
--------------------------------------------------------------------------------
Passed All Tests!
--------------------------------------------------------------------------------
```

In [21]:

```
word2ind_test
```

Out[21]:

```
{'<END>': 0,
 '<START>': 1,
 'All': 2,
 "All's": 3,
 'ends': 4,
 'glitters': 5,
 'gold': 6,
 "isn't": 7,
 'that': 8,
 'well': 9}
```

In [22]:

```
M_test
```

Out[22]:

```
array([[0., 0., 0., 0., 0., 0., 1., 0., 0., 1.],
       [0., 0., 1., 1., 0., 0., 0., 0., 0., 0.],
       [0., 1., 0., 0., 0., 0., 0., 0., 1., 0.],
       [0., 1., 0., 0., 0., 0., 0., 0., 0., 1.],
       [0., 0., 0., 0., 0., 0., 0., 0., 1., 1.],
       [0., 0., 0., 0., 0., 0., 0., 1., 1., 0.],
       [1., 0., 0., 0., 0., 0., 0., 1., 0., 0.],
       [0., 0., 0., 0., 0., 1., 1., 0., 0., 0.],
       [0., 0., 1., 0., 1., 1., 0., 0., 0., 1.],
       [1., 0., 0., 1., 1., 0., 0., 0., 1., 0.]])
```

## Question 1.3: Implement `reduce_to_k_dim` [code] (1 point)

Construct a method that performs dimensionality reduction on the matrix to produce k-dimensional embeddings. Use SVD to take the top k components and produce a new matrix of k-dimensional embeddings.

**Note:** All of numpy, scipy, and scikit-learn (`sklearn`) provide *some* implementation of SVD, but only scipy and sklearn provide an implementation of Truncated SVD, and only sklearn provides an efficient randomized algorithm for calculating large-scale Truncated SVD. So please use [sklearn.decomposition.TruncatedSVD](https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html).

In [23]:

```
from sklearn import decomposition
```

In [24]:

```
def reduce_to_k_dim(M, k=2):
    """ Reduce a co-occurence count matrix of dimensionality (num_corpus_words, num_corpus_words)
        to a matrix of dimensionality (num_corpus_words, k) using the following SVD function from Scikit-Learn:
            - http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html

        Params:
            M (numpy matrix of shape (number of unique words in the corpus , number of unique words in the corpus)): co-occurence matrix of word counts
            k (int): embedding size of each word after dimension reduction
        Return:
            M_reduced (numpy matrix of shape (number of corpus words, k)): matrix of k-dimensioal word embeddings.
                In terms of the SVD from math class, this actually returns U * S
    """
    n_iters = 10     # Use this parameter in your call to `TruncatedSVD`
    M_reduced = None
```

```
        print("Running Truncated SVD over %i words..." % (M.shape[0]))

        # -------------------
        # Write your implementation here.


        svd = decomposition.TruncatedSVD(n_components= k, n_iter= n_iters)
        M_reduced = svd.fit_transform(M)
        print(f"The original size of  embedded matrix is {M.shape}")
        print(f"The reduced size of embedded matrix is {M_reduced.shape}")


        # -------------------

        print("Done.")
        return M_reduced
```

In [25]:

```
# ---------------------
# Run this sanity check
# Note that this is not an exhaustive check for correctness
# In fact we only check that your M_reduced has the right dimensions.
# ---------------------

# Define toy corpus and run student code
test_corpus = ["{} All that glitters isn't gold {}".format(START_TOKEN, END_TOKEN).split(" "), "{} All's well that ends well {}".format(START_TOKEN, END_TOKEN).split(" ")]
M_test, word2ind_test = compute_co_occurrence_matrix(test_corpus, window_size=1)
M_test_reduced = reduce_to_k_dim(M_test, k=2)

# Test proper dimensions
assert (M_test_reduced.shape[0] == 10), "M_reduced has {} rows; should have {}".format(M_test_reduced.shape[0], 10)
assert (M_test_reduced.shape[1] == 2), "M_reduced has {} columns; should have {}".format(M_test_reduced.shape[1], 2)

# Print Success
print ("-" * 80)
print("Passed All Tests!")
print ("-" * 80)
```

```
Running Truncated SVD over 10 words...
The original size of  embedded matrix is (10, 10)
The reduced size of embedded matrix is (10, 2)
Done.
--------------------------------------------------------------------------------
Passed All Tests!
--------------------------------------------------------------------------------
```

In [26]:

```
M_test_reduced
```

Out[26]:

```
array([[ 6.54802087e-01,  7.83221122e-01],
       [ 5.20200324e-01, -4.44089210e-16],
       [ 7.05647176e-01, -4.84057274e-01],
       [ 7.05647176e-01,  4.84057274e-01],
       [ 1.02780472e+00,  0.00000000e+00],
       [ 6.54802087e-01, -7.83221122e-01],
       [ 3.82258491e-01, -6.56224003e-01],
       [ 3.82258491e-01,  6.56224003e-01],
       [ 1.39420808e+00,  1.06179274e+00],
       [ 1.39420808e+00, -1.06179274e+00]])
```

## Question 1.4: Implement `plot_embeddings` [code] (1 point)

Here you will write a function to plot a set of 2D vectors in 2D space. For graphs, we will use Matplotlib ( `plt` ).

In [27]:

```
def plot_embeddings(M_reduced, word2ind, words):
    """ Plot in a scatterplot the embeddings of the words specified in the list "words".
        NOTE: do not plot all the words listed in M_reduced / word2ind.
        Include a label next to each point.

        Params:
            M_reduced (numpy matrix of shape (number of unique words in the corpus , 2)): matrix of 2-dimensioal word embeddings
            word2ind (dict): dictionary that maps word to indices for matrix M
            words (list of strings): words whose embeddings we want to visualize
    """

    # -------------------
    # Write your implementation here.

    # plt.scatter(M_reduced[:, 0], M_reduced[:, 1], c = 'olive', marker = 'x')
    plt.figure(figsize=(8,5))
    distance = 0.001
    for i, word in enumerate(words):
        coordinate = M_reduced[word2ind[words[i]], :]

        print(f" the word {word} lies at {coordinate}")

        plt.scatter(coordinate[0], coordinate[1], c = 'olive', marker = 'x')
        plt.text(coordinate[0] + distance , coordinate[1] + distance, word)

    # -------------------
```

In [28]:

```
# ---------------------
# Run this sanity check
# Note that this is not an exhaustive check for correctness.
# The plot produced should look like the "test solution plot" depicted below.
# ---------------------

print ("-" * 80)
print ("Outputted Plot:")

M_reduced_plot_test = np.array([[1, 1], [-1, -1], [1, -1], [-1, 1], [0, 0]])
word2ind_plot_test = {'test1': 0, 'test2': 1, 'test3': 2, 'test4': 3, 'test5': 4}
words = ['test1', 'test2', 'test3', 'test4', 'test5']
plot_embeddings(M_reduced_plot_test, word2ind_plot_test, words)

print ("-" * 80)
```

```
--------------------------------------------------------------------------------
Outputted Plot:
 the word test1 lies at [1 1]
 the word test2 lies at [-1 -1]
 the word test3 lies at [ 1 -1]
 the word test4 lies at [-1  1]
 the word test5 lies at [0 0]
--------------------------------------------------------------------------------
```

## Question 1.5: Co-Occurrence Plot Analysis [written] (3 points)

Now we will put together all the parts you have written! We will compute the co-occurrence matrix with fixed window of 4 (the default window size), over the Reuters "crude" (oil) corpus. Then we will use TruncatedSVD to compute 2-dimensional embeddings of each word. TruncatedSVD returns U*S, so we need to normalize the returned vectors, so that all the vectors will appear around the unit circle (therefore closeness is directional closeness). **Note**: The line of code below that does the normalizing uses the NumPy concept of *broadcasting*. If you don't know about broadcasting, check out Computation on Arrays: Broadcasting by Jake VanderPlas.

Run the below cell to produce the plot. It'll probably take a few seconds to run. What clusters together in 2-dimensional embedding space? What doesn't cluster together that you might think should have? **Note**: "bpd" stands for "barrels per day" and is a commonly used abbreviation in crude oil topic articles.

In [29]:

```
# ---------------------------
# Run This Cell to Produce Your Plot
# ---------------------------
reuters_corpus = read_corpus()
M_co_occurrence, word2ind_co_occurrence = compute_co_occurrence_matrix(reuters_corpus)
M_reduced_co_occurrence = reduce_to_k_dim(M_co_occurrence, k=2)

# Rescale (normalize) the rows to make them each of unit-length
M_lengths = np.linalg.norm(M_reduced_co_occurrence, axis=1)
M_normalized = M_reduced_co_occurrence / M_lengths[:, np.newaxis] # broadcasting

words = ['barrels', 'bpd', 'ecuador', 'energy', 'industry', 'kuwait', 'oil', 'output', 'petroleum', 'iraq']

plot_embeddings(M_normalized, word2ind_co_occurrence, words)
```

```
Running Truncated SVD over 8185 words...
The original size of  embedded matrix is (8185, 8185)
The reduced size of embedded matrix is (8185, 2)
Done.
 the word barrels lies at [ 0.97614032 -0.2171407 ]
 the word bpd lies at [ 0.95680174 -0.29074118]
 the word ecuador lies at [ 0.95578754 -0.09169068]
 the word energy lies at [ 0.99996827 -0.00796552]
 the word industry lies at [0.998579    0.05329144]
 the word kuwait lies at [ 0.99340325 -0.11467337]
 the word oil lies at [ 0.99884609 -0.04802582]
 the word output lies at [ 0.98371143 -0.1797549 ]
 the word petroleum lies at [0.99714793 0.0754719 ]
 the word iraq lies at [ 0.99671621 -0.08097401]
```



**Write your answer here.**

Two related clusters emerge: **the oil-exporting couties** and the oil-related words.

"iraq", "ecuador" and "kuwait" are closely related as they are the country name, specifically the major oil exporting nations.

"petroleum", "energy", and "oil" is another closely related topic, as it involves around environmental exsahusible resources for consumption. "industry" is also clusterd here as this word often take place with "energy" and/ or "oil".

For words that should have been closer to the oil-related clusters, "bpd" and "barrel" are the units for oil products, thus should be staying closer to "oil". But it turns out some neutral, generic word like "output" is closer.

## Part 2: Prediction-Based Word Vectors (13 points)

As discussed in class, more recently prediction-based word vectors have demonstrated better performance, such as word2vec and GloVe (which also utilizes the benefit of counts). If you're feeling adventurous, challenge yourself and try reading GloVe's original paper.

Then run the following cells to load the GloVe vectors into memory. **Note**: If this is your first time to run these cells, i.e. download the embedding model, it will take a couple minutes to run. If you've run these cells

before, rerunning them will load the model without redownloading it, which will take about 1 to 2 minutes.

In [30]:

```python
def load_embedding_model():
    """ Load GloVe Vectors
        Return:
            word_vectors: All 400000 embeddings, each lengh 100
    """
    import gensim.downloader as api
    word_vectors = api.load("glove-wiki-gigaword-300")
    # print("Loaded vocab size %i" % len(word_vectors))
    return word_vectors
```

In [31]:

```python
import gensim.downloader as api
word_vectors = api.load("glove-wiki-gigaword-300")
```

```
[==================================================] 100.0% 376.1/376.1MB downloaded
```

In [32]:

```python
word_vectors
```

Out[32]:

```
<gensim.models.keyedvectors.Word2VecKeyedVectors at 0x7f23200a2790>
```

In [33]:

```python
# ---------------------------------
# Run Cell to Load Word Vectors
# Note: This will take a couple minutes
# ---------------------------------
# word_vectors = load_embedding_model()
# type(word_vectors)
```

In [34]:

```python
# try some function most_similar
word_vectors.most_similar('queen')
```

Out[34]:

```
[('elizabeth', 0.6771447658538818),
 ('princess', 0.635676383972168),
 ('king', 0.6336469650268555),
 ('monarch', 0.5814188122749329),
 ('royal', 0.543052613735199),
 ('majesty', 0.5350357294082642),
 ('victoria', 0.5239557027816772),
 ('throne', 0.5097099542617798),
 ('lady', 0.5045416355133057),
 ('crown', 0.49980056285858154)]
```

**Note: If you are receiving a "reset by peer" error, rerun the cell to restart the download.**

## Cosine Similarity

Now that we have word vectors, we need a way to quantify the similarity between individual words, according to these vectors. One such metric is cosine-similarity. We will be using this to find words that are "close" and "far" from one another.

We can think of n-dimensional vectors as points in n-dimensional space. If we take this perspective L1 and L2 Distances help quantify the amount of space "we must travel" to get between these two points. Another approach is to examine the angle between two vectors. From trigonometry we know that:



$$\theta = \arccos(x \cdot y / |x| |y|)$$

Instead of computing the actual angle, we can leave the similarity in terms of $similarity = cos(\Theta)$. Formally the Cosine Similarity $s$ between two vectors $p$ and $q$ is defined as:

$$s = \frac{p \cdot q}{||p|| ||q||}, \text{ where } s \in [-1, 1]$$

## Question 2.1: Words with Multiple Meanings (1.5 points) [code + written]

Polysemes and homonyms are words that have more than one meaning (see this wiki page to learn more about the difference between polysemes and homonyms ). Find a word with *at least two different meanings* such that the top-10 most similar words (according to cosine similarity) contain related words from *both* meanings. For example, "leaves" has both "go_away" and "a_structure_of_a_plant" meaning in the top 10, and "scoop" has both "handed_waffle_cone" and "lowdown". You will probably need to try several polysemous or homonymic words before you find one.

Please state the word you discover and the multiple meanings that occur in the top 10. Why do you think many of the polysemous or homonymic words you tried didn't work (i.e. the top-10 most similar words only contain **one** of the meanings of the words)?

**Note**: You should use the `word_vectors.most_similar(word)` function to get the top 10 similar words. This function ranks all other words in the vocabulary with respect to their cosine similarity to the given word. For further assistance, please check the GenSim documentation.

In [35]:

```python
# ------------------
# Write your implementation here.
word_vectors.most_similar('book')
# ------------------
```

Out[35]:

```
[('books', 0.7986249327659607),
 ('author', 0.7123498916625977),
 ('published', 0.6973031759262085),
 ('novel', 0.6966710686683655),
 ('memoir', 0.6465641260147095),
 ('wrote', 0.631791889667511),
 ('biography', 0.6225202083587646),
 ('autobiography', 0.603348970413208),
 ('essay', 0.5995662212371826),
 ('illustrated', 0.5914922952651978)]
```

In [36]:

```python
word_vectors.most_similar('free')
```

Out[36]:

```
[('freedom', 0.4974668025970459),
```

```
('without', 0.48699691891670227),
('available', 0.48415109515190125),
('allowed', 0.4798740744590759),
('give', 0.4775182008743286),
('access', 0.47569626569747925),
('unrestricted', 0.4755344092845917),
('go', 0.47410497069358826),
('allow', 0.473043829202652),
('all', 0.47048091888427734)]
```

In [37]:
```
word_vectors.most_similar('kind')
```

Out[37]:
```
[('sort', 0.9157270789146423),
 ('thing', 0.8076795339584351),
 ('something', 0.8019971251487732),
 ('really', 0.770393503959656),
 ('think', 0.7361322641372681),
 ('what', 0.7342531681060791),
 ('nothing', 0.7274984121322632),
 ('anything', 0.7253236174583435),
 ('you', 0.7192934155464172),
 ('certainly', 0.7054327726364136)]
```

In [38]:
```
word_vectors.most_similar('bright')
```

Out[38]:
```
[('colors', 0.6087764501571655),
 ('blue', 0.595706582069397),
 ('yellow', 0.5886391401290894),
 ('dark', 0.5874146223068237),
 ('colored', 0.5830612778663635),
 ('brightly', 0.5787879824638367),
 ('brighter', 0.5764387845993042),
 ('light', 0.5631998777389526),
 ('shiny', 0.548277735710144),
 ('glow', 0.5474361181259155)]
```

In [39]:
```
word_vectors.most_similar('train')
```

Out[39]:
```
[('trains', 0.823096513748169),
 ('bus', 0.670767068862915),
 ('rail', 0.635927140712738),
 ('commuter', 0.5964051485061646),
 ('freight', 0.5893657803535461),
 ('passenger', 0.5863174200057983),
 ('railway', 0.5833747982978821),
 ('buses', 0.5471792221069336),
 ('subway', 0.5458870530128479),
 ('passengers', 0.5390786528587341)]
```

In [40]:
```
word_vectors.most_similar('let')
```

Out[40]:
```
[('want', 0.7677186727523804),
 ('you', 0.755077600479126),
 ("'ll", 0.7326532006263733),
 ('letting', 0.7323309779167175),
 ('go', 0.7310500741004944),
 ("n't", 0.7278926968574524),
 ('do', 0.7209917902946472),
 ('ca', 0.7142171263694763),
 ('tell', 0.7063310146331787),
 ('ask', 0.700111448764801)]
```

In [41]:
```
word_vectors.most_similar('flat')
```

Out[41]:
```
[('lower', 0.47144830226898193),
 ('thin', 0.4606584906578064),
 ('steep', 0.44942817091941833),
 ('bottom', 0.4260002374649048),
 ('prices', 0.4257771968841553),
 ('relatively', 0.4247564673423767),
 ('slightly', 0.42393964529037476),
 ('fairly', 0.40919122099876404),
 ('plain', 0.40724343061447144),
 ('falling', 0.4068793058395386)]
```

**Write your answer here.**

Most of the polysemes and hypernyms used above seems to have failed to capture its additional or secondary meaning.

For example, the word "book", when searched for 10 most similar words, the results are mostly about the materials related to knowledges, which is the primary meaning of such word. But the alternative meaning of "book" as verb, as in booking a hotel room or making reservation, is entirely overlooked.

The same phenomenon where the primary meaning also reflects in 10 most reflected words are 'train' and 'flat'.

As for train, the related words captures the primary use of a word as locomotive vehicles, not the verb as in 'train the muscles, train your mind, etc. Also, in "flate" where the adjective meaning is shown, without the other meaning as accomodotation, as in "he lives in a flat".

---

## Question 2.3: Synonyms & Antonyms (2 points) [code + written]

When considering Cosine Similarity, it's often more convenient to think of Cosine Distance, which is simply 1 - Cosine Similarity.

Find three words $(w_1, w_2, w_3)$ where $w_1$ and $w_2$ are synonyms and $w_1$ and $w_3$ are antonyms, but Cosine Distance $(w_1, w_3) <$ Cosine Distance $(w_1, w_2)$.

As an example, $w_1$="happy" is closer to $w_3$="sad" than to $w_2$="cheerful". Please find a different example that satisfies the above. Once you have found your example, please give a possible explanation for why this counter-intuitive result may have happened.

You should use the the `word_vectors.distance(w1, w2)` function here in order to compute the cosine distance between two words. Please see the [GenSim documentation](#) for further assistance.

In [42]:
```
# ------------------
# Write your implementation here.

w1 = 'black'
w2 = 'dark'
w3 = 'white'
print(f"The distance between synonym words {w1}, {w2}: {word_vectors.distance(w1, w2)}")
print(f"The distance between antonyms words {w1}, {w3}: {word_vectors.distance(w1, w3)}")

# ------------------
```

```
The distance between synonym words black, dark: 0.4207456707954407
The distance between antonyms words black, white: 0.2864179015159607
```

From thie above example, it shows that the distace between antonyms is smaller than that of synonym, meaning they are more related to one another!

**Write your answer here.**

The plausible explanation for such counterintuitive finding that antonyms are more closely related than synonyms is the frequency they appear together. Synonyms, though sharing the mutual semantics, may not occur adjacent to another synonyms, as it would lead to redundancy. On the other hand, opposite words may be place close to each other, for stark comparison between two things, making the algorithm understand that these antonyms are more interrelated.

.

## Question 2.4: Analogies with Word Vectors [written] (1.5 points)

Word vectors have been shown to *sometimes* exhibit the ability to solve analogies.

As an example, for the analogy "man : king :: woman : x" (read: man is to king as woman is to x), what is x?

In the cell below, we show you how to use word vectors to find x using the `most_similar` function from the [GenSim documentation](#). The function finds words that are most similar to the words in the `positive` list and most dissimilar from the words in the `negative` list (while omitting the input words, which are often the most similar; see [this paper](#)). The answer to the analogy will have the highest cosine similarity (largest returned numerical value).

In [43]:
```
# Run this cell to answer the analogy -- man : king :: woman : x
pprint.pprint(word_vectors.most_similar(positive=['woman', 'king'], negative=['man']))
```

```
[('queen', 0.6713277101516724),
 ('princess', 0.5432624220848083),
 ('throne', 0.5386104583740234),
 ('monarch', 0.5347574949264526),
 ('daughter', 0.498025119304657),
 ('mother', 0.4956442713737488),
 ('elizabeth', 0.4832652509212494),
 ('kingdom', 0.47747087478637695),
 ('prince', 0.4668239951133728),
 ('wife', 0.4647327661514282)]
```

## Question 2.5: Finding Analogies [code + written] (1.5 points)

Find an example of analogy that holds according to these vectors (i.e. the intended word is ranked top).

**Note**: You may have to try many analogies to find one that works!

In [44]:
```
# ------------------
# Write your implementation here.

word_vectors.most_similar(positive=['woman', 'male'], negative=['man'])
# ------------------
```

Out[44]:
```
[('female', 0.8488893508911133),
 ('women', 0.5879524946212769),
 ('males', 0.5201805830001831),
 ('females', 0.5195831656455994),
 ('pregnant', 0.5157788991928101),
 ('adult', 0.4962882995605469),
 ('sex', 0.48539209365844727),
 ('girls', 0.46756434440612793),
 ('gender', 0.45728668570518494),
 ('sexual', 0.44429826736450195)]
```

In [45]:
```
word_vectors.most_similar(positive=['daughter', 'boy'], negative=['son'])
```

Out[45]:
```
[('girl', 0.887003481388092),
 ('girls', 0.6548776030540466),
 ('mother', 0.6525065302848816),
 ('woman', 0.6493716239929199),
 ('child', 0.6013819575309753),
 ('grandmother', 0.5976220965385437),
 ('teenager', 0.589718222618103),
 ('teenage', 0.5784621238708496),
 ('boys', 0.5780875086784363),
 ('girlfriend', 0.5757790803909302)]
```

In [46]:
```
word_vectors.most_similar(positive=['sky', 'green'], negative=['trees'])
```

Out[46]:
```
[('blue', 0.5512756109237671),
 ('bright', 0.45733803510665894),
 ('dark', 0.44471269845962524),
 ('red', 0.4309699535369873),
 ('purple', 0.42887037992477417),
 ('brown', 0.41242167353630066),
 ('yellow', 0.3855505585670471),
 ('black', 0.3697638511657715),
 ('skies', 0.3605387806892395),
 ('colour', 0.3567811846733093)]
```

## Question 2.6: Incorrect Analogy [code + written] (1.5 points)

Find an example of analogy that does *not* hold according to these vectors.

In [47]:
```
# ------------------
# Write your implementation here.

word_vectors.most_similar(positive=['tv', 'read'], negative=['book'])

# ------------------
```

Out[47]:
```
[('television', 0.7062621116638184),
 ('broadcast', 0.6143285632133484),
 ('broadcasts', 0.5504031181335449),
 ('radio', 0.5457763671875),
 ('channel', 0.5182363986968994),
 ('watch', 0.5027680397033691),
 ('watching', 0.501125693321228),
 ('viewers', 0.48809871077537537),
```

```
    ('cnn', 0.48373645544052124),
    ('cable', 0.48029640316963196)]
```

```
word_vectors.most_similar(positive=['yellow', 'apple'], negative=['red'])
```

Out[48]:

```
[('iphone', 0.49975669384002686),
 ('ipad', 0.4942600727081299),
 ('ipod', 0.49260827898979187),
 ('macintosh', 0.4613226652145386),
 ('google', 0.45330381393432617),
 ('microsoft', 0.43000563979148865),
 ('intel', 0.42604002356529236),
 ('imac', 0.41385483741760254),
 ('ibm', 0.3970990777015686),
 ('motorola', 0.3851204514503479)]
```

In [49]:

```
word_vectors.most_similar(positive=['red', 'banana'], negative=['yellow'])
```

Out[49]:

```
[('bananas', 0.5103011131286621),
 ('mango', 0.4726327359676361),
 ('coffee', 0.4201294183731079),
 ('strawberry', 0.4139113128185272),
 ('fruit', 0.4084935784339905),
 ('sugar', 0.406865090113175964),
 ('papaya', 0.40092340111732483),
 ('growers', 0.3974452018737793),
 ('coconut', 0.396612107753753366),
 ('pineapple', 0.39483538269996643)]
```

In [50]:

```
word_vectors.most_similar(positive=['morning', 'sleep'], negative=['night'])
```

Out[50]:

```
[('asleep', 0.513579249382019),
 ('sleeping', 0.5048273801803589),
 ('woke', 0.4851740300655365),
 ('slept', 0.45095255970954895),
 ('bed', 0.44980132579803467),
 ('nap', 0.4347187280654907),
 ('awake', 0.42049551010131836),
 ('waking', 0.4183635115623474),
 ('awoke', 0.395452618598938),
 ('a.m.', 0.3887363076210022)]
```

## Question 2.7: Guided Analysis of Bias in Word Vectors [written] (1 point)

It's important to be cognizant of the biases (gender, race, sexual orientation etc.) implicit in our word embeddings. Bias can be dangerous because it can reinforce stereotypes through applications that employ these models.

Run the cell below, to examine (a) which terms are most similar to "woman" and "worker" and most dissimilar to "man", and (b) which terms are most similar to "man" and "worker" and most dissimilar to "woman". Point out the difference between the list of female-associated words and the list of male-associated words, and explain how it is reflecting gender bias.

In [51]:

```
# Run this cell
# Here `positive` indicates the list of words to be similar to and `negative` indicates the list of words to be
# most dissimilar from.
pprint.pprint(word_vectors.most_similar(positive=['woman', 'worker'], negative=['man']))
pprint.pprint(word_vectors.most_similar(positive=['man', 'worker'], negative=['woman']))
```

```
[('employee', 0.5915157794952393),
 ('workers', 0.5560789108276367),
 ('nurse', 0.514857828617096),
 ('pregnant', 0.4897522032260895),
 ('mother', 0.48388367891311646),
 ('female', 0.46243947744369507),
 ('child', 0.4448588490486145),
 ('teacher', 0.44152435660362244),
 ('waitress', 0.44121503829956055),
 ('employer', 0.4378713071346283)]
[('workers', 0.5640615224838257),
 ('employee', 0.5365462303161621),
 ('laborer', 0.483084499835968),
 ('working', 0.474678635597229),
 ('factory', 0.4493158459663391),
 ('mechanic', 0.43802663683891296),
 ('work', 0.4276600480079651),
 ('unemployed', 0.42742660641670227),
 ('worked', 0.4222966134548187),
 ('job', 0.42074185609817505)]
```

**Write your answer here.**

(a) which terms are most similar to "woman" and "worker" and most dissimilar to "man"

The result shows the stereotypical mentality about women in professional setting. Word related to woman in jobs involves around motherly nature such as 'nurse', 'teacher' and 'waitress', and also 'pregnant'.

(b) which terms are most similar to "man" and "worker" and most dissimilar to "woman"

A stereotype of physical strenghts used in job reflects throught the (b) as the list contains such as 'laborer', 'mechanic' and 'factory'. Some negative biase toward lazinesso of man is also captured as in the word 'unemployed'.

-

## Question 2.8: Independent Analysis of Bias in Word Vectors [code + written] (1 point)

Use the `most_similar` function to find another case where some bias is exhibited by the vectors. Please briefly explain the example of bias that you discover.

In [52]:

```
# -------------------
# Write your implementation here.
pprint.pprint(word_vectors.most_similar(positive=['female', 'food'], negative=['male']))
pprint.pprint(word_vectors.most_similar(positive=['male', 'food'], negative=['female']))
# -------------------
```

```
[('foods', 0.6103315353393555),
 ('supplies', 0.54627762117385864),
 ('products', 0.5149519443511963),
 ('meals', 0.511093258857727),
 ('shortages', 0.50217401981135376),
 ('aid', 0.4981635808944702),
 ('vegetables', 0.49693071842193604),
 ('eat', 0.49676570296287537),
 ('nutrition', 0.4899175763130188),
 ('meat', 0.48773324489593506)]
[('foods', 0.6060277223587036),
 ('eat', 0.5761719346046448),
 ('supplies', 0.5565237998962402),
```

```
 ('eating', 0.5500626564025879),
 ('meat', 0.5441587567329407),
 ('meal', 0.5331789255142212),
 ('meals', 0.5195803642272949),
 ('supply', 0.5031853318214417),
 ('feed', 0.4982718229293823),
 ('medicines', 0.4943087100982666)]
```

In [53]:

```python
# ------------------
# Write your implementation here.

pprint.pprint(word_vectors.most_similar(positive=['son', 'education'], negative=['daughter']))
pprint.pprint(word_vectors.most_similar(positive=['daughter', 'education'], negative=['son']))
# ------------------
```

```
[('educational', 0.6169756054878235),
 ('schools', 0.5731616020202637),
 ('health', 0.5516809225082397),
 ('teaching', 0.5405141711235046),
 ('curriculum', 0.537800133228302),
 ('school', 0.5329912900924683),
 ('vocational', 0.5234990119934082),
 ('programs', 0.509280800819397),
 ('teachers', 0.503345787525177),
 ('students', 0.49978548288345337)]
[('educational', 0.6250971555709839),
 ('schools', 0.570708692073822),
 ('teaching', 0.54721999168396),
 ('curriculum', 0.5404050946235657),
 ('teacher', 0.5312469005584717),
 ('school', 0.5282274484634399),
 ('health', 0.5276145339012146),
 ('programs', 0.5196062326431274),
 ('students', 0.5155026912689209),
 ('social', 0.5053571462631226)]
```

**Write your answer here.**

The first example of gender and food shows the bias about how women and men are supposed to eat. In words related to women, we see words like "vegetables", and "nutrituion" to suggest the bias that women should take care of their food regime, while for men-related words about food, we have "meat" instead to suggest that men are more carnivorous.

The second example shows the bias of educational choices based on gender. For "son", the more technical fields are more related as we have "vocational" listed, contrary to the "daughters" that has "social" implying that women are more generally exposed to less marketable fields.

.

---

**Question 2.9: Thinking About Bias [written] (2 points)**

Give one explanation of how bias gets into the word vectors. Argue whether this can be lead to problems into the society. Last, how do we address this.

**Write your answer here.**

The word vectors that we are implemented has been trained on the expansive amount of real world datasets such as from newspaper or social media. The bias simply reflects how the society percieve the characteristics of each topic that has been circulated through comments, articles, journal. etc.

These stereotypical idea, conveyed through coocurrence of words, is the source of model training, thus keeping the bias in its word vectors. Solutions about bias is a challenging topic, if not impossible to settle, as it involves the value, beliefe and perception that has been embedded at both personal, organizaitonal, national and even global levels.

One thing we can do about the bias, however, is to practice critical thinking. Bias or stereotype may give the general idea of how such things are percieved by the external parties, but it is not necessarily the universal truth that reflects the individual virtues.

In [53]: