# ASSIGNMENT 5 - Similarity Analysis on Tweets

**Matthew Dunne**

You have ~57K news articles loaded into Linux directory on RCC:
/project/msca/kadochnikov/webhose/news_university_2018_02.json

Your goal is to determine how many of these English articles are unique vs. how many are "near-duplicate".

## Load Data and Functions

In [1]:

```python
import pandas as pd
import re
import json
from itertools import combinations, takewhile
import collections
from matplotlib import pyplot as plt
from simhash import Simhash, SimhashIndex
```

In [2]:

```python
def get_features(s):
    width = 3
    s = s.lower()
    s = re.sub(r'[^\w]+', '', s)
    return [s[i:i + width] for i in range(max(len(s) - width + 1, 1))]
```

In [3]:

```python
data = pd.read_json('news_university_2018_02.json')
```

In [4]:

```python
data.head()
```

Out[4]:

| | author | crawled | entities | external_links | highlightText | highlightTitle | language | ord_ |
|---|---|---|---|---|---|---|---|---|
| 0 | Jackie Harper (noreply@blogger.com) | 2018-01-22T06:31:32.009+02:00 | {'persons': [], 'organizations': [{'name': 'ws... | [] | | | english | 0 |
| 1 | n0lifeismylife | 2018-01-22T06:31:35.026+02:00 | {'persons': [], 'organizations': [], 'location... | [] | | | english | 0 |
| 2 | redaintded | 2018-01-22T06:31:43.018+02:00 | {'persons': [], 'organizations': [], 'location... | [] | | | english | 0 |
| 3 | Manoj Pandey (noreply@blogger.com) | 2018-01-22T06:31:55.016+02:00 | {'persons': [{'name': 'vivek dubey', 'sentimen... | [] | | | english | 0 |

| | author | crawled | entities | external_links | highlightText | highlightTitle | language | ord_ |
|---|---|---|---|---|---|---|---|---|
| 4 | vision Arts (noreply@blogger.com) | 2018-01-22T06:32:24.091+02:00 | {'persons': [{'name': 'sharon wagner', 'sentim... | [] | | | english | 0 |

In [5]:

```
data.shape
```

Out[5]:

```
(57500, 15)
```

## Similarity Analysis

1. Run similarity analysis on the "title" variable
2. Explain how you selected a similarity threshold for "near-duplicate"

To do this we shall use the following steps:

- Determine how many unique values there are in the title variable (those that are not unique are perfect duplicates) and then run the analysis only the unique values.
- Run the SimHash on a given threshold (start at 1).
- Cycle through each title and assess whether it is unique or a near duplicate at that threshold.
- Sample some titles (actually groups of titles) that are marked as near duplicate to see how the given threshold performs in labelling titles as unique or near duplicate. We will take several samples.
- Once we increase the threshold to a point that the titles marked as near duplicate are not really near duplicates we stop increasing the threshold.

In [6]:

```
#make sure it is just the English tweets
data=data[data['language']=='english']
```

How many unique values are there in the title variable?

In [7]:

```
len(data['title'].unique())
```

Out[7]:

```
44527
```

This means there are about 13,000 perfect duplicates.

In [8]:

```
len(data['title'])-len(data['title'].unique())
```

Out[8]:

```
12973
```

We will take these out of consideration.

In [9]:

```
#just take the title column as a list
titles = data['title'].unique()
#convert to series so you can make into dictionary
titles = pd.Series(titles)
```

```
                                    ...
#create a dataframe from it
titles_df = pd.DataFrame(titles)
#create a dictionary from it
titles_dict = titles.to_dict()
```

In [10]:

```
titles_df.rename(columns={0: 'Title'}, inplace=True)
titles_df.head()
```

Out[10]:

| | Title |
|---|---|
| 0 | FREE public planetarium shows at Wayne State U... |
| 1 | A computer Science professor at my university ... |
| 2 | Whos the Greatest College Football Team Ever |
| 3 | FULL HD VIDEO || ओढ़नी से मुह बाध के || College... |
| 4 | Tips For Newborn Photography College Station TX |

### *First a threshold of 1*

In [11]:

```
#NEED TO MAKE SURE THE INDEX IS RESET FOR THE DICTIONARY IN ORDER FOR THE SIMHASH TO WORK
#get the features for all titles. We can re-use this across thresholds
objs = [(str(k), Simhash(get_features(v))) for k, v in titles_dict.items()]
#run the simhash on a threshold of one on the dictionary
index_one = SimhashIndex(objs, k=1)
```

In [41]:

```
objs[0]
```

Out[41]:

```
('0', <simhash.Simhash at 0x23d3cd80ba8>)
```

In [13]:

```
#create a list, based on the dictionary, of whether a title is unique or near duplicate
unique_or_dup_one = []
#cycle through all titles, get features and Simhash. Then get near duplicates.
#If more than one record is returned it's a near duplicate (if it returns one it is just the origi
nal title)
#mark it a unique or near duplicate
for i in range(len(titles_dict)):
        s1 = Simhash(get_features(titles_dict[i]))
        s1_dups = index_one.get_near_dups(s1)
        if len(s1_dups)>1:
            unique_or_dup_one.append('near-duplicate')
        else:
            unique_or_dup_one.append('unique')
```

How many titles are deemed to be near duplicate?

In [14]:

```
#add the list you just created to the dataframe to mark the title as unique or near duplicate base
d on the threshold
se1=pd.Series(unique_or_dup_one)
titles_df['Threshold_1']=se1.values
#limit to just the near dups
near_dup_one= titles_df[titles_df['Threshold_1']=='near-duplicate']
len(near_dup_one)
```

1785

Now we look at how similar the documents marked as near duplicate actually are. We shall take three random samples:

In [20]:

```
#select three random entries from the near dups
sample_one=near_dup_one.sample(n=3, random_state=80129111)
sample_one
```

Out[20]:

| | Title | Threshold_1 |
|---|---|---|
| **5842** | Eldest son of allegedly captive siblings described by peers at college as 'pale' and 'depressive' | near-duplicate |
| **17361** | 'Deplorable' NYU Professor Sues University, Colleagues for Defamation | near-duplicate |
| **29931** | BREAKING News: Buhari Renames Federal University Ebonyi, 'Alex Ekwueme University | near-duplicate |

In [23]:

```
#look at the duplicates of each
test = Simhash(get_features(titles_dict[sample_one.index[0]]))
test_dups = index_one.get_near_dups(test)
pd.set_option('max_colwidth',100)
titles_df.iloc[test_dups]
```

Out[23]:

| | Title | Threshold_1 |
|---|---|---|
| **5842** | Eldest son of allegedly captive siblings described by peers at college as 'pale' and 'depressive' | near-duplicate |
| **6199** | Eldest son of allegedly captive siblings described by peers at college as 'pale' and 'depressive' | near-duplicate |

In [24]:

```
test = Simhash(get_features(titles_dict[sample_one.index[1]]))
test_dups = index_one.get_near_dups(test)
pd.set_option('max_colwidth',100)
titles_df.iloc[test_dups]
```

Out[24]:

| | Title | Threshold_1 |
|---|---|---|
| **18839** | 'Deplorable' NYU Professor Sues University, Colleagues for Defamation | near-duplicate |
| **17361** | 'Deplorable' NYU Professor Sues University, Colleagues for Defamation | near-duplicate |

In [25]:

```
test = Simhash(get_features(titles_dict[sample_one.index[2]]))
test_dups = index_one.get_near_dups(test)
titles_df.iloc[test_dups]
```

Out[25]:

| | Title | Threshold_1 |
|---|---|---|
| **29521** | BREAKING News: Buhari Renames Federal University Ebonyi, "Alex Ekwueme University" | near-duplicate |
| **26281** | BREAKING News: Buhari Renames Federal University Ebonyi, "Alex Ekwueme University" | near-duplicate |
| **29931** | BREAKING News: Buhari Renames Federal University Ebonyi, 'Alex Ekwueme University | near-duplicate |

At threshold = 1 the near duplicates are actually quite close. In fact it seems that the reason that some of these were not caught as perfect duplicates only because the apostrophes have a slightly different format. We shall increase the threshold.

### *Threshold = 2*

In [26]:

```python
#set the threshold = 2
index_two = SimhashIndex(objs, k=2)

unique_or_dup_two = []
for i in range(len(titles_dict)):
        #this is different than what you did for Threshold 1 because you are not re-doing get feat
ures
        #so no equivalent of s1 = Simhash(get_features(titles_dict[i])) is necessary
        s2_dups = index_two.get_near_dups(objs[i][1])
        if len(s2_dups)>1:
            unique_or_dup_two.append('near-duplicate')
        else:
            unique_or_dup_two.append('unique')
```

In [27]:

```python
#add the list you just created to the dataframe to mark the title as unique or near duplicate base
d on the threshold
se2=pd.Series(unique_or_dup_two)
titles_df['Threshold_2']=se2.values
```

Let us look at the titles that were unique under a threshold of 1 but are near duplicate under a threshold of 2.

In [28]:

```python
near_dup_two= titles_df[(titles_df['Threshold_2']=='near-duplicate') & (titles_df['Threshold_1']==
'unique')]
```

In [29]:

```python
#select three random entries
sample_two = near_dup_two.sample(n=3, random_state=1234)
sample_two
```

Out[29]:

| | | Title | Threshold_1 | Threshold_2 |
|---|---|---|---|---|
| 32844 | Women's Track vs Queens University Multis | | unique | near-duplicate |
| 31941 | Anthony Chukwudi Mazeli, Soldier Tops His Class In University Of Lancaster, UK (1) | | unique | near-duplicate |
| 5256 | Skylight damaged when ice flies off wind turbine at Mount Wachusett Community College #MA | | unique | near-duplicate |

In [30]:

```python
test = Simhash(get_features(titles_dict[sample_two.index[0]]))
test_dups = index_two.get_near_dups(test)
pd.set_option('max_colwidth',150)
titles_df.iloc[test_dups]
```

Out[30]:

| | | Title | Threshold_1 | Threshold_2 |
|---|---|---|---|---|

| | Title | Threshold_1 | Threshold_2 |
|---|---|---|---|
| 32843 | Men's Track vs Queens University Multis | unique | near-duplicate |
| 32844 | Women's Track vs Queens University Multis | unique | near-duplicate |

In [31]:

```python
test = Simhash(get_features(titles_dict[sample_two.index[1]]))
test_dups = index_two.get_near_dups(test)
pd.set_option('max_colwidth',100)
titles_df.iloc[test_dups]
```

Out[31]:

| | Title | Threshold_1 | Threshold_2 |
|---|---|---|---|
| 32549 | Anthony Chukwudi Mazeli, Soldier Tops His Class In University Of Lancaster, UK | unique | near-duplicate |
| 31941 | Anthony Chukwudi Mazeli, Soldier Tops His Class In University Of Lancaster, UK (1) | unique | near-duplicate |

In [32]:

```python
test = Simhash(get_features(titles_dict[sample_two.index[2]]))
test_dups = index_two.get_near_dups(test)
pd.set_option('max_colwidth',100)
titles_df.iloc[test_dups]
```

Out[32]:

| | Title | Threshold_1 | Threshold_2 |
|---|---|---|---|
| 3860 | Skylight damaged when ice flies off wind turbine at Mount Wachusett Community College | unique | near-duplicate |
| 5256 | Skylight damaged when ice flies off wind turbine at Mount Wachusett Community College #MA | unique | near-duplicate |

After taking a few samples I found that some of the "near duplicates" can start to breakdown. The first example under this threshold says a title on men's track and women's track are near duplicate, though for the most part looks like threshold = 2 holds up pretty well.

### Threshold = 3

In [33]:

```python
#set the threshold = 3
index_three = SimhashIndex(objs, k=3)

unique_or_dup_three = []
for i in range(len(titles_dict)):
        s3_dups = index_three.get_near_dups(objs[i][1])
        if len(s3_dups)>1:
            unique_or_dup_three.append('near-duplicate')
        else:
            unique_or_dup_three.append('unique')
```

In [34]:

```python
#add the list you just created to the dataframe to mark the title as unique or near duplicate based on the threshold
se3=pd.Series(unique_or_dup_three)
titles_df['Threshold_3']=se3.values
```

Let us look at the titles that were unique under a threshold of 2 but are near duplicate under a threshold of 3.

In [35]:

```python
near_dup_three= titles_df[(titles_df['Threshold_3']=='near-duplicate') & (titles_df['Threshold_2']=='unique')]
#select three random entries
```

```
#select three random entries
sample_three = near_dup_three.sample(n=3, random_state=80129111)
```

In [36]:

```
test = Simhash(get_features(titles_dict[sample_three.index[0]]))
test_dups = index_three.get_near_dups(test)
pd.set_option('max_colwidth',150)
titles_df.iloc[test_dups]
```

Out[36]:

|  | Title | Threshold_1 | Threshold_2 | Threshold_3 |
|---|---|---|---|---|
| **26638** | Chitwan: Gaindakot FC Enters QFs Of Pacific College 4th Chitwan Championship | unique | unique | near-duplicate |
| **39732** | Chitwan:Gaindakot FC Enters SFs Of Pacific College 4th Chitwan Championship | unique | unique | near-duplicate |

In [37]:

```
test = Simhash(get_features(titles_dict[sample_three.index[1]]))
test_dups = index_three.get_near_dups(test)
pd.set_option('max_colwidth',150)
titles_df.iloc[test_dups]
```

Out[37]:

|  | Title | Threshold_1 | Threshold_2 | Threshold_3 |
|---|---|---|---|---|
| **14549** | Doctoral student can predict where football prospects will go to college | unique | unique | near-duplicate |
| **14511** | How an Iowa doctoral student can predict where football prospects will go to college | near-duplicate | near-duplicate | near-duplicate |

In [38]:

```
test = Simhash(get_features(titles_dict[sample_three.index[2]]))
test_dups = index_three.get_near_dups(test)
pd.set_option('max_colwidth',150)
titles_df.iloc[test_dups]
```

Out[38]:

|  | Title | Threshold_1 | Threshold_2 | Threshold_3 |
|---|---|---|---|---|
| **10126** | BAGHDAD \| Al-Mashriq University \| U/C | unique | unique | near-duplicate |
| **10082** | BAGHDAD \| Al-Mashriq University | unique | unique | near-duplicate |

Here we see that using a threshold of 3 we see some instances of classification of near duplication as correct, but others are not, as in the first example where it classifies the article about a team reaching the Quarterfinals (QF) to be a near duplicate of an article about that team reaching the Semifinals (SF).

We also see in the second example how a title (index 14511) was classified as a near duplicate under threshold=2 but then under threshold=3 is classified as a near duplicate of a different title.

**For these reasons I will stop at a threshold of 3. If we go any further we risk making clearly erroneous decisions in classifying documents at near duplicates.**

We see this when we go to a threshold of 4.

## *Threshold = 4*

```python
#set the threshold = 4
index_four = SimhashIndex(objs, k=4)

unique_or_dup_four = []
for i in range(len(titles_dict)):
        s4_dups = index_four.get_near_dups(objs[i][1])
        if len(s4_dups)>1:
            unique_or_dup_four.append('near-duplicate')
        else:
            unique_or_dup_four.append('unique')
```

```
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
```

```
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
```

```
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
```

```
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
```

```
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
```

```
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
```

```
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
```

```
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
```

```
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
```

```
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
```

```
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
```

```
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
```

```
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
```

```
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
```

```
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
```

```
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
```

```
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
```

```
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
```

```
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:4c9:1, len:206
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
```

```
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:cc9:1, len:254
Big bucket found. key:ad6:3, len:230
Big bucket found. key:4c9:1, len:206
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:4c9:1, len:206
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:8d6:3, len:282
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:2d6:3, len:205
Big bucket found. key:d6:3, len:326
Big bucket found. key:d6:3, len:326
Big bucket found. key:ad6:3, len:230
Big bucket found. key:cc9:1, len:254
Big bucket found. key:4c9:1, len:206
Big bucket found. key:ad6:3, len:230
Big bucket found. key:8d6:3, len:282
Big bucket found. key:d6:3, len:326
Big bucket found. key:cc9:1, len:254
```

In [346]:

```python
#add the list you just created to the dataframe to mark the title as unique or near duplicate based on the threshold
se4=pd.Series(unique_or_dup_four)
titles_df['Threshold_4']=se4.values
```

In [348]:

```python
near_dup_four= titles_df[(titles_df['Threshold_4']=='near-duplicate') & (titles_df['Threshold_3']=='unique')]
#select three random entries
sample_four = near_dup_four.sample(n=3, random_state=80129111)
```

In [349]:

```python
test = Simhash(get_features(titles_dict[sample_four.index[0]]))
test_dups = index_four.get_near_dups(test)
pd.set_option('max_colwidth',150)
```

```
titles_df.iloc[test_dups]
```

Out[349]:

| | Title | Threshold_1 | Threshold_2 | Threshold_3 | Threshold_4 |
|---|---|---|---|---|---|
| 28649 | Manasuku Nachindi Movie Team At Kasturba Gandhi College for Women is an institute in west Marredpally- 2 | unique | unique | unique | near-duplicate |
| 28650 | Manasuku Nachindi Movie Team At Kasturba Gandhi College for Women is an institute in west Marredpally- 3 | unique | unique | unique | near-duplicate |
| 28644 | Manasuku Nachindi Movie Team At Kasturba Gandhi College for Women is an institute in west Marredpally- 1 | unique | unique | unique | near-duplicate |

In [350]:
```
test = Simhash(get_features(titles_dict[sample_four.index[1]]))
test_dups = index_four.get_near_dups(test)
pd.set_option('max_colwidth',150)
titles_df.iloc[test_dups]
```

Out[350]:

| | Title | Threshold_1 | Threshold_2 | Threshold_3 | Threshold_4 |
|---|---|---|---|---|---|
| 41407 | iwantmyvices comments on I recently landed a well paying job. I'm still in college, and need some good advice going forward. | unique | unique | unique | near-duplicate |
| 41319 | DejectedUnicorn comments on I recently landed a well paying job. I'm still in college, and need some good advice going forward. | unique | unique | unique | near-duplicate |
| 41465 | NewbSaysRawr comments on I recently landed a well paying job. I'm still in college, and need some good advice going forward. | unique | unique | unique | near-duplicate |

In [351]:
```
test = Simhash(get_features(titles_dict[sample_four.index[2]]))
test_dups = index_four.get_near_dups(test)
pd.set_option('max_colwidth',150)
titles_df.iloc[test_dups]
```

Big bucket found. key:ad6:3, len:230

Out[351]:

| | Title | Threshold_1 | Threshold_2 | Threshold_3 | Threshold_4 |
|---|---|---|---|---|---|
| 1554 | Oxford University gives women more time to pass exams (Tony Diver/Telegraph) | unique | unique | unique | near-duplicate |
| 726 | Oxford University gives women more time to pass exams | unique | unique | unique | near-duplicate |

At a threshold of 4 we start to see titles marked as near duplicate that are somewhat similar but cannot really be considered near duplicate. That is why we will stick with a threshold of 3.

# Visualization

**Build a bar-chart visualization for two variables (count of unique and count of "near-duplicate")**

I will assume this means under the threshold that I chose.

**note from Igor's comments: "Question 6) meant that once you select a threshold, create a histogram of bucket sizes"**

In [1]:

```
#THIS IS HOW YOU GET NUMBER OF BUCKETS (WHICH IS WHAT HE MEANS)
#objs = [(str(k), Simhash(get_features(v))) for k, v in data.items()]
#index = SimhashIndex(objs, k=3) #`k` is the tolerance
#print (index.bucket_size())
```

In [57]:

```
counts=titles_df.Threshold_3.value_counts()
counts=pd.Series(counts)
counts
```
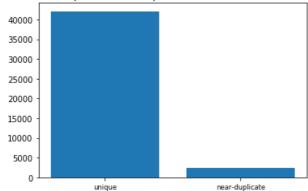
Out[57]:

```
unique           42077
near-duplicate    2450
Name: Threshold_3, dtype: int64
```

In [58]:

```
plt.bar(range(len(counts)), counts.values, align='center')
plt.title('Distribution of Unique and Near Duplicate Titles Under Simhash Threshold = 3')
plt.xticks(range(len(counts)), counts.index.values, size='small')
plt.show()
```

Distribution of Unique and Near Duplicate Titles Under Simhash Threshold = 3



Keep in mind there are about 13,000 exact duplicates which are neither unique or near-duplicates. They are not included in this chart.

**Build a histogram showing the overall distribution of "near-duplication"**

I will assume this means for all thresholds up to and including the threshold that I chose.

In [62]:

```
#get the number of near-duplicates under Thresholds 1 and 2
counts_one=titles_df[titles_df['Threshold_1']=='near-duplicate']
counts_one=len(counts_one)
counts_two=titles_df[titles_df['Threshold_2']=='near-duplicate']
counts_two=len(counts_two)
```
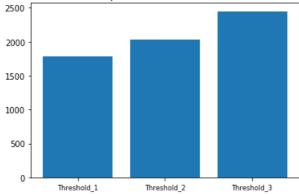
Out[62]:

```
1785
```

In [67]:

```
Near_Dup = {'Threshold_1': counts_one, 'Threshold_2': counts_two, 'Threshold_3': counts.loc['near-d
uplicate']}
df = pd.Series(Near_Dup)
plt.bar(range(len(df)), df.values, align='center')
plt.title('Distribution of Near Duplicate Titles Under Simhash Threshold 1 to 3')
plt.xticks(range(len(df)), df.index.values, size='small')
```

```
plt.xticks(range(len(df)), df.index.values, size='small')
plt.show()
```

Distribution of Near Duplicate Titles Under Simhash Threshold 1 to 3



We see the number of near duplicates slowly but steadily increasing as the Simhash threshold increases.