# Tokenization

```python
import nltk, string

def tokenize(text):
    tokens = []

    for word in nltk.word_tokenize(text):
        if word \
            not in nltk.corpus.stopwords.words('english') \
            and word not in string.punctuation \
            and word != '``':
                tokens.append(word)

    return tokens
```

```python
tokenized_rdd = essay_rdd_repartition.filter(lambda row: row['essay'] and row['essay'] != '') \
                    .map(lambda row: row['essay']) \
                    .map(lambda text: text.replace('\\n', '').replace('\r', '')) \
                    .map(lambda text: tokenize(text))
```

```python
tokenized_rdd.cache()
```

```
PythonRDD[31] at RDD at PythonRDD.scala:43
```

# Vectorization

```python
vocab = tokenized_rdd.flatMap(lambda words: words).distinct()
vocab.collect()
```
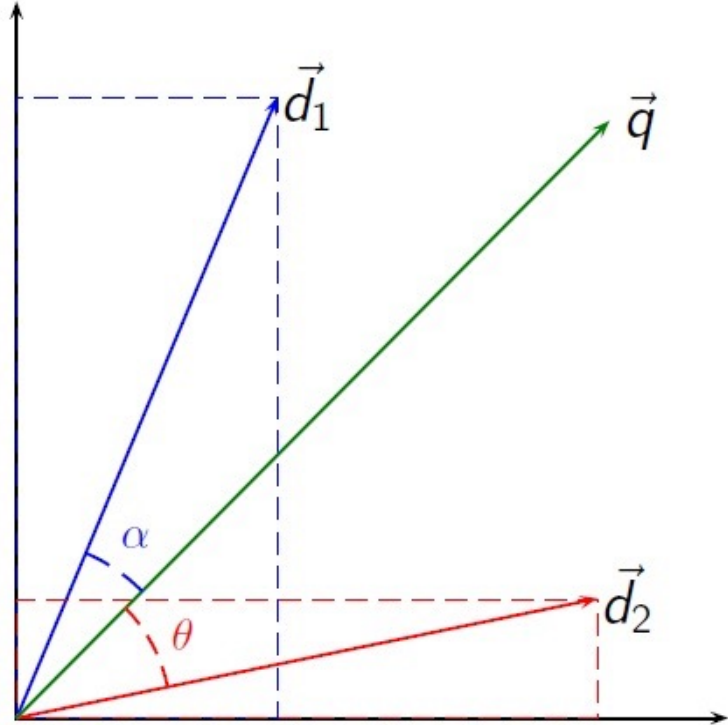
```python
from collections import Counter
import numpy as np

broadcastVocab = sc.broadcast(vocab.collect())

def bow_vectorize(tokens):
    word_counts = Counter(tokens)
    vector = [word_counts[v] if v in word_counts else 0 for v in broadcastVocab.value]
    return np.array(vector)
```

```python
tokenized_rdd.map(bow_vectorize).collect()
```

# Vector Space Model

Similarity is a measure of "distance"

©2016 Pearson, Inc.

# Interlude: How to Scale