# Handout 4: Wordlists and Frequencies

## Texts and wordlists

**1.** A text behaves just like a list.

```
1    >>> text1[:6]
2    ['[', 'Moby', 'Dick', 'by', 'Herman', 'Melville']
3    >>> text1[1]
4    'Moby'
5    >>> len(text1)
6    260819
```

**2.** But it has some additional **methods**.

```
1    >>> text1.name
2    'Moby Dick by Herman Melville 1851'
3    >>> text1.count('Ishmael')
4    19
5    >>> text1.concordance('monstrous')
```

**3.** You can create your own text from a list.

```
1    >>> mytext = Text(['Can', 'a', 'wood', 'chuck', 'chuck', 'wood', '?'])
2    >>> len(mytext)
3    7
4    >>> mytext.count('wood')
5    2
```

**4.** *Review.*

    **a.** To eliminate duplicates, convert to a set:  `vocab = set(text1)`

    **b.** Elements of the vocabulary are **types**, elements of the text are **tokens**

    **c.** Alphabetizing:

```
1    >>> v1 = set(mytext)
2    >>> sorted(v1)
3    ['?', 'Can', 'a', 'chuck', 'wood']
```

    **d.** Downcasing, eliminating punctuation:  `set(w.lower() for w in text1 if w.isalpha())`

## Sorting

**5.** Use the function `sorted`

    **a.** Works with any iterable; produces a list

```
1    >>> words = {'hi', 'bobby', 'bye'}
2    >>> sorted(words)
3    ['bobby', 'bye', 'hi']
4    >>> sorted([5, 0, 3])
5    [0, 3, 5]
```

   **b.** What if you want biggest to smallest?

```
1    >>> sorted([5, 0, 3], reverse=True)
2    [5, 3, 0]
```

   **c.** "reverse=True" is a **keyword argument**

**6.** What if you want to sort words by length?

   **a.** Keyword argument "key"

```
1    >>> sorted(words, key=len)
2    ['hi', 'bye', 'bobby']
```

   **b.** Notice: the value of `len` is a *function.* The function is called on each word, and the words are sorted according to the values:

| **Word:** | 'hi' | 'bobby' | 'bye' |
|-----------|------|---------|-------|
| **Len:**  | 2    | 5       | 3     |

   **c.** Combining `key` and `reverse`:

```
1    >>> sorted(words, key=len, reverse=True)
2    ['bobby', 'bye', 'hi']
```

**7.** Stability

   **a.** Sorting is guaranteed to be **stable**: words that tie are kept in their *original order*

   **b.** So—if we want tying words to be sorted alphabetically, do an alphabetic sort first:

```
1    >>> tmp = sorted({'hi', 'cat', 'bye'})
2    >>> tmp
3    ['bye', 'cat', 'hi']
4    >>> sorted(tmp, key=len)
5    ['hi', 'bye', 'cat']
```

**8.** *Exercises.*

   **a.** Sort `text1` ... `text9` by their length (number of tokens).

   **b.** Sort them by their vocabulary size.

**9.** Note: one can create an "anonymous" function with `lambda`:

```
1    >>> sorted(x, key=lambda text: text.count('you')/len(text))
```

# A little data exploration

10. Lexical diversity

    a. Contrary to the book, let us define the lexical diversity of a text to be the average number of types per 1000 tokens

    b. *Moby Dick* has 19,317 types per 260,819 tokens

    ```
    1    >>> 19317 / 260819
    2    0.07406285585022564
    3    >>> _ * 1000
    4    74.06285585022563
    ```

    c. General function

    ```
    1    >>> def diversity (text):
    2    ...      k = len(set(text))
    3    ...      n = len(text)
    4    ...      return 1000 * k / n
    5    ...
    6    >>> diversity(text1)
    7    74.06285585022563
    ```

11. Note: k and n are **local variables**; created by assignment. Completely contained inside the function.

    ```
    1    >>> k = 'hi'
    2    >>> diversity(text1)
    3    74.06285585022563
    4    >>> k
    5    'hi'
    ```

12. Diversity of personals

    a. Consider the following. Is this surprising?

    ```
    1    >>> diversity(text8)
    2    227.65564002465584
    3    >>> text8
    4    <Text: Personals Corpus>
    ```

    b. What does the personals corpus look like?

    ```
    1    >>> ' '.join(text8[:8])
    2    '25 SEXY MALE , seeks attrac older single'
    ```

    c. Seems like a pretty limited vocabulary. Why is the diversity so high?

    d. Is it any different somewhere in the middle?

    ```
    1    >>> ' '.join(text8[2000:2010])
    2    'rship . WLTM sincere , caring Lady to share life'
    ```

3

**13.** How else do the texts differ?

    **a.** In length: `len(text1)` is 260,819; `len(text8)` is 4867.

    **b.** What if we shorten *Moby Dick*? Aha!

```
1    >>> diversity(text1[:4867])
2    335.319498664475
```

**14.** *Homework.* Redefine `diversity()` to count the number of types in the first thousand words of the text. Which text is most diverse now? Sort the texts from most diverse to least diverse.

## Frequency dists

**15.** A **dict** is a table mapping keys to values

```
1    >>> tab = dict()
2    >>> tab['moby'] = len(text1)
3    >>> tab['moby']
4    260819
```

**16.** A frequency distribution is a specialized dict that maps items to counts.

```
1    >>> tokens = list('abbdabdbbd')
2    >>> fd = FreqDist(tokens)
3    >>> fd.tabulate()
4       b    d    a
5       5    3    2
```

**17.** Access

    **a.** Original items (tokens) are **keys**; access counts by key

```
1    >>> fd['a']
2    2
```

    **b.** Relative frequency (probability):

```
1    >>> fd.freq('a')
2    0.2
3    >>> fd['a'] / fd.N()
4    0.2
```

**18.** Methods

    **a.** How many types? How many tokens?

```
1    >>> len(fd)
2    3
3    >>> fd.N()
4    10
```

**b.** Dist behaves like list of keys (random order)

```
1    >>> for k in fd:
2    ...     print(k, fd[k])
3    ...
4    d 3
5    b 5
6    a 2
```

**c.** Sorting by frequency

```
1    >>> fd.most_common()
2    [('b', 5), ('d', 3), ('a', 2)]
3    >>> fd.most_common(1)
4    [('b', 5)]
```

**19.** ('b', 5) is a **tuple**.

    **a.** Tuples are just like lists, except they cannot be modified.

```
1    >>> x = ('b', 5)
2    >>> len(x)
3    2
4    >>> x[0]
5    'b'
```

    **b.** A useful trick

```
1    >>> (w, ct) = x
2    >>> w
3    'b'
4    >>> ct
5    5
```

**20.** *Exercises.*

    **a.** How do we find the five most-frequent tokens in *Moby Dick*?

    **b.** What if we want just the words, without counts?

    **c.** What if we want only the real words, not punctuation?

    **d.** What are the commonest word lengths in *Moby Dick*?

**21.** Zipf's Law

    **a.** Let $w_1$ be the most-frequent word, $w_2$ the second most-frequent, etc. Zipf's law states:
$$f(w_r) = K/r$$

    **b.** It implies that the most-frequent words account for most of the text

    **c.** It also implies that the commonest rate of occurrence is 1.

    **d.** *Homework.* Determine whether these statements are true for our texts

22. **Cumulative frequency**: add up the frequencies of the first $n$ elements

    a. Aggregation function `sum`

```
1    >>> freqs = [.4, .3, .2, .1]
2    >>> sum(freqs)
3    0.9999999999999999
```

    b. How do we get the sum of the first $n$ freqs?

23. Joint distributions.

    a. A distribution over pairs of items is a **joint** distribution.

    b. Pairs of adjacent words are called **bigrams**. For our corpus:

|   | $a$ | $b$ | $d$ |     |
|---|-----|-----|-----|-----|
| $a$ |     | 2/9 |     | 2/9 |
| $b$ |     | 2/9 | 3/9 | 5/9 |
| $d$ | 1/9 | 1/9 |     | 2/9 |
|   | 1/9 | 5/9 | 3/9 |     |

24. Marginal probability.

    a. The distribution over the individual items is called the **marginal** distribution.

    b. In the case of bigrams, the marginal distribution is called the **unigram distribution**.

    c. There are actually two unigram distributions. Explain. How could we fix that?

25. Bigrams

```
1    >>> bigrams(tokens)
2    <generator object bigrams at 0x1075e0990>
3    >>> list(bigrams(tokens))
4    [('a', 'b'), ('b', 'b'), ('b', 'd'), ('d', 'a'), ...]
5    >>> bd = FreqDist(bigrams(tokens))
6    >>> bd.freq(('a','b'))
7    0.2222222222222222
8    >>> bd.most_common(3)
9    [(('b', 'd'), 3), (('b', 'b'), 2), (('a', 'b'), 2)]
```

26. *Exercises.*

    a. What are the most common pairs of words in *Moby Dick*?

```
1    [(',', 'and'), ('of', 'the'), ("'", 's'), ('in', 'the')]
```
    (How did I get that result?)

    b. What are the most common words preceding *whale* in *Moby Dick*?

```
1    [(',', 18713), ('the', 13721), ('.', 6861), ('of', 6536)]
```
    c. How do we get the *most significant* pairs? $\rightarrow$ collocations