

Motivation

Suppose I had Melville's *MOBY DICK* stored in a text file called `moby.txt`. What if I was interested in finding the most frequent word used in the text? It's easy enough to hold all of *MOBY DICK* in memory, so I can read the entire text into a string, split the words using whitespace as my delimiter and produce a list of words, which we call tokens.

```

1 def file_to_tokens(filename):
2     with open(filename) as fin:
3         return fin.read().split()

```

Now I'm left with the task of counting how many times each token occurs in the list. I could use list operations to first find the set of unique tokens, and then count the occurrences of those tokens.

```

1 def wc_list(tokens):
2     uniq = []
3     for token in tokens:
4         if token not in uniq:
5             uniq.append(token)
6     return [(t, tokens.count(t)) for t in uniq]

```

Let's think about this. Suppose there are n unique tokens. When considering a new token, we had to scan the list of unique tokens seen so far. In the worst case, our set of unique tokens has size that's linear in n , which means that finding the unique tokens takes around n^2 operations. After that, we still have to count the number of occurrences of each unique token in the original token list, which is at least another n^2 operations because n is bounded above by the number of words in *MOBY DICK*.

Let's do a back-of-the-envelope calculation. There are around 212,000 tokens in *MOBY DICK*. Suppose that we only consider the first 5,000 tokens:

```

>>> tokens = file_to_tokens("moby.txt")
>>> first5000 = tokens[:5000]

```

There are 2285 unique tokens in the first 5,000 words of *MOBY DICK* and it takes around 0.2 seconds to both identify these words and construct their counts.

```

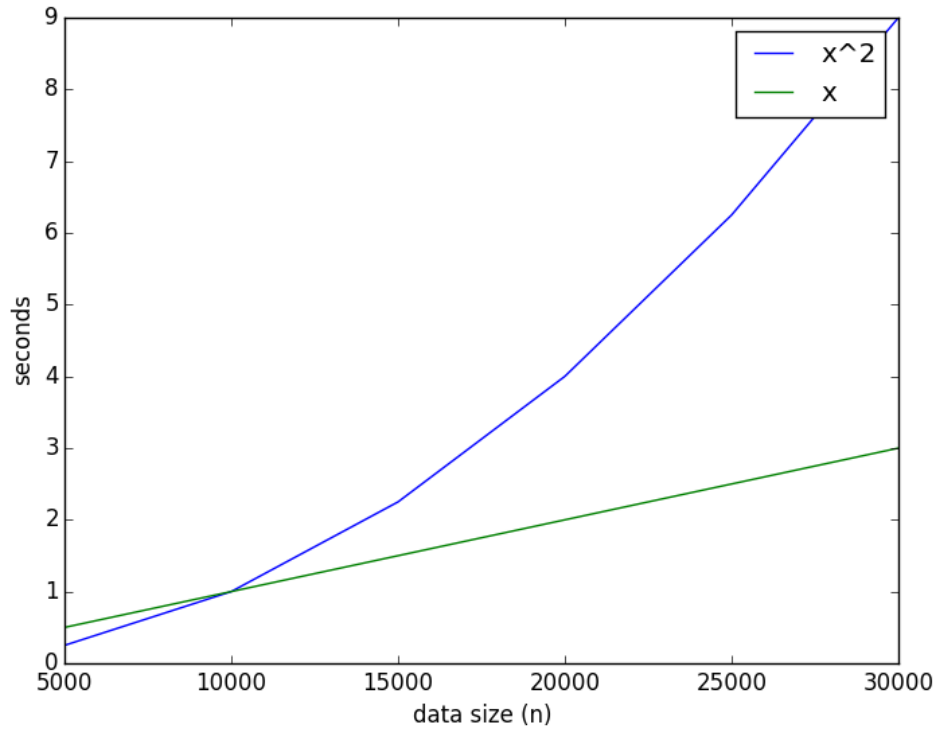
>>> cProfile.run('wc_list(first5000)')
4575 function calls in 0.238 seconds

```

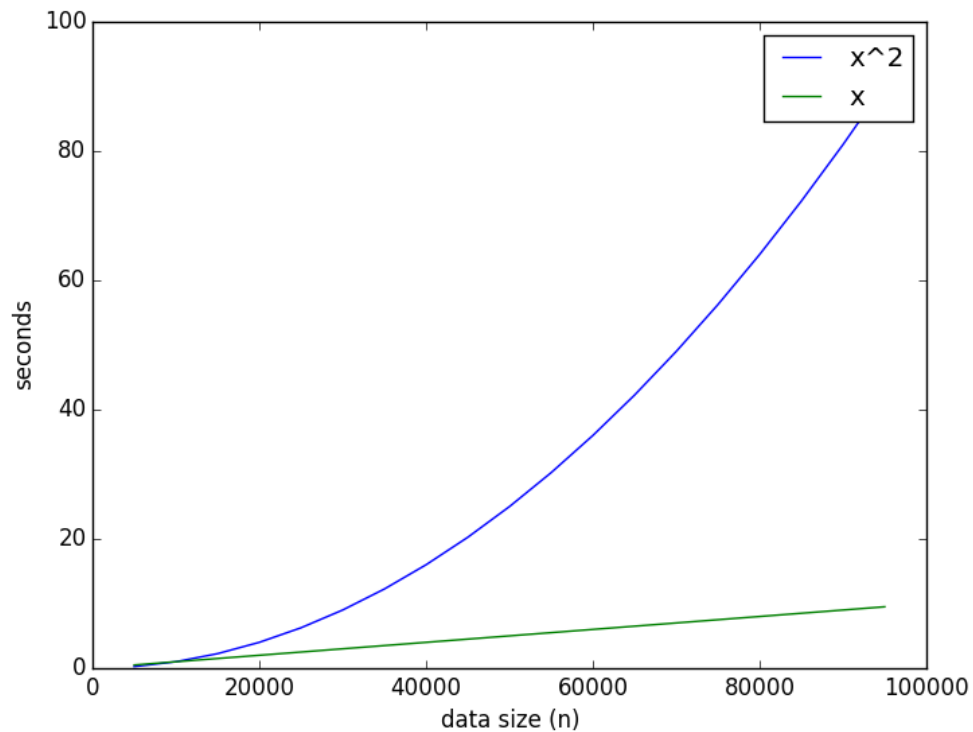
Ordered by: standard name

ncalls	tottime	percall	cumtime	percall	filename:lineno(function)
1	0.000	0.000	0.238	0.238	<string>:1(<module>)
1	0.060	0.060	0.238	0.238	freq.py:12(wc_list)
1	0.001	0.001	0.177	0.177	freq.py:18(<listcomp>)
1	0.000	0.000	0.238	0.238	{built-in method builtins.exec}
2285	0.000	0.000	0.000	0.000	{method 'append' of 'list' objects}
2285	0.176	0.000	0.176	0.000	{method 'count' of 'list' objects}
1	0.000	0.000	0.000	0.000	{method 'disable' of '_lsprof.Profiler' objects}

We know, however, that this searching grows quadratically, as the number of unique tokens increases. *MOBY DICK* contains around 35,000 unique tokens so it will take at least 9 seconds just to identify the unique tokens.



But for a text with 100,000 unique tokens, it will take close to two minutes to find the counts. Thus, lists don't seem viable for even moderately sized texts.



Dictionaries

Ideally, we could scan the list of tokens once and keep a count associated with each token. If we see a token for the first time, we associate a 1 with it. If we see it again, we add 1 to its current count. To only scan the tokens once requires the ability to look up the token count immediately. Dictionaries, more or less, allow us to do just that. Contrast this with our list implementation: determining the count for a single token required scanning the entire list. But dictionaries aren't magic. For them to work requires the ability to take a Python object like a string, integer, or tuple and *hash* it to an integer so that different Python objects rarely collide. We won't worry about these details now, but by choosing a hash function intelligently, one can keep the number of collisions minimal.

Syntax

Dictionaries maintain a set of *key/value pairs*. In Python, the *key* can be anything that is immutable—strings, tuples, numbers, etc. The *value* can be any Python object, including mutable data structures like lists or other hash tables.

Here is the syntax for creating an empty dictionary, adding keys, retrieving the values associated with keys, and checking if a dictionary has a given key.

```
>>> d = {}
>>> d["brent"] = 40
>>> d["courtney"] = 41
>>> d["oscar"] = 6
>>> d["george"] = 3
>>> d
{'oscar': 6, 'courtney': 41, 'brent': 40, 'george': 3}
>>> d["oscar"]
6
>>> "brent" in d
True
>>> "amy" in d
False
```

There are three ways of iterating over the data in a dictionary *d*:

by key `d.keys()` returns a view of the keys of the dictionary that is iterable. In other words `list(d.keys())` will give you a list of the keys in the dictionary.

by value `d.values()` returns a view of the values of the dictionary that is iterable. Similarly, `list(d.values())` returns a list of the values of the dictionary.

by key/value pairs `d.items()` returns a view of the dictionary yielding tuples of the form `(key, value)`. Calling `list(d.items())` returns a list of tuples.

```
>>> list(d.keys())
['oscar', 'courtney', 'brent', 'george']
>>> list(d.values())
[5, 40, 38, 1]
>>> list(d.items())
[('oscar', 5), ('courtney', 40), ('brent', 38), ('george', 1)]
```

With these operations in hand, we can now construct an efficient version of word count.

```
1 def wc_dict(tokens):
2     counts = {}
```

```
3 for token in tokens:
4     if token in counts:
5         counts[token] += 1
6     else:
7         counts[token] = 1
8 return counts.items()
```

Here are some other methods that capture idioms often encountered with a dictionary `d`:

- `setdefault(key, default=None)` If key is in `d`, then return `d[key]`, otherwise insert key with value `default` and return `default`.
- `get(key, default=None)` If key is in `d`, then return `d[key]`, otherwise, return `default`.

Notice that `get` makes our method above more compact.

```
1 def wc_dict(tokens):
2     counts = {}
3     for token in tokens:
4         counts[token] = counts.get(token,0) + 1
5     return counts.items()
```

Practice

Suppose we wanted to create an index of the positions of each token in the original text. Write a function called `token_locations` that, when given a list of tokens, returns a dictionary where each key is a token and each value is list of indices where that token appears.

```
>>> l = "brent sucks big rocks through a big straw".split()
>>> print(token_locations(l))
{'big': [2, 6], 'straw': [7], 'brent': [0], 'a': [5], 'through': [4], 'sucks': [1], 'rocks': [3]}
```

Sets

Sets are like dictionaries without keys—you're interested in a collection of unique objects, but you're not interested in their order. Besides adding, removing, and testing membership, sets support the following set operations:

- union,
- intersection, and
- symmetric difference.

By default these operation return new sets, but there are also version that are side-effecting.

```
>>> s = set(range(10))
>>> len(s)
10
>>> s.add(8)
>>> s
{0, 1, 2, 3, 4, 5, 6, 7, 8, 9}
>>> len(s)
```

```
10
>>> s2 = s.union(set(range(20)))
>>> s2
{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19}
>>> s
{0, 1, 2, 3, 4, 5, 6, 7, 8, 9}
>>> s3 = s2.intersection(set(range(10,20)))
>>> s3
{10, 11, 12, 13, 14, 15, 16, 17, 18, 19}
```