# Access the slides and files here:

https://github.com/j-berg/bioinformatics_bootcamp

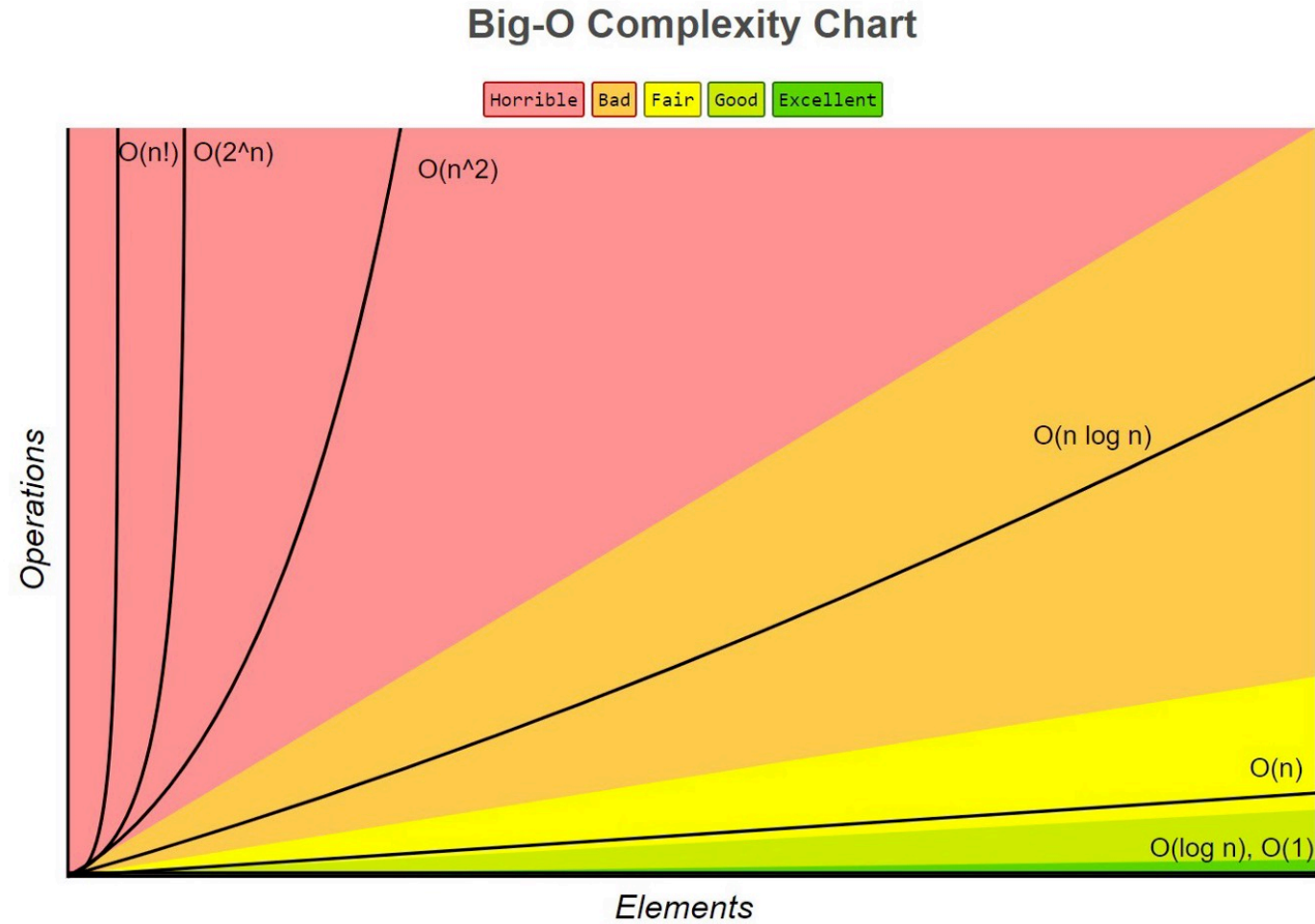# #4.2

Complex tasks

Performance

Debugging

# Questions from writing a script?

# Performance



Big-O Complexity Chart

# Linear Time

- for x in my_list:
    if x == 100000:
        print("I'm here")
        break


    Best case: item is at beginning of the list
    Worst case: item is at the end of the list

# Constant Time

- if x == y:
	print("True")
  else:
	print("False")


- if value in my_set:
	print("True")


- Not dependent on the input

```python
>>> import time

>>> def one(l):
...     start_time = time.time()
...     for x in l:
...         if x == 900000000:
...             print("found it")
...             break
... ...     seconds = time.time() - start_time
    print("Time:", seconds)

>>> def two(s):
...     start_time = time.time()
...     if 900000000 in s:
...         print("found it")
...     seconds = time.time() - start_time
...     print("Time:", seconds, 'seconds')

>>> l = [x for x in range(900000010)]
>>> s = set(l)

>>> one(l)
found it
Time: 10.97479 seconds

>>> two(s)
found it
Time: 0.00194 seconds
```
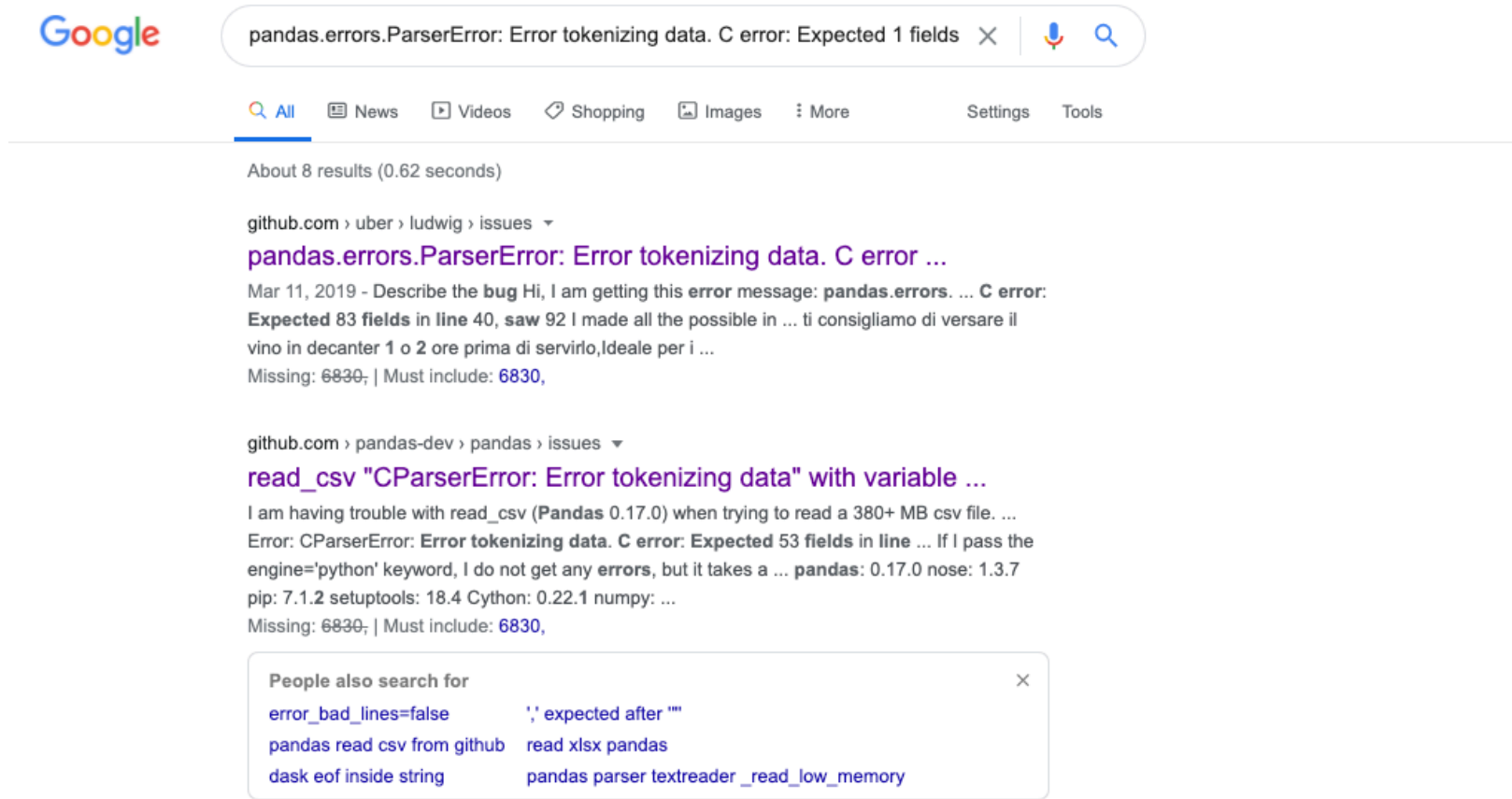
# Troubleshooting

```
[>>> d = pd.read_csv("~/Desktop/SCE_data_table.tsv", sep="t")
```

# Troubleshooting

```
[>>> d = pd.read_csv("~/Desktop/SCE_data_table.tsv", sep="t")
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
  File "/Users/jordan/miniconda3/lib/python3.8/site-packages/pandas/io/parsers.py", line 676, in parser_f
    return _read(filepath_or_buffer, kwds)
  File "/Users/jordan/miniconda3/lib/python3.8/site-packages/pandas/io/parsers.py", line 454, in _read
    data = parser.read(nrows)
  File "/Users/jordan/miniconda3/lib/python3.8/site-packages/pandas/io/parsers.py", line 1133, in read
    ret = self._engine.read(nrows)
  File "/Users/jordan/miniconda3/lib/python3.8/site-packages/pandas/io/parsers.py", line 2037, in read
    data = self._reader.read(nrows)
  File "pandas/_libs/parsers.pyx", line 860, in pandas._libs.parsers.TextReader.read
  File "pandas/_libs/parsers.pyx", line 875, in pandas._libs.parsers.TextReader._read_low_memory
  File "pandas/_libs/parsers.pyx", line 929, in pandas._libs.parsers.TextReader._read_rows
  File "pandas/_libs/parsers.pyx", line 916, in pandas._libs.parsers.TextReader._tokenize_rows
  File "pandas/_libs/parsers.pyx", line 2071, in pandas._libs.parsers.raise_parser_error
pandas.errors.ParserError: Error tokenizing data. C error: Expected 1 fields in line 6830, saw 2
```

# Copy the most informational error line

# Not an exact answer, but we can work with it

File "/Library/Frameworks/Python.framework/Versions/3.6/lib/python3.6/site-packages/ludwig/data/preprocessing.py", line 54, in build_dataset
dataset_df = read_csv(dataset_csv)
File "/Library/Frameworks/Python.framework/Versions/3.6/lib/python3.6/site-packages/ludwig/utils/data_utils.py", line 48, in read_csv
logging.WARNING('Failed to parse the CSV with pandas default way,'
TypeError: 'int' object is not callable

**w4nderlust** commented on Mar 11, 2019                                    Collaborator  ☺  ...

@IzzyHibbert thanks for posting this. In the docs we suggest to escape the commas within the text with \\, , so first thing I would try to do that.
Let me know if this solves your problem.

Personally I prefer to be a bit more strict on the data side rather than letting things pass or being filtered out, because those could become problems down the line.

🏷  😺 **w4nderlust** added the `waiting for answer` label on Mar 11, 2019

# Troubleshooting

```
[>>> d = pd.read_csv("~/Desktop/SCE_data_table.tsv", sep="\t")
```

# Troubleshooting



```
[>>> d = pd.read_csv("~/Desktop/SCE_data_table.tsv", sep="\t")
[>>> d.head()
   Unnamed: 0  ...  14251X9_170420_D00294_0314_BCB1VVANXX_6_Aligned
0      ETS1-1  ...                                                0
1      ETS1-2  ...                                                0
2      ETS2-1  ...                                                0
3      ETS2-2  ...                                                0
4        HRA1  ...                                                2

[5 rows x 25 columns]
[>>>
[>>> d = pd.read_csv("~/Desktop/SCE_data_table.tsv", sep="\t", index_col=0)
[>>> d.head()
        14251X10_170420_D00294_0314_BCB1VVANXX_6_Aligned  ...  14251X9_170420_D00294_0314_BCB1VVANXX_6_Aligned
ETS1-1                                                 0  ...                                                0
ETS1-2                                                 0  ...                                                0
ETS2-1                                                 0  ...                                                0
ETS2-2                                                 0  ...                                                0
HRA1                                                   0  ...                                                2

[5 rows x 24 columns]
>>>
```

# Learn to describe the problem to Google

# Again, not an exact answer, but we can work with it

# Again, not an exact answer, but we can work with it



You could also optionally tell `read_csv` that the first column is the index column by passing `index_col=0`:

```
In [40]:
pd.read_csv(io.StringIO(df.to_csv()), index_col=0)

Out[40]:
          a          b          c
0  0.109066  -1.112704  -0.545209
1  0.447114   1.525341   0.317252
2  0.507495   0.137863   0.886283
3  1.452867   1.888363   1.168101
4  0.901371  -0.704805   0.088335
```

share   edit   follow                     edited Apr 9 '16 at 16:16                     answered Apr 9 '16 at 15:50

# Homework

- Review concepts from Python classes
- Finish gene dictionary project