# 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**

Horrible | Bad | Fair | Good | Excellent

O(n!) | O(2^n) | O(n^2)

O(n log n)

O(n)

O(log n), O(1)

Operations

Elements

# 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

```
1    >>> import time
2
3    >>> def one(l):
4    ...      start_time = time.time()
5    ...      for x in l:
6    ...          if x == 900000000:
7    ...              print("found it")
8    ...              break
9    ... ...      seconds = time.time() - start_time
10       print("Time:", seconds)
11
12   >>> def two(s):
13   ...      start_time = time.time()
14   ...      if 900000000 in s:
15   ...          print("found it")
16   ...      seconds = time.time() - start_time
17   ...      print("Time:", seconds, 'seconds')
18
19   >>> l = [x for x in range(900000010)]
20   >>> s = set(l)
21
22   >>> one(l)
23   found it
24   Time: 10.97479 seconds
25
26   >>> two(s)
27   found it
28   Time: 0.00194 seconds
29
30
31
32
```

# Troubleshooting

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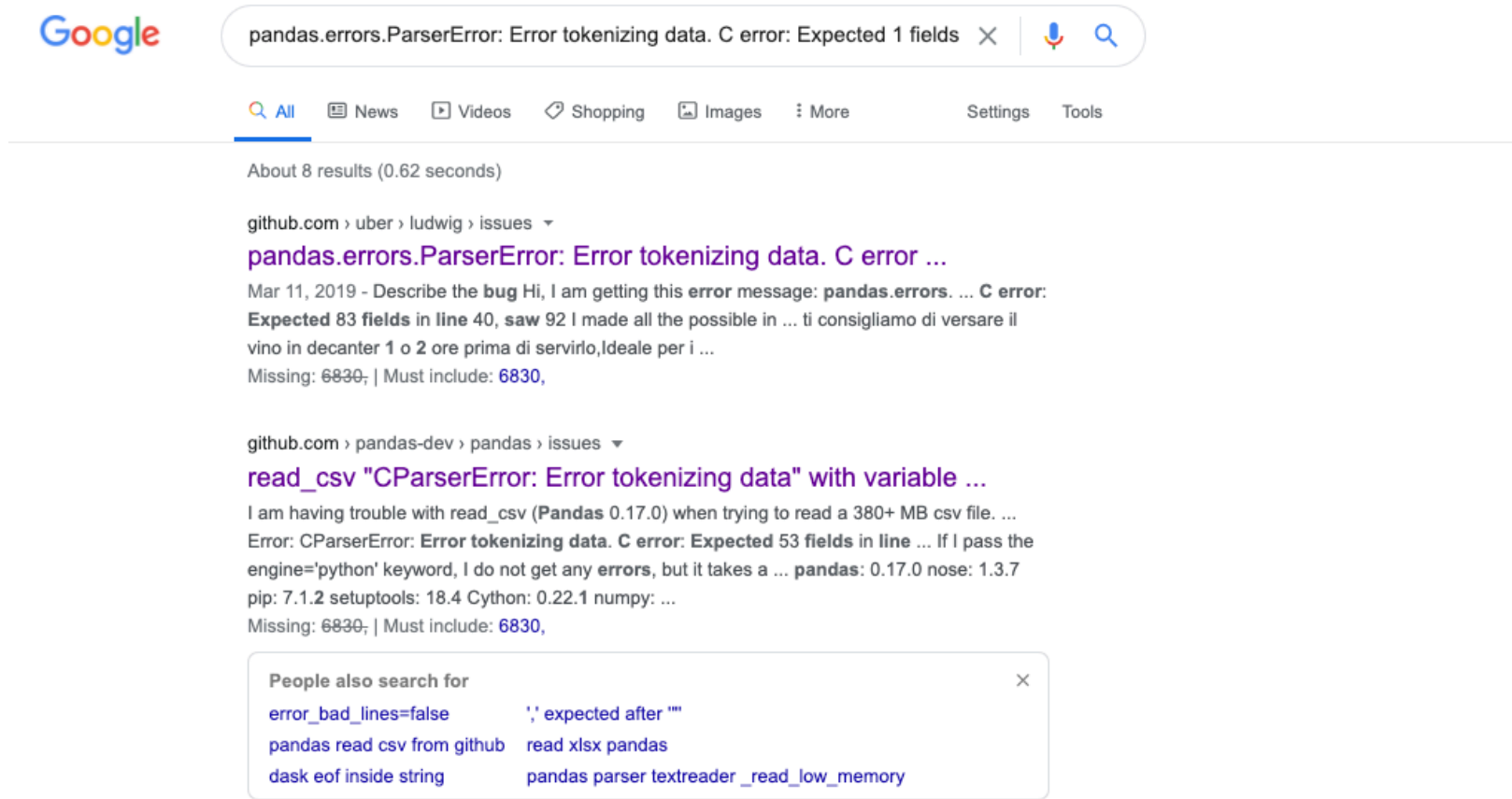
# Troubleshooting



```
[>>> 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]
>>> ▮
```

# 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