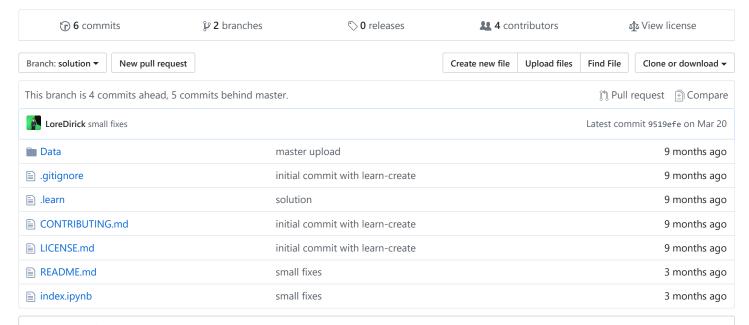
learn-co-curriculum / dsc-importing-data-using-pandas-lab

No description, website, or topics provided.



■ README.md

Importing Data Using Pandas - Lab

Introduction

In this lab, you'll get some practice with loading files with summary or metadata, and if you find that easy, the optional "level up" content covers loading data from a corrupted csv file!

Objectives

You will be able to:

- Import data from csv files and Excel files
- Understand and explain key arguments for imports
- Save information to csv and Excel files
- Access data within a Pandas DataFrame (print() and .head())

Loading Files with Summary or Meta Data

Load either of the files Zipcode_Demos.csv or Zipcode_Demos.xlsx. What's going on with this dataset? Clean it up into a useable format and describe the nuances with how the data is currently formatted.

All data files are stored in a folder titled 'Data'.

```
import pandas as pd

df = pd.read_csv('Data/Zipcode_Demos.csv')
df.head()
```

```
<style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }
   .dataframe tbody tr th {
       vertical-align: top;
   }
   .dataframe thead th {
       text-align: right;
   }
```

	0	Average Statistics	Unnamed: 2	Unnamed: 3	Unnamed: 4	Unnamed: 5	Unnamed: 6	Unnamed: 7	Unnam 8
0	1	NaN	0	NaN	NaN	NaN	NaN	NaN	NaN
1	2	JURISDICTION NAME	10005.8	NaN	NaN	NaN	NaN	NaN	NaN
2	3	COUNT PARTICIPANTS	9.4	NaN	NaN	NaN	NaN	NaN	NaN
3	4	COUNT FEMALE	4.8	NaN	NaN	NaN	NaN	NaN	NaN
4	5	PERCENT FEMALE	0.404	NaN	NaN	NaN	NaN	NaN	NaN

5 rows × 47 columns

```
df.tail()
```

<style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }

```
.dataframe tbody tr th {
    vertical-align: top;
}
.dataframe thead th {
    text-align: right;
}
```

</style>

	0	Average Statistics	Unnamed: 2	Unnamed:	Unnamed:	Unnamed: 5	Unnamed:	Unnamed:	Unnamed
52	53	10006	6	2	0.33	4	0.67	0	0
53	54	10007	1	0	0	1	1	0	0
54	55	10009	2	0	0	2	1	0	0
55	56	10010	0	0	0	0	0	0	0
56	57	10011	3	2	0.67	1	0.33	0	0

5 rows × 47 columns

Commentary:

Dataframe is really two table views, one on top of the other. The first is a summary view of the raw data below. There is also a blank row at row 1 in the file.

```
prev_count = 10**3
 for row in df.index:
     count = 0
     for entry in df.iloc[row].isnull():
         if entry:
             count += 1
     if count != prev_count and row!=0:
         print('On row {} there are {} null values. The previous row had {} null values.'.format(row, count, prev_cou
     prev_count = count
 On row 1 there are 44 null values. The previous row had 45 null values.
 On row 46 there are 0 null values. The previous row had 44 null values.
 df1 = pd.read_csv('Data/Zipcode_Demos.csv', skiprows=[1], nrows=45, usecols=[0,1,2])
 df1.head()
<style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }
  .dataframe tbody \operatorname{tr} th \{
     vertical-align: top;
  .dataframe thead th \{
     text-align: right;
```

1/ Jty			
	0	Average Statistics	Unnamed: 2
0	2	JURISDICTION NAME	10005.800
1	3	COUNT PARTICIPANTS	9.400
2	4	COUNT FEMALE	4.800
3	5	PERCENT FEMALE	0.404
4	6	COUNT MALE	4.600

```
df1.tail()

<style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }

.dataframe tbody tr th {
    vertical-align: top;
  }

.dataframe thead th {
    text-align: right;
  }

</style>
```

	0	Average Statistics	Unnamed: 2
40	42	COUNT NRECEIVES PUBLIC ASSISTANCE	7.100
41	43	PERCENT NRECEIVES PUBLIC ASSISTANCE	0.649
42	44	COUNT PUBLIC ASSISTANCE UNKNOWN	0.000
43	45	PERCENT PUBLIC ASSISTANCE UNKNOWN	0.000
44	46	COUNT PUBLIC ASSISTANCE TOTAL	9.400

```
df2 = pd.read_csv('Data/Zipcode_Demos.csv', skiprows=47)
df2.head()

<style scoped > .dataframe tbody tr th:only-of-type { vertical-align: middle; }
    .dataframe tbody tr th {
        vertical-align: top;
    }
    .dataframe thead th {
        text-align: right;
    }
```

	47	JURISDICTION NAME	COUNT PARTICIPANTS	COUNT	PERCENT FEMALE	COUNT MALE	PERCENT MALE	COUNT GENDER UNKNOWN	PERCE GEND UNKNO
0	48	10001	44	22	0.50	22	0.50	0	0
1	49	10002	35	19	0.54	16	0.46	0	0
2	50	10003	1	1	1.00	0	0.00	0	0
3	51	10004	0	0	0.00	0	0.00	0	0
4	52	10005	2	2	1.00	0	0.00	0	0

5 rows × 47 columns

Level Up (Optional) - Loading Corrupt CSV files

Occasionally, you encountered some really ill formatted data. One example of this can be data that has strings containing commas in a csv file. Under the standard protocol, when this occurs, one is supposed to use quotes to differentiate between the commas denoting fields and commas within those fields themselves. For example, we could have a table like this:

ReviewerID,Rating,N_reviews,Review,VenueID 123456,4,137,This restaurant was pretty good, we had a great time.,98765

Which should be saved like this if it were a csv (to avoid confusion with the commas in the Review text): "ReviewerID", "Rating", "N_reviews", "Review", "VenueID" "123456", "4", "137", "This restaurant was pretty good, we had a great time.", "98765"

Attempt to import the corrupt file, or at least a small preview of it. It is appropriately titled Yelp_Reviews_corrupt.csv. Investigate some of the intricacies of skipping rows to then pass over this error and comment on what you think is going on.

```
#Hint: here's a useful programming pattern to use.
try:
```

```
#do something
 except Exception as e:
     #handle your exception e
   File "<ipython-input-8-13f6e15364f1>", line 4
     except Exception as e:
 IndentationError: expected an indented block
 #Your code here
     df = pd.read_csv('Data/Yelp_Reviews_Corrupt.csv')
 except Exception as e:
     print(e)
 Error tokenizing data. C error: Expected 10 fields in line 2331, saw 11
 # # Iteration 1
 for i in range(1500,2000):
         df = pd.read_csv('Data/Yelp_Reviews_Corrupt.csv', nrows=i)
     except:
          print('First failure at: {}'.format(i))
 df1 = pd.read_csv('Data/Yelp_Reviews_Corrupt.csv', nrows=i-1)
 print(len(df))
 df1.head()
 First failure at: 1962
 1961
<style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }
  .dataframe tbody tr th {
     vertical-align: top;
 }
 .dataframe thead th {
     text-align: right;
 }
```

	Unnamed:	business_id	cool	date	funny	review_id	stars	
0	1	pomGBqfbxcqPv14c3XH-ZQ	0	2012- 11-13	0.0	dDl8zu1vWPdKGihJrwQbpw	5.0	l p N fi A

	Unnamed:	business_id	cool	date	funny	review_id	stars
1	2	jtQARsP6P-LbkyjbO1qNGg	1	2014- 10-23	1.0	LZp4UX5zK3e-c5ZGSeo3kA	1.0
2	4	Ums3gaP2qM3W1XcA5r6SsQ	0	2014- 09-05	0.0	jsDu6QEJHbwP2Blom1PLCA	5.0
3	5	vgfcTvK81oD4r50NMjU2Ag	0	2011- 02-25	0.0	pfavA0hr3nyqO61oupj-IA	1.0
4	10	yFumR3CWzpfvTH2FCthvVw	0	2016- 06-15	0.0	STiFMww2z31siPY7BWNC2g	5.0

```
df1.tail()

<style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }
   .dataframe tbody tr th {
      vertical-align: top;
   }
   .dataframe thead th {
      text-align: right;
```

\/3tyle/	ystyles										
	Unnamed:	business_id	cool	date	funny	review_i					
1956	4993	u8C8pRvaHXg3PgDrsUHJHQ	0	2016- 08-08	0.0	gXmHGBSBBz2 uHdvGf4lZQ					

	Unnamed: 0	business_id	cool	date	funny	review_i
1957	4998	-9nai28tnoylwViuJVrYEQ	0	2015- 03-22	0.0	u- zqCN_IXfypJIUz
1958	I had an awesome great time with friends.	NaN	NaN	NaN	NaN	NaN
1959	I loved the tapas and the excellent paella.	NaN	NaN	NaN	NaN	NaN
1960	I can't wait to come back soon.	0	otDVyX37h61WEbqPLEjCmQ	NaN	NaN	NaN

Comments:

Be careful, even prior to the error, the last few entries look faulty here; these could very well be the spillovers of unencapsulated commas!

```
# # Iteration 2
for i in range(0,500):
    try:
        temp = pd.read_csv('Data/Yelp_Reviews_Corrupt.csv', skiprows=1962, nrows=i, names=df1.columns)
    except:
        print('First failure at: {}'.format(i))
        break

df2 = pd.read_csv('Data/Yelp_Reviews_Corrupt.csv', skiprows=1962, nrows=i-1, names=df1.columns)
print(len(df2))
df2.head()

498

<style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }
    .dataframe tbody tr th {
        vertical-align: top;
    }
    .dataframe thead th {
        text-align: right;
    }
```

</th <th>'styl</th> <th>le</th> <th>></th>	'styl	le	>
---	-------	----	---

Unnamed: 0	business_id	cool	date	funny	revie
---------------	-------------	------	------	-------	-------

		icam ce camedam/ace importing data doing particle lab at collater.										
Unnamed: 0	business_id	cool	date	funny	revie							
STAY AWAY FROM THIS PLACE!!!!!!	5	sDofYImMQQmu4Le5G9zmpQ	NaN	NaN	NaN							
3948	GAKFx4jFUtTOTpp_jDJnuA	0	2017- 09-01	0	OUZWMw7EgC							
3949	0QzCeORfF8EY34UODWRV9A	0	2017- 09-03	0	7lbykaWFD8YB							
3950	tlt8zNrZ6_A3DmXiM-cnBA	0	2016- 06-12	0	Nd_soHwCYi8a							
3952	XD0LjNuPPwJPsTAHecUh7A	0	2015- 08-23	0	FUUTAr5CECrkf							
	O STAY AWAY FROM THIS PLACE!!!!!! 3948	STAY AWAY FROM THIS PLACE!!!!!! 3948 GAKFx4jFUtTOTpp_jDJnuA 3949 OQzCeORfF8EY34UODWRV9A tlt8zNrZ6_A3DmXiM-cnBA	STAY AWAY FROM THIS PLACE!!!!!! 3948 GAKFx4jFUtTOTpp_jDJnuA 0 0 1950 tlt8zNrZ6_A3DmXiM-cnBA 0	o business_id cool date STAY AWAY FROM THIS PLACE!!!!!! 5 sDofYImMQQmu4Le5G9zmpQ NaN 3948 GAKFx4jFUtTOTpp_jDJnuA 0 2017- 09-01 3949 0QzCeORfF8EY34UODWRV9A 0 2017- 09-03 3950 tlt8zNrZ6_A3DmXiM-cnBA 0 2016- 06-12 3952 XDQLiNuPPw/IPsTAHect IIb7A 0 2015-	0 business_id cool date funny STAY AWAY FROM THIS PLACE!!!!!! 5 sDofYlmMQQmu4Le5G9zmpQ NaN NaN 3948 GAKFx4jFUtTOTpp_jDJnuA 0 2017- 09-01 0 3949 0QzCeORfF8EY34UODWRV9A 0 2017- 09-03 0 3950 tlt8zNrZ6_A3DmXiM-cnBA 0 2016- 06-12 0 3952 XD0LiNupPpw/PsTAHecl lh7A 0 2015- 0 0							

```
temp = pd.read_csv('Data/Yelp_Reviews_Corrupt.csv')
print(len(temp))
temp.head()
ParserError
                                          Traceback (most recent call last)
<ipython-input-13-6c179a8f3c47> in <module>()
----> 1 temp = pd.read_csv('Data/Yelp_Reviews_Corrupt.csv')
      2 print(len(temp))
      3 temp.head()
/usr/local/lib/python3.6/site-packages/pandas/io/parsers.py in parser_f(filepath_or_buffer, sep, delimiter,
header, names, index_col, usecols, squeeze, prefix, mangle_dupe_cols, dtype, engine, converters, true_values,
false_values, skipinitialspace, skiprows, nrows, na_values, keep_default_na, na_filter, verbose,
skip_blank_lines, parse_dates, infer_datetime_format, keep_date_col, date_parser, dayfirst, iterator, chunksize,
compression, thousands, decimal, lineterminator, quotechar, quoting, escapechar, comment, encoding, dialect,
tupleize_cols, error_bad_lines, warn_bad_lines, skipfooter, skip_footer, doublequote, delim_whitespace,
as_recarray, compact_ints, use_unsigned, low_memory, buffer_lines, memory_map, float_precision)
    707
                            skip_blank_lines=skip_blank_lines)
    708
--> 709
                return _read(filepath_or_buffer, kwds)
    710
    711
            parser_f.__name__ = name
```

```
/usr/local/lib/python3.6/site-packages/pandas/io/parsers.py in _read(filepath_or_buffer, kwds)
     454
  --> 455
                 data = parser.read(nrows)
     456
            finally:
     457
                 parser.close()
 /usr/local/lib/python3.6/site-packages/pandas/io/parsers.py in read(self, nrows)
                          raise ValueError('skipfooter not supported for iteration')
    1068
  -> 1069
                ret = self._engine.read(nrows)
    1070
    1071
                 if self.options.get('as_recarray'):
 /usr/local/lib/python3.6/site-packages/pandas/io/parsers.py in read(self, nrows)
            def read(self, nrows=None):
    1838
                try:
  -> 1839
                     data = self._reader.read(nrows)
    1840
                 except StopIteration:
    1841
                     if self._first_chunk:
 pandas/_libs/parsers.pyx in pandas._libs.parsers.TextReader.read()
 pandas/_libs/parsers.pyx in pandas._libs.parsers.TextReader._read_low_memory()
 pandas/_libs/parsers.pyx in pandas._libs.parsers.TextReader._read_rows()
 pandas/_libs/parsers.pyx in pandas._libs.parsers.TextReader._tokenize_rows()
 pandas/ libs/parsers.pyx in pandas. libs.parsers.raise parser error()
 ParserError: Error tokenizing data. C error: Expected 10 fields in line 2331, saw 11
 temp = pd.read_csv('Data/Yelp_Reviews_Corrupt.csv', names=df1.columns, skiprows=1)
 print(len(temp))
 temp.head()
 4651
<style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }
  .dataframe tbody tr th {
     vertical-align: top;
 }
  .dataframe thead th {
     text-align: right;
 }
</style>
```

Unnamed:	business_id	cool	date	funny	review_id	stars

	Unnamed:	business_id	cool	date	funny	review_id	stars	
0	1	pomGBqfbxcqPv14c3XH-ZQ	0	2012- 11-13	0	dDl8zu1vWPdKGihJrwQbpw	5	t p M f
1	2	jtQARsP6P-LbkyjbO1qNGg	1	2014- 10-23	1	LZp4UX5zK3e-c5ZGSeo3kA	1	T E k F V
2	4	Ums3gaP2qM3W1XcA5r6SsQ	0	2014- 09-05	0	jsDu6QEJHbwP2Blom1PLCA	5	h fo T s a
3	5	vgfcTvK81oD4r50NMjU2Ag	0	2011- 02-25	0	pfavA0hr3nyqO61oupj-IA	1	T p s T c c s h
4	10	yFumR3CWzpfvTH2FCthvVw	0	2016- 06-15	0	STiFMww2z31siPY7BWNC2g	5	l b C n fo

```
pd.read_csv('Data/Yelp_Reviews_Corrupt.csv', skiprows=len(df1)+len(df2), names=df1.columns)

<style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }

.dataframe tbody tr th {
    vertical-align: top;
}

.dataframe thead th {
    text-align: right;
}
```

Uni	named: 0 busine	ss_id cool	date	
-----	-----------------	------------	------	--

	Unnamed: 0	business_id	cool	date
0	Cons:	NaN	NaN	NaN
1	- Dusty! Not sure if it's all of Vegas but I	NaN	NaN	NaN
2	- Valet parking: kinda inconvenient when you	NaN	NaN	NaN
3	- Sofabed is extremely flimsy	if you have more than 2 people	insist on 2 queen beds. the sofa cushions ar	NaN
4	Other points:	NaN	NaN	NaN
5	* Should call ahead of time to make sure your	NaN	NaN	NaN
6	* Hotel lobby is extremely small!	NaN	NaN	NaN
7	* In-room food service was overpriced (and fo	NaN	NaN	NaN
8	* Don't go to the 7-11	it's shady. You can shop at the am/pm or the	NaN	NaN
9	Overall	it was a good experience for the price we pai	3	DZYGeWwBRKHgLUSk12sCvA
10	4058	WPCgtEG-bJt0cZtnM-x7yw	0	2012-02-28
11	1624	T5R6alLLDBnHQvfejY7dgA	1	2012-07-09
12	4938	53BSdnhzcCBfBH_6TgX63Q	0	2014-08-31
13	2897	hOB3NHuF-iVFdEkrA-PUlg	1	2012-02-17

	Unnamed: 0	business_id	cool	date
14	The room was lovely	we had a loft basic package and really enjoye	NaN	NaN
15	We took part in the spa package offered throug	NaN	NaN	NaN
16	I emailed the concierge ahead of time to reque	NaN	NaN	NaN
17	I've given 4 stars instead of 5 for a couple o	the breakfast is not included in the room rat	and the gym is adequate but quite small.	3
18	2422	vwQvDIb_F7AqwCPaQhHrwg	0	2012-10-29
19	1527	04u-szAykldu-caSDHQaKA	0	2012-02-09
20	4080	EDcZRvERC22Cvw1yi4-VKg	1	2017-12-05
21	4237	ZM-ljL_Y6bR4qEYsGHws5A	0	2016-11-05
22	4498	VRTfAP2DjvUYxRY3dw37hA	0	2014-01-08
23	Chi	my pedicurist	was wonderful. Super sweet and very attentive	it was heavenly. When it was all over
24	Four stars instead of five because:	NaN	NaN	NaN

	Unnamed: 0	business_id	cool	date
25	- The ladies in reception were a bit rude	both over the phone and more so in person.	NaN	NaN
26	- My pedicure only lasted two days! A pedicure	even 3 weeks without a single chip.	7	ETmpBain2s02PqHGwSr7hQ
27	2716	4VHp2gei1bpY68ZzEZE9Bg	2	2013-05-22
28	3426	8g3u6g7J93nIOF8owARxew	0	2014-08-13
29	Had the traditional chicken shawarma. one of m	NaN	NaN	NaN
•••				
2552	Starting off with drinks	Bamburger serves beer	wine	and old-fashioned Stewart's soda
2553	The menu is varied	offering up soup	salads	sandwiches and desserts
2554	We went with the Bambamburger (\$11.50)	which is 2/3 of a pound of prime ground chuck	and the chicken burger (\$9.95) for myself	on whole-wheat buns. Both of us outfitted our
2555	Bamburger serves up great burgers	fries and shakes at a fairly good price	although if you go a little overboard with th	you might quickly end up with a \$20 burger
2556	If you are hunting for a real deal on a burger	Bamburger might not be what you are looking for	but if you are more on the adventurous side	and want to have fun creating your own burger

	Unnamed: 0	business_id	cool	date
2557	3225	iyyWYpWm8X-6i7kBR3JHuw	0	2014-01-27
2558	4674	4KfDcE9iU2isFpoaKeDpgw	0	2012-06-14
2559	4719	P4Plzlfm4uJjNmH3wY4W1Q	0	2014-01-14
2560	3440	nW45ez1L6U4PsYhV1BTrGQ	0	2012-05-19
2561	I had given my husband some ideas of engagemen	a friend of ours recommended H&F.	NaN	NaN
2562	The customer service here is great - they are	attentive	honest and the prices are very reasonable! My	NaN
2563	We went back for our wedding bands and I worke	NaN	NaN	NaN
2564	We've recommended 5 other friends to H&F - the	NaN	NaN	NaN
2565	Go to H&F for all your jewellery needs - you w	1	PkRFSQgSfca9Tamq7b2LdQ	NaN
2566	4812	e13SEvJud_vgeDR_doL4sQ	0	2013-03-01

	Unnamed: 0	business_id	cool	date
2567	1970	9NBkIExYYz3w9O5JdzDOMA	0	2013-11-03
2568	689	BTcY04QFiS1uh-RpkR7rAg	1	2013-06-02
2569	4874	t0T_4MM4EUHbCzBTF11FHA	0	2016-08-14
2570	564	5XYR6doRa5Nj1JMfSDei6A	1	2016-06-14
2571	Highly recommend the custard cakes they are th	NaN	NaN	NaN
2572	The rice flour cake is also really good and a	NaN	NaN	NaN
2573	The bean cakes are great here too! orange	almond	and a few others I have tried are all good.	NaN
2574	Can't go wrong with Nova Era	0	kBNFdviedCPFWyR- wVaAzw	NaN
2575	3458	aLcFhMe6DDJ430zelCpd2A	0	2013-10-02
2576	This was disappointing.	NaN	NaN	NaN
2577	First off	it was really awkward sitting on the benches	as people walked past us while to wait for ou	NaN
2578	Second	when we were seated	it was so loud. It felt like we were in a hig	NaN

	Unnamed: 0	business_id	cool	date
2579	Finally - Food was mediocre. I was extremely d	but it wasn't flavourful.	NaN	NaN
2580	Wasn't worth the hype	unfortunately.	1	PkRFSQgSfca9Tamq7b2LdQ
2581	4206	WdBWhGe4Siqg3IYTc4_K4A	0	2016-08-15

2582 rows × 10 columns

Summary

Congratulations, you now practiced your pandas-importing skills!