

tu15_WranglingIII_HW

March 9, 2023

1 Wrangling III homework

1.1 Deleting invalid rows from the cancer data

It's been decided by committee that duplicate data in the (smaller version of) the cancer data set is going to go as follows. If a row is **identical to the one immediately above it**, we'll **consider it an accidental entry** due to fatigue or whatever. But a row that is *not* identical to the one above it will be considered valid, even it has a duplicate somewhere else in the data set; we'll assume such duplicates represent separate visits.

So our rule is: **if a row is identical to the one above it, we drop it.**

A small pre-cleaned version of the data with only 4 columns is in **small_cancer_data.csv**, so you can read it in directly without having to clean it up.

Spend a minute or two thinking about how you would approach this problem.

If you are ready to go on your own, then go!

Once you have working code – once you can take **small_cancer_data.csv** and trim out the unwanted rows – then wrap your code into a function, so all you have to do drop unwanted rows is call your function!

Spend some time thinking about and working on the problem. If you get to an impass and you'd like some hints, read on.

1.2 Preliminaries

As usual, we'll load some libraries we'll be likely to use.

```
[1]: import pandas as pd
```

1.3 Make a mini data set for testing

Rather than taking a crack at the whole set, make a small data frame named **tiny** with 10 rows and two columns. Put successive repeated rows in two places (like rows 2 and 3 could repeat, as could rows 6 and 7). Put an additional repeated row on its own.

Something like this:

```
[2]: tiny = pd.DataFrame(dict(a = [1, 2, 3, 3, 4, 5, 5, 5, 6, 3],
                               b = ['a', 'b', 'c', 'c', 'd', 'e', 'e', 'f', 'g',
                                     ↪ 'c'])))
```

Check our tiny data frame.

```
[3]: # check your test data frame
tiny
```

```
[3]:   a  b
0   1  a
1   2  b
2   3  c
3   3  c
4   4  d
5   5  e
6   5  e
7   5  f
8   6  g
9   3  c
```

There should be two rows that need to be dropped, and one (the last) that should be kept even though it's a duplicate.

Just to be sure, check the output of `.duplicated(keep=False)` – it should show back-to-back `True` values in 2 places, and one solo `True` at the end.

```
[5]: # check .duplicated output
tiny.duplicated(keep = False)
```

```
[5]: 0    False
1    False
2     True
3     True
4    False
5     True
6     True
7    False
8    False
9     True
dtype: bool
```

1.4 Make a plan

There are probably 100 ways to solve this problem. Many are probably very clever and involve using fancy pandas functions.

A straightforward plan using things we already know about might be something like

- go through the rows of the data frame with a **for** loop, starting with the second row
 - at each row, compare the current row with the previous one
 - if they're the same, save the index of the current row
 - after the for loop, delete the unwanted rows using the saved indexes
-

1.5 Test the parts of the plan

Now that we've got a plan, let's get the pieces of the plan to work before putting the whole plan together.

1.5.1 Make sure we can get rows

We should be able to get rows of a data frame in a couple of ways. These are

- using `.loc[]` with the value of rows index (it's name)
- using `.iloc[]` and indexing into the data like it were a numpy array

Let's try the `.loc[]` way.

```
[12]: tiny.loc[0]
```

```
[12]: a    1
      b    a
      Name: 0, dtype: object
```

And let's try the `.iloc[]` method.

```
[11]: tiny.iloc[1]
```

```
[11]: a    2
      b    b
      Name: 1, dtype: object
```

Looks like either will work!

1.5.2 Figure out how to compare rows

We are going to need to compare rows. Let's see how that is going to work.

compare the first and second rows - these *should not* match

```
[15]: # which things in the rows match?
      tiny.loc[0:1]
```

```
[15]:      a  b
      0  1  a
      1  2  b
```

compare the third and fourth rows - these *should* match

```
[17]: # which things in the rows match?
      tiny.loc[2:3]
```

```
[17]:      a  b
      2  3  c
      3  3  c
```

The rows only match if **all** the columns match, so we can see if this is the case with the `all()` function.

```
[20]: # do all the column match?
      tiny.all()
```

```
[20]: a      True
      b      True
      dtype: bool
```

Now we have a way to compare rows and get a single `True` if the rows are identical, and a `False` if they're not.

And now that we know how to do the row comparison, let's get a `for` loop working.

1.5.3 Confirm we can get rows with a `for` loop

Loop through the first few rows Let's make sure we can index into rows with a `for` loop. Let's try to get the first few using `.loc[]` and print them. Like

```
for ... :
    print(...)
```

```
[39]: # loop through the first few rows
      for i in range(2):
          print(tiny.loc[i])
```

```
a      1
b      a
Name: 0, dtype: object
a      2
b      b
Name: 1, dtype: object
```

Loop through the all rows To loop through all the rows, we first need to get the number of rows. We can do this using the `shape` attribute.

```
[30]: # get the number of rows using shape
tiny.shape
```

```
[30]: (10, 2)
```

```
[37]: # loop through all the rows
for i in range(len(tiny)):
    print(tiny.loc[i])
```

```
a    1
b    a
Name: 0, dtype: object
a    2
b    b
Name: 1, dtype: object
a    3
b    c
Name: 2, dtype: object
a    3
b    c
Name: 3, dtype: object
a    4
b    d
Name: 4, dtype: object
a    5
b    e
Name: 5, dtype: object
a    5
b    e
Name: 6, dtype: object
a    5
b    f
Name: 7, dtype: object
a    6
b    g
Name: 8, dtype: object
a    3
b    c
Name: 9, dtype: object
```

1.6 Putting it all together

Get the number of rows

```
[31]: # get the number of rows using shape
tiny.shape
```

```
[31]: (10, 2)
```

```
[41]: len(tiny)
```

```
[41]: 10
```

Make an empty list to hold the indexes of the columns we're going to drop

```
[40]: empty_list = []
```

Make a for loop that

- goes from 1 (i.e. the second row) to the end
- tests the current row against previous
- stores index for dropping

```
[43]: tiny['duplicated'] = tiny.duplicated(keep = False)
tiny
```

```
[43]:
```

	a	b	duplicated
0	1	a	False
1	2	b	False
2	3	c	True
3	3	c	True
4	4	d	False
5	5	e	True
6	5	e	True
7	5	f	False
8	6	g	False
9	3	c	True

```
[45]: tiny[tiny['duplicated'] == True]
```

```
[45]:
```

	a	b	duplicated
2	3	c	True
3	3	c	True
5	5	e	True
6	5	e	True
9	3	c	True

```
[46]: for row in range(1, len(tiny)):
        if((tiny.iloc[row - 1] == tiny.iloc[row]).all()):
            empty_list.append(row)
```

Check that we got the correct indexes.

```
[47]: empty_list
```

```
[47]: [3, 6]
```

Make a new data frame with the unwanted rows `.dropped`.

```
[49]: tiny_drop = tiny.drop(empty_list)
tiny_drop
```

```
[49]:   a  b  duplicated
0  1  a         False
1  2  b         False
2  3  c          True
4  4  d         False
5  5  e          True
7  5  f         False
8  6  g         False
9  3  c          True
```

Use `.reset_index()` to make a new sequential index for our data frame.

```
[53]: tiny_drop = tiny_drop.reset_index()
```

Marvel at your work!

```
[54]: tiny_drop
```

```
[54]:   index  a  b  duplicated
0      0  1  a         False
1      1  2  b         False
2      2  3  c          True
3      4  4  d         False
4      5  5  e          True
5      7  5  f         False
6      8  6  g         False
7      9  3  c          True
```

If you don't like the "index" column with old indexes (sometimes it's useful to have the old indexes – here it's just annoying), you can set `drop=True` when you call `.reset_index()` above.

1.7 Run your code on the cancer data

Try our code on the (small version of the) cancer data!

1.7.1 Load the data

```
[33]: cancer = pd.read_csv('./data/small_cancer_data.csv')
cancer.head()
```

```
[33]:   id  thick  chrom  class
0  1000025   5.0    3.0  benign
```

```

1  1002945    5.0    3.0  benign
2  1015425    3.0    3.0  benign
3  1016277    6.0    3.0  benign
4  1017023    4.0    3.0  benign

```

1.7.2 Get the number of rows

```
[55]: # get the number of rows using shape
      cancer.shape
```

```
[55]: (699, 4)
```

```
[56]: len(cancer)
```

```
[56]: 699
```

1.7.3 Make an empty list for indexes

```
[58]: cancer_empty = []
```

1.7.4 Run your for loop!

```
[59]: for cancer_row in range(1, len(cancer)):
      if((cancer.iloc[cancer_row - 1] == cancer.iloc[cancer_row]).all()):
          cancer_empty.append(cancer_row)
```

1.7.5 Check the indexes you found

```
[60]: cancer_empty
```

```
[60]: [208, 322, 443, 561, 684, 690, 698]
```

1.7.6 Drop the unwanted rows

```
[62]: cancer_drop = cancer.drop(cancer_empty)
      cancer_drop
```

```
[62]:
```

	id	thick	chrom	class
0	1000025	5.0	3.0	benign
1	1002945	5.0	3.0	benign
2	1015425	3.0	3.0	benign
3	1016277	6.0	3.0	benign
4	1017023	4.0	3.0	benign
..
693	763235	3.0	2.0	benign
694	776715	3.0	1.0	benign


```

695    841769    2.0    1.0    benign
696    888820    5.0    8.0  malignant
697    897471    4.0   10.0  malignant

```

[692 rows x 4 columns]

1.7.7 Reset the row indexes

```
[64]: cancer_drop = cancer_drop.reset_index()
      cancer_drop
```

```
[64]:
```

	level_0	index	id	thick	chrom	class
0	0	0	1000025	5.0	3.0	benign
1	1	1	1002945	5.0	3.0	benign
2	2	2	1015425	3.0	3.0	benign
3	3	3	1016277	6.0	3.0	benign
4	4	4	1017023	4.0	3.0	benign
..
687	687	693	763235	3.0	2.0	benign
688	688	694	776715	3.0	1.0	benign
689	689	695	841769	2.0	1.0	benign
690	690	696	888820	5.0	8.0	malignant
691	691	697	897471	4.0	10.0	malignant

[692 rows x 6 columns]

1.7.8 Check the shape to confirm the rows were dropped!

```
[66]: # Old row = 699
      # Drop = 7
      699 - 7
```

[66]: 692

```
[65]: len(cancer_drop)
```

[65]: 692

1.8 Wrapping it all in a function

Once you've got your code running, put it all in a function so it's reusable!

```
[80]: # Reset index function
      def drop_duplicated(file_name):
```

```

empty_list = []

for row in range(1, len(file_name)):
    if((file_name.iloc[row - 1] == file_name.iloc[row]).all()):
        empty_list.append(row)

file_name_drop = file_name.drop(empty_list)

file_name_reset_index = file_name_drop.reset_index()

file_name = file_name_reset_index

return file_name

```

Run your function!

```
[81]: cancer_small = pd.read_csv('./data/small_cancer_data.csv')
```

```
[82]: cancer_small.shape
```

```
[82]: (699, 4)
```

```
[83]: cancer_small_test = drop_duplicated(cancer_small)
```

Check the shape to confirm your function worked!

```
[85]: cancer_small_test.shape
```

```
[85]: (692, 5)
```

1.8.1 High-five the person closest to you!

Because you deserve a high-five right now.
