

Wrangling III

In this tutorial, we'll round out our focus on data wrangling by looking

- handling duplicate values
- data transformations

Preliminaries

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

```
In [4]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

Now we'll get set up to work by

- loading the cancer data and cleaning it up (as before)
- trim out some columns so we can look at the data frame more easily
- shorten up some of the column names to save ourselves some typing

Let's reuse our function to do the loading and cleaning.

```
In [5]: def bcd_load_clean():
bcd = pd.read_csv('C:\\Users\\wgero\\Downloads\\PSY 341K\\data\\breast_cancer_data.csv')
bcd['patient_id'] = bcd['patient_id'].astype('string')
bcd['doctor_name'] = bcd['doctor_name'].str.split().str[1]
bcd['bare_nuclei'] = bcd['bare_nuclei'].replace('?', '')
bcd['bare_nuclei'] = pd.to_numeric(bcd['bare_nuclei'])

return bcd
```

```
In [11]: bcd = bcd_load_clean()
bcd
```

```
Out[11]:
```

	patient_id	clump_thickness	cell_size_uniformity	cell_shape_uniformity	marginal_adhesion	single
0	1000025	5.0	1.0	1	1	
1	1002945	5.0	4.0	4	5	
2	1015425	3.0	1.0	1	1	
3	1016277	6.0	8.0	8	1	
4	1017023	4.0	1.0	1	3	
...
694	776715	3.0	1.0	1	1	
695	841769	2.0	1.0	1	1	
696	888820	5.0	10.0	10	3	
697	897471	4.0	8.0	6	4	
698	897471	4.0	8.0	8	5	

699 rows × 12 columns

Make a little version with just two numeric columns to play with.

```
In [12]: bcd2 = bcd[['patient_id', 'clump_thickness', 'bland_chromatin', 'class']].copy()
```

Let's give the columns shorter names to save some typing.

```
In [14]: bcd2 = bcd2.rename(columns={'clump_thickness': 'thick',
                                     'bland_chromatin': 'chrom',
                                     'patient_id': 'id'})
```

Duplicate entries

As we have already seen, datasets can contain strange things that we have to overcome prior to analysis. One of the most common issues in a dataset are duplicate entries. These are common with large datasets that have been transcribed by humans at some point. Humans get bored, lose their place, etc.

Let's look at the shape of our cancer data frame (remember data frames have a `shape` attribute).

```
In [16]: bcd2.shape
```

```
Out[16]: (699, 4)
```

Now let's look at the number of unique entries using the `nunique()` data frame method; this will return the number of distinct values in each column.

```
In [17]: bcd2.nunique()
```

```
Out[17]: id          645  
thick         10  
chrom         10  
class          2  
dtype: int64
```

So we can see that, while there are 699 observations in our data, there are only 645 unique patient ids. This tells us that several patients have multiple entries. These could be from patients making multiple visits to the doctor, or they could be a mistakes, or some combination thereof.

We can find out which rows – which entire observations – are identical with the `duplicated()` method.

```
In [18]: bcd2.duplicated()
```

```
Out[18]: 0          False  
1          False  
2          False  
3          False  
4          False  
...  
694        False  
695        False  
696        False  
697        False  
698         True  
Length: 699, dtype: bool
```

That's not terribly helpful by itself, but...

In the cell below, count the number of duplicated rows (remember a True is a 1).

```
In [21]: bcd2.duplicated().sum()
```

```
Out[21]: 12
```

We can also use the output of `.duplicated()` to do logical indexing to see the observations that have duplicates. Do that in the cell below.

```
In [36]: duplicates = bcd2.duplicated()  
duplicate_index = bcd2[duplicates]  
duplicate_index
```

Out[36]:

	id	thick	chrom	class
208	1218860	1.0	3.0	benign
253	1100524	6.0	7.0	malignant
254	1116116	9.0	3.0	malignant
258	1198641	3.0	3.0	benign
272	320675	3.0	7.0	malignant
322	733639	3.0	3.0	benign
338	704097	1.0	2.0	benign
443	734111	1.0	1.0	benign
561	1321942	5.0	3.0	benign
684	466906	1.0	1.0	benign
690	654546	1.0	1.0	benign
698	897471	4.0	10.0	malignant

This is promising but, if we look at what is listed, we don't actually see any duplicates. So what is `duplicated()` doing?

Use the cell below to get help on `duplicated()` using `help()` or `?`.

In [27]: `help(bcd2.duplicated)`

Help on method duplicated in module pandas.core.frame:

duplicated(subset: 'Hashable | Sequence[Hashable] | None' = None, keep: "Literal['first'] | Literal['last'] | Literal[False]" = 'first') -> 'Series' method of pandas.core.frame.DataFrame instance

Return boolean Series denoting duplicate rows.

Considering certain columns is optional.

Parameters

subset : column label or sequence of labels, optional

Only consider certain columns for identifying duplicates, by default use all of the columns.

keep : {'first', 'last', False}, default 'first'

Determines which duplicates (if any) to mark.

- ``first`` : Mark duplicates as ``True`` except for the first occurrence.
- ``last`` : Mark duplicates as ``True`` except for the last occurrence.
- False : Mark all duplicates as ``True``.

Returns

Series

Boolean series for each duplicated rows.

See Also

Index.duplicated : Equivalent method on index.

Series.duplicated : Equivalent method on Series.

Series.drop_duplicates : Remove duplicate values from Series.

DataFrame.drop_duplicates : Remove duplicate values from DataFrame.

Examples

Consider dataset containing ramen rating.

```
>>> df = pd.DataFrame({
...     'brand': ['Yum Yum', 'Yum Yum', 'Indomie', 'Indomie', 'Indomie'],
...     'style': ['cup', 'cup', 'cup', 'pack', 'pack'],
...     'rating': [4, 4, 3.5, 15, 5]
... })
>>> df
   brand style  rating
0  Yum Yum   cup    4.0
1  Yum Yum   cup    4.0
2  Indomie   cup    3.5
3  Indomie  pack   15.0
4  Indomie  pack    5.0
```

By default, for each set of duplicated values, the first occurrence is set on False and all others on True.

```
>>> df.duplicated()
0    False
1     True
2    False
3    False
4    False
dtype: bool
```

By using 'last', the last occurrence of each set of duplicated values is set on False and all others on True.

```
>>> df.duplicated(keep='last')
0    True
1    False
2    False
3    False
4    False
dtype: bool
```

By setting ``keep`` on False, all duplicates are True.

```
>>> df.duplicated(keep=False)
0    True
1    True
2    False
3    False
4    False
dtype: bool
```

To find duplicates on specific column(s), use ``subset``.

```
>>> df.duplicated(subset=['brand'])
0    False
1     True
2    False
3     True
4     True
dtype: bool
```

... we can see that it has a "keep" argument. By default, `duplicated()` it gives us the *first* instance of any duplicated rows. We can make it show all the rows with `keep=False`.

Go ahead and do that in the cell below.

```
In [34]: dupes = bcd2.duplicated(keep=False)
         dupe_index = bcd2[dupes]
         dupe_index
```

Out[34]:

	id	thick	chrom	class
42	1100524	6.0	7.0	malignant
62	1116116	9.0	3.0	malignant
168	1198641	3.0	3.0	benign
207	1218860	1.0	3.0	benign
208	1218860	1.0	3.0	benign
253	1100524	6.0	7.0	malignant
254	1116116	9.0	3.0	malignant
258	1198641	3.0	3.0	benign
267	320675	3.0	7.0	malignant
272	320675	3.0	7.0	malignant
314	704097	1.0	2.0	benign
321	733639	3.0	3.0	benign
322	733639	3.0	3.0	benign
338	704097	1.0	2.0	benign
442	734111	1.0	1.0	benign
443	734111	1.0	1.0	benign
560	1321942	5.0	3.0	benign
561	1321942	5.0	3.0	benign
683	466906	1.0	1.0	benign
684	466906	1.0	1.0	benign
689	654546	1.0	1.0	benign
690	654546	1.0	1.0	benign
697	897471	4.0	10.0	malignant
698	897471	4.0	10.0	malignant

Hm. That's somewhat helpful. If we look near the bottom, we see that the last 5 or so duplicates occur in successive rows, perhaps indicating a data entry mistake. Perhaps looking at the data sorted by patient ID would be more helpful.

In the cell below, use the the `.sort_values()` method to look at our duplicates sorted by ID.

In [33]: `bcd2.sort_values(by = 'id')`

Out[33]:

	id	thick	chrom	class
0	1000025	5.0	3.0	benign
485	1002025	1.0	1.0	benign
382	1002504	3.0	3.0	benign
1	1002945	5.0	3.0	benign
2	1015425	3.0	3.0	benign
...
364	896404	2.0	3.0	benign
365	897172	2.0	NaN	benign
697	897471	4.0	10.0	malignant
698	897471	4.0	10.0	malignant
366	95719	6.0	7.0	malignant

699 rows × 4 columns

So most of the duplicates occur in adjacent rows, but others do not. Perhaps we should check and see if the same patients occur multiple times with different measurements, indicating multiple visits to the doctor.

Use the cell below and the `subset` argument to `duplicated()` to look at multiple entries for any patients that have them.

In [110...

```
multiple_entries = bcd2[bcd2.duplicated(subset=['id'], keep = False)]
multiple_entries
```



```
Out[110]:
```

	chrom	class	id	thick
0	-0.183305	NaN	<NA>	0.206942
1	-0.183305	NaN	<NA>	0.206942
2	-0.183305	NaN	<NA>	-0.502864
3	-0.183305	NaN	<NA>	0.561845
4	-0.183305	NaN	<NA>	-0.147961
...
694	-1.002577	NaN	<NA>	-0.502864
695	-1.002577	NaN	<NA>	-0.857766
696	1.864876	NaN	<NA>	0.206942
697	2.684148	NaN	<NA>	-0.147961
698	2.684148	NaN	<NA>	-0.147961

699 rows × 4 columns

Now, in the cell below, do the same thing but sort the output by patient ID.

```
In [111]: multi_entries = bcd2[bcd2.duplicated(subset=['id'])]
multi_entries
```

```
Out[111]:
```

	chrom	class	id	thick
1	-0.183305	NaN	<NA>	0.206942
2	-0.183305	NaN	<NA>	-0.502864
3	-0.183305	NaN	<NA>	0.561845
4	-0.183305	NaN	<NA>	-0.147961
5	2.274512	NaN	<NA>	1.271650
...
694	-1.002577	NaN	<NA>	-0.502864
695	-1.002577	NaN	<NA>	-0.857766
696	1.864876	NaN	<NA>	0.206942
697	2.684148	NaN	<NA>	-0.147961
698	2.684148	NaN	<NA>	-0.147961

698 rows × 4 columns

So it looks like patients do come in multiple times and the values can change between visits.

We can look at repeat patient's number of visits directly if we want. We'll take advantage of the fact that the `.size` of a `groupby()` object returns the number of rows for each group.

```
In [49]: repeat_patients = bcd2.groupby('id').size().sort_values(ascending=False)
```

```
In [50]: repeat_patients
```

```
Out[50]: id
1182404    6
1276091    5
1198641    3
1299596    2
1158247    2
..
1200892    1
1200952    1
1201834    1
1201870    1
95719      1
Length: 645, dtype: int64
```

So one patient came in 6 times.

Use the cell below look at the data for the patient with 6 visits.

```
In [20]: bcd2[bcd2['id'] == '1182404']
```

```
Out[20]:
```

	id	thick	chrom	class
136	1182404	4.0	2.0	benign
256	1182404	3.0	1.0	benign
257	1182404	3.0	2.0	benign
265	1182404	5.0	3.0	benign
448	1182404	1.0	1.0	benign
497	1182404	4.0	1.0	benign

So it appears that some patients have multiple legitimate entries in the data frame.

If you were put in charge of analyzing these data, what would you do with duplicate observations in this data frame, and why?

I would ask the ones who collected the data for context. Maybe there were patients who had multiple visits. If they tell me that there was, I would consider keeping it in the dataset. If not, I

would treat the duplicates as data entry errors and remove them from the dataset.

Transforming data

Sometimes we wish to apply a transform to data by pushing each data value through some function. Common transformations are unit conversions (miles to kilometers, for example), log or power transformations, and normalizing data (for example, converting data to z-scores).

Transforming data with a built-in function

Consider the following data...

```
In [53]: df = pd.DataFrame({'x': range(6),  
                           'y': [0.1, 0.9, 4.2, 8.7, 15.9, 26]})
```

```
In [52]: df
```

```
Out[52]:
```

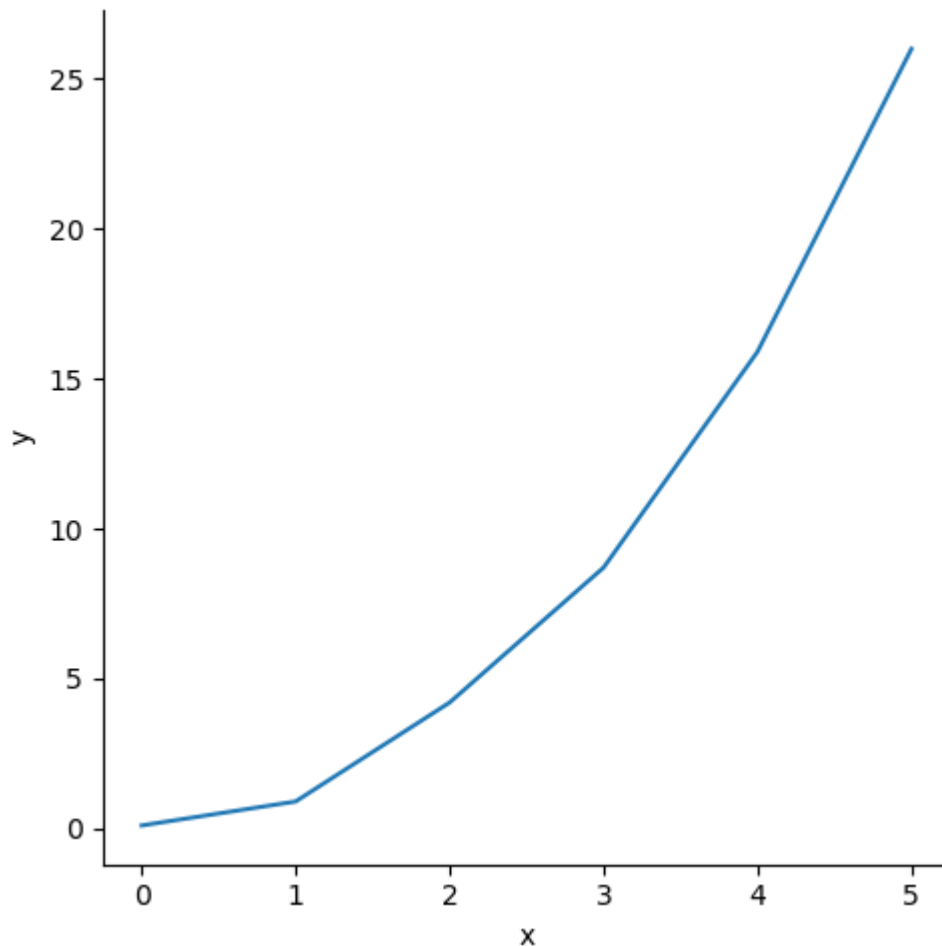
	x	y
0	0	0.1
1	1	0.9
2	2	4.2
3	3	8.7
4	4	15.9
5	5	26.0

Plot the data (y vs. x) (seaborn's `relplot()` is handy).

```
In [54]: %matplotlib inline
```

```
In [57]: # plot y vs. x  
sns.relplot(x = 'x', y = 'y', data = df, kind = 'line')
```

```
Out[57]: <seaborn.axisgrid.FacetGrid at 0x181e3ed99a0>
```

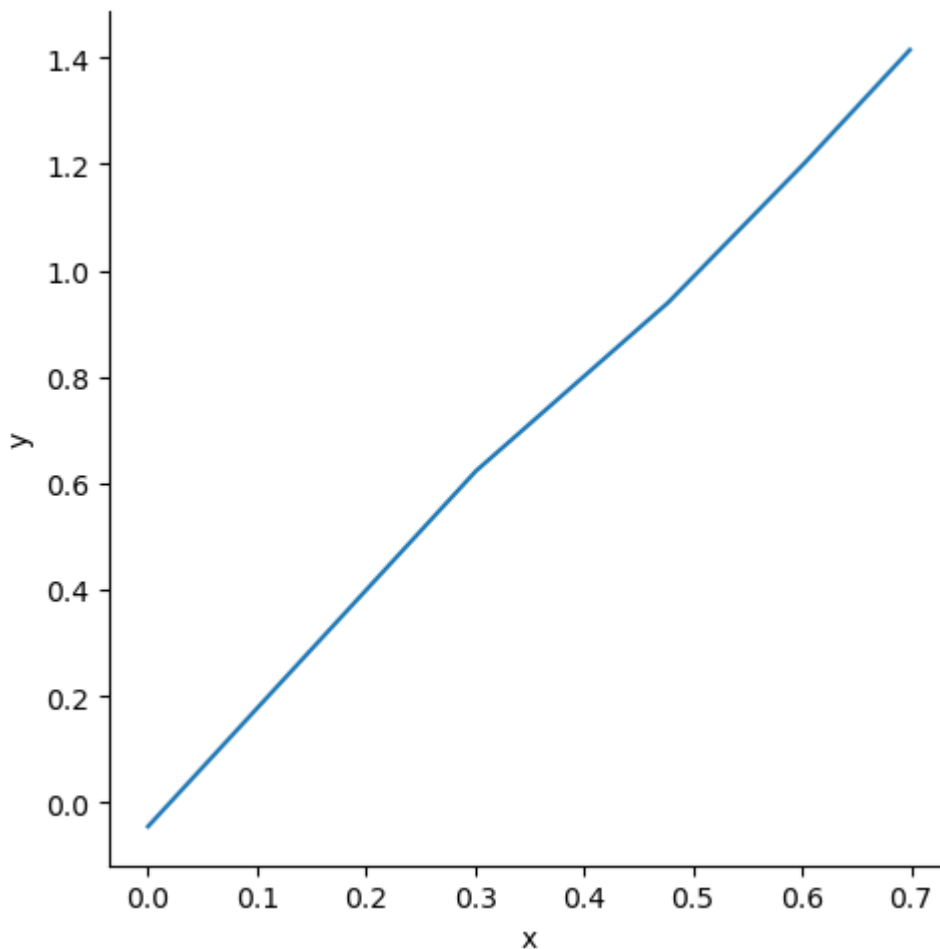


These data look non-linear, like they are following a power law. If that's true, we should get a straight line if we plot the log of the values against one another. In order to get these values, we will use the `transform()` method to convert the values into their logs.

```
In [56]: df_trans = df.copy()
df_trans['y'] = df['y'].transform(np.log10)
df_trans['x'] = df['x'].transform(np.log10)
```

```
In [58]: # plot new y vs. new x
sns.relplot(x = 'x', y = 'y', data = df_trans, kind = 'line')
```

```
Out[58]: <seaborn.axisgrid.FacetGrid at 0x181e3fb3be0>
```



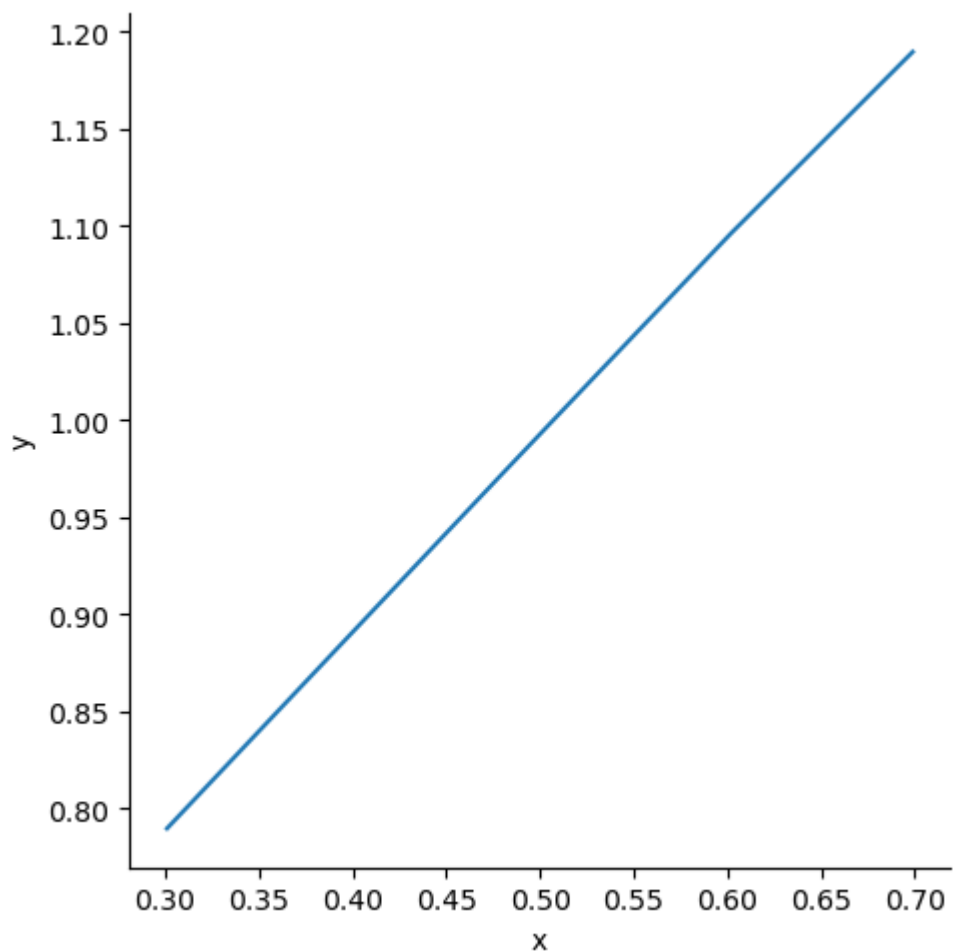
Sure enough. The slope of the line should tell us the exponent of the power law, and it looks to be about 2. If that's the case, then transforming the original y-values with a square-root function should also produce a straight line.

In the cells below, use `transform()` to get the square root of the original y values, and plot them against the x values.

```
In [60]: # get sqrts
df_sqrt = df_trans.copy()
df_sqrt['y'] = df_trans['y'].transform(np.sqrt)
```

```
In [61]: #plot
sns.relplot(x='x', y='y', data=df_sqrt, kind='line')
```

```
Out[61]: <seaborn.axisgrid.FacetGrid at 0x181e47dcac0>
```



We could also transform our cancer data. In the cell below, create a new data frame in which the numeric values are the natural log of the original values.

```
In [113... # compute log vals
bcd2_log = bcd2.copy()
numeric_cols = bcd2.select_dtypes(include = np.number).columns
bcd2_log[numeric_cols] = bcd2[numeric_cols].apply(np.log)

bcd2_log.head()
```

```
Out[113]:
```

	chrom	class	id	thick
0	NaN	NaN	<NA>	-1.575317
1	NaN	NaN	<NA>	-1.575317
2	NaN	NaN	<NA>	NaN
3	NaN	NaN	<NA>	-0.576530
4	NaN	NaN	<NA>	NaN

Applying a custom function to data

A great thing about `transform()` (and some other data frame methods) is you can use your own functions, not just built in ones.

For `transform()`, the only requirement is that your function

- be able to take a data frame as input
- produce output the same size as the input, or
- produce a single value

Here's a function to "center" data by subtracting the mean from each value.

```
In [66]: def center_data(grp):  
        grp_mean = grp.mean(numeric_only = True)  
        grp = (grp - grp_mean)  
        return grp
```

In the cell below, use our new function to create a new version of our cancer data frame with the mean removed from each group of data. The `.transform()` method works column-by-column, so you don't need to worry about grouping the data.

```
In [74]: bcd2_no_mean = bcd2.transform(center_data)  
bcd2_no_mean
```

```
Out[74]:
```

	chrom	class	id	thick
0	-0.447482	NaN	<NA>	0.583095
1	-0.447482	NaN	<NA>	0.583095
2	-0.447482	NaN	<NA>	-1.416905
3	-0.447482	NaN	<NA>	1.583095
4	-0.447482	NaN	<NA>	-0.416905
...
694	-2.447482	NaN	<NA>	-1.416905
695	-2.447482	NaN	<NA>	-2.416905
696	4.552518	NaN	<NA>	0.583095
697	6.552518	NaN	<NA>	-0.416905
698	6.552518	NaN	<NA>	-0.416905

699 rows × 4 columns

Confirm this worked by computing the mean for each column of your transformed data.

```
In [80]: num_cols = bcd2.select_dtypes(include = np.number)
bcd2_mean = num_cols.mean()
bcd2_mean
```

```
Out[80]: thick    4.416905
chrom    3.447482
dtype: float64
```

In the cells below, write a function to convert the cancer data to z-scores, and use your new function to convert the numeric columns of our cancer data frame.

```
In [97]: # my z-score function!
def zscore_conversion(grp):
    num_cols = grp.select_dtypes(include = np.number)
    grp_mean = num_cols.mean()
    grp_std = num_cols.std()
    grp_zscore = (num_cols - grp_mean) / grp_std
    return grp_zscore
```

```
In [98]: # run transform() with my function
bcd2_zscore = bcd2.transform(zscore_conversion)
```

```
In [99]: # Look at the transformed data
bcd2_zscore
```

```
Out[99]:
```

	chrom	thick
0	-0.183305	0.206942
1	-0.183305	0.206942
2	-0.183305	-0.502864
3	-0.183305	0.561845
4	-0.183305	-0.147961
...
694	-1.002577	-0.502864
695	-1.002577	-0.857766
696	1.864876	0.206942
697	2.684148	-0.147961
698	2.684148	-0.147961

699 rows × 2 columns

```
In [101]: # see what the means are
means = bcd2_zscore.mean()
means
```



```
Out[101]: chrom    4.600636e-17  
         thick    -6.131676e-17  
         dtype: float64
```

```
In [103... # see what the ... are  
st_devs = bcd2_zscore.std()  
st_devs
```

```
Out[103]: chrom    1.0  
         thick    1.0  
         dtype: float64
```

lambda functions

Lambda functions, also know as anonymous functions, are short, one-off functions that are often used in situation in which **all** you need the function for is get passed to a method such as `transform()`

While the structure of a normal function is:

```
In [ ]: def func_name(input_arg) :  
        calculations  
        ret_val = more calculations  
  
        return ret_val
```

The structure of a lambda function is:

```
In [ ]: lambda input_arg : calculation of ret_val
```

Here's how we would compute z-scores using a lambda function:

```
In [104... trans_data = bcd2[['thick', 'chrom']].transform(  
            lambda col_vals: (col_vals - col_vals.mean()) / col_vals.std()  
            )
```

Note that the entire lambda function is the one and only input to `transform()` .

In the cell below, confirm that the lambda function method worked.

```
In [105... trans_data
```

```
Out[105]:
```

	thick	chrom
0	0.206942	-0.183305
1	0.206942	-0.183305
2	-0.502864	-0.183305
3	0.561845	-0.183305
4	-0.147961	-0.183305
...
694	-0.502864	-1.002577
695	-0.857766	-1.002577
696	0.206942	1.864876
697	-0.147961	2.684148
698	-0.147961	2.684148

699 rows × 2 columns

For very simple transformations, using a lambda function makes a lot of sense. For more complicated transformations, we'd probably want to just create a regular function, or the code could become unreadable.

How complicated is too complicated? That's up to you, but anything more complicated than applying an offset and a scale factor (like computing a z-score), probably deserves its own function.

In the cell below, transform the numeric cancer data so the values range from 0 to 1 using a lambda function. You can assume that the maximum value is 10 and the minimum value is 1.

```
In [108... lambda_fn = lambda x: (x - 1) / (10 - 1) # max - min
normalized_bcd2 = bcd2.select_dtypes(include = np.number).apply(lambda_fn)
normalized_bcd2
```

```
Out[108]:
```

	chrom	thick
0	-0.131478	-0.088118
1	-0.131478	-0.088118
2	-0.131478	-0.166985
3	-0.131478	-0.048684
4	-0.131478	-0.127551
...
694	-0.222509	-0.166985
695	-0.222509	-0.206418
696	0.096097	-0.088118
697	0.187128	-0.127551
698	0.187128	-0.127551

699 rows × 2 columns

In the cell below, use a regular function to rescale the values from 0 to 1. In this case, however, do not assume you know the minimum and maximum values ahead of time.

```
In [114... # rescaling function
def rescale_0_to_1(column):
    min_val = column.min()
    max_val = column.max()
    return (column - min_val) / (max_val - min_val)

rescaled_bcd2 = bcd2.select_dtypes(include = np.number).apply(rescale_0_to_1)
rescaled_bcd2
```

Out[114]:

	chrom	thick
0	0.222222	0.444444
1	0.222222	0.444444
2	0.222222	0.222222
3	0.222222	0.555556
4	0.222222	0.333333
...
694	0.000000	0.222222
695	0.000000	0.111111
696	0.777778	0.444444
697	1.000000	0.333333
698	1.000000	0.333333

699 rows × 2 columns
