

Homework 1: Exploring & Visualizing Data

Part 1: Data Exploration and Preparation

Setup

Make sure you have seaborn and missingno installed. Run `pip3 install seaborn` and `pip3 install missingno` in your container/shell if you don't.

In this homework, we will more rigorously explore data visualization and data manipulation with a couple datasets. Please fill in the cells with `## YOUR CODE HERE` following the appropriate directions.

```
In [78]: # removes the need to call plt.show() every time
%matplotlib inline
```

Seaborn is a powerful data visualization library built on top of matplotlib. We will be using seaborn for this homework (since it is a better tool and you should know it well). Plus seaborn comes default with *much* better aesthetics (invoked with the `set()` function call).

Please import seaborn as `sns`

```
In [79]: import seaborn as sns
sns.set() # This sets aesthetic parameters in one step
```

First load the `titanic` dataset directly from seaborn. The `load_dataset` function will return a pandas dataframe. Documentation: https://seaborn.pydata.org/generated/seaborn.load_dataset.html (https://seaborn.pydata.org/generated/seaborn.load_dataset.html)

```
In [80]: titanic = sns.load_dataset('titanic')
```

Preparing a new dataset (Pandas)

Import `numpy` and `pandas` (remember to abbreviate them accordingly!)

```
In [81]: import numpy as np
import pandas as pd
```

Now use some pandas functions to get a quick overview/statistics on the dataset. Take a quick glance at the overview you create.

```
In [82]: titanic.describe()
```

```
Out[82]:
```

	survived	pclass	age	sibsp	parch	fare
count	891.000000	891.000000	714.000000	891.000000	891.000000	891.000000
mean	0.383838	2.308642	29.699118	0.523008	0.381594	32.204208
std	0.486592	0.836071	14.526497	1.102743	0.806057	49.693429
min	0.000000	1.000000	0.420000	0.000000	0.000000	0.000000
25%	0.000000	2.000000	20.125000	0.000000	0.000000	7.910400
50%	0.000000	3.000000	28.000000	0.000000	0.000000	14.454200
75%	1.000000	3.000000	38.000000	1.000000	0.000000	31.000000
max	1.000000	3.000000	80.000000	8.000000	6.000000	512.329200

With your created overview, you should be able to answer these questions:

- What was the age of the oldest person on board? 80
- What was the survival rate of people on board? 38.38%
- What was the average fare of people on board? 32.204

By the way, for getting overviews, pandas also has a `groupby` function that is quite nice to use. example:

```
In [83]: titanic.groupby(['sex', 'embark_town'])['survived'].mean()
```

```
Out[83]: sex      embark_town
female  Cherbourg      0.876712
         Queenstown     0.750000
         Southampton     0.689655
male     Cherbourg      0.305263
         Queenstown     0.073171
         Southampton     0.174603
Name: survived, dtype: float64
```

Now we have an overview of our dataset. The next thing we should do is clean it.

One quick observation is that there are a couple of repetitive columns. One example is the 'embarked_town' and 'embarked' columns which are conveying the same information. As we are preparing this dataset to be used in a model we only want to keep one of those columns. **Find the other repetitive pair of columns and drop the non-numerical one.**

Hint: is there a function that could help us preview the first few items in the dataset?

Note: adult_male differs from sex since there could be children males

```
In [84]: titanic.drop('embark_town', axis=1, inplace=True)
titanic.drop('who', axis=1, inplace=True)
titanic.drop('alive', axis=1, inplace=True)
titanic.head()
```

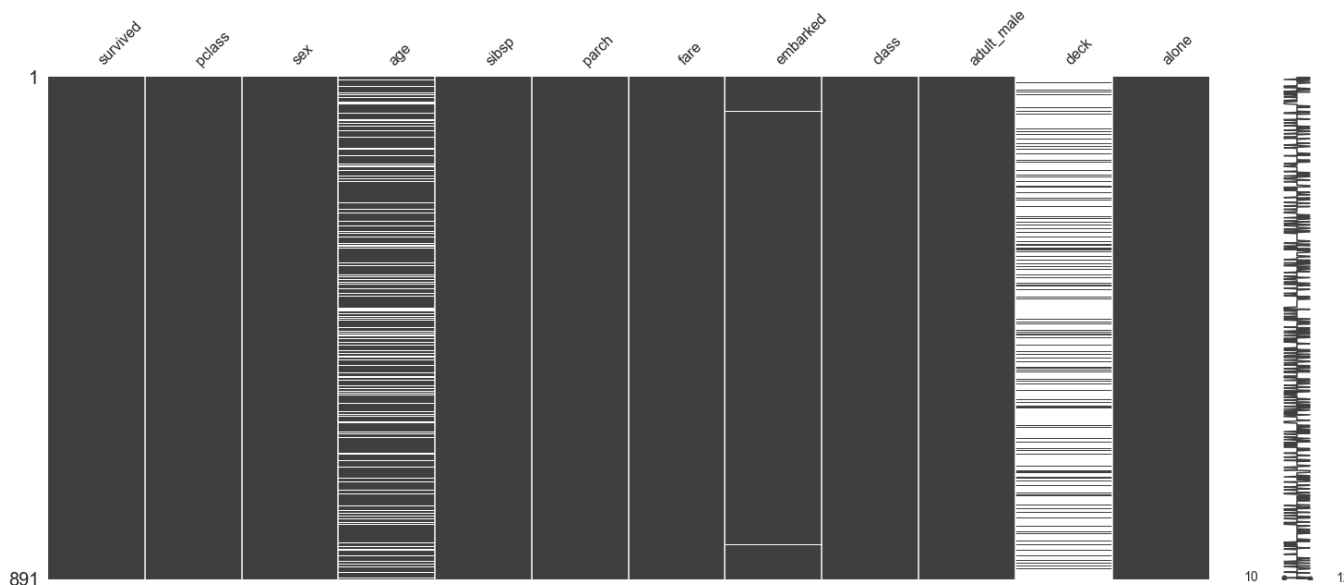
Out[84]:

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	adult_male	deck	alone
0	0	3	male	22.0	1	0	7.2500	S	Third	True	NaN	False
1	1	1	female	38.0	1	0	71.2833	C	First	False	C	False
2	1	3	female	26.0	0	0	7.9250	S	Third	False	NaN	True
3	1	1	female	35.0	1	0	53.1000	S	First	False	C	False
4	0	3	male	35.0	0	0	8.0500	S	Third	True	NaN	True

check for missing values and deal with them appropriately. `missingno` allows us to really easily see where missing values are in our dataset. It's a simple command. **Import missingno as msno and use its matrix function on titanic.** <https://www.residentmar.io/2016/03/28/missingno.html> (<https://www.residentmar.io/2016/03/28/missingno.html>)

```
In [85]: import missingno as msno
msno.matrix(titanic)
```

Out[85]: <AxesSubplot:>



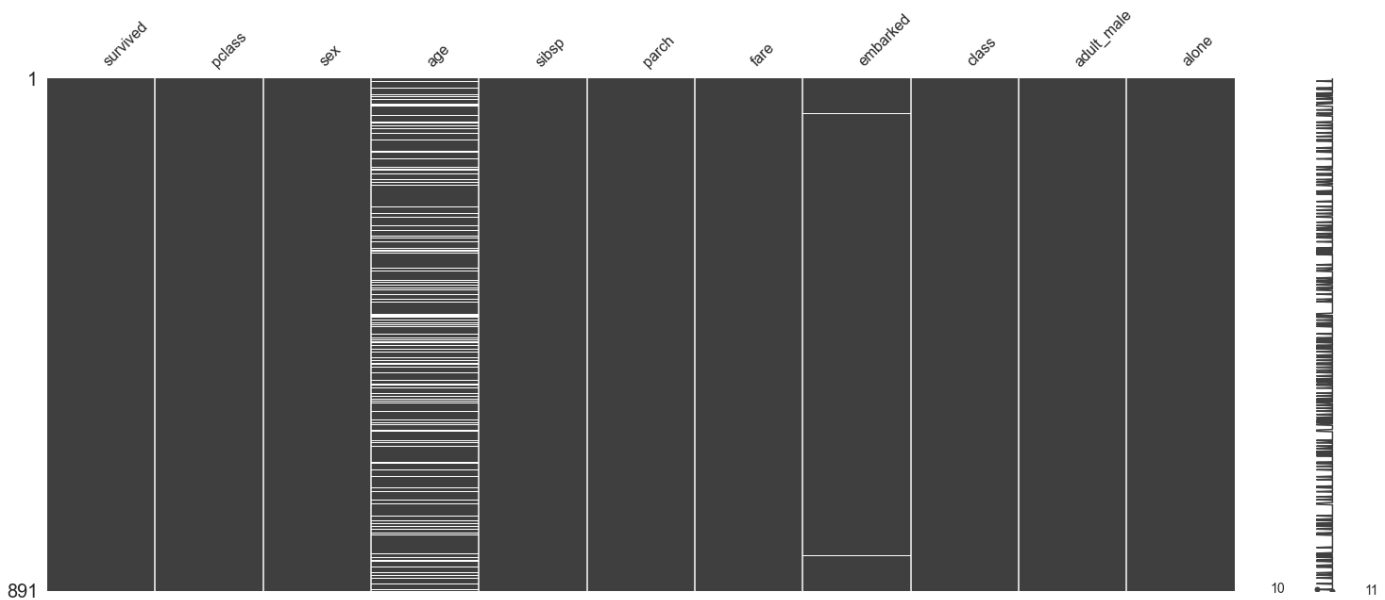
The white lines show us the missing data. One quick observation is the `deck` has a lot of missing data. Let's just go ahead and **drop the deck column from the dataset** since it's not that relevant. Make sure to set `inplace=True`. Documentation: <https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.drop.html> (<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.drop.html>).

```
In [86]: titanic.drop('deck', axis=1, inplace=True)
```

Now let's rerun the matrix and see. All that white is gone! Nice.

```
In [87]: msno.matrix(titanic)
```

```
Out[87]: <AxesSubplot:>
```



We still have a bunch of **missing values for the age field**. We can't just drop the age column since it is a pretty important datapoint. One way to deal with this is simply to just remove the records with missing information with `dropna()`, but this would end up removing out a significant amount of our data.

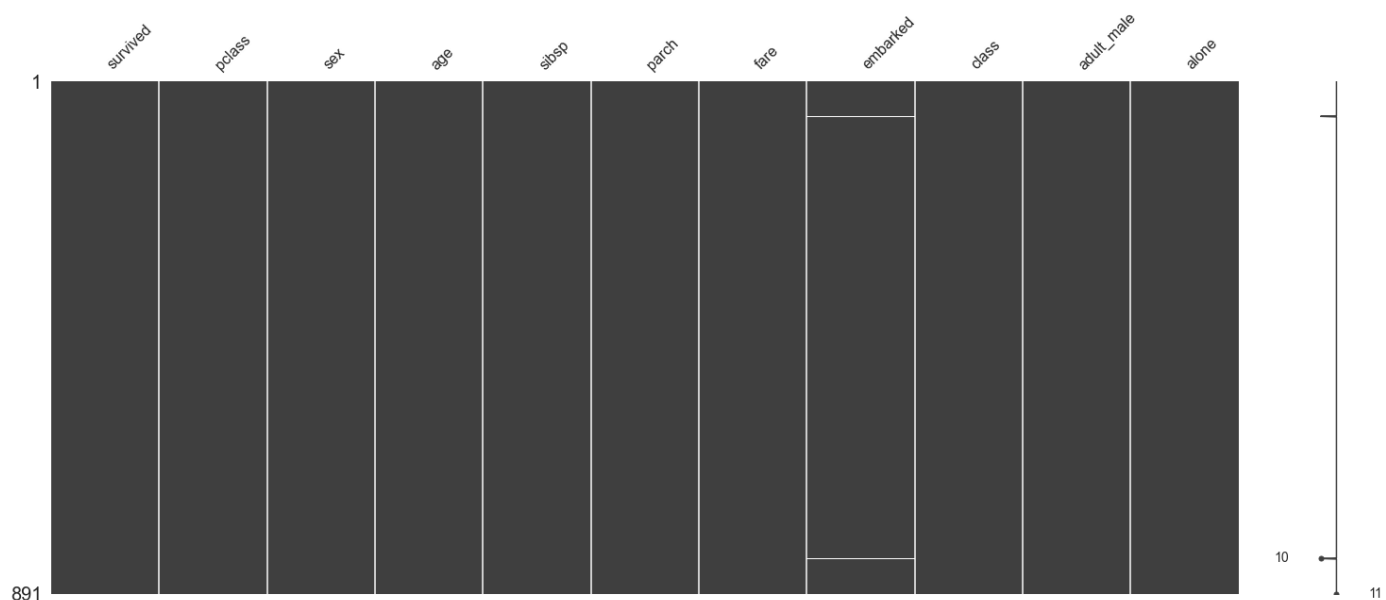
What do we do now? We can now explore a technique called `missing value imputation`. What this means is basically we find a reasonable way to *replace* the unknown data with workable values.

There's a lot of theory regarding how to do this properly, ([for the curious look here \(http://www.stat.columbia.edu/~gelman/arm/missing.pdf\)](http://www.stat.columbia.edu/~gelman/arm/missing.pdf)). We can simply put in the average age value for the missing ages. But this really isn't so great, and will skew our stats.

If we assume that the data is missing *at random* (which actually is rarely the case and very hard to prove), we can just fit a model to predict the missing value based on the other available factors. One popular way to do this is to use KNN (where you look at the nearest datapoints to a certain point to conclude the missing value), but we can also use deep neural networks to achieve this task.

You must now make you own decision on how to deal with the missing data. You may choose any of the methods discussed above. Easiest would be to fill in with average value (but this will skew our visualizations) (if you use pandas correctly, you can do this in one line - try looking at pandas documentation!). After writing your code, verify the result by rerunning the matrix - you should not see any white lines.

```
In [88]: ## YOUR CODE HERE
mean_age = titanic["age"].mean()
titanic['age'] = titanic['age'].fillna(mean_age)
msno.matrix(titanic)
titanic.dropna(inplace=True)
```



If we were to feed our dataset into a model, the model cannot understand features such as 'True', 'Yes', 'male'. You have to turn the data into binary format, 1 for positive and 0 negative. ** Fill in toBinary and then use it to binarize the 'alive', 'sex', 'alone', and 'adult_male' columns. **

```
In [89]: def toBinary(data, positive):
          data = (data == positive).astype(int)
          return data

          toBinary(titanic['sex'], 'male')
          toBinary(titanic['alone'], 'True')
          toBinary(titanic['adult_male'], 'True')
```

```
Out[89]: 0      0
          1      0
          2      0
          3      0
          4      0
          ..
          886    0
          887    0
          888    0
          889    0
          890    0
          Name: adult_male, Length: 889, dtype: int64
```

As we learned in lecture, we have to Normalize our numerical data points. Now you must **define the function 'Standardize' that standardizes the column passed in**, and then pass in the 'age' and 'fare' columns.

Hint: There is a sklearn function that can help with normalization

```
In [90]: from sklearn import preprocessing
          def Standardize(data):
              numpied = data.to_numpy().reshape(-1,1)
              scaler = preprocessing.StandardScaler().fit(numpied)
              return scaler.transform(numpied)

          titanic['age'] = Standardize(titanic['age'])
          titanic['fare'] = Standardize(titanic['fare'])
```

Finally we have to deal with the categorical columns that have more than two categories. For this situation we will use One-Hot encoding <http://queirozf.com/entries/one-hot-encoding-a-feature-on-a-pandas-dataframe-an-example> (<http://queirozf.com/entries/one-hot-encoding-a-feature-on-a-pandas-dataframe-an-example>).

Pandas has a good method 'get_dummies' that makes this task much easier. Read the documentation and fill in the code to **** One-Hot encode the 'pclass' and 'embarked' columns. .concat the output from get_dummies and drop the original columns. ****

```
In [91]: one_hot = pd.get_dummies(titanic["embarked"],prefix='embarked',drop_first=True)
titanic = pd.concat([titanic, one_hot])
titanic.drop('embarked', axis=1, inplace=True)

## YOUR CODE HERE
one_hot = pd.get_dummies(titanic["pclass"],prefix='pclass',drop_first=True)
titanic = pd.concat([titanic, one_hot])
titanic.drop('pclass', axis=1, inplace=True)
```

Out[91]:

	survived	sex	age	sibsp	parch	fare	class	adult_male	alone	embarked_Q	embarked_S	pclass
0	0.0	male	-0.590495	1.0	0.0	-0.500240	Third	True	False	NaN	NaN	3
1	1.0	female	0.643971	1.0	0.0	0.788947	First	False	False	NaN	NaN	1
2	1.0	female	-0.281878	0.0	0.0	-0.486650	Third	False	True	NaN	NaN	3
3	1.0	female	0.412509	1.0	0.0	0.422861	First	False	False	NaN	NaN	1
4	0.0	male	0.412509	0.0	0.0	-0.484133	Third	True	True	NaN	NaN	3

```
In [94]: titanic.head()
```

Out[94]:

	survived	sex	age	sibsp	parch	fare	class	adult_male	alone	embarked_Q	embarked_S	pclass
0	0.0	male	-0.590495	1.0	0.0	-0.500240	Third	True	False	NaN	NaN	3
1	1.0	female	0.643971	1.0	0.0	0.788947	First	False	False	NaN	NaN	1
2	1.0	female	-0.281878	0.0	0.0	-0.486650	Third	False	True	NaN	NaN	3
3	1.0	female	0.412509	1.0	0.0	0.422861	First	False	False	NaN	NaN	1
4	0.0	male	0.412509	0.0	0.0	-0.484133	Third	True	True	NaN	NaN	3

Intro to Seaborn

Seaborn can handle categorical data and NaN directly, **reload the titanic dataset to its original version**

```
In [95]: titanic = sns.load_dataset('titanic')
```

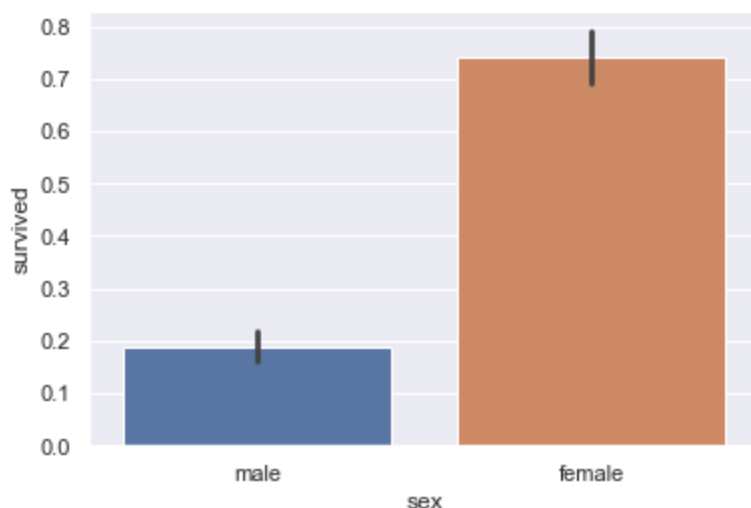
There are 2 types of data in any dataset: categorial and numerical data. We will first explore categorical data.

One really easy way to show categorical data is through bar plots. Let's explore how to make some in seaborn. We want to investigate the difference in rates at which males vs females survived the accident. Using the [documentation here \(https://seaborn.pydata.org/generated/seaborn.barplot.html\)](https://seaborn.pydata.org/generated/seaborn.barplot.html) and [example here \(http://seaborn.pydata.org/examples/color_palettes.html\)](http://seaborn.pydata.org/examples/color_palettes.html), create a `barplot` to depict this. It should be a really simple one-liner.

We will show you how to do this so you can get an idea of how to use the API.

```
In [96]: sns.barplot(x='sex', y='survived', data=titanic)
```

```
Out[96]: <AxesSubplot:xlabel='sex', ylabel='survived'>
```

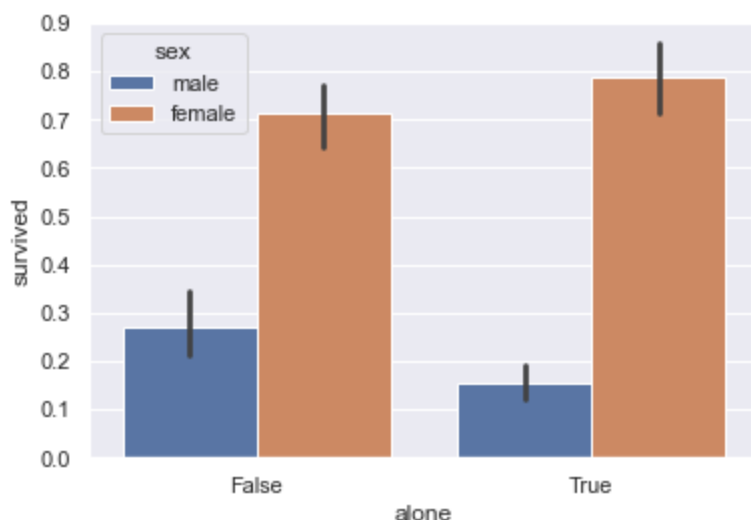


Notice how it was so easy to create the plot! You simply passed in the entire dataset, and just specified the `x` and `y` fields that you wanted exposed for the barplot. Behind the scenes seaborn ignored `NaN` values for you and automatically calculated the survival rate to plot. Also, that black tick is a 95% confidence interval that seaborn plots.

So we see that females were much more likely to make it out alive. What other factors do you think could have an impact on survival rate? ** Plot a couple more barplots below. ** Make sure to use *categorical* values, not something numerical like age or fare.*

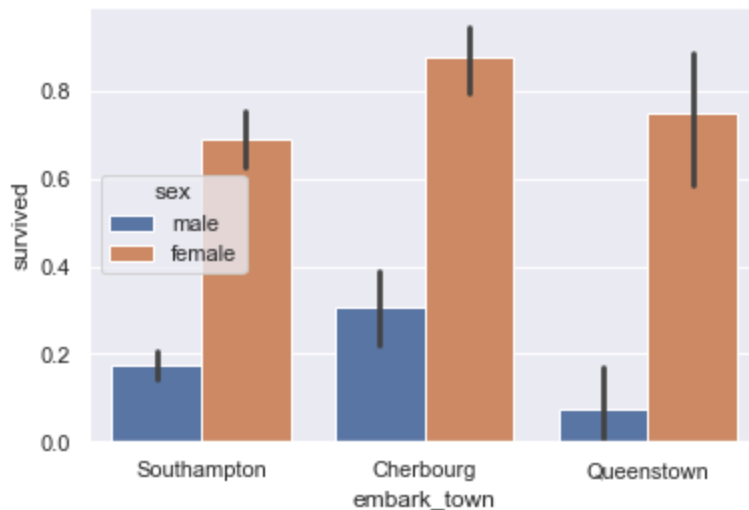
```
In [101]: sns.barplot(x='alone', y='survived', data=titanic, hue='sex')
```

```
Out[101]: <AxesSubplot:xlabel='alone', ylabel='survived'>
```




```
In [102]: sns.barplot(x='embark_town', y='survived', data=titanic, hue='sex')
```

```
Out[102]: <AxesSubplot:xlabel='embark_town', ylabel='survived'>
```



What if we wanted to add a further sex breakdown for the categories chosen above? Go back and add a `hue='sex'` parameter for the couple plots you just created, and seaborn will split each bar into a male/female comparison.

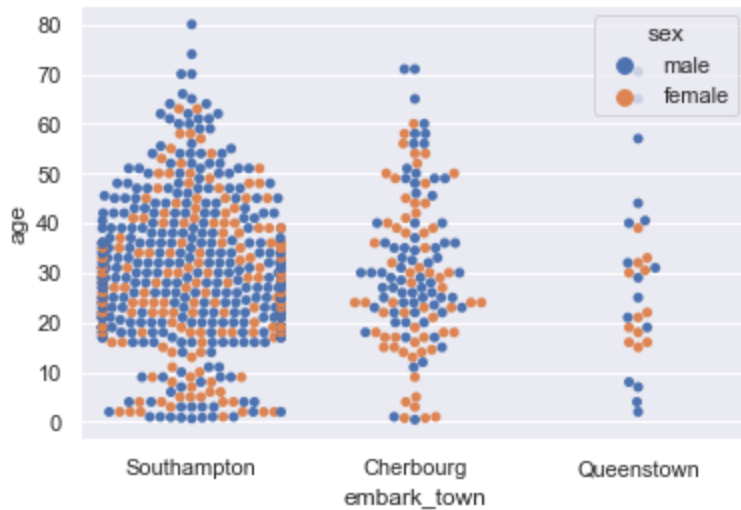
Now we want to compare the embarking town vs the age of the individuals. We don't simply want to use a barplot, since that will just give the average age; rather, we would like more insight into the relative and numeric *distribution* of ages.

A good tool to help us here is `swarmplot` (<https://seaborn.pydata.org/generated/seaborn.swarmplot.html>). Use this function to view `embark_town` vs `age`, again using `sex` as the `hue`.

```
In [103]: ax = sns.swarmplot(x="embark_town", y="age", data=titanic, hue='sex')
```

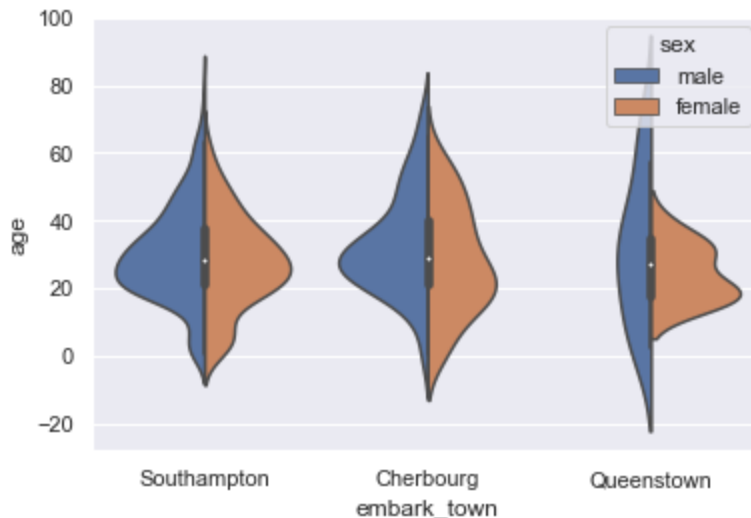
/usr/local/lib/python3.8/site-packages/seaborn/categorical.py:1296: UserWarning: 26.6% of the points cannot be placed; you may want to decrease the size of the markers or use stripplot.

```
warnings.warn(msg, UserWarning)
```



Cool! This gives us much more information. What if we didn't care about the number of individuals in each category at all, but rather just wanted to see the *distribution* in each category? [violinplot](https://seaborn.pydata.org/generated/seaborn.violinplot.html) (<https://seaborn.pydata.org/generated/seaborn.violinplot.html>) plots a density distribution. Plot that. Keep the hue .

```
In [106]: ax = sns.violinplot(x="embark_town", y="age", data=titanic, hue='sex', split='True')
```



Go back and clean up the violinplot by adding `split='True'` parameter.

Now take a few seconds to look at the graphs you've created of this data. What are some observations? Jot a couple down here.

Your observations Here

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/serif/regular/

• The distribution of surviving men is much wider than that of women.

- Southampton seems to have comparatively more survivors.

- Queenstown seems to have many less survivors.

As I mentioned, data is categorical or numeric. We already started getting into numerical data with the swarmplot and violinplot. We will now explore a couple more examples.

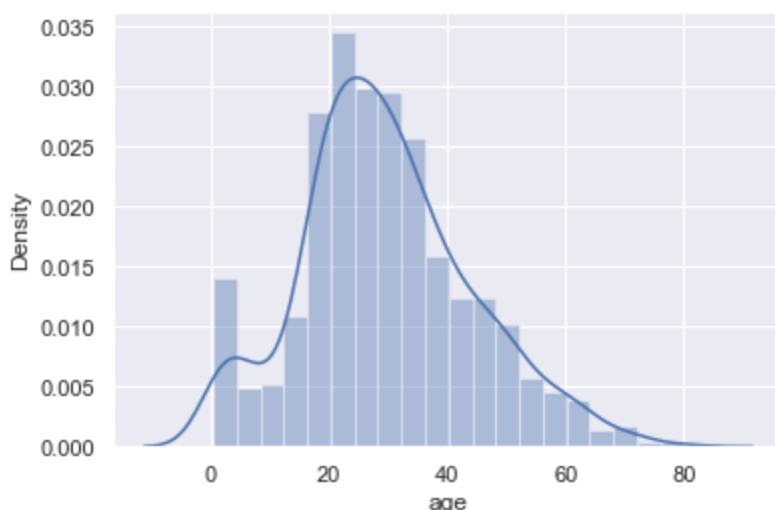
Let's look at the distribution of ages. Use `distplot` (<https://seaborn.pydata.org/generated/seaborn.distplot.html>) to make a histogram of just the ages.

```
In [108]: sns.distplot(titanic['age'])
```

```
/usr/local/lib/python3.8/site-packages/seaborn/distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).
```

```
warnings.warn(msg, FutureWarning)
```

```
Out[108]: <AxesSubplot:xlabel='age', ylabel='Density'>
```



If you did your missing value imputation by average value (If we had not reloaded the data set), your results will look very skewed. This is why we don't normally just fill in an average. As a quick fix for now, though, you can filter out the age values that equal the mean before passing it in to `displot`.

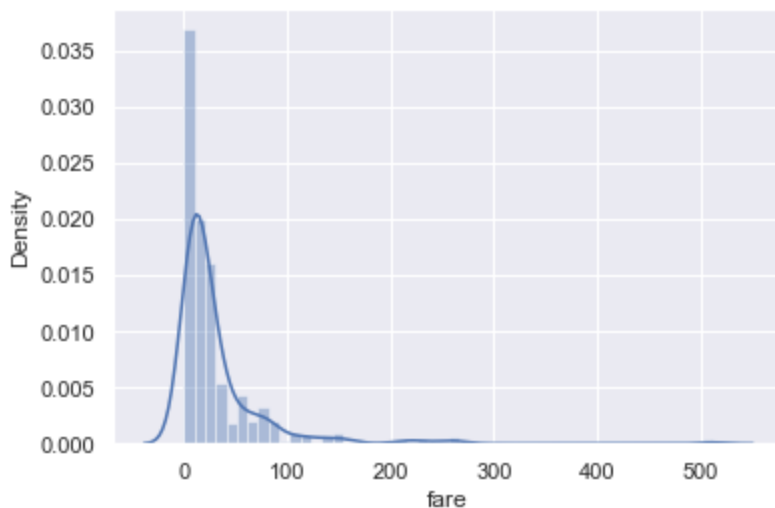
A histogram can nicely represent numerical data by breaking up numerical ranges into chunks so that it is easier to visualize. As you might notice from above, seaborn also automatically plots a gaussian kernel density estimate.

Do the same thing for fares - do you notice something odd about that histogram? What does that skew mean?

```
In [110]: sns.distplot(titanic['fare'])
```

```
/usr/local/lib/python3.8/site-packages/seaborn/distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).  
  warnings.warn(msg, FutureWarning)
```

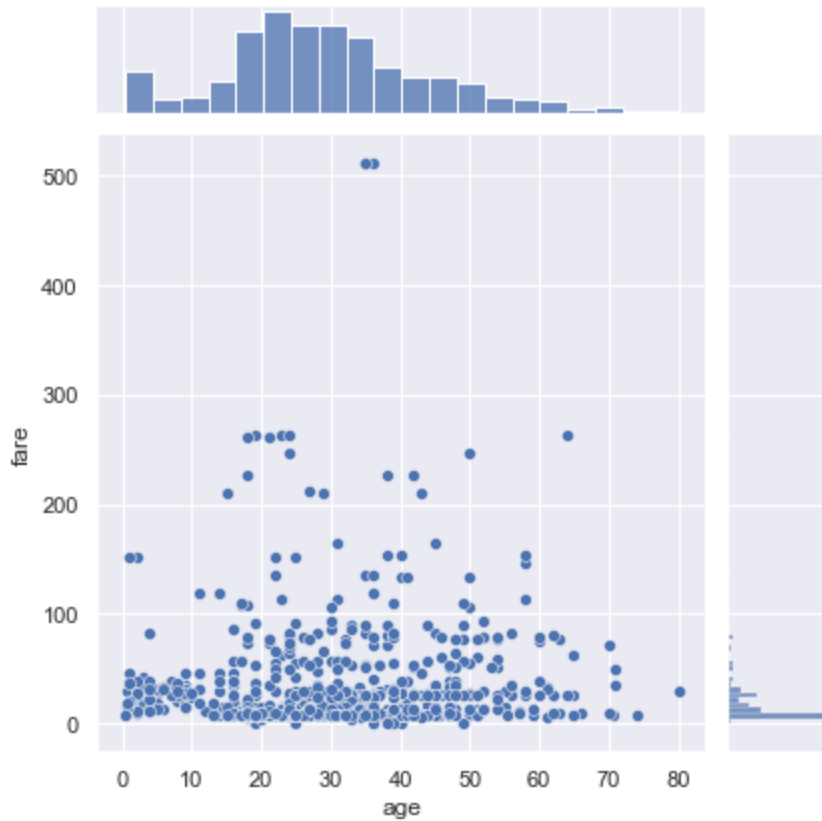
```
Out[110]: <AxesSubplot:xlabel='fare', ylabel='Density'>
```



Now, using the `jointplot` [function](https://seaborn.pydata.org/generated/seaborn.jointplot.html#seaborn.jointplot), make a scatterplot of the `age` and `fare` variables to see if there is any relationship between the two.

```
In [111]: sns.jointplot(x="age", y="fare", data=titanic)
```

```
Out[111]: <seaborn.axisgrid.JointGrid at 0x136b33b50>
```

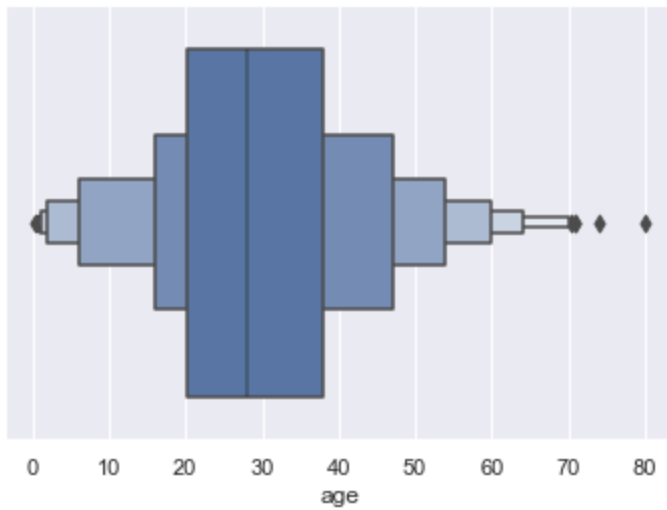


Scatterplots allow one to easily see trends/coorelations in data. As you can see here, there seems to be very little correlation. Also observe that seaborn automatically plots histograms.

Now, use a seaborn function we haven't used yet to plot something. The [API](http://seaborn.pydata.org/api.html) (<http://seaborn.pydata.org/api.html>) has a list of all the methods.

```
In [113]: sns.boxenplot(x = titanic["age"])
```

```
Out[113]: <AxesSubplot:xlabel='age'>
```



```
In [ ]:
```