

Visualization principles

1. Show the data - Jitter
2. Consider the scale - logs
3. Baselines
4. Ease comparisons
5. Set the scale

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import matplotlib as mpl
```

1. Show the data - Jitter

Google it: [Jitter in python](#)

Documentation contains such a good example we'll just [follow it](#)

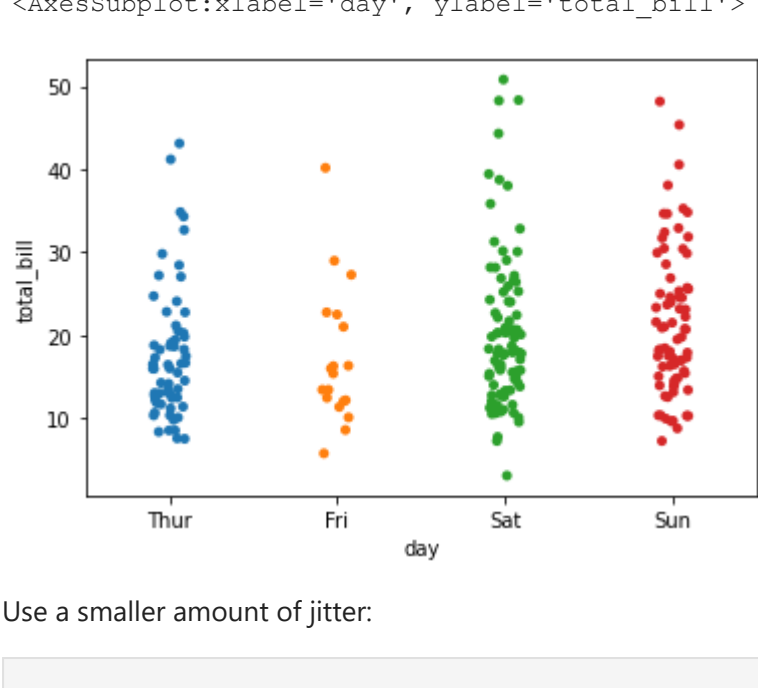
```
In [2]: tips = sns.load_dataset("tips")
```

```
In [3]: tips.head()
```

```
Out[3]:   total_bill  tip    sex  smoker  day  time  size
0      16.99  1.01  Female    No  Sun  Dinner     2
1      10.34  1.66    Male    No  Sun  Dinner     3
2      21.01  3.50    Male    No  Sun  Dinner     3
3      23.68  3.31    Male    No  Sun  Dinner     2
4      24.59  3.61  Female    No  Sun  Dinner     4
```

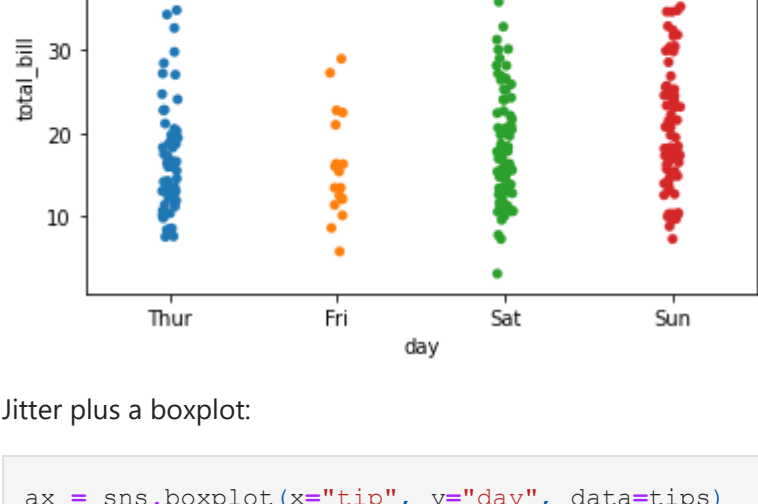
```
In [4]: sns.barplot(x = "smoker", y = "size", hue = "sex", data = tips)
```

```
Out[4]: <AxesSubplot: xlabel='smoker', ylabel='size'>
```



```
In [5]: sns.stripplot(x="day", y="total_bill", data=tips)
```

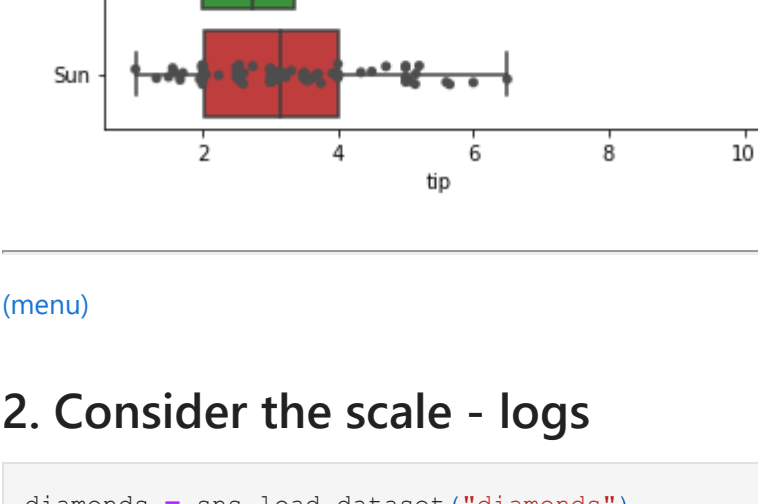
```
Out[5]: <AxesSubplot: xlabel='day', ylabel='total_bill'>
```



Use a smaller amount of jitter:

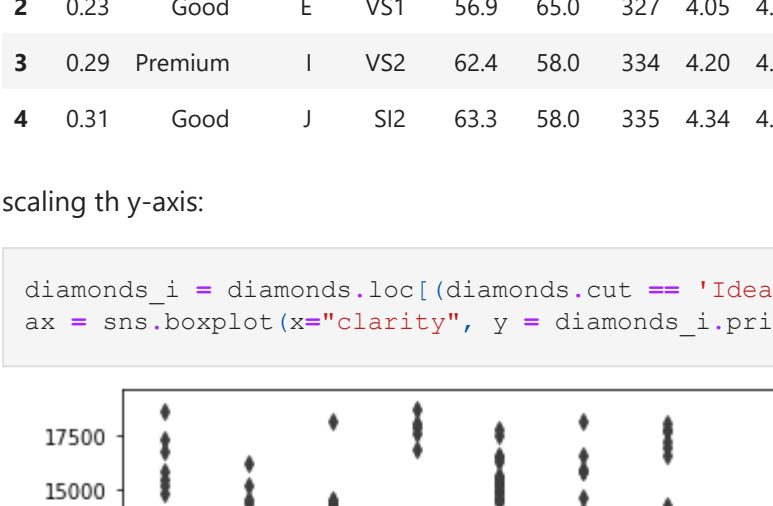
```
In [6]: sns.stripplot(x="day", y="total_bill", data=tips, jitter=0.05)
```

```
Out[6]: <AxesSubplot: xlabel='day', ylabel='total_bill'>
```



Jitter plus a boxplot:

```
In [7]: ax = sns.boxplot(x="tip", y="day", data=tips)
ax = sns.stripplot(x="tip", y="day", data=tips, color=".3")
```



(menu)

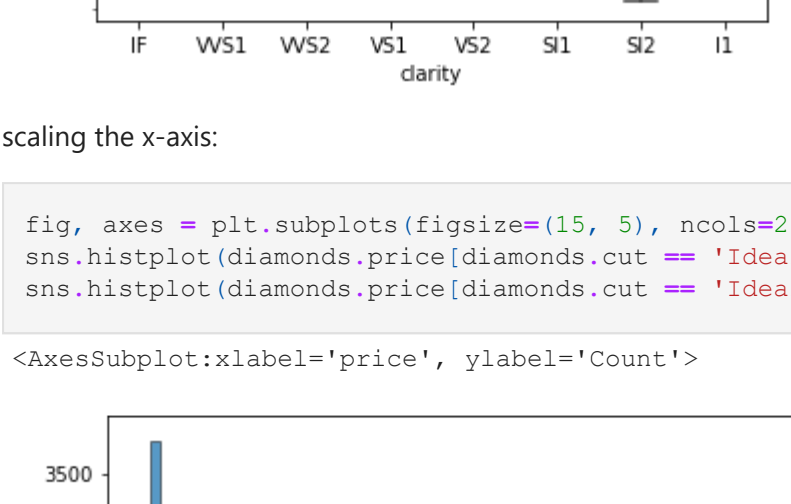
2. Consider the scale - logs

```
In [8]: diamonds = sns.load_dataset("diamonds")
diamonds.head()
```

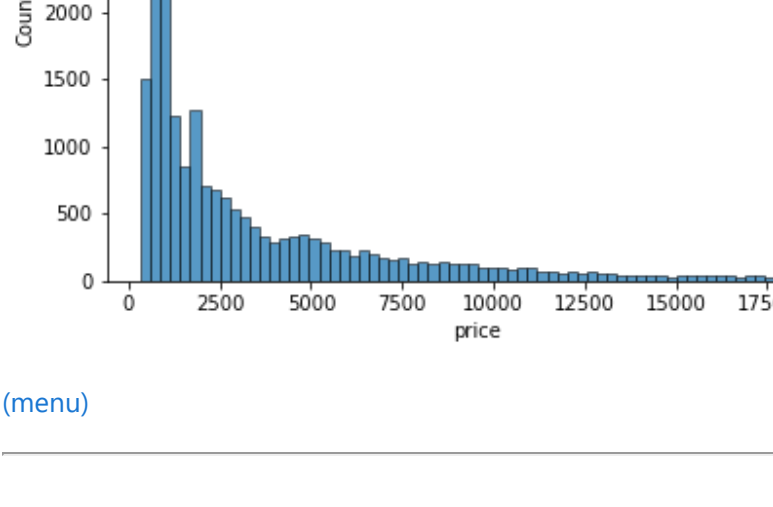
```
Out[8]:   carat    cut  color  clarity  depth  table  price     x     y     z
0   0.23  Ideal     E   SI2     61.5   55.0   326   3.95   3.98   2.43
1   0.21  Premium     E   SI1     59.8   61.0   312   3.89   3.84   2.31
2   0.23   Good     E   VS1     56.9   65.0   327   4.05   4.07   2.31
3   0.29  Premium     I   VS2     62.4   58.0   334   4.20   4.23   2.63
4   0.31   Good     J   SI2     63.3   58.0   335   4.34   4.35   2.75
```

scaling the y-axis:

```
In [9]: diamonds_i = diamonds.loc[(diamonds.cut == 'Ideal') & (diamonds.color == 'E')]
ax = sns.boxplot(x="clarity", y=diamonds_i.price, data=diamonds_i)
```



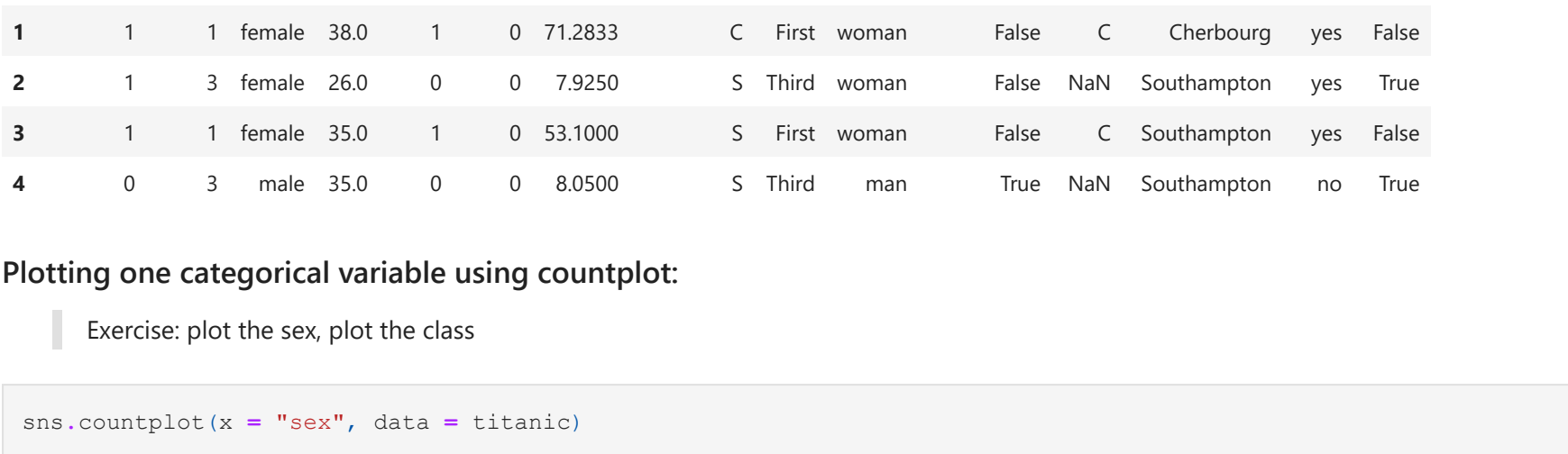
```
In [10]: diamonds_i = diamonds.loc[(diamonds.cut == 'Ideal') & (diamonds.color == 'E')]
ax = sns.boxplot(x="clarity", y=diamonds_i.price, data=diamonds_i)
ax.set_yscale("log")
plt.show()
```



scaling the x-axis:

```
In [11]: fig, axes = plt.subplots(figsize=(15, 5), ncols=2)
sns.histplot(diamonds.price[diamonds.cut == 'Ideal'], ax = axes[0])
sns.histplot(diamonds.price[diamonds.cut == 'Ideal'], log_scale = True, ax = axes[1])
```

```
Out[11]: <AxesSubplot: xlabel='price', ylabel='Count'>
```



(menu)

3. Baselines

Let look at a case where we have 2 binary variables: 'sex' and 'survived'

```
In [12]: titanic = sns.load_dataset("titanic")
titanic.head()
```

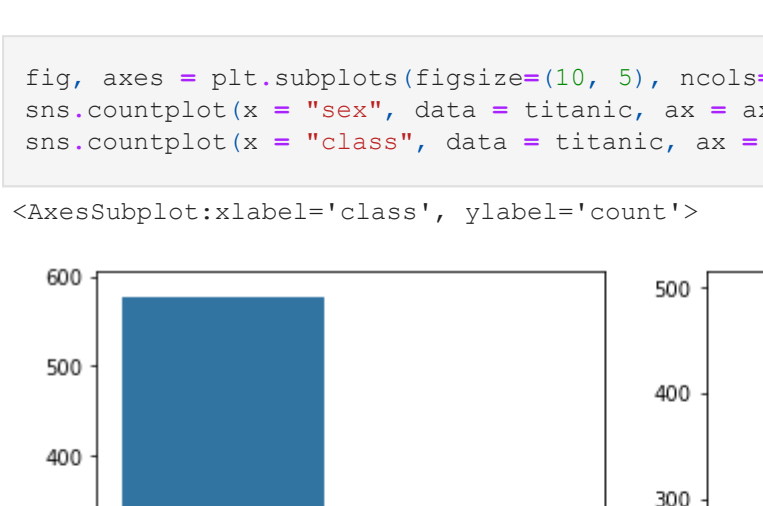
```
Out[12]:   survived  pclass    sex  age  sibsp  parch    fare  embarked  class  who  adult_male  deck  embark_town  alive  alone
0         0         3    male  22.0      1      0   7.2500         S  Third   man           True   NaN  Southampton    no   False
1         1         1  female  38.0      0      0  53.1000         C  First  woman        False   C    Cherbourg    yes   True
2         1         3  female  26.0      0      0   7.9250         S  Third  woman        False  NaN  Southampton    yes   True
3         1         1  female  35.0      1      0  53.1000         S  First  woman        False   C    Southampton    yes   False
4         0         3    male  35.0      0      0  8.0500         S  Third   man           True  NaN  Southampton    no   True
```

Plotting one categorical variable using countplot:

Exercise: plot the sex, plot the class

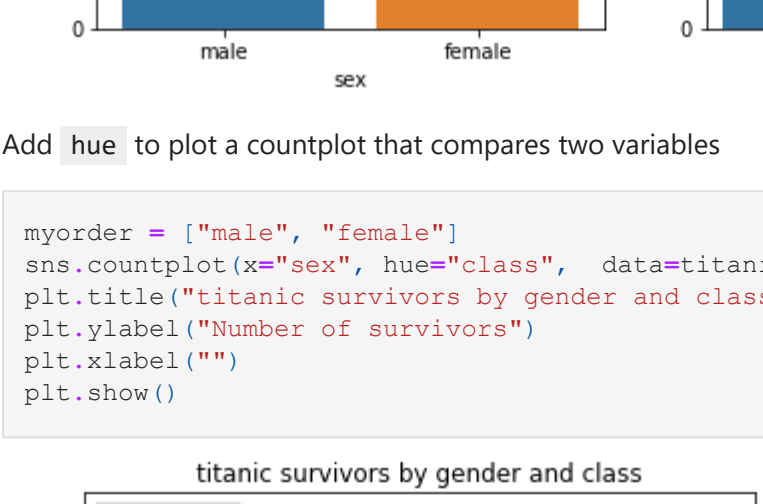
```
In [13]: sns.countplot(x = "sex", data = titanic)
```

```
Out[13]: <AxesSubplot: xlabel='sex', ylabel='count'>
```



```
In [14]: sns.countplot(x = "class", data = titanic)
```

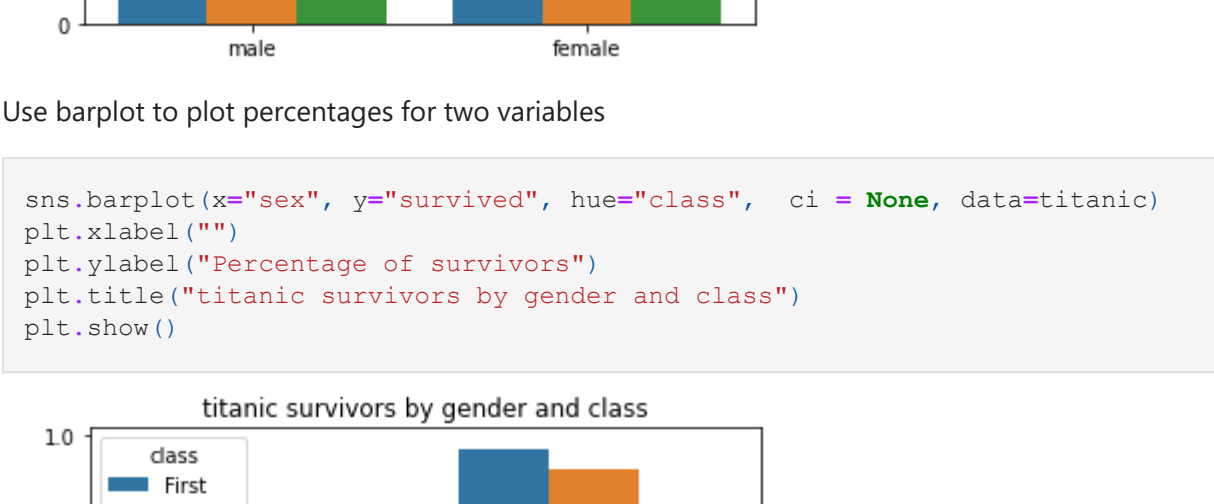
```
Out[14]: <AxesSubplot: xlabel='class', ylabel='count'>
```



advanced: plot them side by side

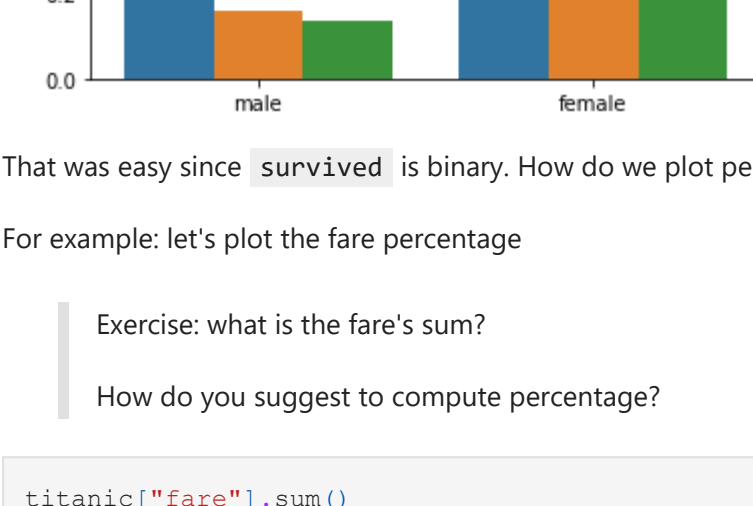
```
In [15]: fig, axes = plt.subplots(figsize=(10, 5), ncols=2)
sns.countplot(x = "sex", data = titanic, ax = axes[0])
sns.countplot(x = "class", data = titanic, ax = axes[1])
```

```
Out[15]: <AxesSubplot: xlabel='class', ylabel='count'>
```



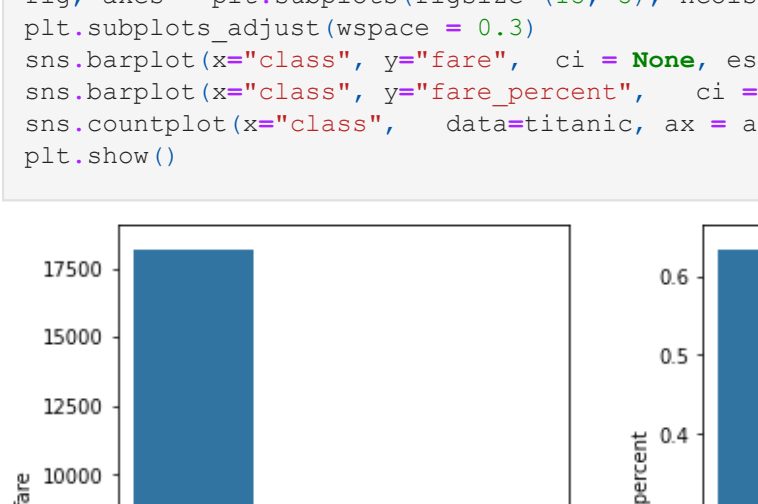
Add 'hue' to plot a countplot that compares two variables

```
In [16]: myorder = ["male", "female"]
sns.countplot(x="sex", hue="class", data=titanic.loc[titanic["survived"] == 1], order = myorder)
plt.title("titanic survivors by gender and class")
plt.xlabel("number of survivors")
plt.ylabel("")
plt.show()
```



Use barplot to plot percentages for two variables

```
In [17]: sns.barplot(x="sex", y="survived", hue="class", ci = None, data=titanic)
plt.xlabel("")
plt.ylabel("Percentage of survivors")
plt.title("titanic survivors by gender and class")
plt.show()
```



That was easy since 'survived' is binary. How do we plot percentages for non-binary data?

For example: let's plot the fare percentage

Exercise: what is the fare's sum?

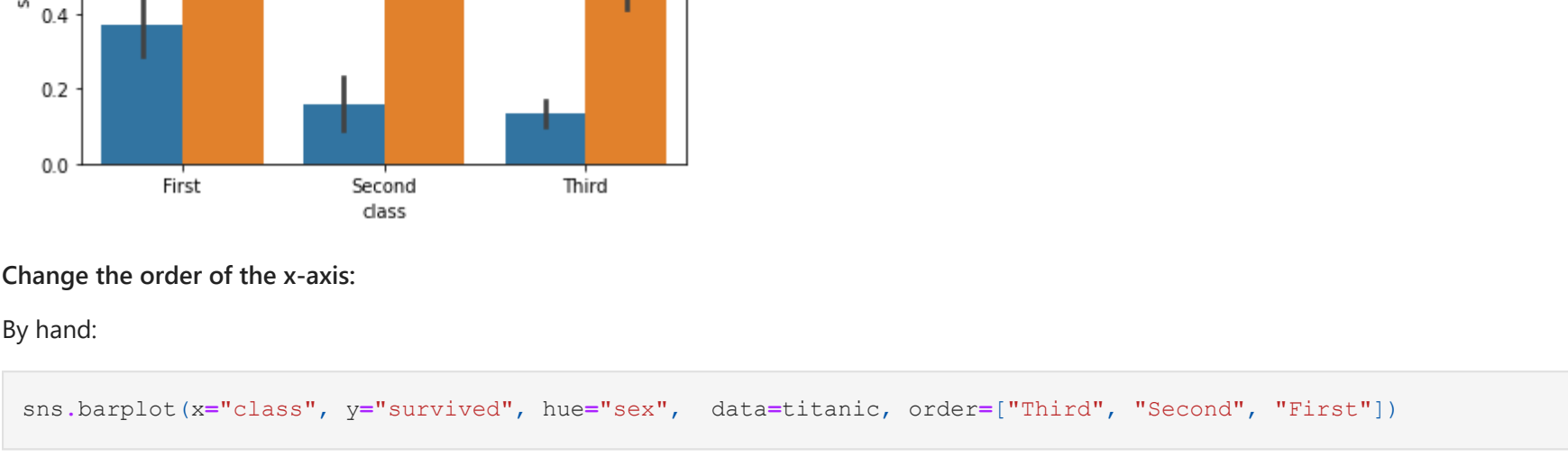
How do you suggest to compute percentage?

```
In [18]: titanic["fare"].sum()
```

```
Out[18]: 28693.9493
```

```
In [19]: titanic["fare_percent"] = titanic["fare"]/titanic["fare"].sum()
```

```
In [20]: fig, axes = plt.subplots(figsize=(15, 5), ncols=3)
plt.subplots_adjust(wspace = 0.3)
sns.countplot(x="sex", data=titanic, ax=axes[0])
sns.barplot(x="class", y="fare", ci = None, estimator = sum, data=titanic, ax=axes[1])
sns.barplot(x="class", y="fare_percent", ci = None, estimator = sum, data=titanic, ax = axes[2])
plt.show()
```

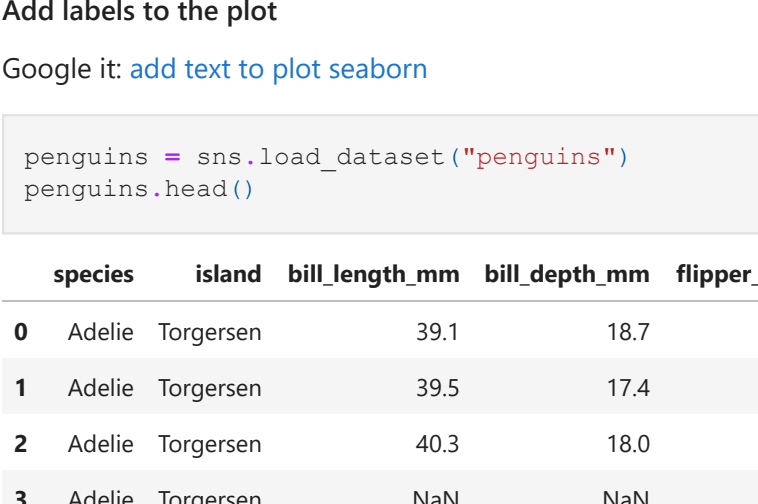


(menu)

4. Ease comparisons

Note what happens when we swap the hue and the x-axis. You can also remove the hue altogether (check what happens).

```
In [21]: sns.barplot(x="class", y="survived", hue="sex", data=titanic)
plt.show()
```

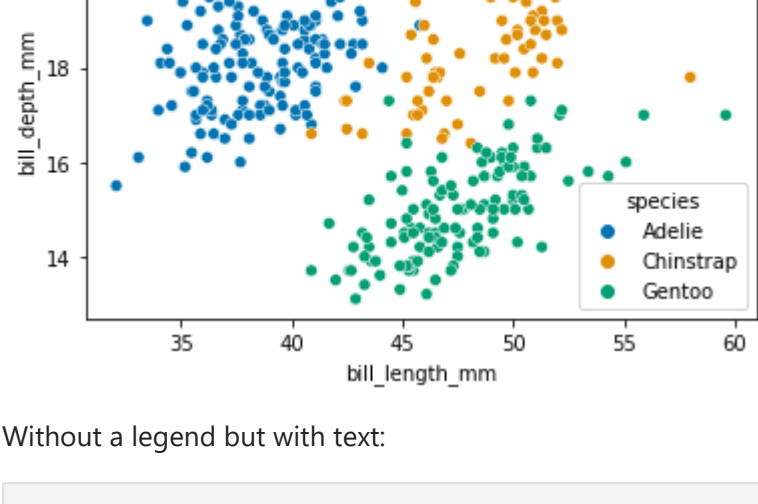


Change the order of the x-axis:

By hand:

```
In [22]: sns.barplot(x="class", y="survived", hue="sex", data=titanic, order=["Third", "Second", "First"])
```

```
Out[22]: <AxesSubplot: xlabel='class', ylabel='survived'>
```



Add labels to the plot

Google it: [add text to plot seaborn](#)

```
In [23]: penguins = sns.load_dataset("penguins")
penguins.head()
```

```
Out[23]:   species  island  bill_length_mm  bill_depth_mm  flipper_length_mm  body_mass_g  sex
0  Adeline  Torgersen      39.1           18.7           181.0         3750.0  Male
1  Adeline  Torgersen      39.5           17.4           186.0         3800.0  Female
2  Adeline  Torgersen      40.3           18.0           195.0         3250.0  Female
3  Adeline  Torgersen      NaN            NaN            NaN            NaN    NaN
4  Adeline  Torgersen      36.7           19.3           193.0         3450.0  Female
```

Exercise: plot the relationship between 'bill_length_mm' and 'bill_depth_mm'

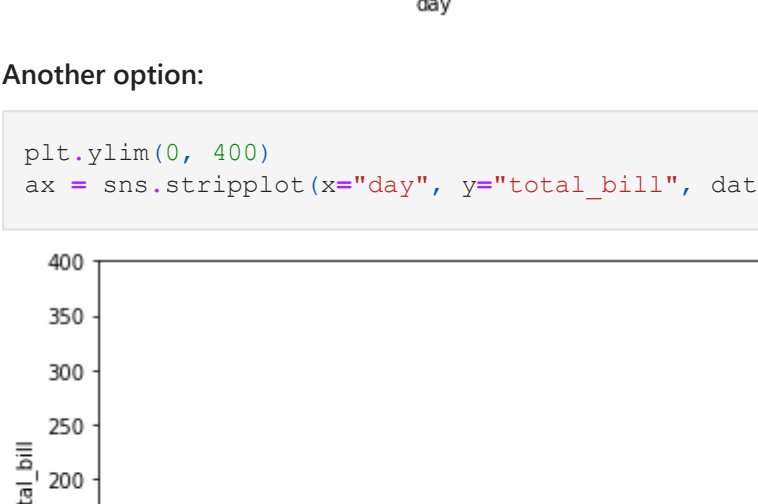
advanced: color the different species

```
In [24]: ax = sns.scatterplot(data=penguins, x="bill_length_mm", y="bill_depth_mm", hue = 'species', palette = 'colorblind')
```



Without a legend but with text:

```
In [25]: ax = sns.scatterplot(data=penguins, x="bill_length_mm", y="bill_depth_mm", hue = 'species', palette = 'colorblind')
style = dict(size=12, color='black')
ax.text(35, 15, "Adeline", **style)
ax.text(55, 20, "Chinstrap", **style)
ax.text(52, 14, "Gentoo", **style)
plt.show()
```



(menu)

5. Set the scale

Google it: [set scale seaborn](#)

One option:

```
In [26]: ax = sns.stripplot(x="day", y="total_bill", data=tips, jitter=0.05)
ax.set(ylim=(0, 100))
```

```
Out[26]: [(0.0, 100.0)]
```


Another option:

```
In [27]: plt.ylim(0, 400)
ax = sns.stripplot(x="day", y="total_bill", data=tips, jitter=0.05)
```