

Week 3

Exploratory Data Analysis

1. Descriptive stats
2. GroupBy
3. ANOVA
4. Correlation

1. Desc. statistics describe the basic features of the data (summary)

- `df.describe()` sum's stats (~~inc.~~ NaN's)
 - Shows count, mean, stdev, min, 25%, 50%, 75%, and max of data
- For categorical vars, `df[var].value_counts()` value counts can describe categorical data *
- To create a box plot, use `sns.boxplot(x=var, y=var2, data=df)`
- To create a scatterplot, first def the vars and use the matplotlib function `plt.scatter(x, y)`
 - `plt.title("Title"); plt.xlabel("var"); plt.ylabel("var2")`

2. GroupBy

With a scatterplot, you can find correlations b/c both vars are #s. What if you want to see how a categorical var affects a numerical var?

1. Assign a var to encompass any indep + dep vars:
`var0 = df['var1', 'var0']`
2. Group categorical indep vars with dep var:
`var1 = var0.groupby(['var11', 'var12'].as_index=False).mean()`
3. You can make a pivot table, which is easier to visualize:
`var2 = var1.pivot(index='var11', columns='var12')`

* To convert to a dataframe:

`df[var].valuecounts().to_frame()`

* To just find correlation, use:

```
df[['var1', 'var2']].corr()
```

3+4. Correlation

```
import matplotlib.pyplot as plt  
import seaborn as sns
```

↓ To plot a regression line with a scatterplot, use:

```
sns.regplot(x=var1, y=var2, data=df)  
plt.ylim(0,)
```



Stat Methods

and 'r'

- To take the p-value (certainty), use:

```
Pearson-coef, p-value = stats.ppersonr[['var1'], df[['var2']]]
```

- ANOVA

- Correlation bet diff groups of a categorical var

- F-test: $\frac{\text{var of means}}{\text{var of group}}$

- 3 lines of code:

1. `df_anova = df[['make', 'price']]` ← initialize cat's
2. `grouped_anova = df_anova.groupby(['make'])` ← group by dep. var.
3. `result = stats.f_oneway(grouped_anova.get_group('honda')['price'],
grouped_anova.get_group('subaru')['price'])`

```
1 import pandas as pd  
2 import numpy as np  
3 import matplotlib.pyplot as plt  
4 import seaborn as sns  
5 path = _____  
6 df = pd.read_csv(path)  
7 ...  
:
```