

IMD0033 - Probabilidade

Lesson 22 - Comparing Frequency Distributions

Ivanovitch Silva
November, 2018



Agenda

- Grouped bar plots
- Comparing histograms
- Kernel density estimates
- Strip plots
- Box plots
- Outliers

Atualizar o repositório

```
git clone https://github.com/ivanovitchm/imd0033_2018_2.git
```

Ou

```
git pull
```

PREVIOUSLY ON...

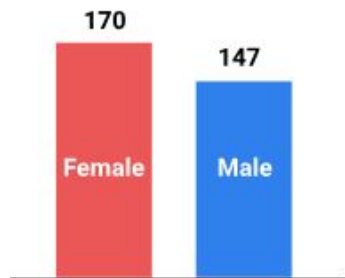


Getting **good** data



Id	Name	Salary	...	Gender
1	Mary Ann	\$35 000	...	Female
2	Marc Downey	\$55 000	...	Male
..
51	Juliet Ali	\$45 000	...	Female
...
317	Jane Ace	\$95 000	...	Female

Understand how the data is **structured** and **measured**



Visualize the patterns

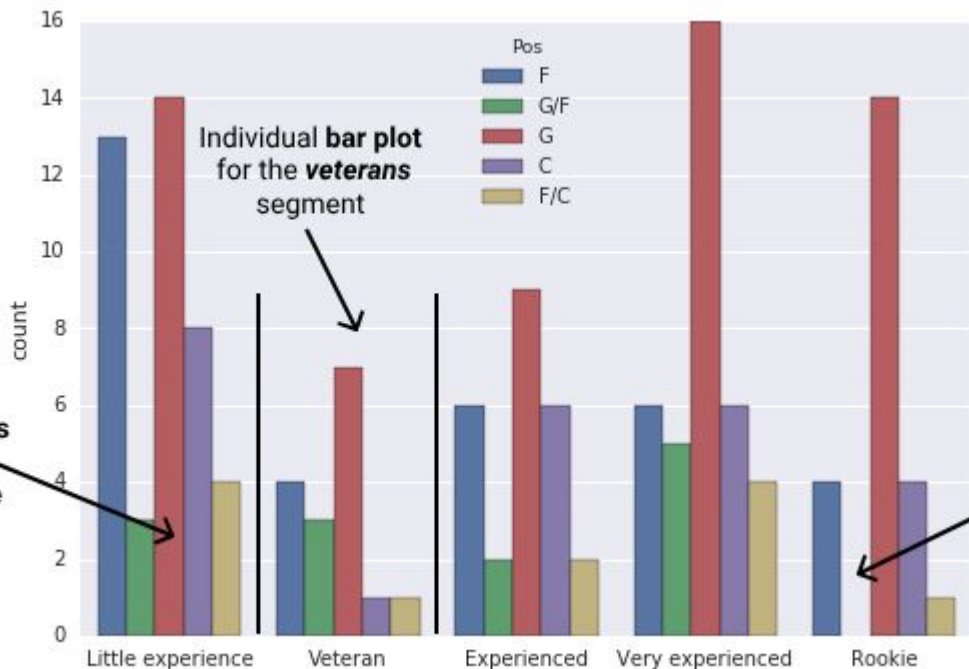
50 %

Gender	Frequency
Male	147
Female	170

Organize the data in **comprehensible** forms to find patterns

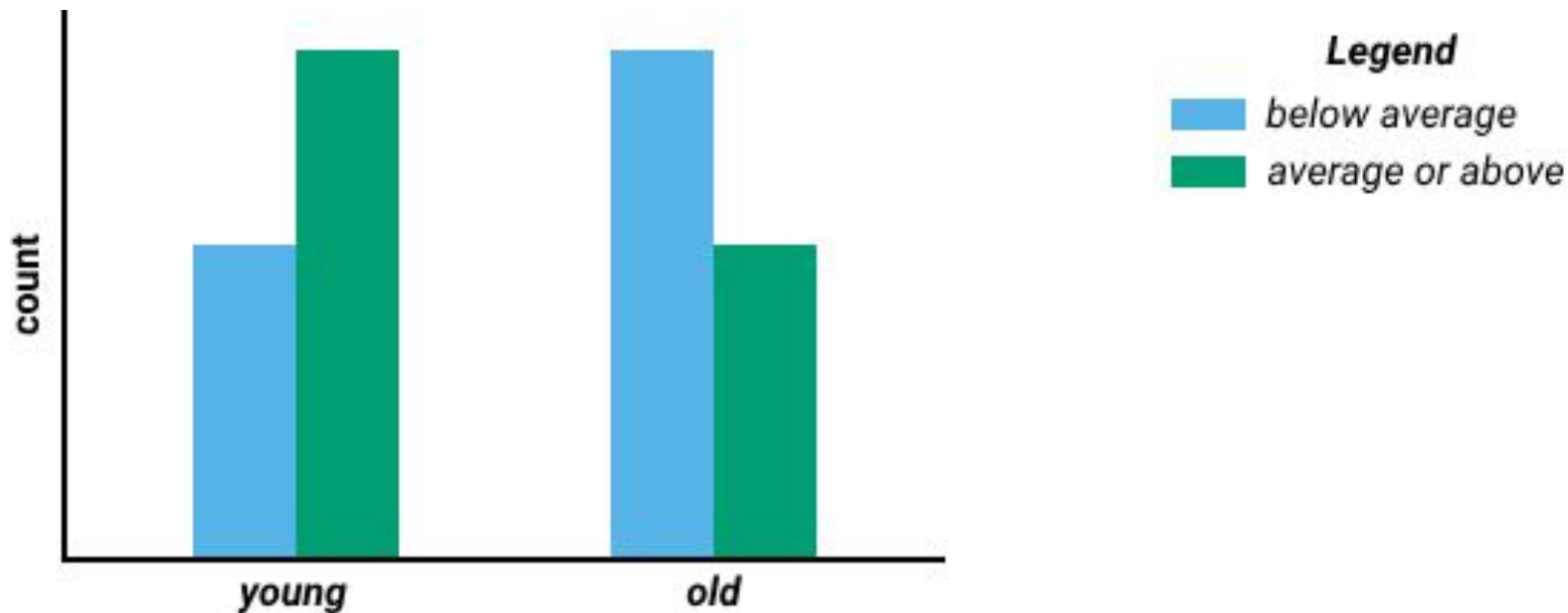


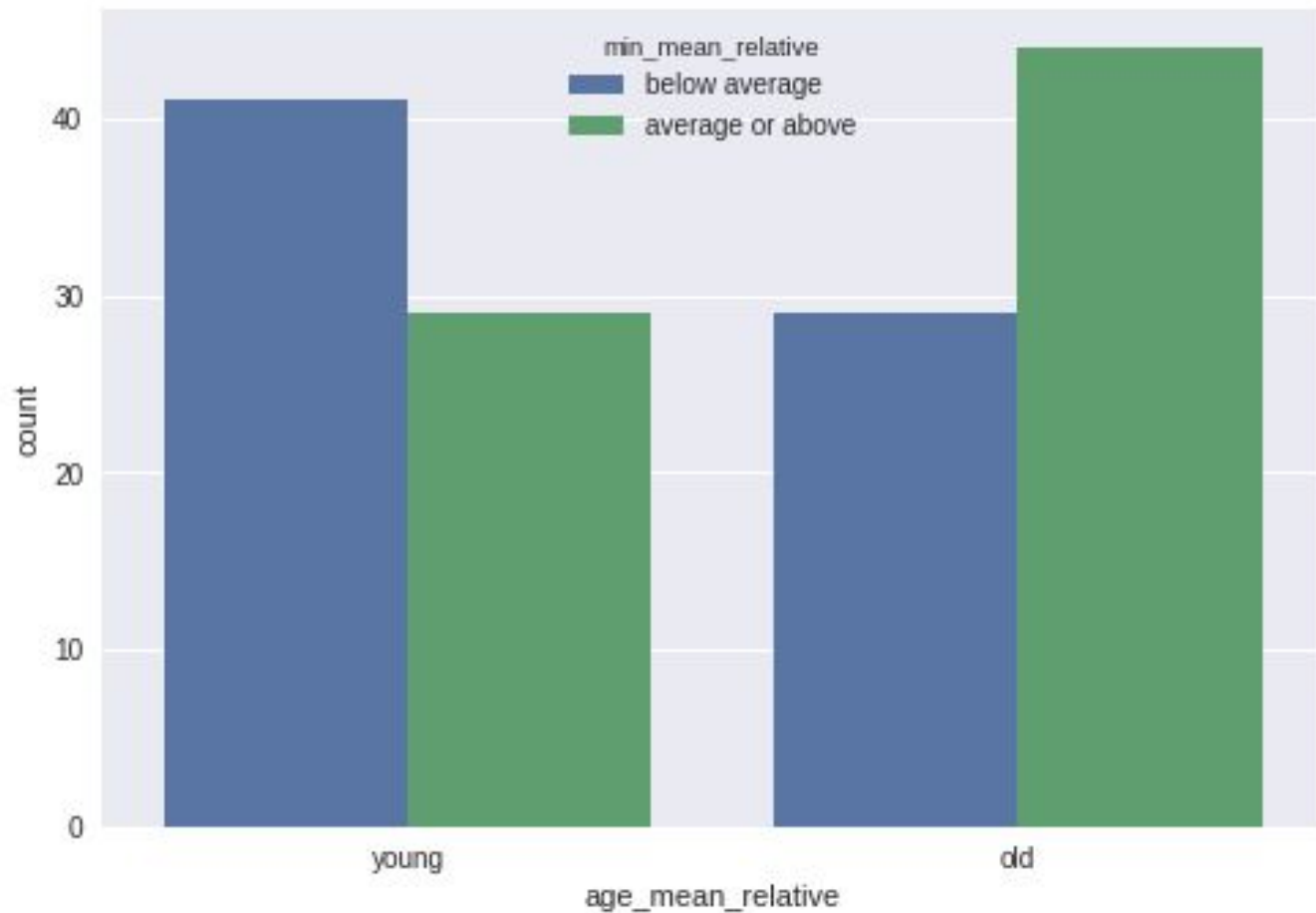
Comparing Frequency Distribution



Years in WNBA	Label
0	Rookie
1-3	Little experience
4-5	Experienced
5-10	Very experienced
>10	Veteran

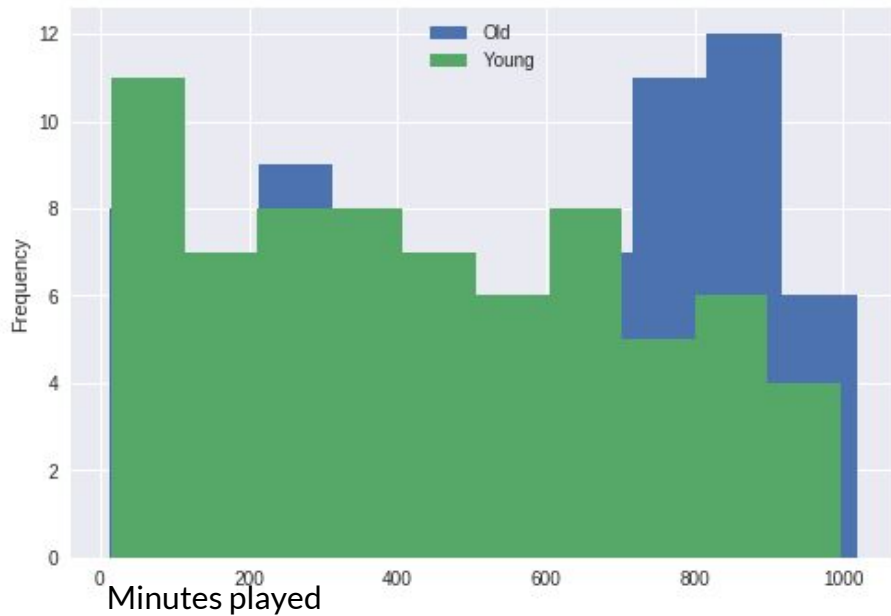
Challenge: Do older players play less?



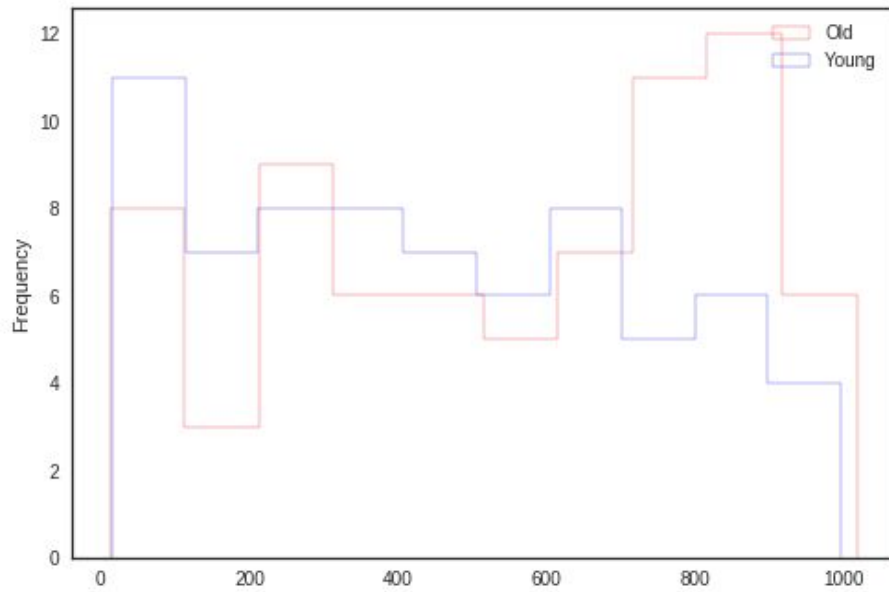


```
sns.countplot(x = 'age_mean_relative', hue = 'min_mean_relative', data = wnba)
```

Comparing Histograms

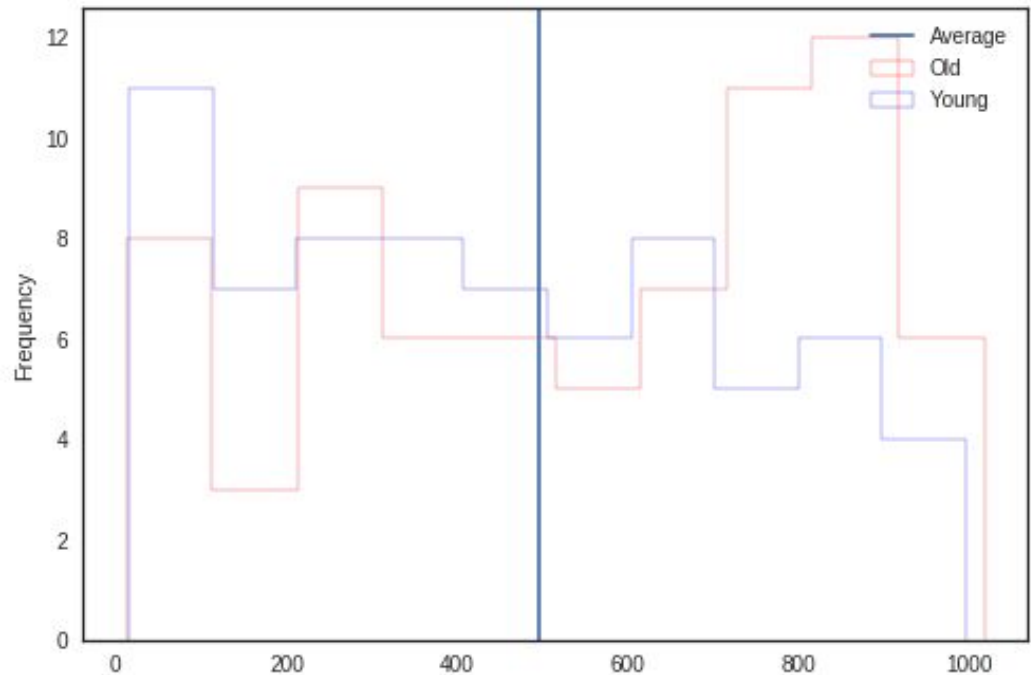


```
wnba[wnba.Age >= 27]['MIN'].plot.hist(label = 'Old', legend = True)
wnba[wnba.Age < 27]['MIN'].plot.hist(label = 'Young', legend = True)
```



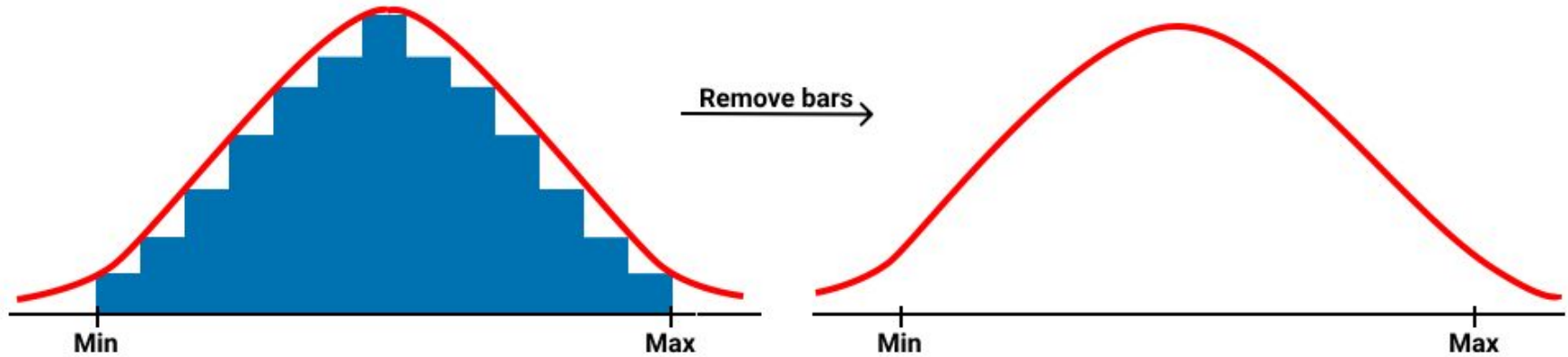
```
sns.set_style("white")
wnba[wnba.Age >= 27]['MIN'].plot.hist(histtype = 'step',
label = 'Old',
legend = True,color="red")
```


Comparing Histograms

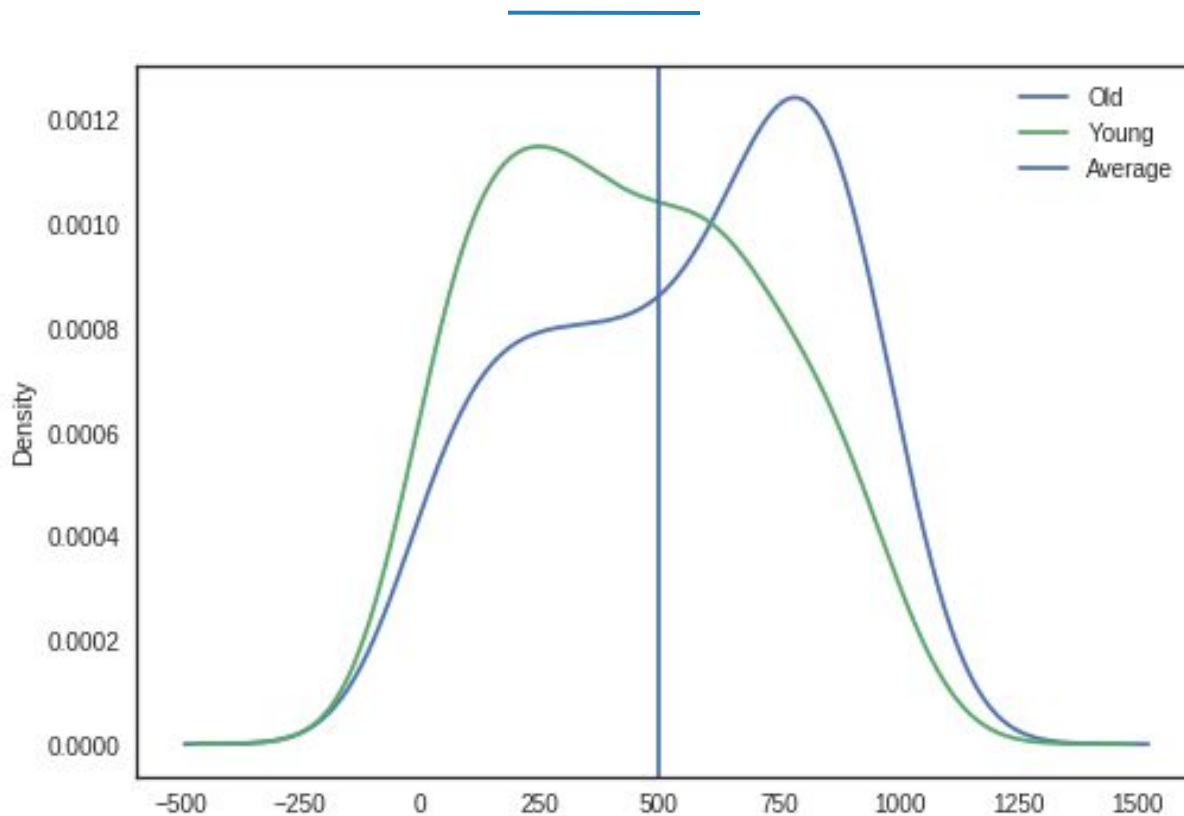


```
plt.axvline(497, label = 'Average')
```

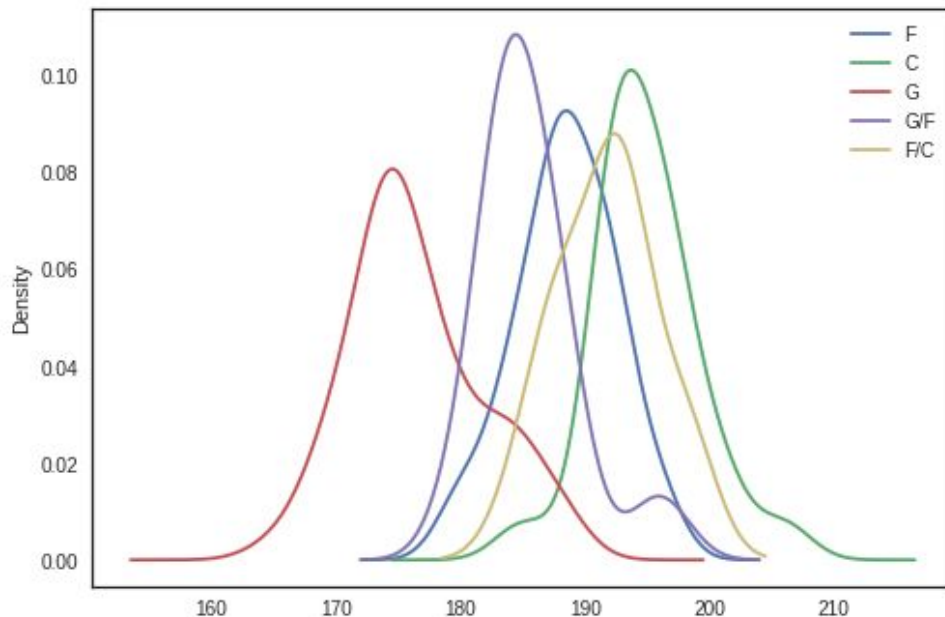
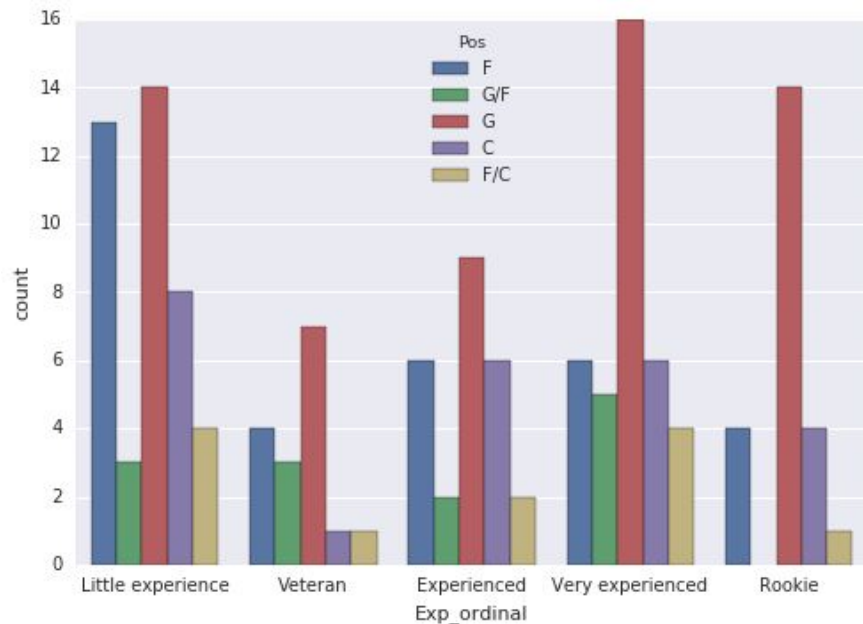
Kernel Density Estimate (KDE) Plots



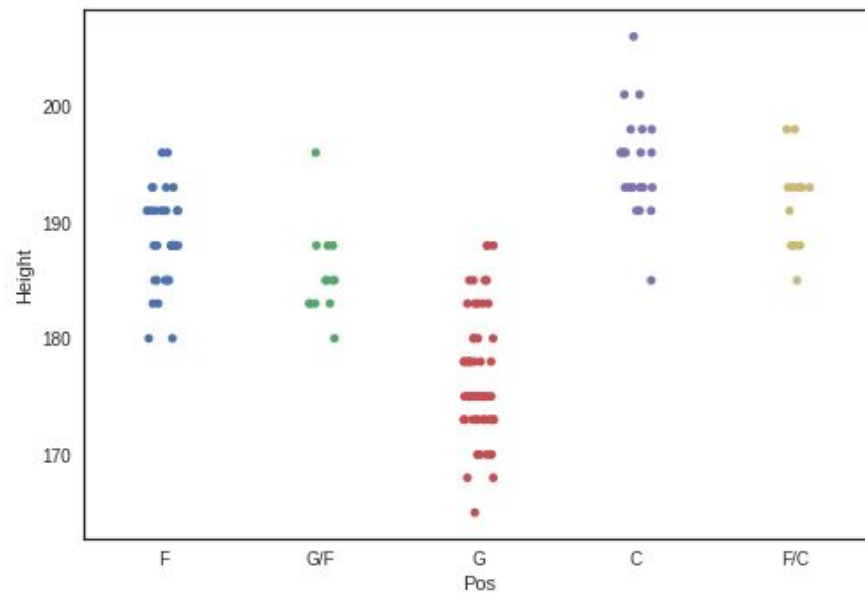
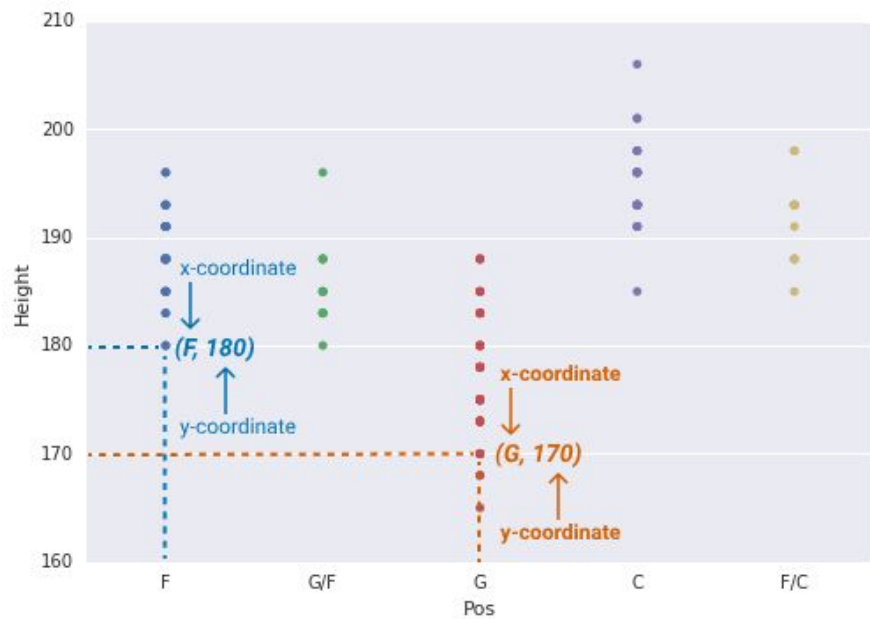
Kernel Density Estimate Plots



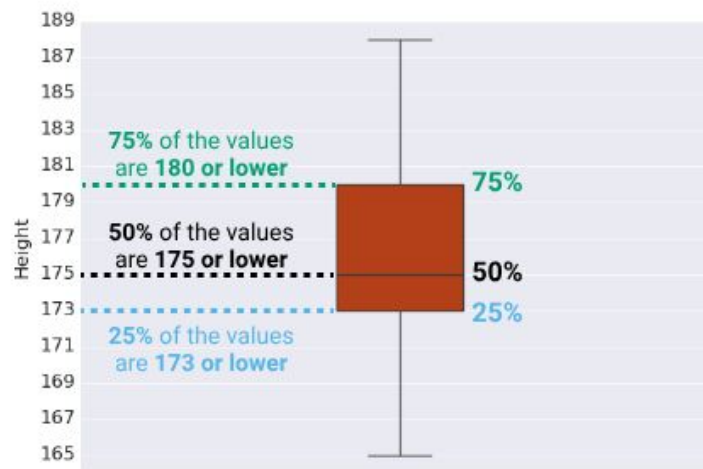
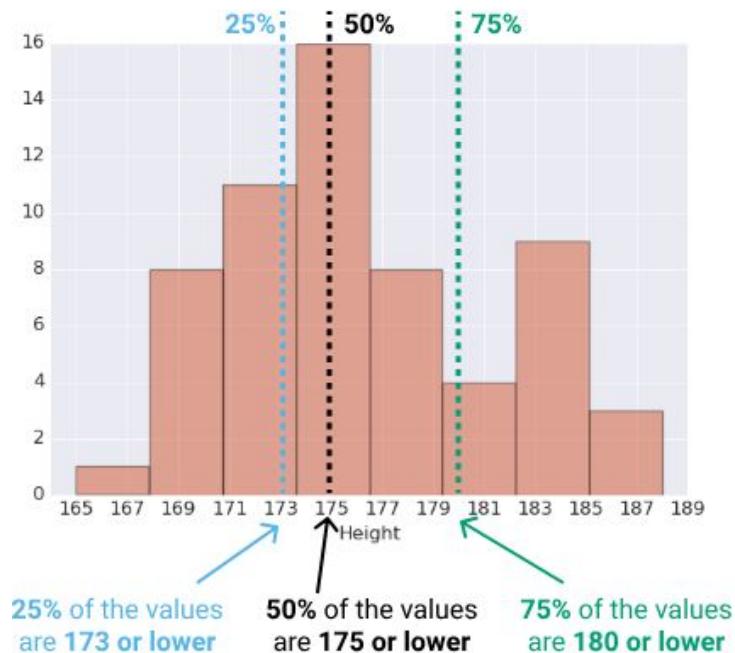
Drawbacks of Kernel Density Plots



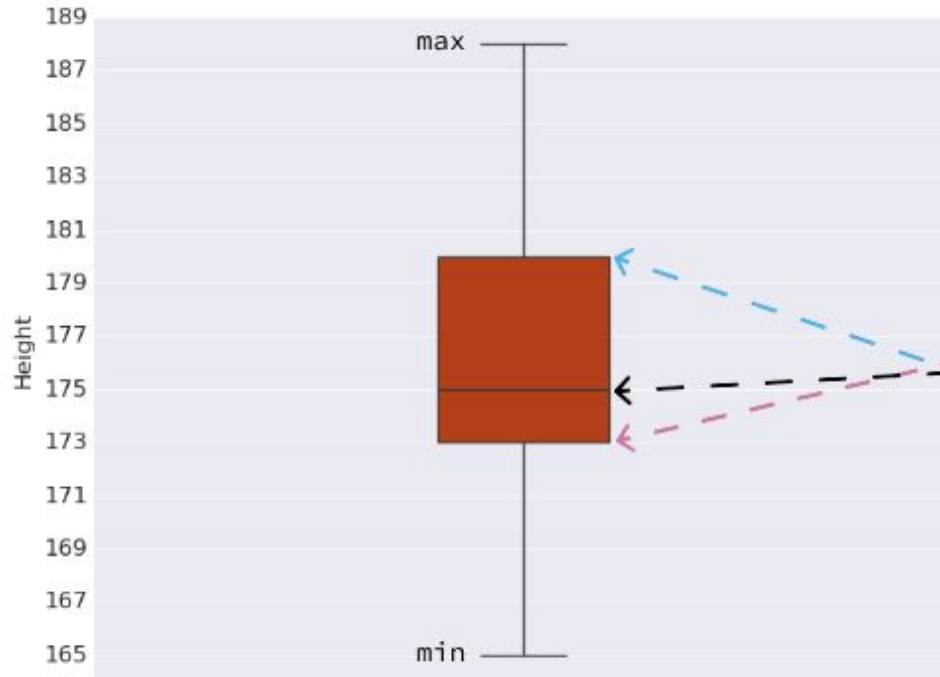
Strip Plots



Box Plots

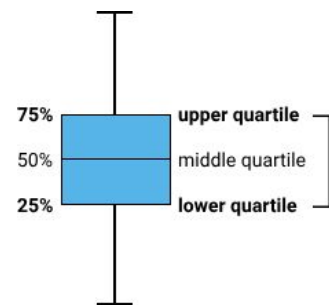
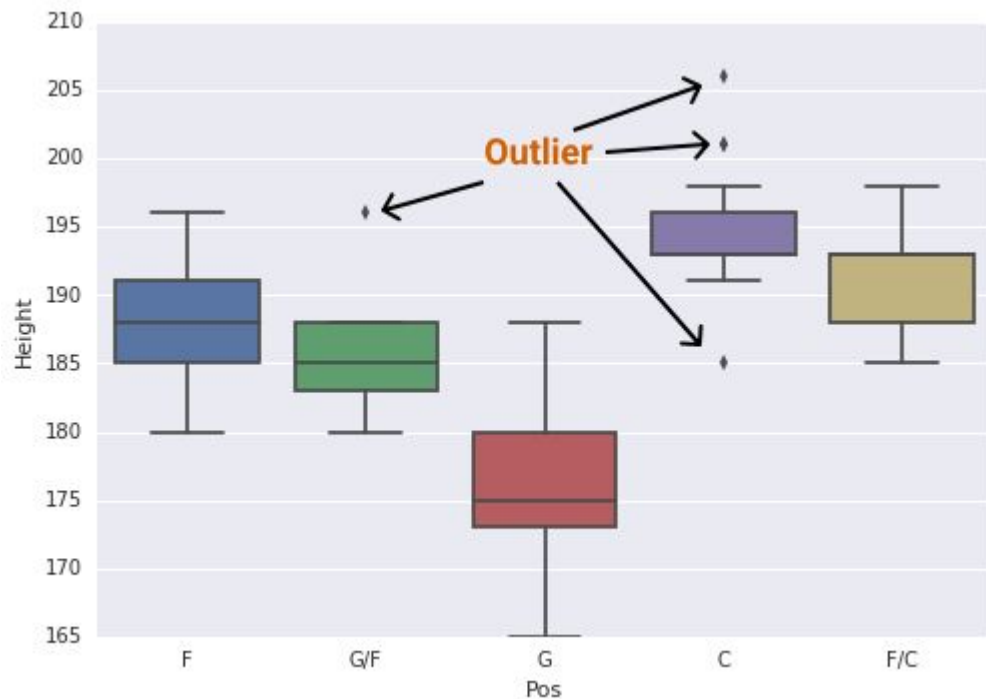


Box Plots



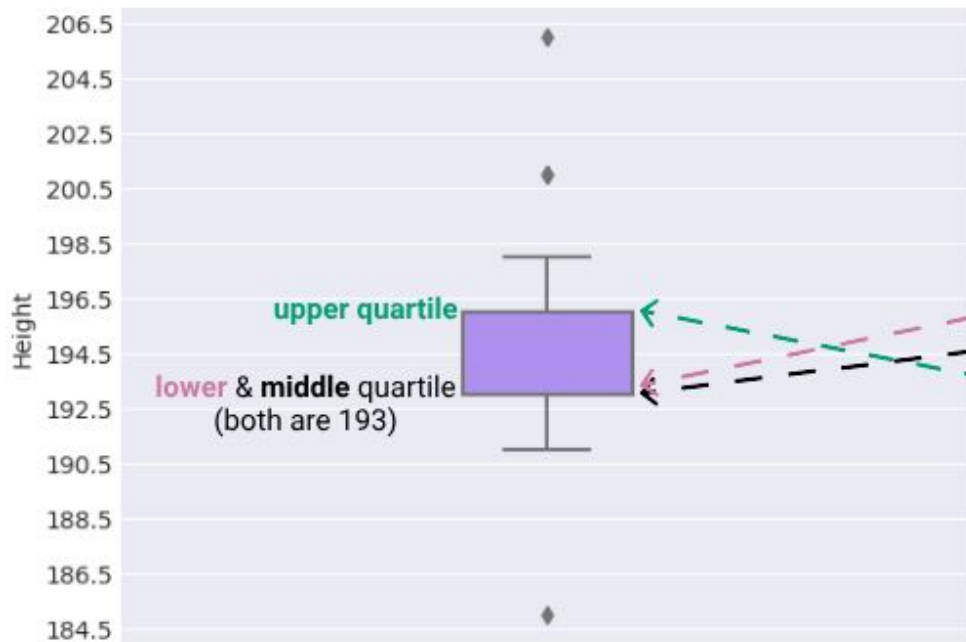
```
>> wnba[wnba['Pos'] == 'G']['Height'].describe()
count      60.000000
mean       176.450000
std         5.312043
min        165.000000
25%        173.000000
50%        175.000000
75%        180.000000
max        188.000000
Name: Height, dtype: float64
```

Outliers



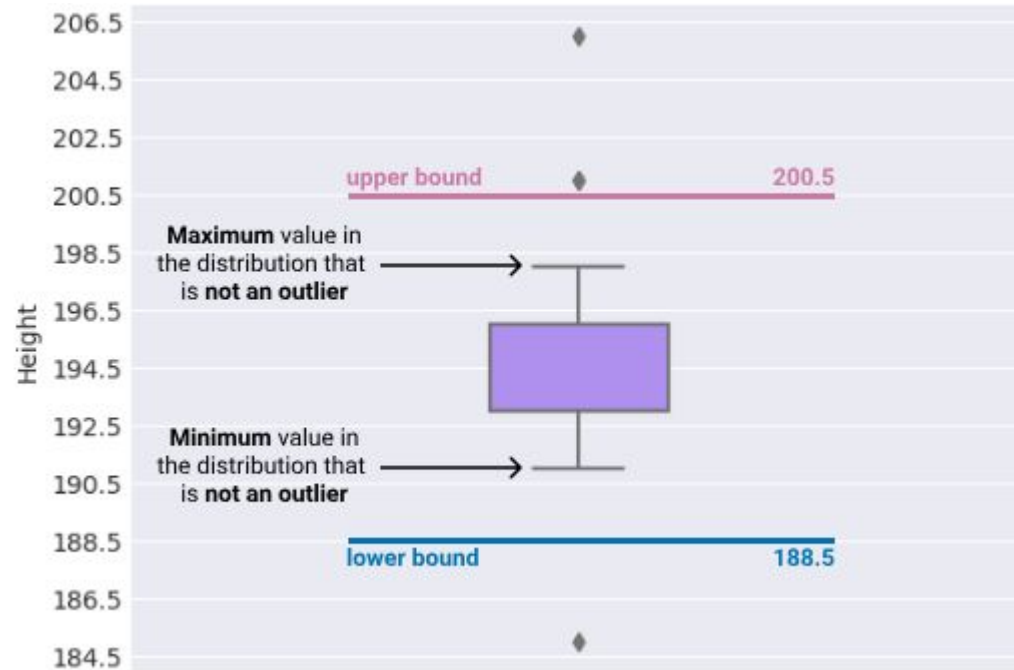
upper quartile - lower quartile = interquartile range

Outliers



```
>> wnba[wnba['Pos'] == 'C']['Height'].describe()
count      25.000000
mean       194.920000
std         4.132392
min        185.000000
25%        193.000000
50%        193.000000
75%        196.000000
max        206.000000
Name: Height, dtype: float64
```

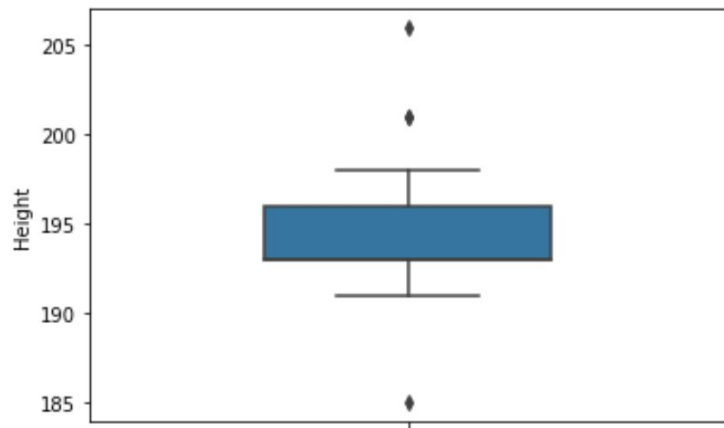
Outliers



Outliers

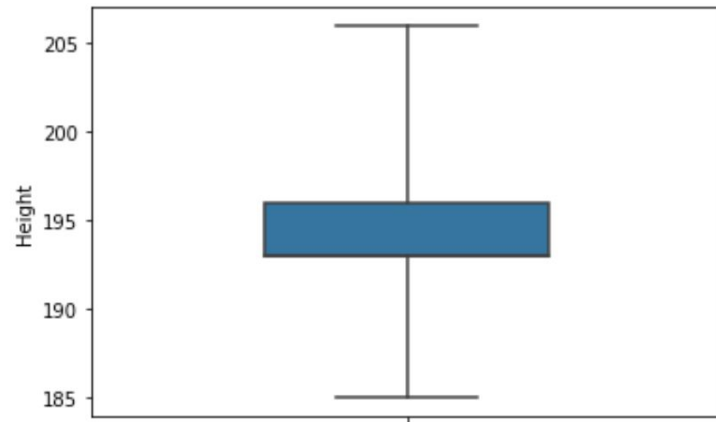
```
sns.boxplot(wnba[wnba['Pos'] == 'C']['Height'], whis = 1.5,  
            orient = 'vertical', width = .45)
```



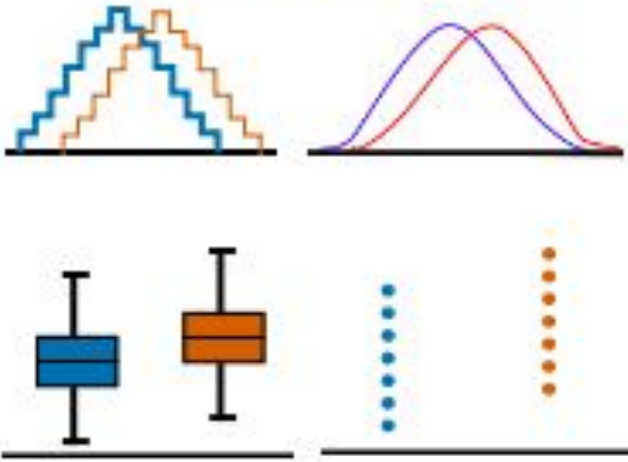
<matplotlib.axes._subplots.AxesSubplot at 0x1a180c4518>



```
sns.boxplot(wnba[wnba['Pos'] == 'C']['Height'], whis = 4,  
            orient = 'vertical', width = .45)
```

<matplotlib.axes._subplots.AxesSubplot at 0x1a18180208>



Scale of measurement	Graphs we can use to compare distributions
Nominal	 A bar chart with six groups of bars. Each group contains three bars of different colors (yellow, pink, blue). The heights of the bars vary across the groups, representing categorical data.
Ordinal	 A bar chart identical to the one for nominal data, showing six groups of three colored bars (yellow, pink, blue) with varying heights.
Interval & Ratio	 Two types of graphs are shown. The top row features two bell-shaped curves: a stepped histogram on the left (blue and orange) and two smooth normal distribution curves on the right (purple and red). The bottom row features two box plots on the left (blue and orange) and two vertical dot plots on the right (blue and orange).

