# Agenda (Parte I)

- Visualizando distribuições
- Gráficos de barra, pizza e histogramas
- Assimetria
- Distribuições simétricas

# Atualizar o repositório
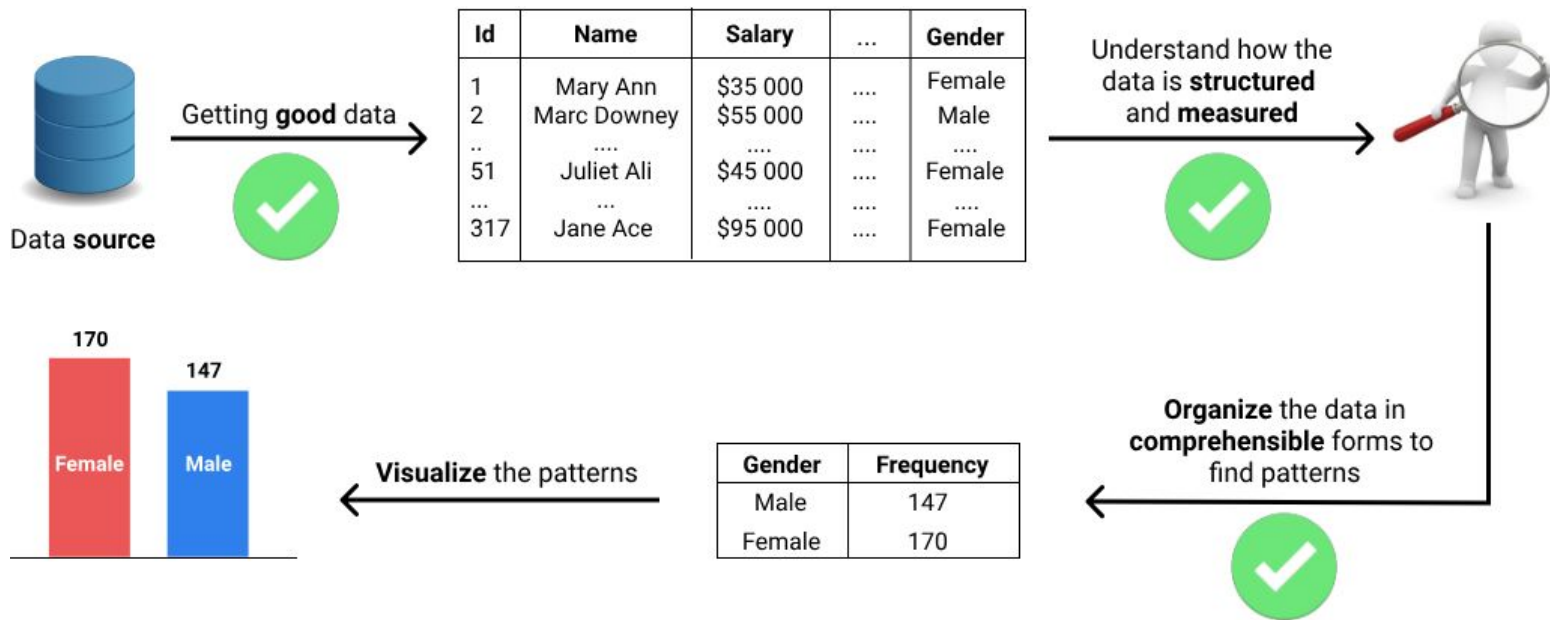
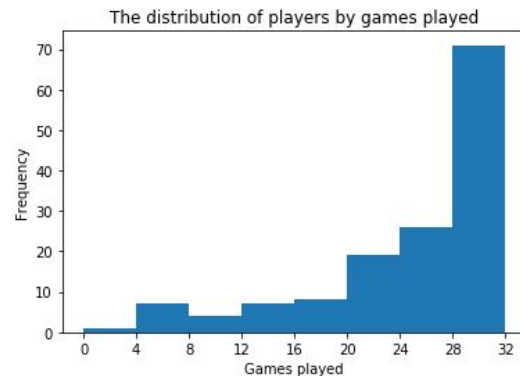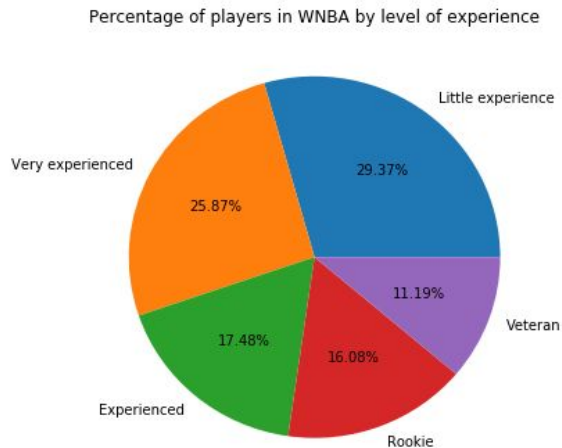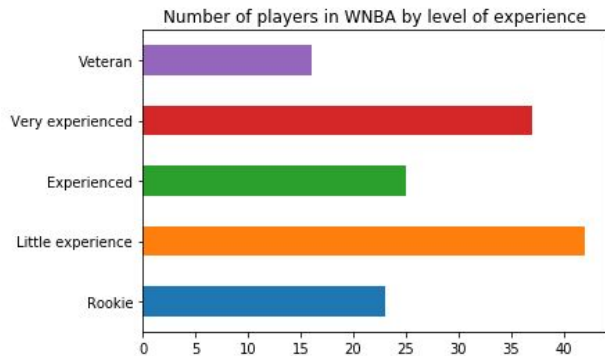git clone https://github.com/ivanovitchm/imd0033_2019_1.git

Ou ....

git pull

# PREVIOUSLY ON...



Data **source**

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**

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

**Visualize** the patterns

| Gender | Frequency |
|--------|-----------|
| Male | 147 |
| Female | 170 |

170 Female
147 Male

# Visualizing Distributions



Number of players in WNBA by level of experience



Percentage of players in WNBA by level of experience



The distribution of players by games played
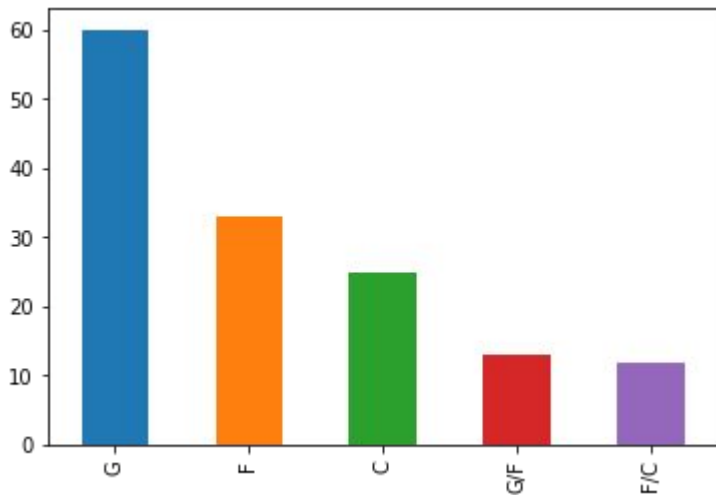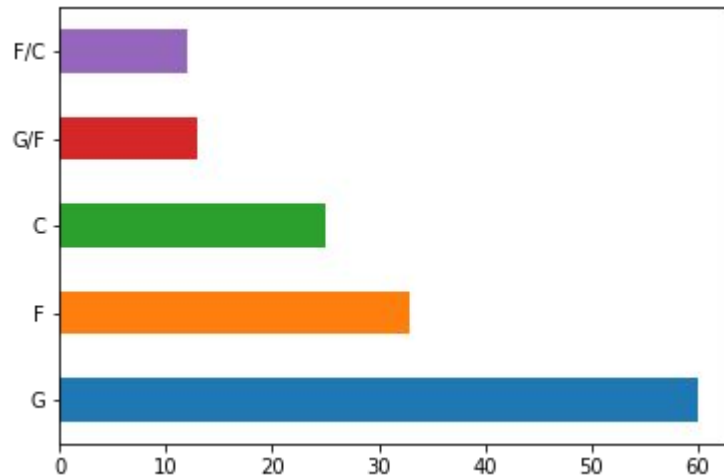
Graphs make easy to scan and compare frequencies, providing us with a single picture of the entire distribution of a variable (**nominal** or **ordinal scale**)

# Bar Plots

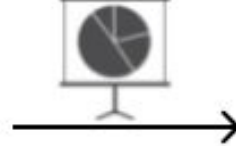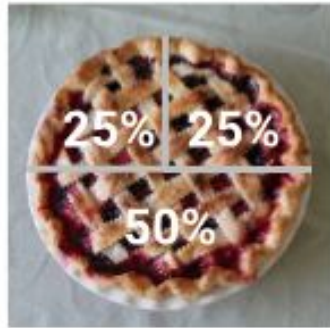horizontal bar plots are ideal to use when the labels of the unique values are long
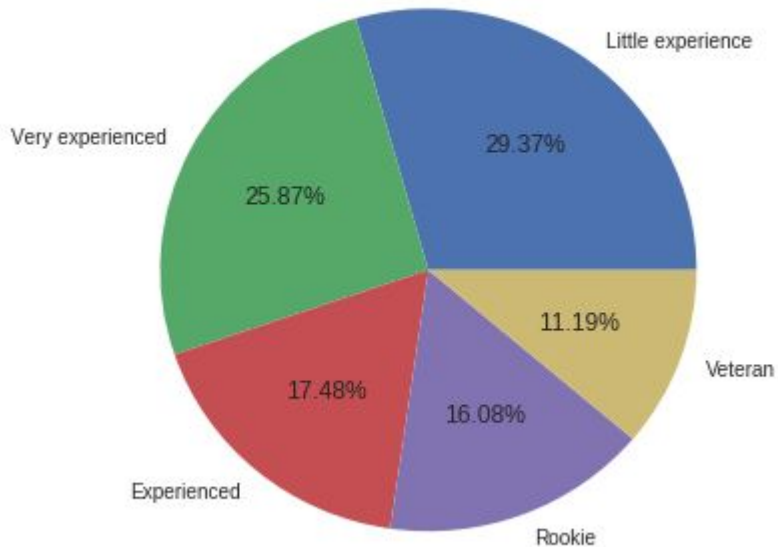


```
wnba['Pos'].value_counts().plot.bar()
```

```
wnba['Pos'].value_counts().plot.barh()
```
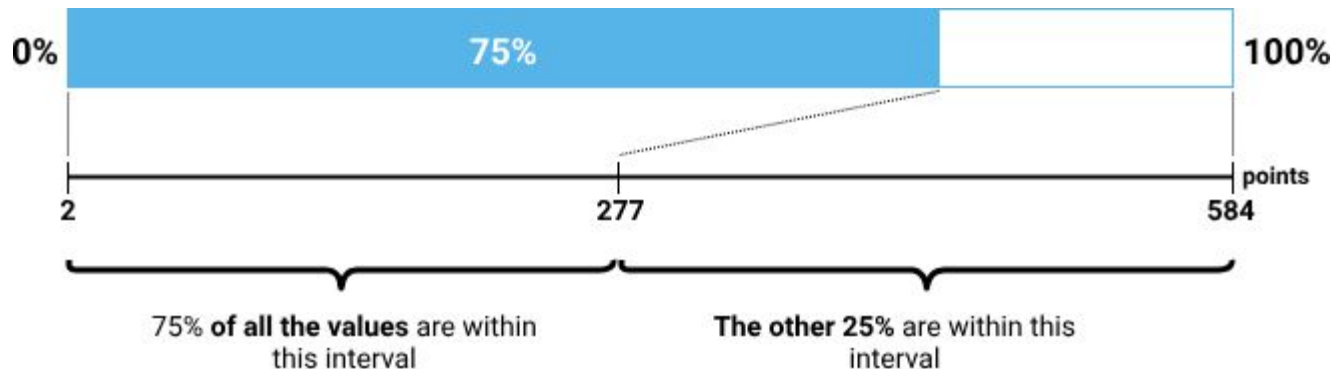
# Pie Charts

# Pie Charts

Percentage of players in WNBA by level of experience



```
wnba['Exp_ordinal'].value_counts().\
plot.pie(figsize = (6,6),
         autopct = '%.2f%%',
         title = 'Percentage of players in \
         WNBA by level of experience')
plt.ylabel('')
```
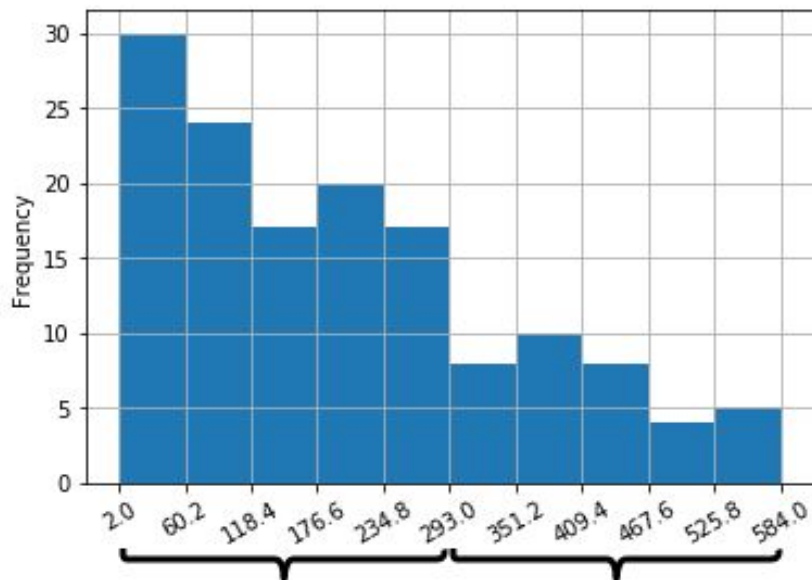
# Histograms



| | |
|---|---|
| 0% | 75% | 100% |

```
2                 277                 584     points
```

75% **of all the values** are within this interval

The other **25%** are within this interval

We can see that 75% of the values are distributed within a relatively narrow interval (between 2 and 277), while the remaining 25% are distributed in an interval that's slightly larger.

```
>> wnba['PTS'].describe()
count        143.000000
mean         201.790210
std          153.381548
min            2.000000
25%           75.000000
50%          177.000000
75%          277.500000
max          584.000000
```
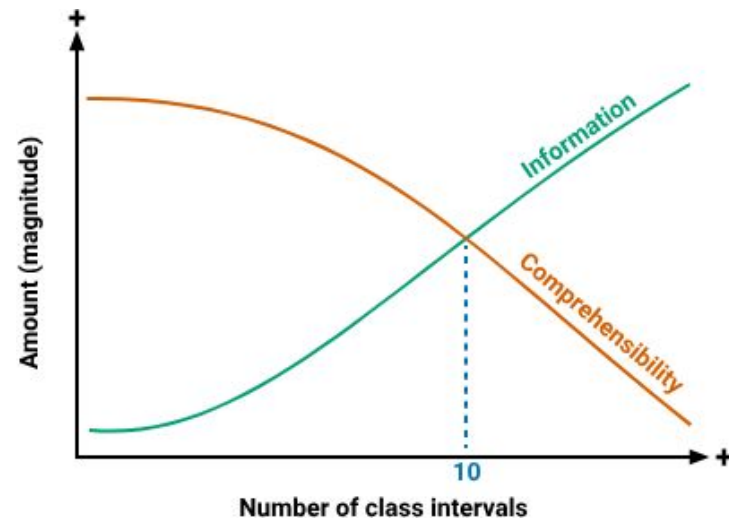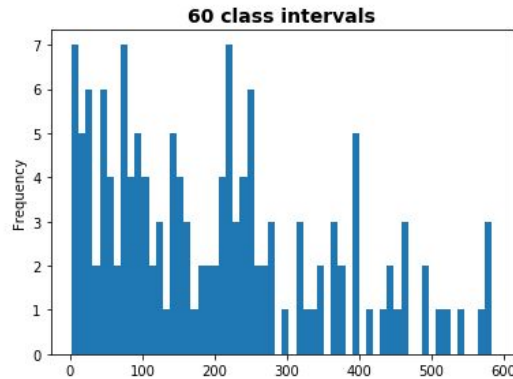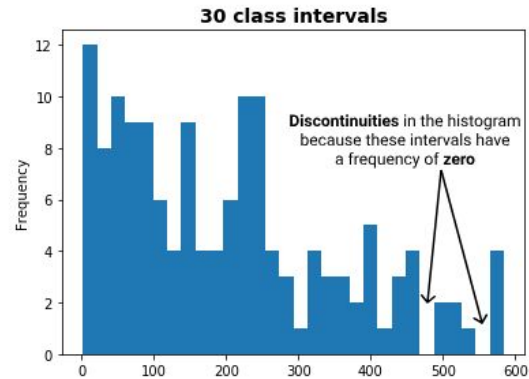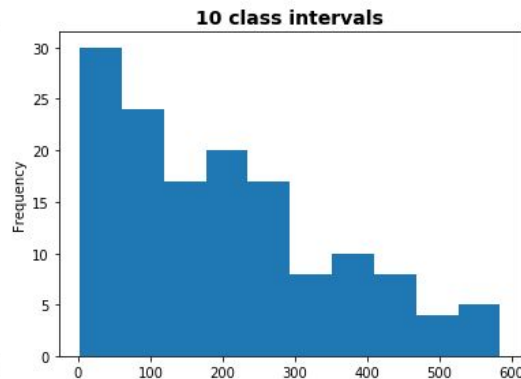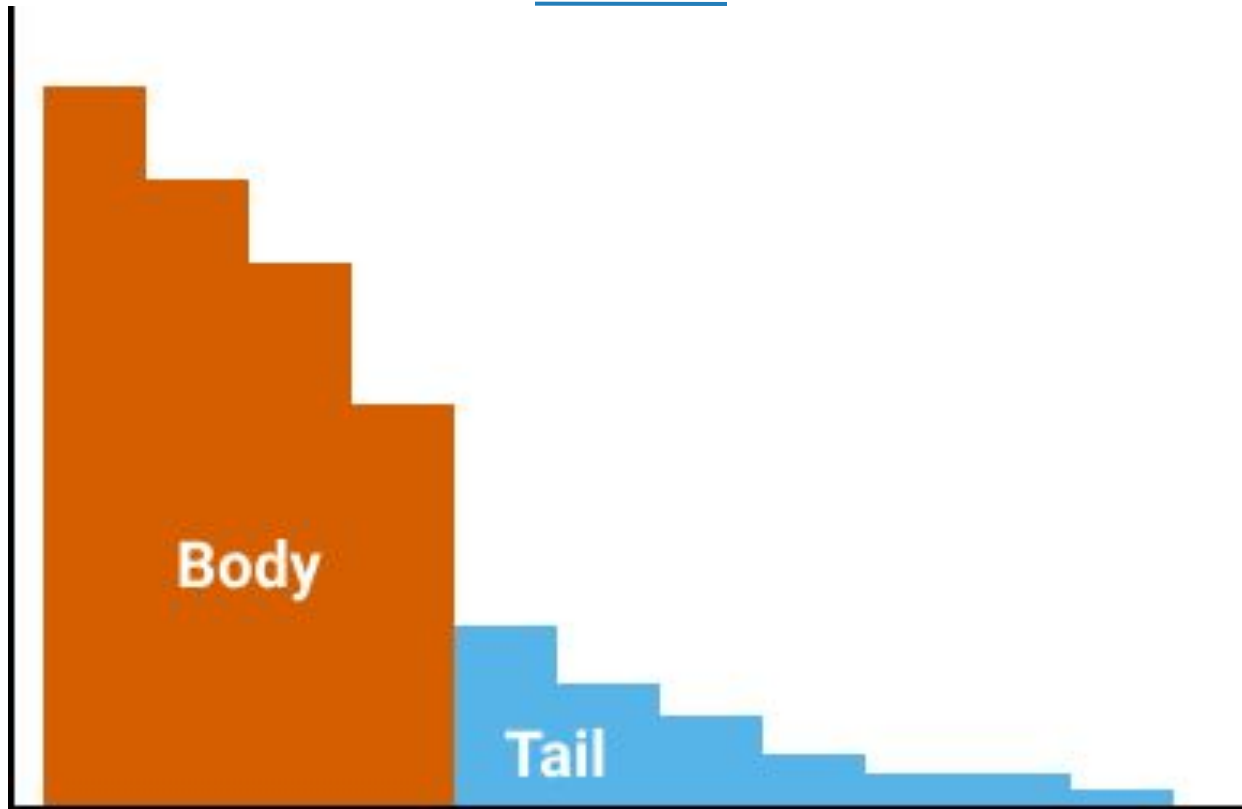
# The Statistics Behind Histograms



```
>> wnba['PTS'].describe()
count      143.000000
mean       201.790210
std        153.381548
min          2.000000
25%         75.000000
50%        177.000000
75%        277.500000
max        584.000000
Name: PTS, dtype: float64
```

```
>> wnba['PTS'].plot.hist()
```
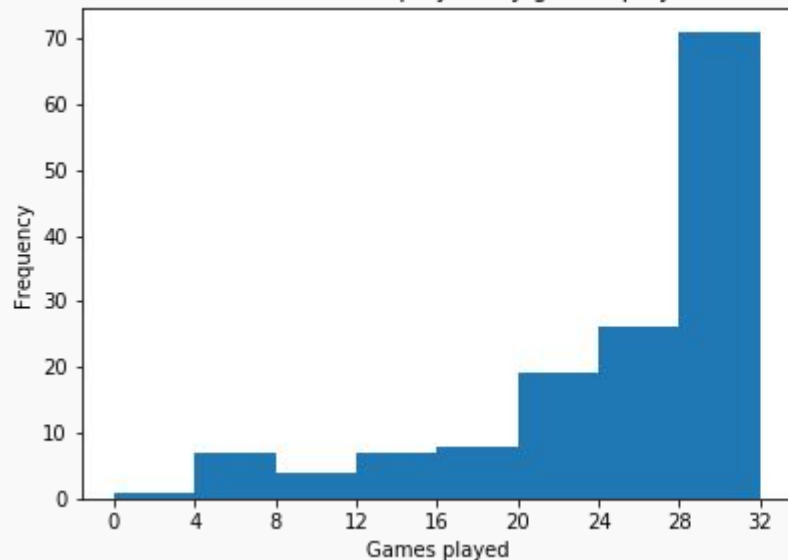
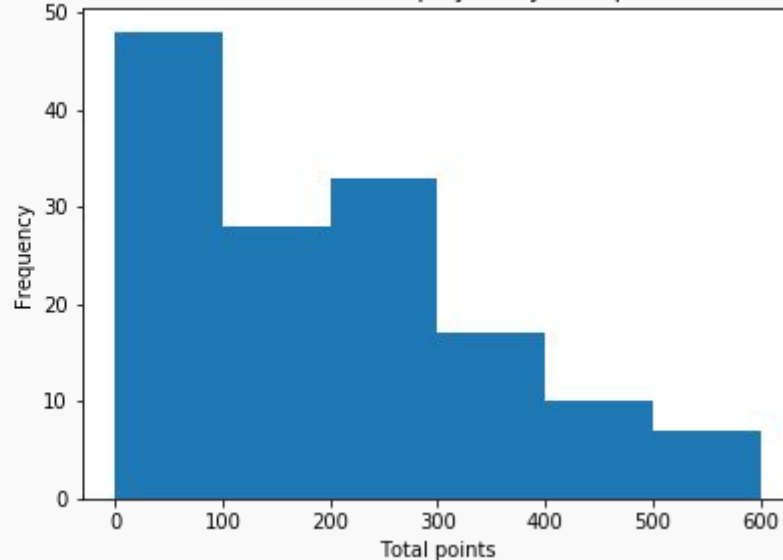# Binning for Histograms
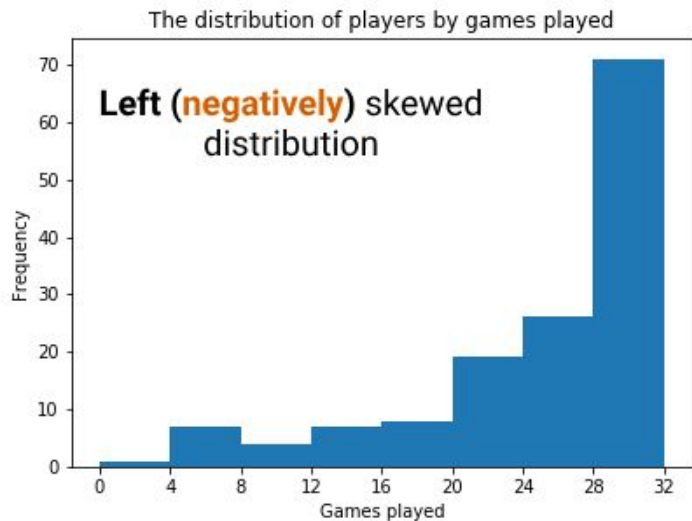
# Skewed Distributions

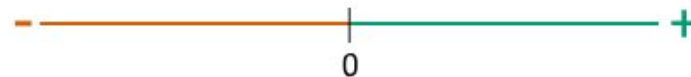# Skewed Distributions



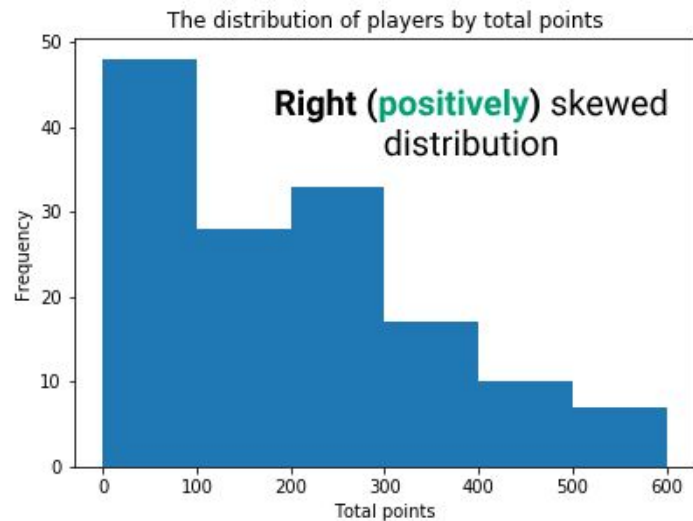The distribution of players by games played

The distribution of players by total points

# Skewed Distributions



The distribution of players by games played

**Left (negatively) skewed** distribution

The distribution of players by total points
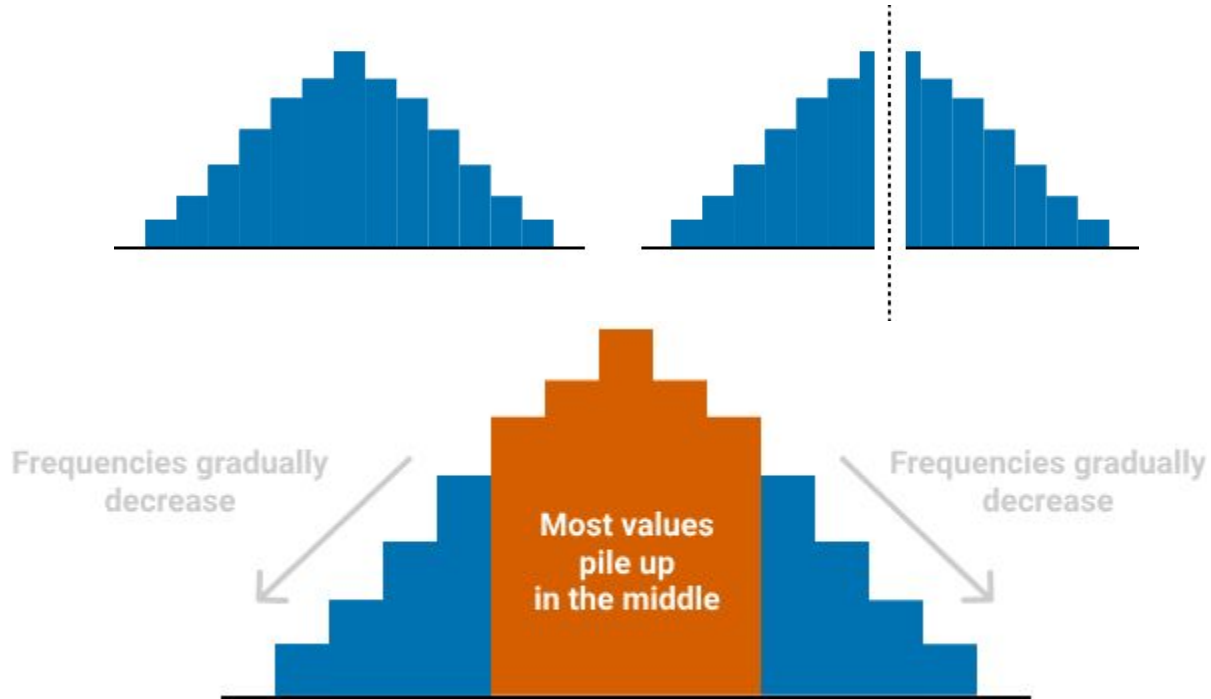
**Right (positively) skewed** distribution

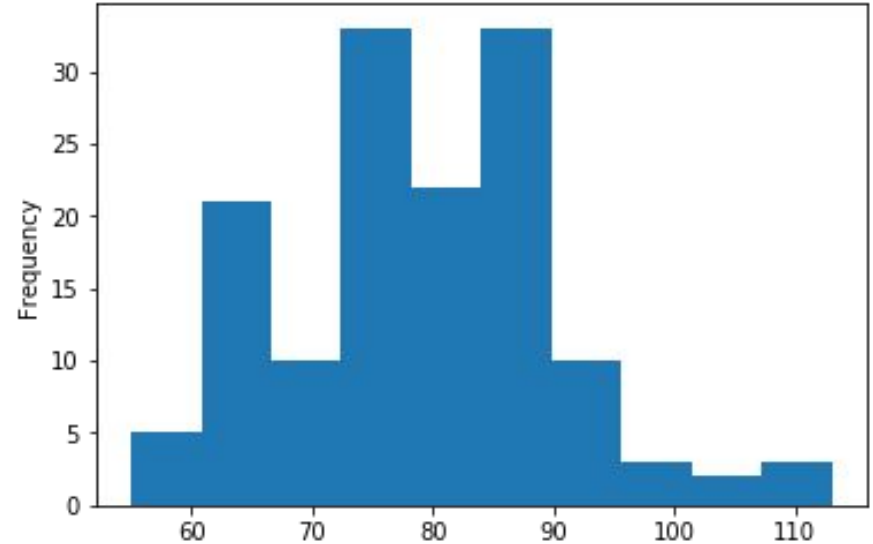If the **tail** points in the direction of **negative numbers** then the distribution is **negatively skewed**

If the **tail** points in the direction of **positive numbers** then the distribution is **positively skewed**

# Symmetrical Distributions



Frequencies gradually decrease

Most values pile up in the middle

Frequencies gradually decrease
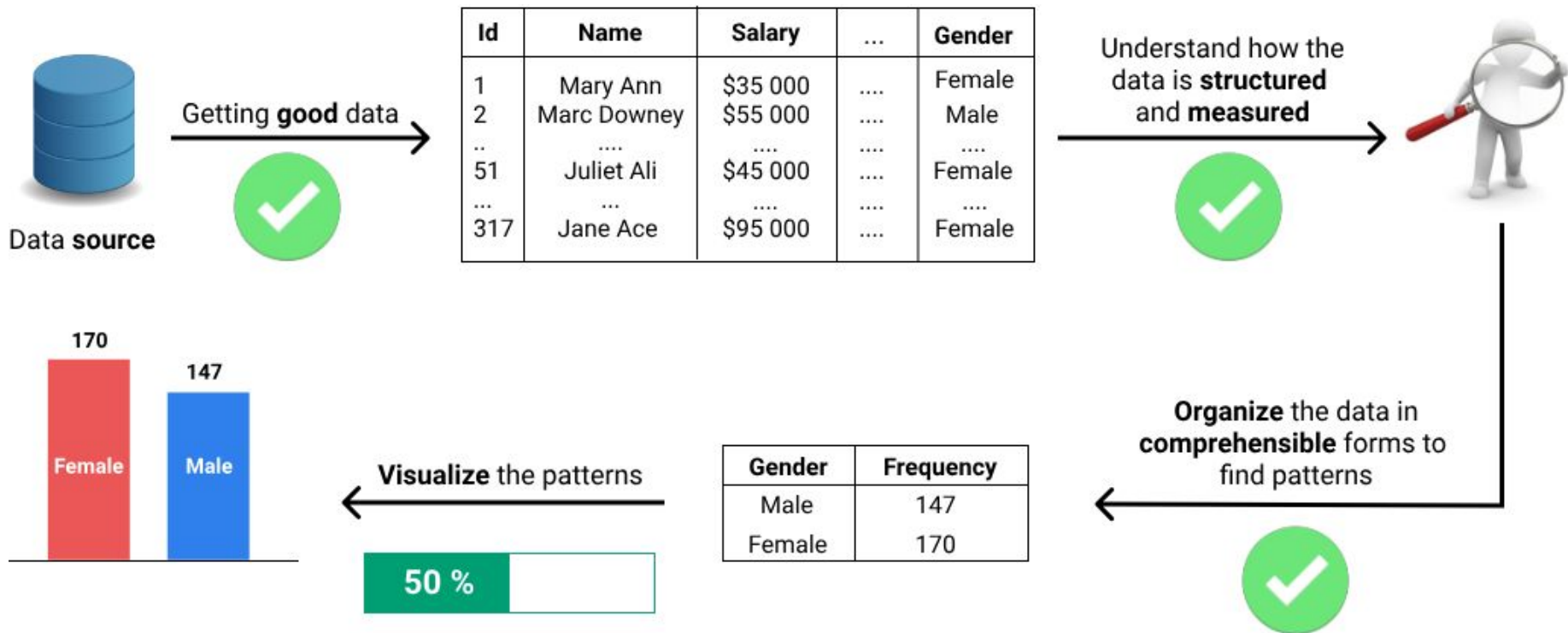
# Symmetrical Distribution (uniform)



The values are distributed uniformly

| Scale of measurement | Graphs we can use to show the distribution |
|---|---|
| Nominal |  |
| Ordinal |  |
| Interval |  |
| Ratio |  |

Data **source**

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**

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

| Gender | Frequency |
|--------|-----------|
| Male | 147 |
| Female | 170 |

**Visualize** the patterns
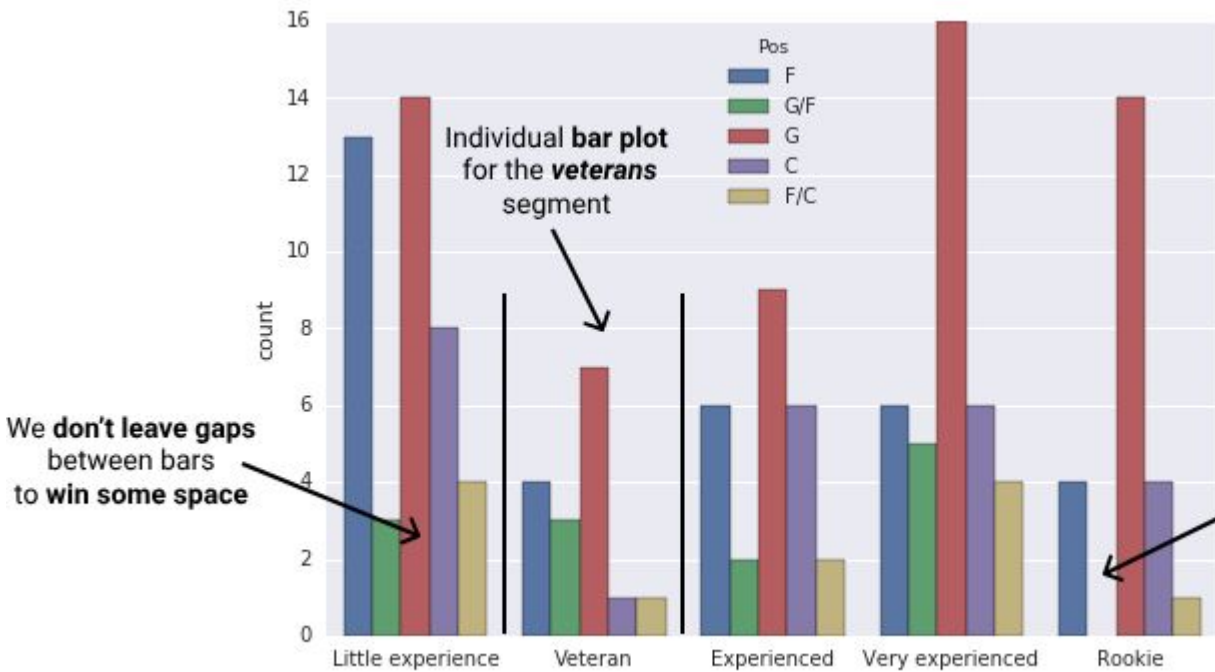
50 %

170

147

Female     Male

# Agenda (Parte II)

- Agrupamentos de gráficos de barras
- Comparando histogramas
- Estimativa de densidade kernel
- Gráficos de faixa e caixa
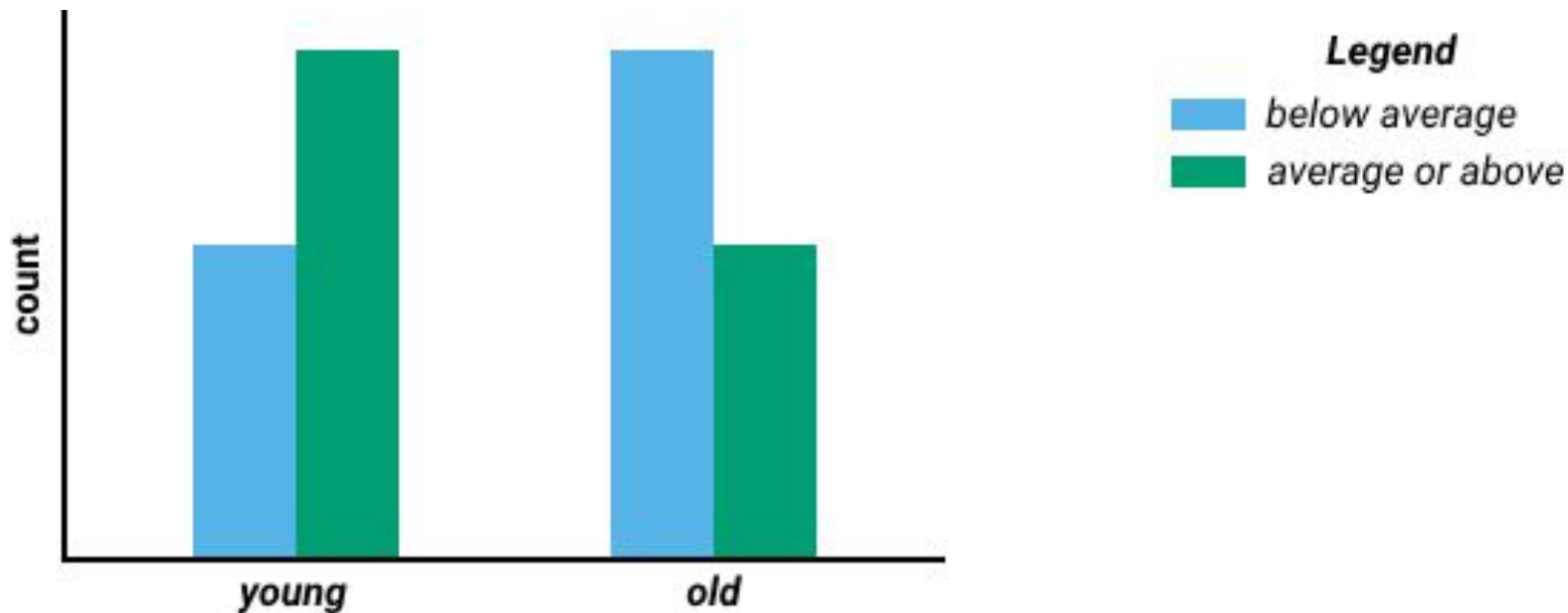- Pontos fora da curva

# Seaborn

# Comparing Frequency Distribution



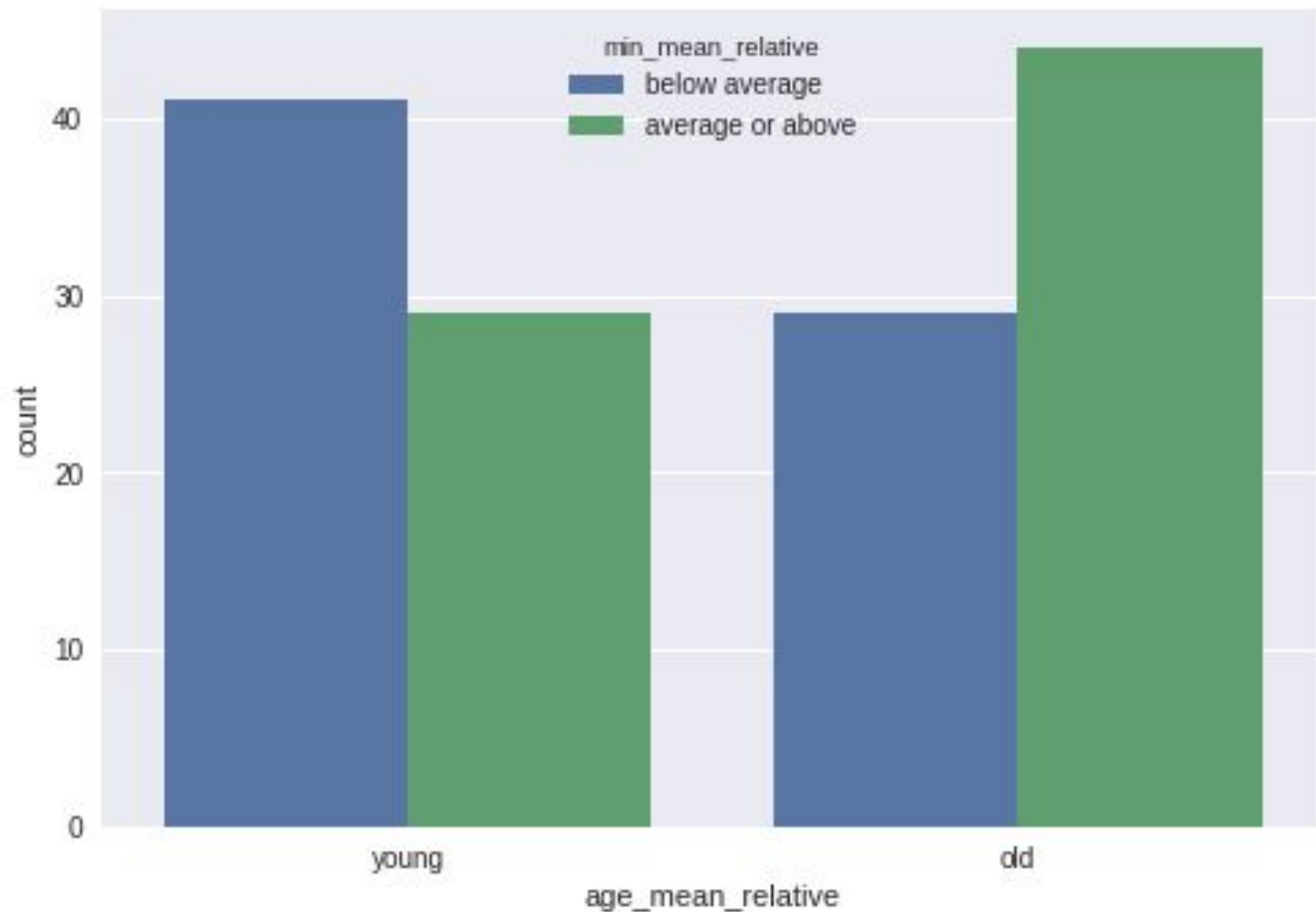| Years in WNBA | Label |
|---|---|
| 0 | Rookie |
| 1-3 | Little experience |
| 4-5 | Experienced |
| 5-10 | Very experienced |
| >10 | Veteran |

```
sns.countplot(x = 'Exp_ordinal', hue = 'Pos', data = wnba)
```
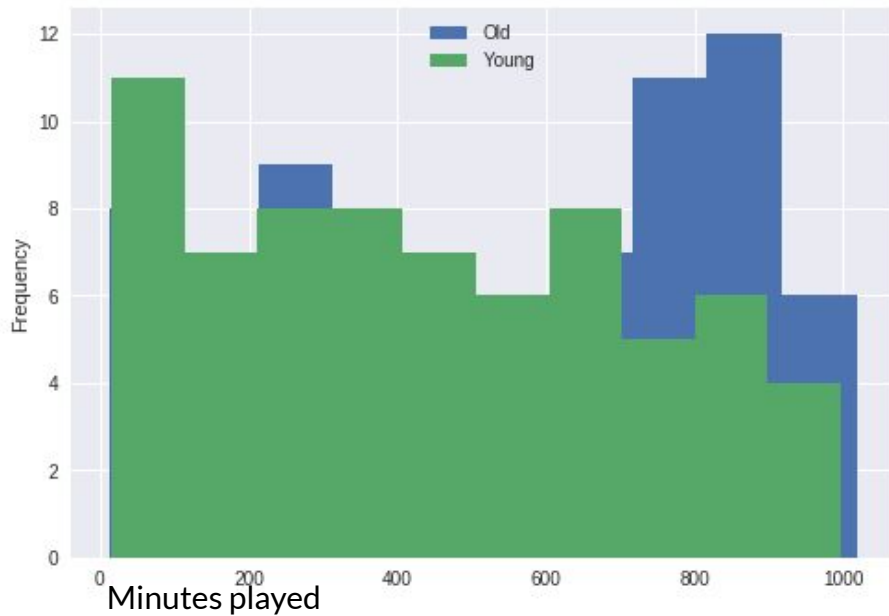
# Challenge: Do older players play less?

```
wnba['age_mean_relative'] = wnba['Age'].apply(lambda x: 'old' if x >= 27 else 'young')
wnba['min_mean_relative'] = wnba['MIN'].apply(lambda x: 'average or above' if x >= 497 else
                                              'below average')
cols = ["Name","Age","age_mean_relative","MIN","min_mean_relative"]
```

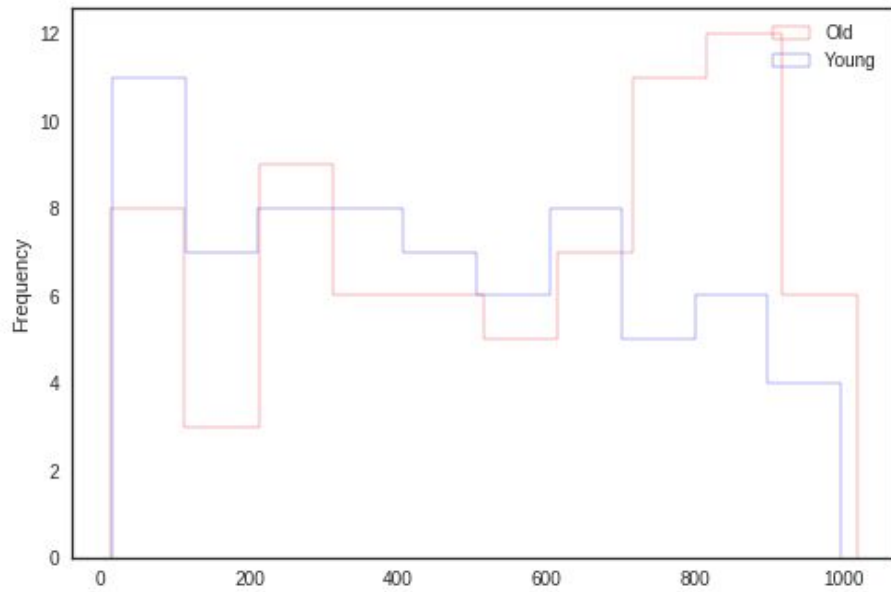| | Name | Age | age_mean_relative | MIN | min_mean_relative |
|---|---|---|---|---|---|
| 0 | Aerial Powers | 23 | young | 173 | below average |
| 1 | Alana Beard | 35 | old | 947 | average or above |
| 2 | Alex Bentley | 26 | young | 617 | average or above |
| 3 | Alex Montgomery | 28 | old | 721 | average or above |
| 4 | Alexis Jones | 23 | young | 137 | below average |

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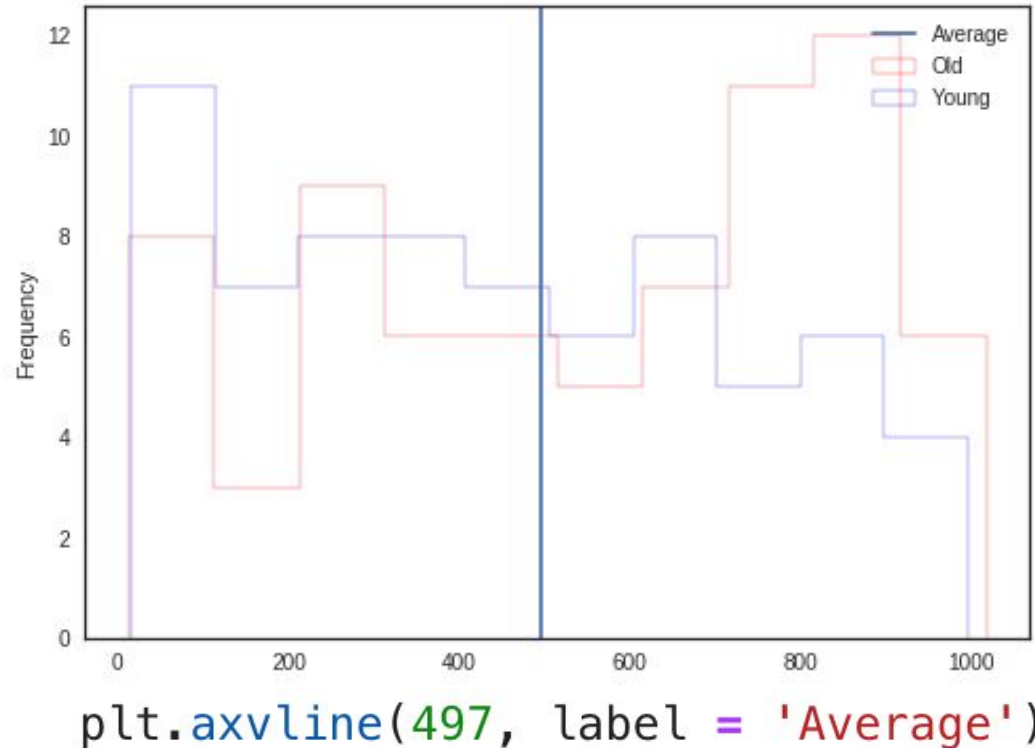
# Comparing Histograms



Minutes played

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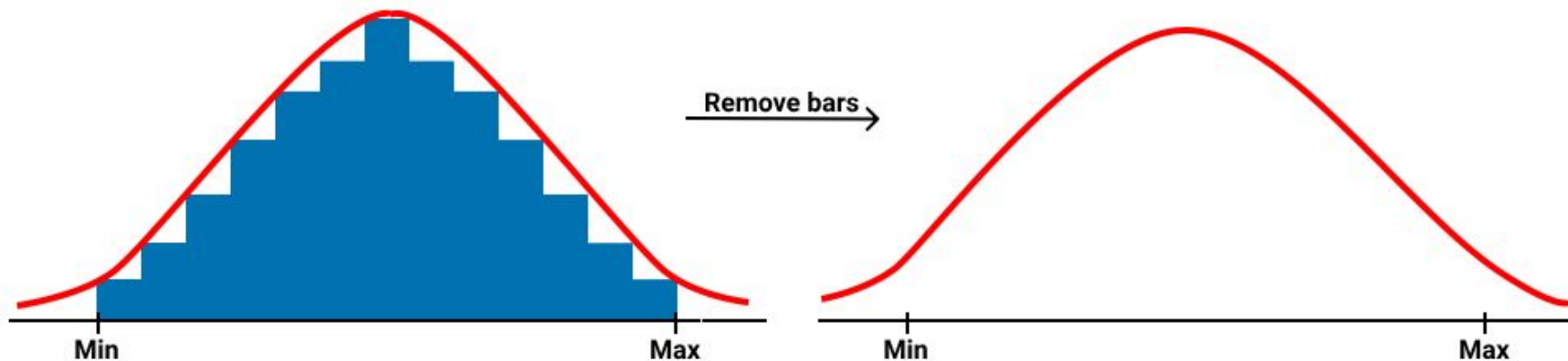
```
sns.set_style("white")
wnba[wnba.Age >= 27]['MIN'].plot.hist(histtype = 'step',
                                       label = 'Old',
                                       legend = True,color="red")
```
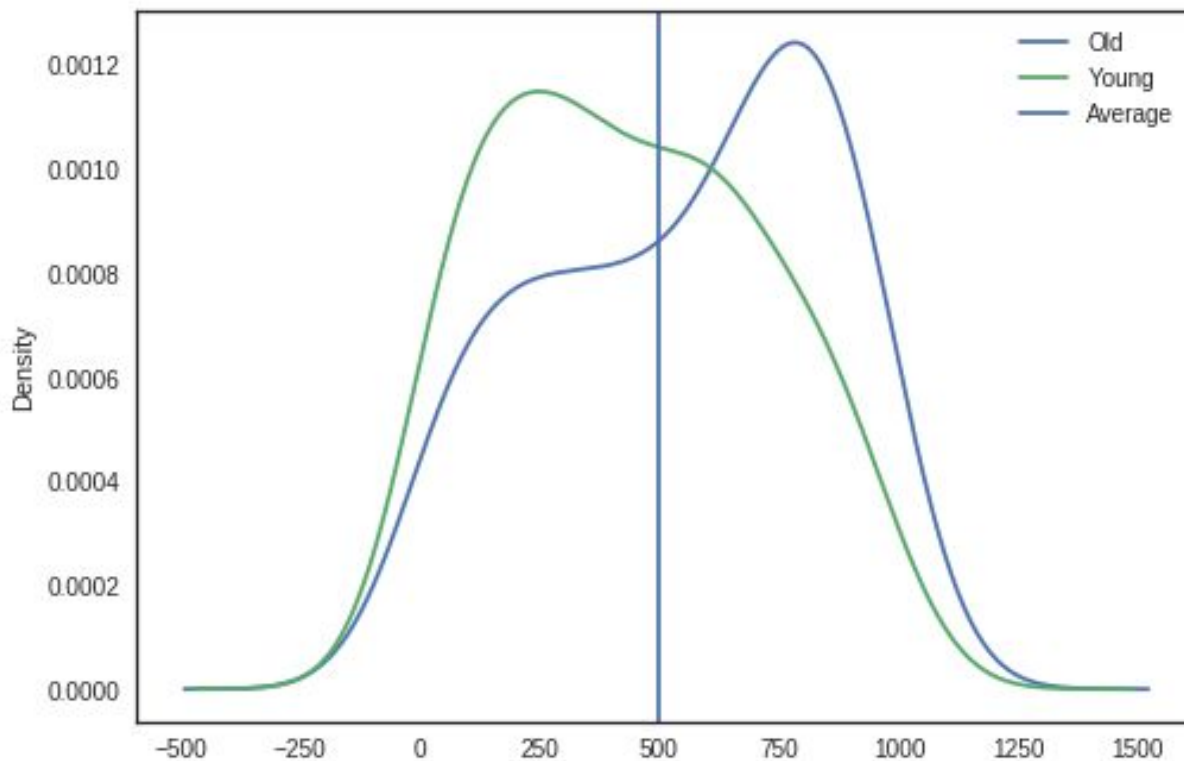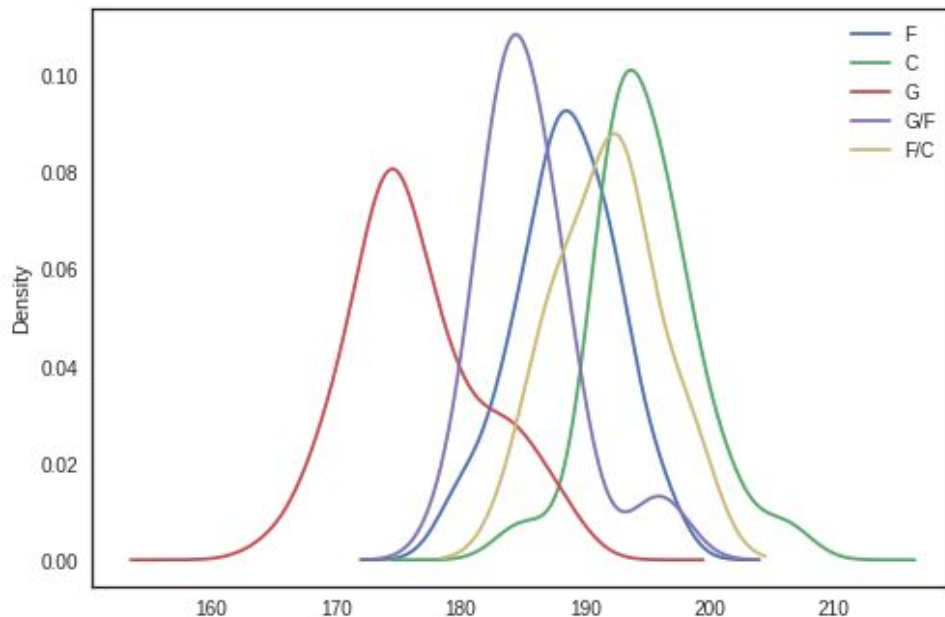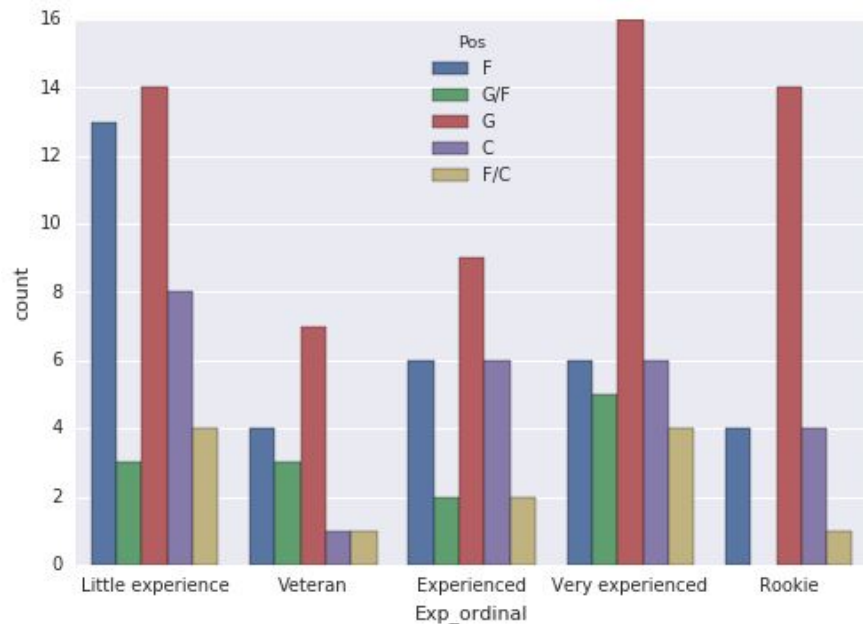
# Comparing Histograms



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

# Kernel Density Estimate (KDE) Plots
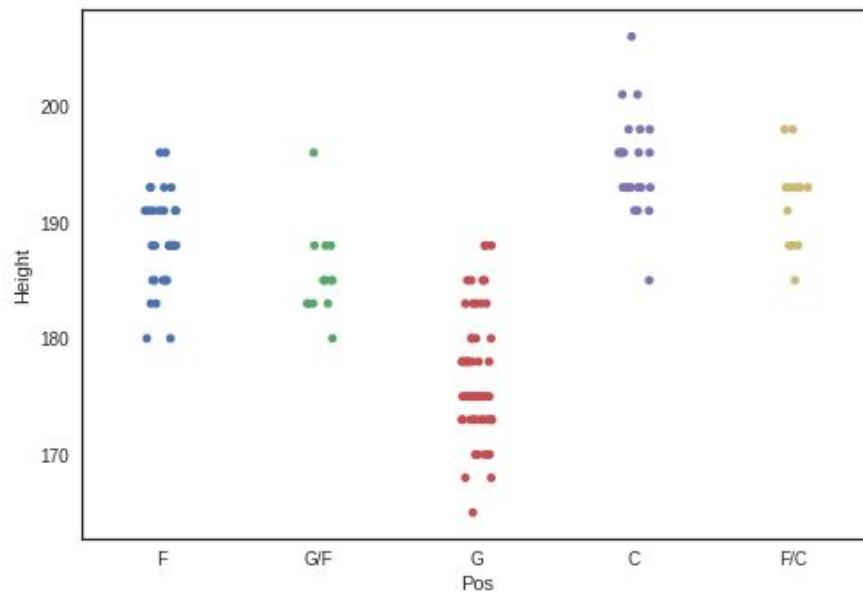
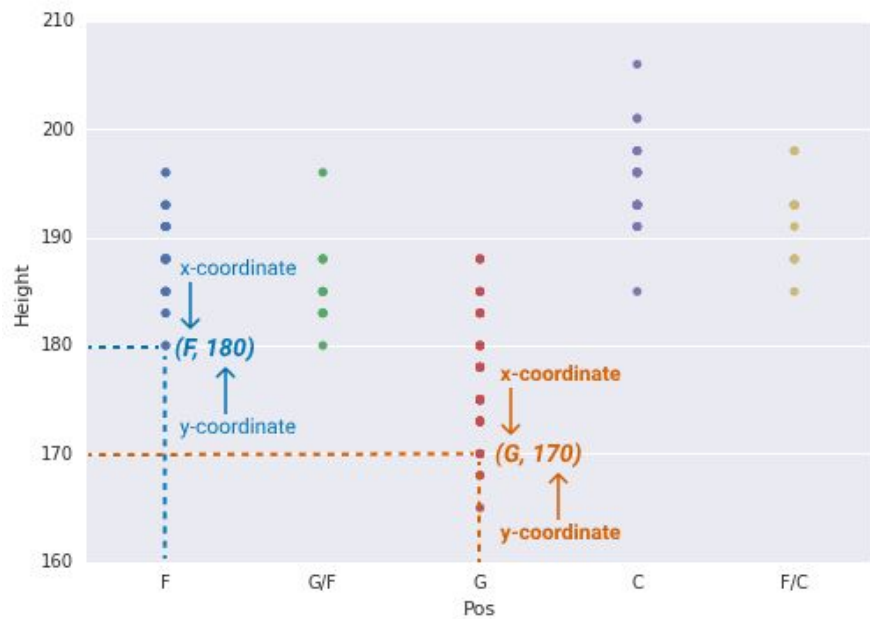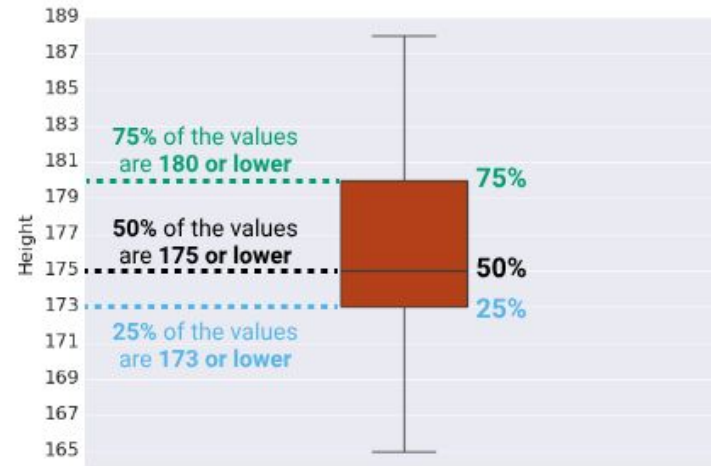Remove bars →
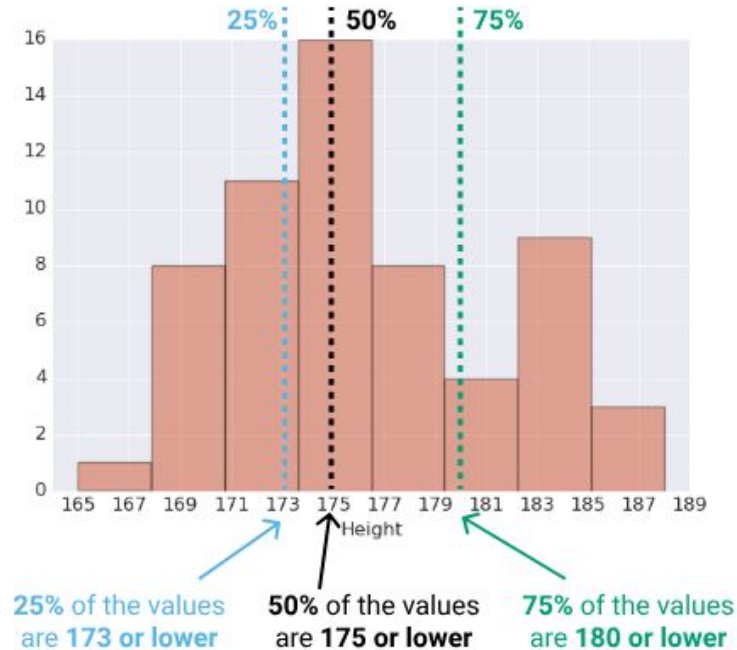
Min    Max

Min    Max

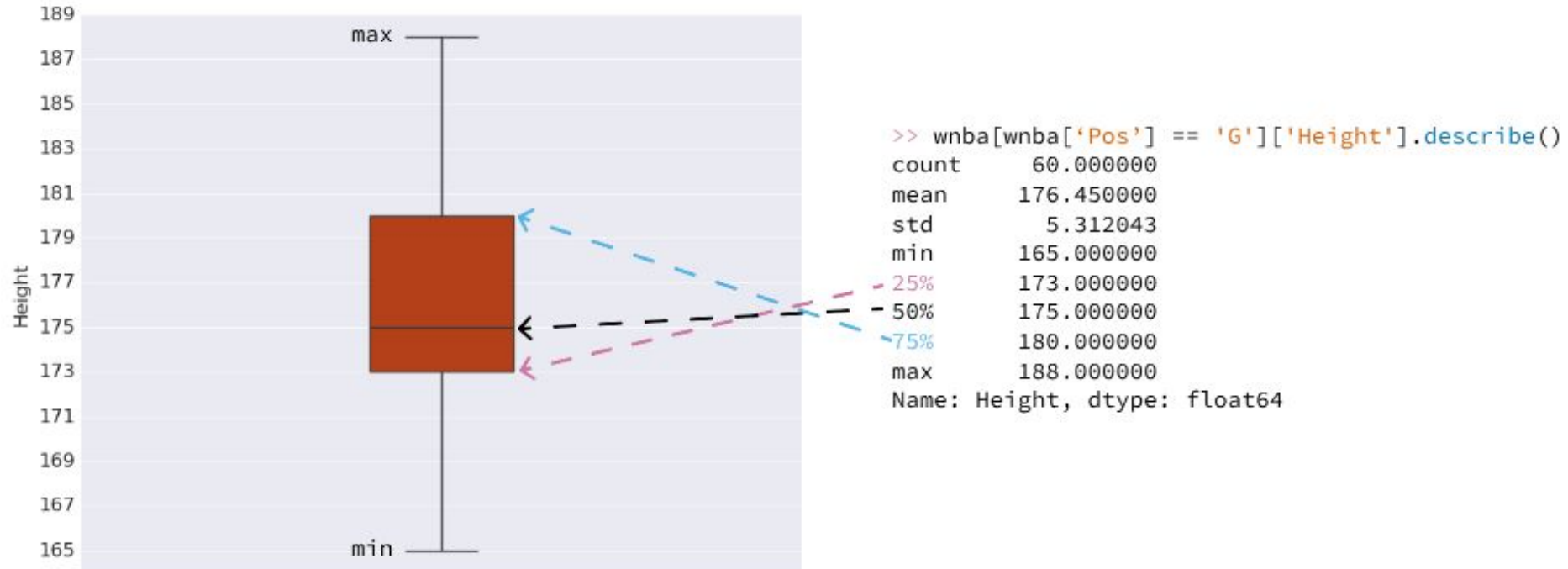# Kernel Density Estimate Plots

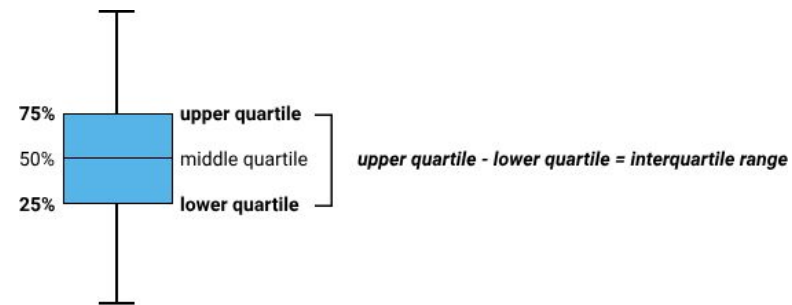# Drawbacks of Kernel Density Plots

# 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
```
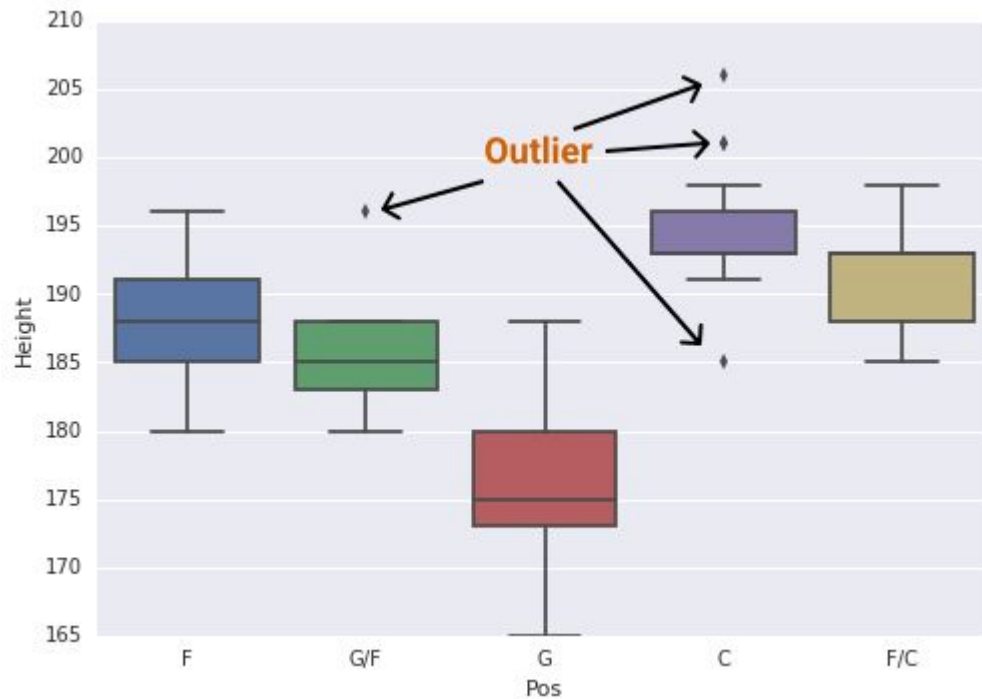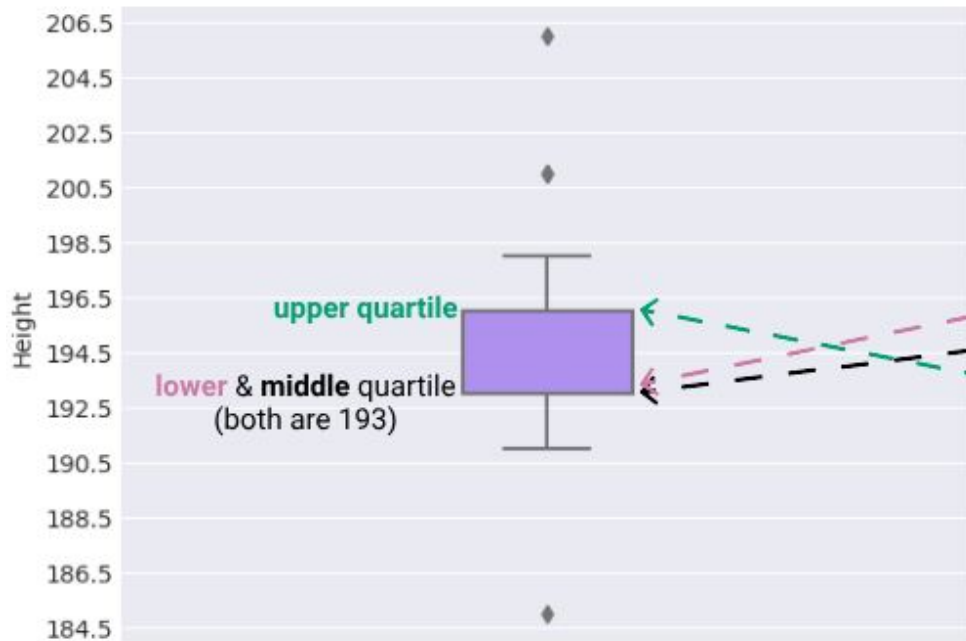
# Outliers
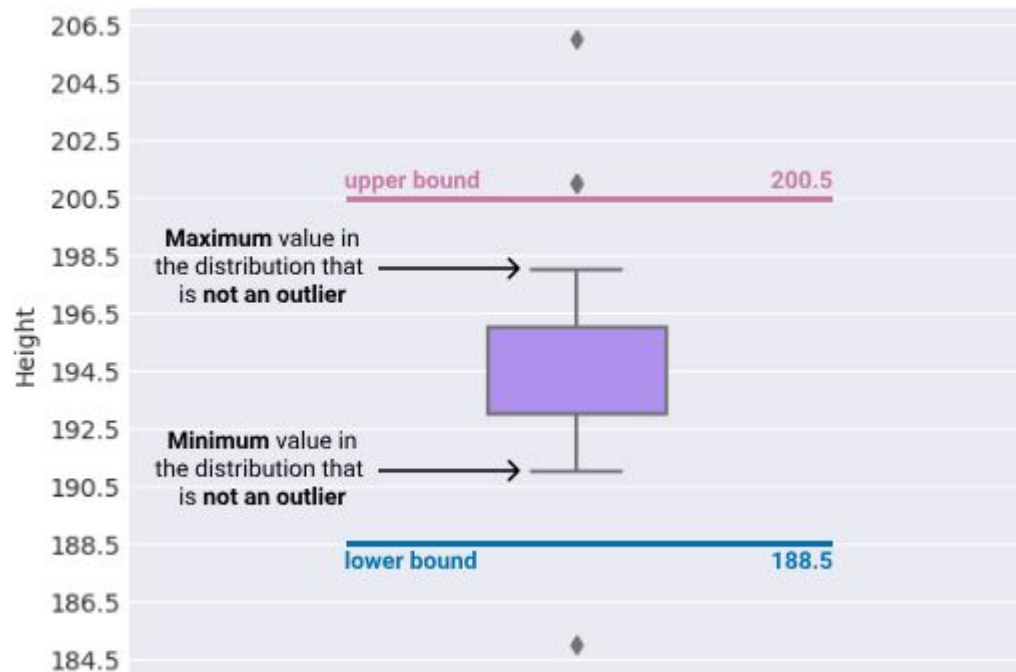
# Outliers



```
>> 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
```
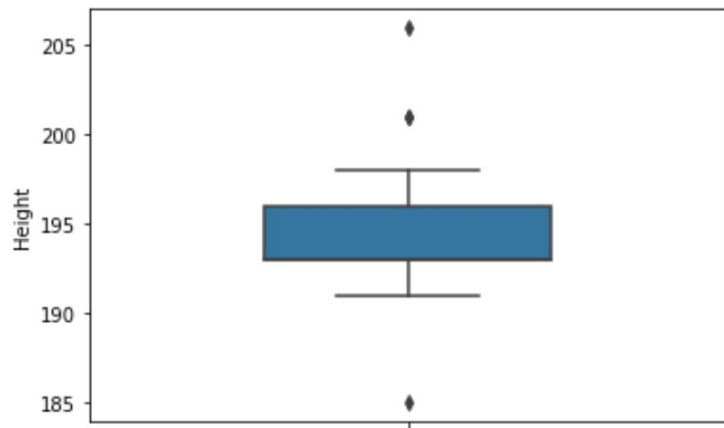
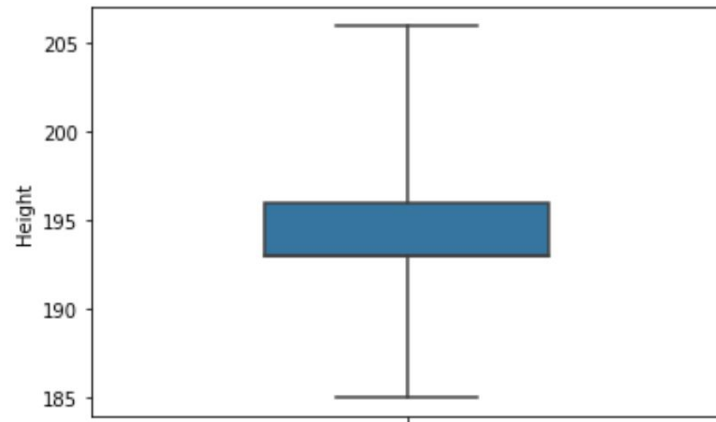# Outliers

# 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)
```
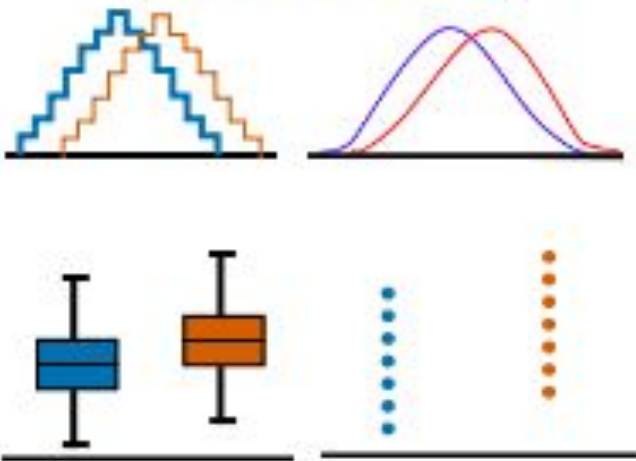
<matplotlib.axes._subplots.AxesSubplot at 0x1a18180208>

| Scale of measurement | Graphs we can use to compare distributions |
|---|---|
| Nominal |  |
| Ordinal |  |
| Interval & Ratio |  |

| 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 |

Data **source**

Getting **good** data

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

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

| Gender | Frequency |
|--------|-----------|
| Male | 147 |
| Female | 170 |

**Visualize** the patterns

170

147

Female   Male