

Session 03: Statistische Messgrößen

Dominic Schmitz & Janina Esser

Verein für Diversität in der Linguistik

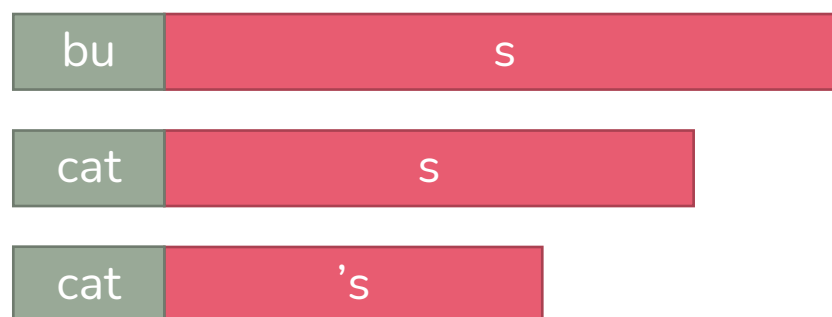
Statistische Messgrößen



- Für die folgenden Beispiele werden wir Daten folgender Studie nutzen:

The duration of word-final /s/ differs across morphological categories in English: Evidence from pseudowords¹

- Wort-finales /s/ zeigt je nach Bedeutung unterschiedliche Dauern



¹ Schmitz, D., Baer-Henney, D., & Plag, I. (2021). The duration of word-final /s/ differs across morphological categories in English: Evidence from pseudowords. *Phonetica*, 78(5-6), 571-616. doi: 10.1515/phon-2021-2013

Statistische Messgrößen



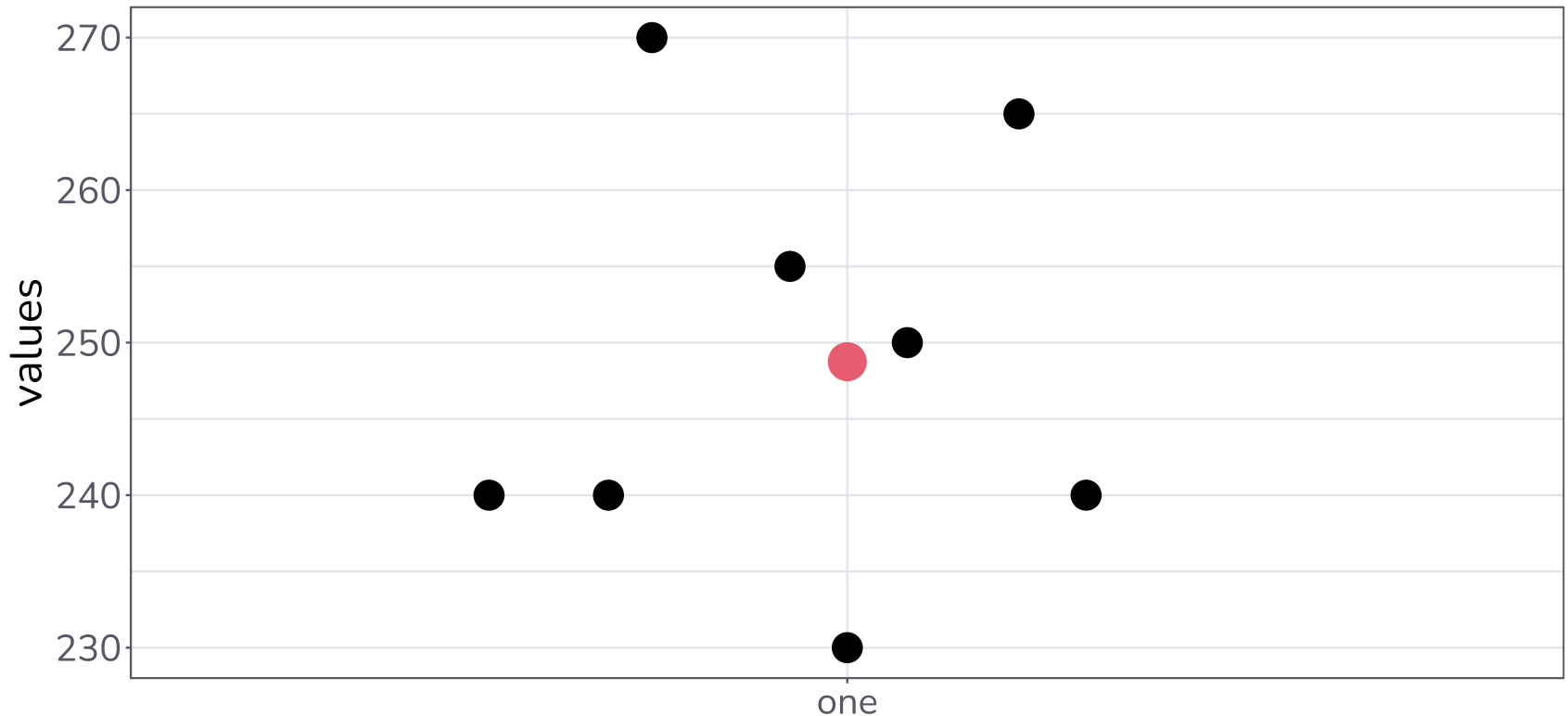
- Measures of Central Tendency / Lagemaße
 - **MEAN / DURCHSCHNITT, ARITHMETISCHES MITTEL**
 - **MEDIAN**
 - **MODE / MODUS**
- Measures of Dispersion / Streuungsmaße
 - **RANGE / SPANNWEITE**
 - **INTERQUARTILE RANGE / INTERQUARTILSPANNWEITE**
 - **SAMPLE COVARIANCE / STICHPROBENVARIANZ**
 - **STANDARD DEVIATION / STANDARDABWEICHUNG**
 - **STANDARD ERROR / STANDARDFEHLER**
- Shape of Distribution / Verteilungsform
 - **SKEWNESS / SCHIEFE**

Lagemaße



- **MEAN / DURCHSCHNITT, ARITHMETISCHES MITTEL**

Die Summe aller Werte geteilt durch die Anzahl der Werte



Lagemaße



- **MEAN / DURCHSCHNITT, ARITHMETISCHES MITTEL**

Die Summe aller Werte geteilt durch die Anzahl der Werte

$$A = \frac{1}{n} \sum_{i=1}^n a_i = \frac{a_1 + a_2 + \dots + a_n}{n}$$

Example:

$$A = \frac{270 + 240 + 240 + 255 + 250 + 265 + 230 + 240}{8} = 248.75$$

Lagemaße



- **MEAN / DURCHSCHNITT, ARITHMETISCHES MITTEL**

Die Summe aller Werte geteilt durch die Anzahl der Werte

```
mean (data$sDur)
```

```
## [1] 0.1315305
```

```
mean (data$baseDur)
```

```
## [1] 0.3190967
```

```
mean (data$speakingRate)
```

```
## [1] 3.449667
```

Lagemaße



- **MEAN / DURCHSCHNITT, ARITHMETISCHES MITTEL**

Die Summe aller Werte geteilt durch die Anzahl der Werte

```
mean (data$sDur[data$typeOfS == "nm"])
```

```
## [1] 0.156608
```

```
mean (data$sDur[data$typeOfS == "pl"])
```

```
## [1] 0.1317052
```

```
mean (data$sDur[data$typeOfS == "is"])
```

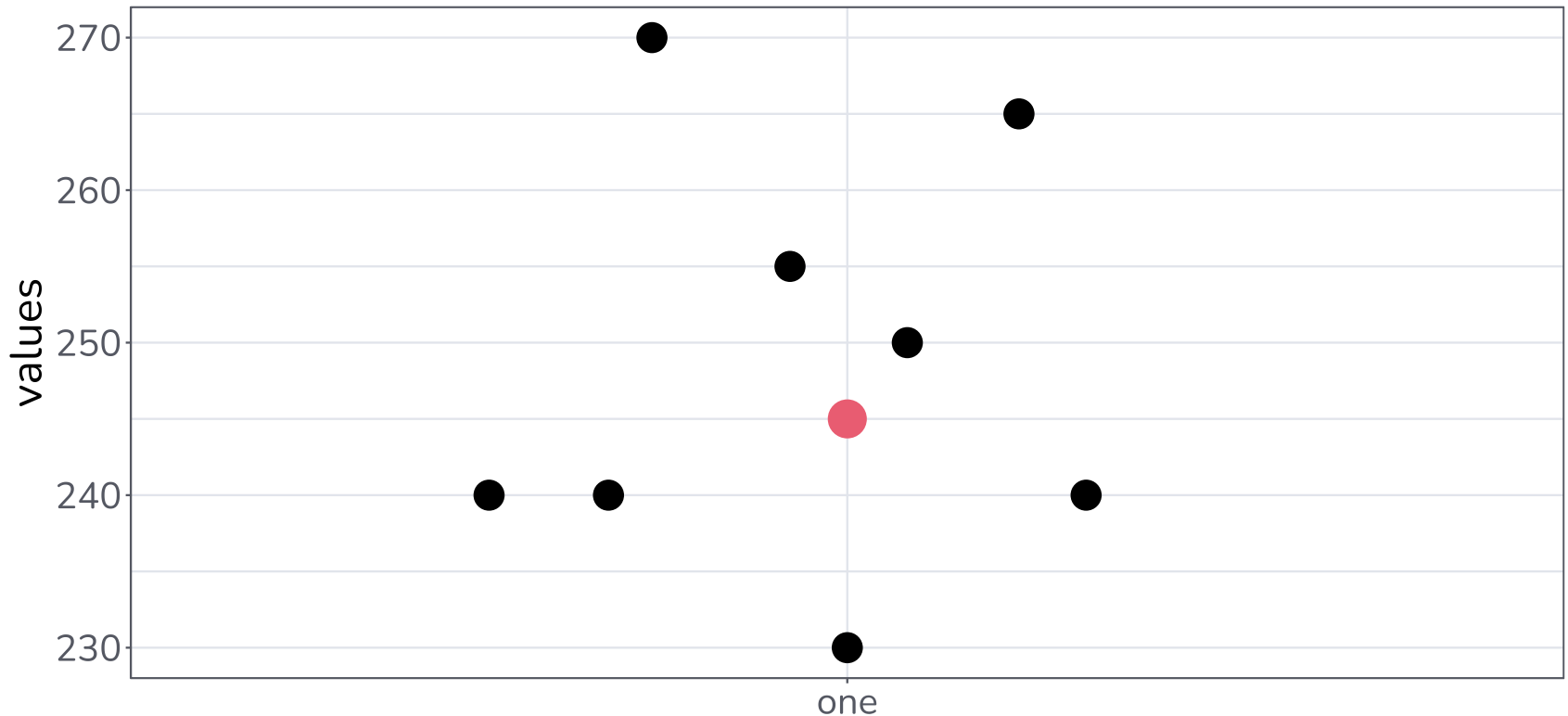
```
## [1] 0.1062782
```

Lagemaße



- **MEDIAN**

Der mittlere Wert in einer Reihe von Werten, geordnet vom Kleinsten zum Größten



Lagemaße



- **MEDIAN**

Der mittlere Wert in einer Reihe von Werten, geordnet vom Kleinsten zum Größten

$$\text{median}(a) = \frac{a_{\lfloor \#x \div 2 \rfloor} + a_{\lfloor \#x \div 2 + 0.5 \rfloor}}{2}$$

Example:

↓ 245
230, 240, 240, 240, 250, 255, 265, 270

↓ 240
230, 240, 240, 240, 250, 255, 265,

Lagemaße



- **MEDIAN**

Der mittlere Wert in einer Reihe von Werten, geordnet vom Kleinsten zum Größten

```
median (data$sDur)
```

```
## [1] 0.118175
```

```
median (data$baseDur)
```

```
## [1] 0.306315
```

```
median (data$speakingRate)
```

```
## [1] 3.355
```

Lagemaße



- **MEDIAN**

Der mittlere Wert in einer Reihe von Werten, geordnet vom Kleinsten zum Größten

```
median (data$sDur [data$typeOfS == "nm"] )
```

```
## [1] 0.15425
```

```
median (data$sDur [data$typeOfS == "pl"] )
```

```
## [1] 0.121815
```

```
median (data$sDur [data$typeOfS == "is"] )
```

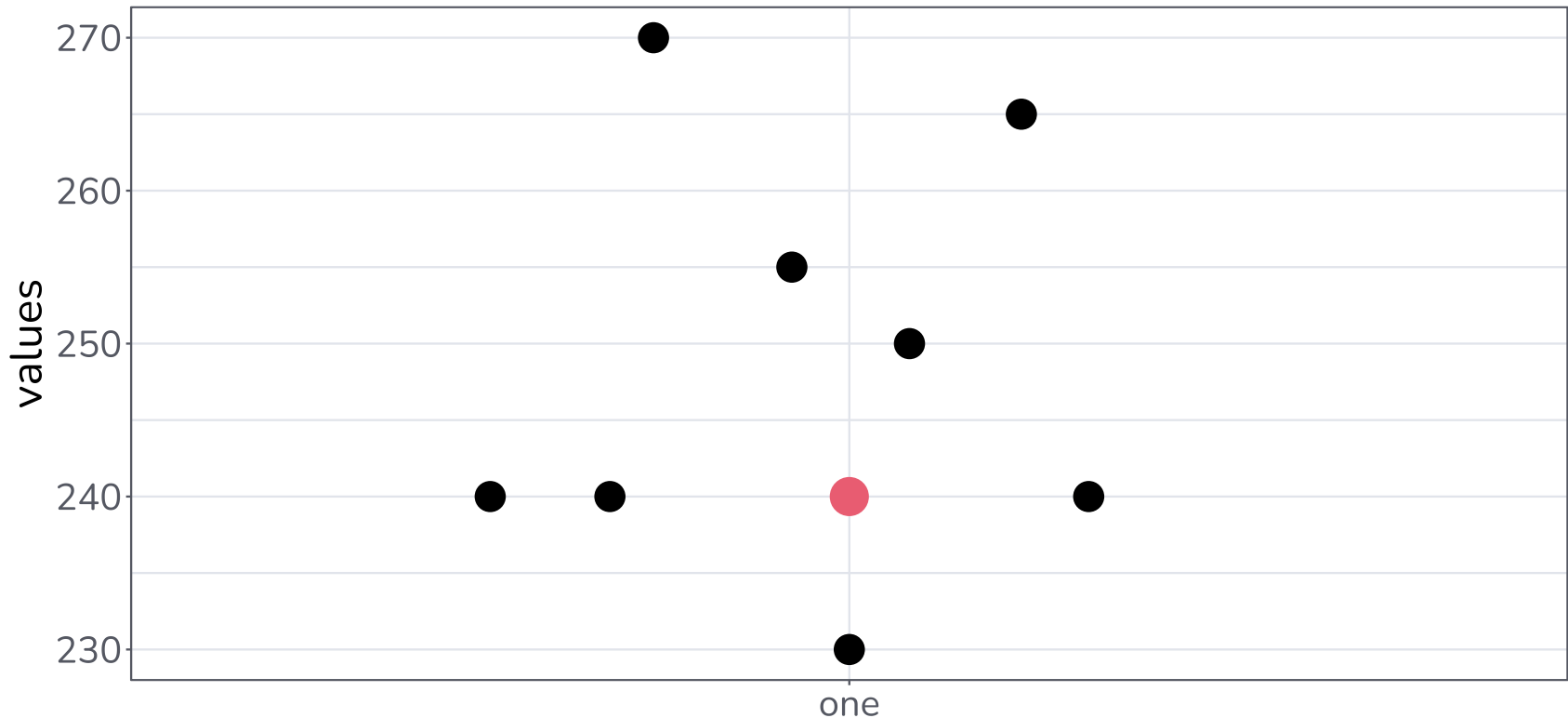
```
## [1] 0.101505
```

Lagemaße



- **MODE**

Der Wert, der am häufigsten in einer Gruppe von Werten vorkommt



Lagemaße



- **MODE**

Der Wert, der am häufigsten in einer Gruppe von Werten vorkommt

$$Modus = L + \frac{(f_m - f_1)h}{2f_m - f_1 - f_2}$$

Example:

↓ ↓ ↓
270, 240, 240, 255, 250, 265, 230, 240

Lagemaße



- **MODE**

Der Wert, der am häufigsten in einer Gruppe von Werten vorkommt

```
mode_stat(data$sDur)
```

```
## [1] 0.1311
```

```
mode_stat(data$baseDur)
```

```
## [1] 0.25162
```

```
mode_stat(data$speakingRate)
```

```
## [1] 2.94
```

Lagemaße



- **MODE**

Der Wert, der am häufigsten in einer Gruppe von Werten vorkommt

```
mode_stat(data$sDur[data$typeOfS == "nm"])\n## [1] 0.096
```

```
mode_stat(data$sDur[data$typeOfS == "pl"])\n## [1] 0.04176
```

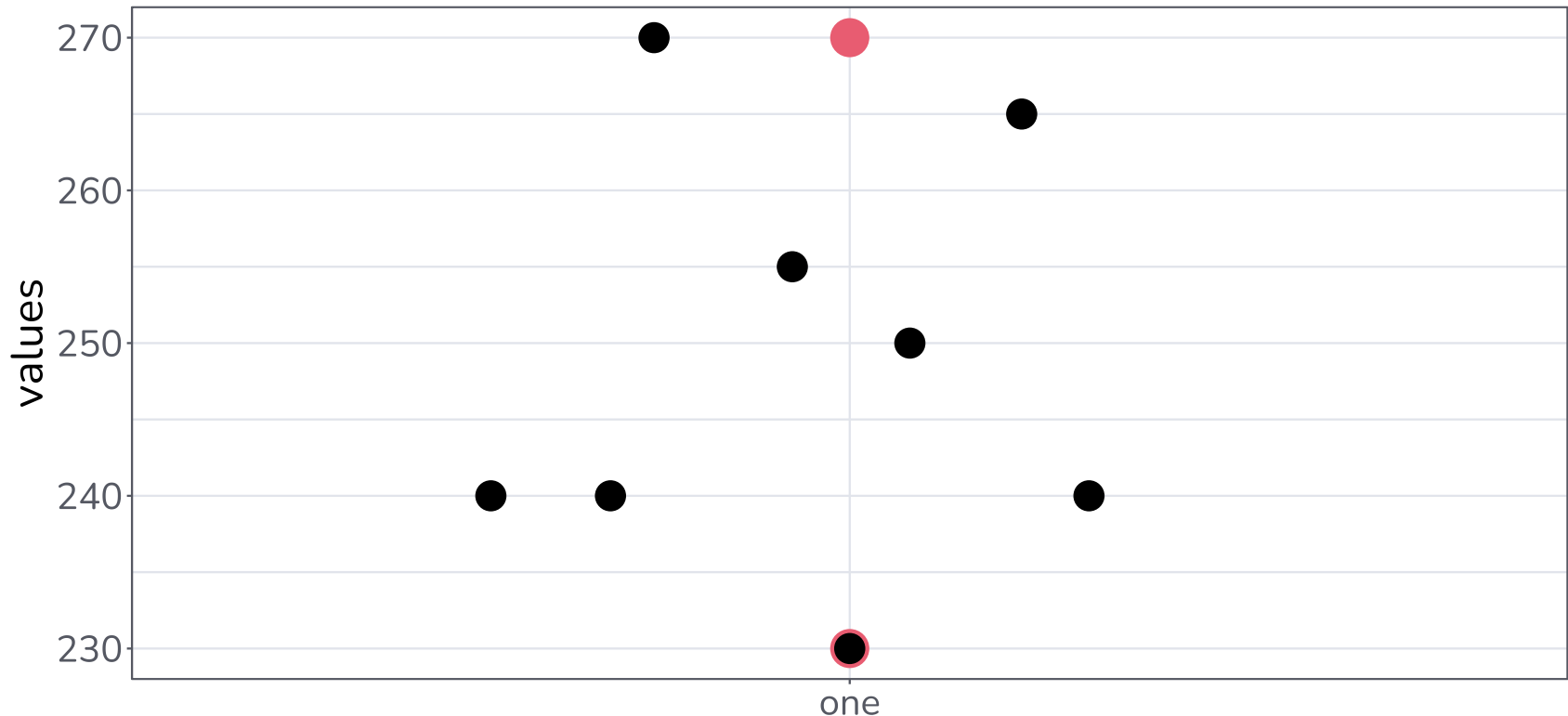
```
mode_stat(data$sDur[data$typeOfS == "is"])\n## [1] 0.1605
```

Streuungsmaße



- **RANGE / SPANNWEITE**

Die Differenz zwischen dem kleinsten und dem größten Wert in einer Gruppe von Werten



Streuungsmaße



- **RANGE / SPANNWEITE**

Die Differenz zwischen dem kleinsten und dem größten Wert in einer Gruppe von Werten

$$R = x_{max} - x_{min}$$

Example:

230, 240, 240, 240, 250, 255, 265, 270

$$R = 280 - 230 = 50$$

Streuungsmaße



- **RANGE / SPANNWEITE**

Die Differenz zwischen dem kleinsten und dem größten Wert in einer Gruppe von Werten

```
range (data$sDur)
```

```
## [1] 0.04176 0.32750
```

```
range (data$baseDur)
```

```
## [1] 0.17995 0.68749
```

```
range (data$speakingRate)
```

```
## [1] 1.52 6.94
```

Streuungsmaße



- **RANGE / SPANNWEITE**

Die Differenz zwischen dem kleinsten und dem größten Wert in einer Gruppe von Werten

```
range (data$sDur[data$typeOfS == "nm"])
```

```
## [1] 0.05202 0.32750
```

```
range (data$sDur[data$typeOfS == "pl"])
```

```
## [1] 0.04176 0.25289
```

```
range (data$sDur[data$typeOfS == "is"])
```

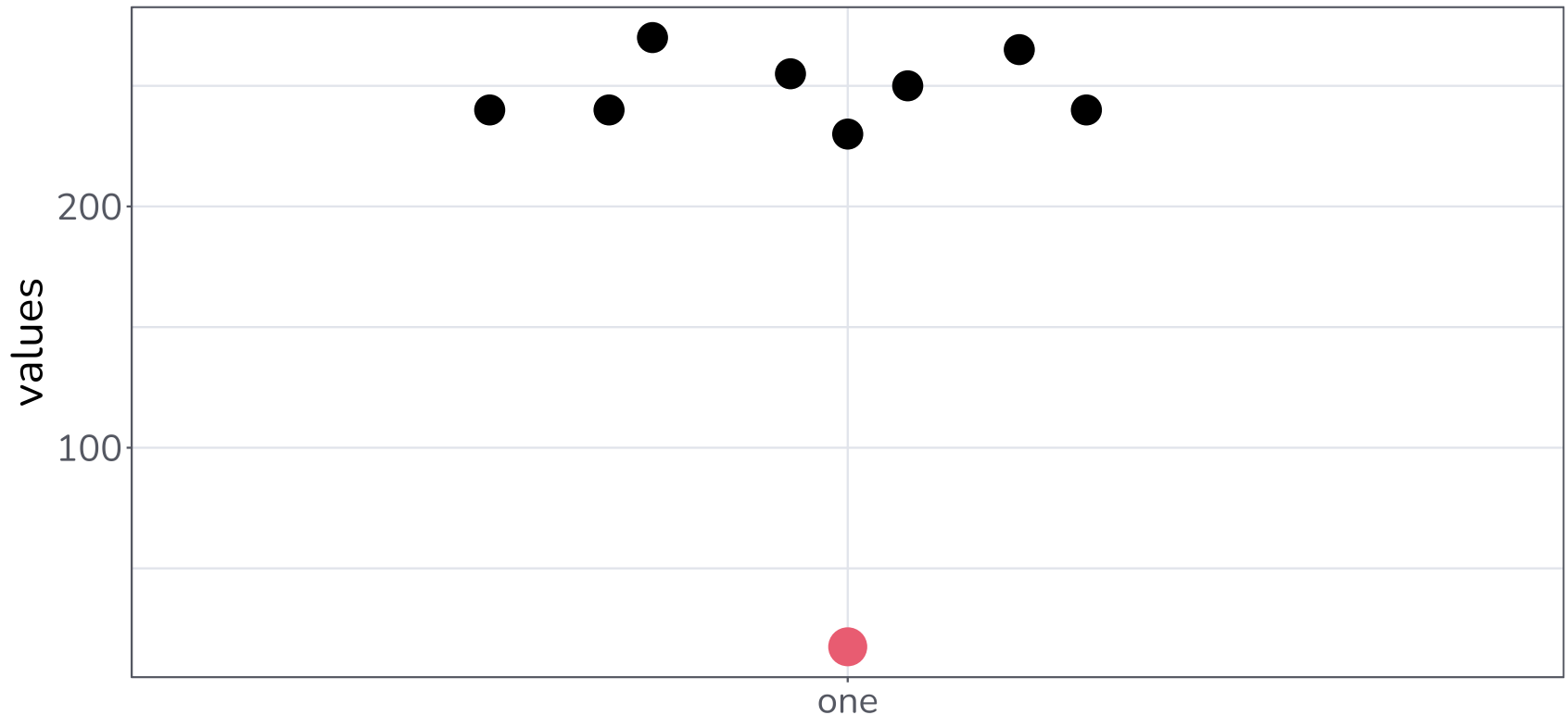
```
## [1] 0.04435 0.22428
```

Streuungsmaße



- **INTERQUARTILE RANGE / INTERQUARTILSPANNWEITE**

Die Spanne des Intervalls zwischen dem unteren und dem oberen Quartil



Streuungsmaße



- **INTERQUARTILE RANGE / INTERQUARTILSPANNWEITE**

Die Spanne des Intervalls zwischen dem unteren und dem oberen Quartil

$$x_{IQM} = \frac{2}{n} \sum_{i=\frac{n}{4}+1}^{\frac{3n}{4}} x_i$$

Example:

1. 270, 240, 240, 255, 250, 265, 230, 240 > sort
2. 230, 240, 240, 240, 250, 255, 265, 270 > quartiles
3. ~~230, 240, 240, 240, 250, 255, 265, 270~~ > remove 1st + 4th
4. $R = 255 - 240 = 15$ > range

Streuungsmaße



- **INTERQUARTILE RANGE / INTERQUARTILSPANNWEITE**

Die Spanne des Intervalls zwischen dem unteren und dem oberen Quartil

```
IQR(data$sDur)
```

```
## [1] 0.06783
```

```
IQR(data$baseDur)
```

```
## [1] 0.1067575
```

```
IQR(data$speakingRate)
```

```
## [1] 1.125
```

Streuungsmaße



- **INTERQUARTILE RANGE / INTERQUARTILSPANNWEITE**

Die Spanne des Intervalls zwischen dem unteren und dem oberen Quartil

```
IQR(data$sDur[data$typeOfS == "nm"])\n## [1] 0.0910275
```

```
IQR(data$sDur[data$typeOfS == "pl"])\n## [1] 0.072535
```

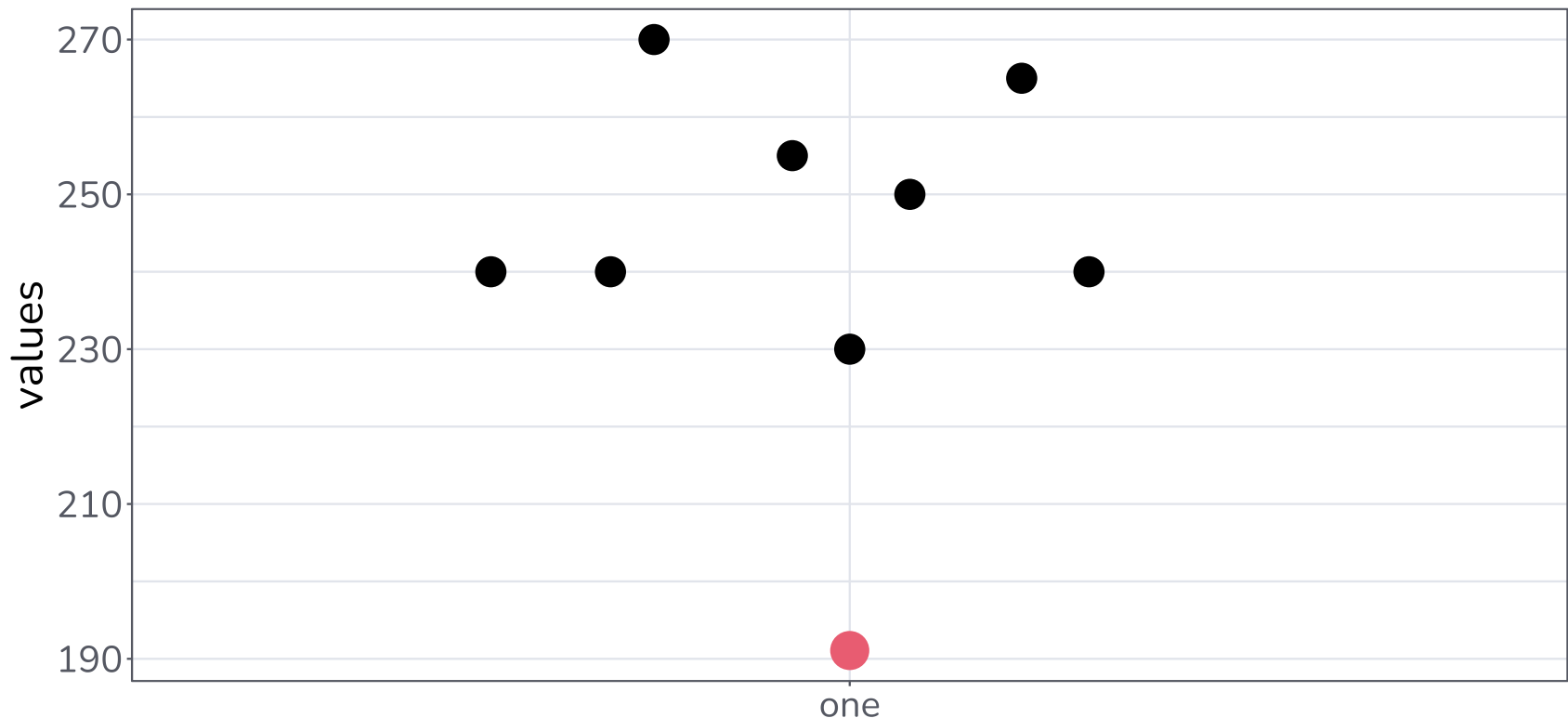
```
IQR(data$sDur[data$typeOfS == "is"])\n## [1] 0.0363475
```

Streuungsmaße



- **SAMPLE COVARIANCE / STICHPROBENVARIANZ**

Ein numerisches Maß dafür, wie die Datenpunkte um den Mittelwert gestreut sind



Streuungsmaße



- **SAMPLE COVARIANCE / STICHPROBENVARIANZ**

Ein numerisches Maß dafür, wie die Datenpunkte um den Mittelwert gestreut sind

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Example:

230, 240, 240, 240, 250, 255, 265, 270

$$s^2 = \frac{1}{8-1} \sum_{i=1}^8 (x_i - \bar{x})^2 = \frac{1337.5}{7} \approx 191.07$$

Streuungsmaße



- **SAMPLE COVARIANCE / STICHPROBENVARIANZ**

Ein numerisches Maß dafür, wie die Datenpunkte um den Mittelwert gestreut sind

```
var(data$sDur)
```

```
## [1] 0.002990366
```

```
var(data$baseDur)
```

```
## [1] 0.007913081
```

```
var(data$speakingRate)
```

```
## [1] 0.8649482
```

Streuungsmaße



- **SAMPLE COVARIANCE / STICHPROBENVARIANZ**

Ein numerisches Maß dafür, wie die Datenpunkte um den Mittelwert gestreut sind

```
var (data$sDur[data$typeOfS == "nm"])
```

```
## [1] 0.003943441
```

```
var (data$sDur[data$typeOfS == "pl"])
```

```
## [1] 0.002601761
```

```
var (data$sDur[data$typeOfS == "is"])
```

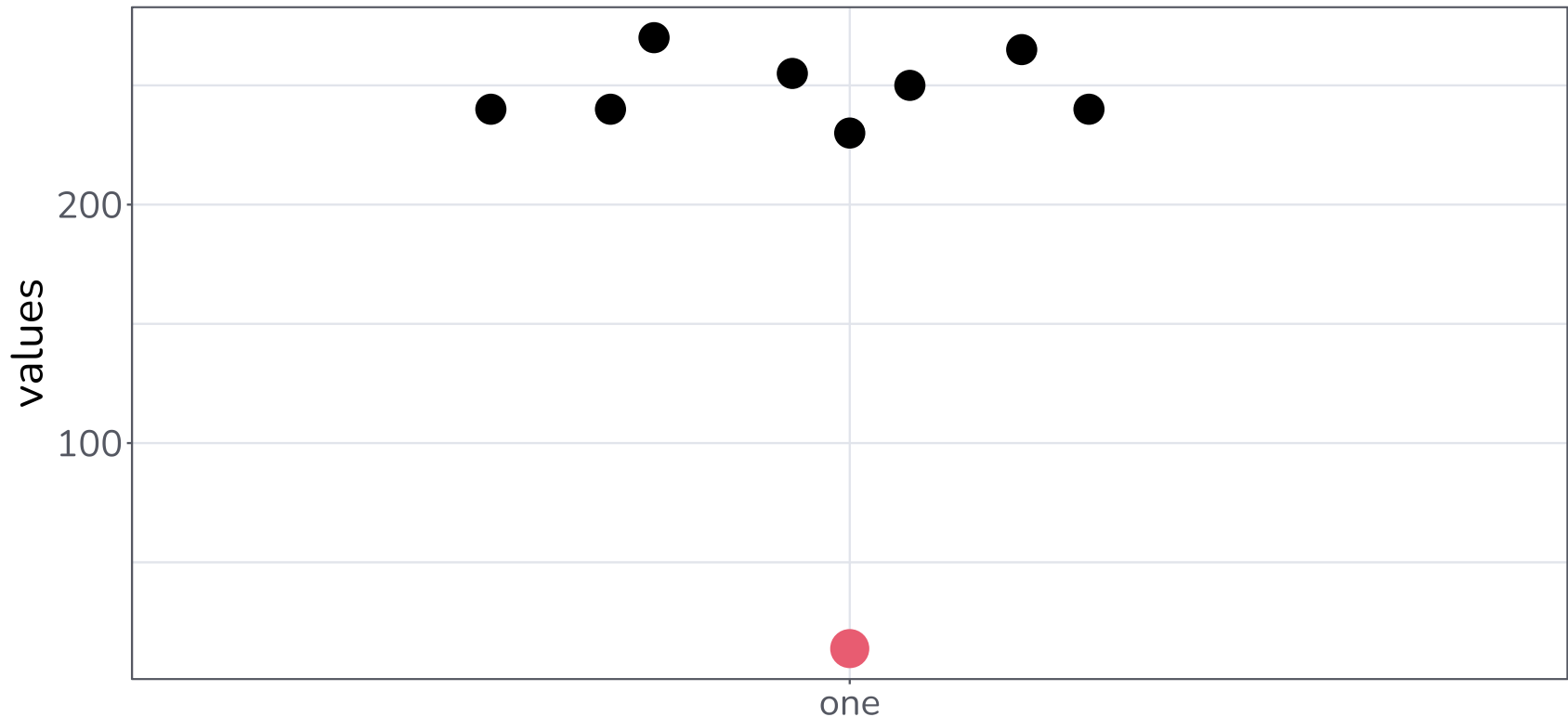
```
## [1] 0.001255514
```

Streuungsmaße



- **STANDARD DEVIATION / STANDARDABWEICHUNG**

Ein Indiz des Gesamtabstands der einzelnen Werte vom Mittelwert



Streuungsmaße



- **STANDARD DEVIATION / STANDARDABWEICHUNG**

Ein Indiz des Gesamtabstands der einzelnen Werte vom Mittelwert

$$s := + \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

Example:

230, 240, 250, 255, 265, 270

Wurzel der
Varianz

$$s = \sqrt{\frac{1337.5}{7}} \approx 13.82$$

Streuungsmaße



- **STANDARD DEVIATION / STANDARDABWEICHUNG**

Ein Indiz des Gesamtabstands der einzelnen Werte vom Mittelwert

```
sd(data$sDur)
```

```
## [1] 0.05468424
```

```
sd(data$baseDur)
```

```
## [1] 0.0889555
```

```
sd(data$speakingRate)
```

```
## [1] 0.9300259
```

Streuungsmaße



- **STANDARD DEVIATION / STANDARDABWEICHUNG**

Ein Indiz des Gesamtabstands der einzelnen Werte vom Mittelwert

```
sd(data$sDur[data$typeOfS == "nm"])
```

```
## [1] 0.06279683
```

```
sd(data$sDur[data$typeOfS == "pl"])
```

```
## [1] 0.05100746
```

```
sd(data$sDur[data$typeOfS == "is"])
```

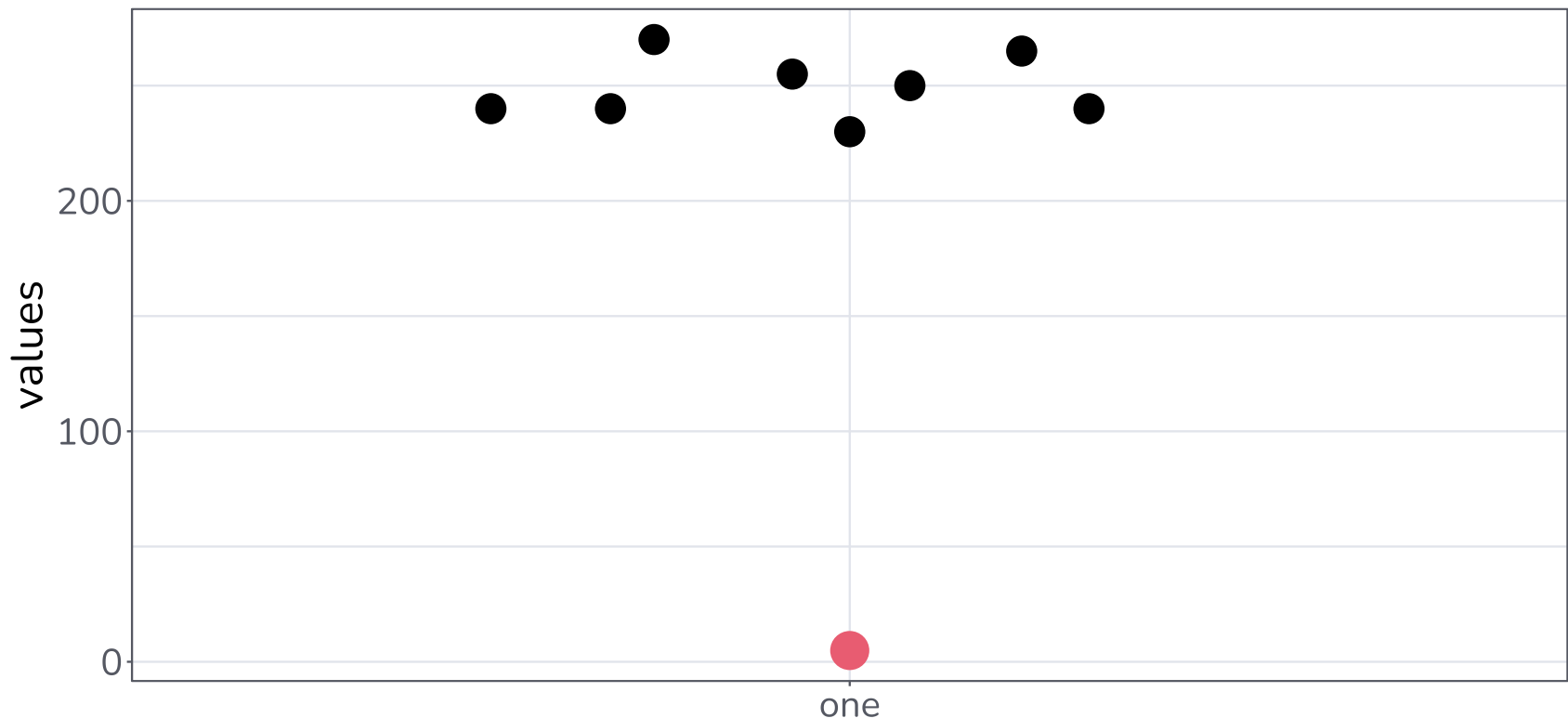
```
## [1] 0.03543323
```

Streuungsmaße



- **STANDARD ERROR / STANDARDFEHLER**

Ein statistischer Begriff, der die Genauigkeit misst, mit der eine Stichprobe eine Grundgesamtheit repräsentiert



Streuungsmaße



- **STANDARD ERROR / STANDARDFEHLER**

Ein statistischer Begriff, der die Genauigkeit misst, mit der eine Stichprobe eine Grundgesamtheit repräsentiert

$$\sigma(\bar{X}) = \frac{\sigma}{\sqrt{n}}$$

σ being the standard deviation of the population

**Standardabweichung geteilt
durch die Wurzel des
Stichprobenumfangs**

Example:

2, 230, 240, 250, 255, 265, 270

$$\sigma(\bar{X}) = \frac{\frac{1}{8-1} \sum_{i=1}^8 (x_i - \bar{x})^2}{\sqrt{8}} \approx 4.89$$

Streuungsmaße



- **STANDARD ERROR / STANDARDFEHLER**

Ein statistischer Begriff, der die Genauigkeit misst, mit der eine Stichprobe eine Grundgesamtheit repräsentiert

```
se (data$sDur)
```

```
## [1] 0.004464949
```

```
se (data$baseDur)
```

```
## [1] 0.007263186
```

```
se (data$speakingRate)
```

```
## [1] 0.0759363
```

Streuungsmaße



- **STANDARD ERROR / STANDARDFEHLER**

Ein statistischer Begriff, der die Genauigkeit misst, mit der eine Stichprobe eine Grundgesamtheit repräsentiert

```
se (data$sDur[data$typeOfS == "nm"] )  
## [1] 0.008880812
```

```
se (data$sDur[data$typeOfS == "pl"] )  
## [1] 0.007213545
```

```
se (data$sDur[data$typeOfS == "is"] )  
## [1] 0.005011015
```

Verteilungsform



- **SKEWNESS / SCHIEFE**

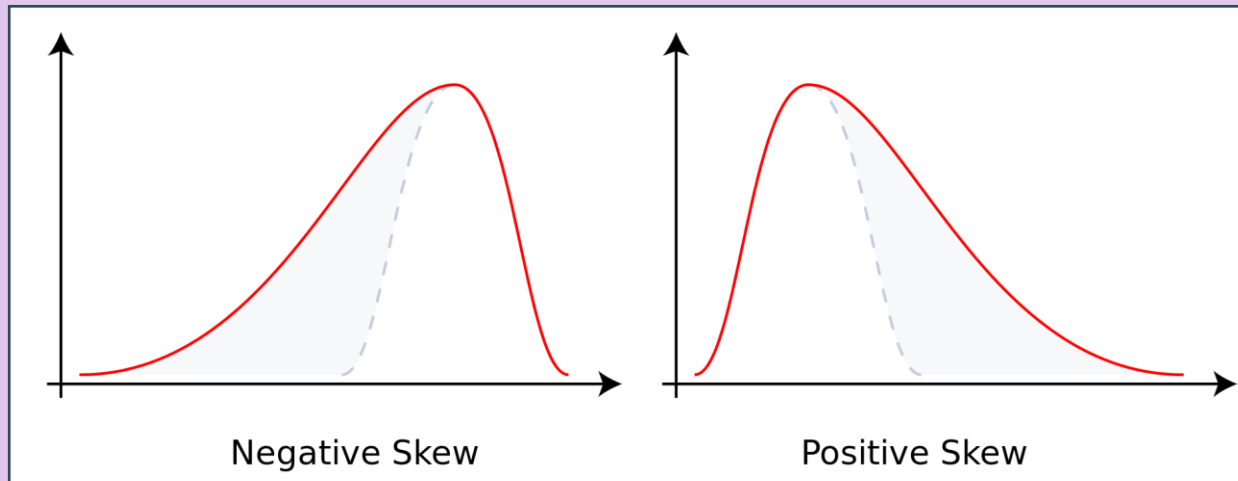
Asymmetrie in einer statistischen Verteilung, bei der die Kurve entweder nach links oder nach rechts verzerrt ist oder schief erscheint

$$v = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s} \right)^3$$

\bar{x} = mean

s = deviation

Example:



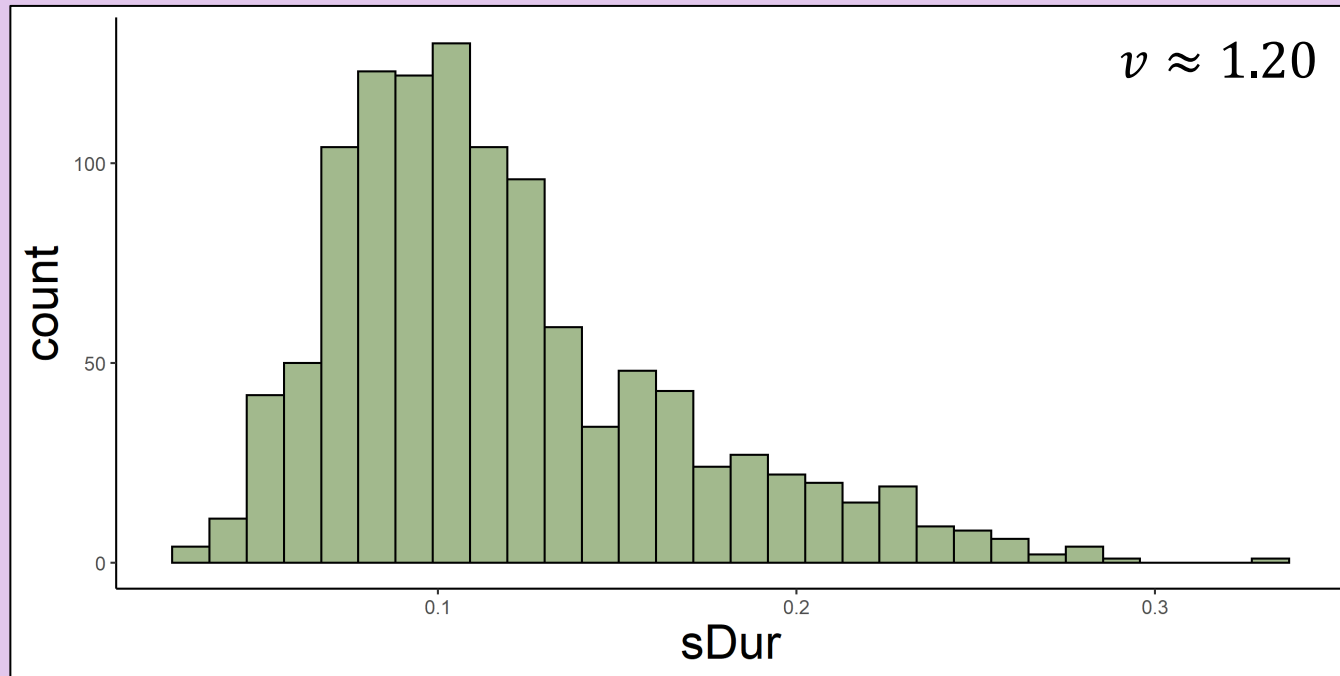
Verteilungsform



- **SKEWNESS / SCHIEFE**

Asymmetrie in einer statistischen Verteilung, bei der die Kurve entweder nach links oder nach rechts verzerrt ist oder schief erscheint

Example:



Verteilungsform



- **SKEWNESS / SCHIEFE**

Asymmetrie in einer statistischen Verteilung, bei der die Kurve entweder nach links oder nach rechts verzerrt ist oder schief erscheint

```
skewness (data$sDur)
```

```
## [1] 0.9483159
```

```
skewness (data$baseDur)
```

```
## [1] 1.360664
```

```
skewness (data$speakingRate)
```

```
## [1] 0.8348821
```

Verteilungsform



- **SKEWNESS / SCHIEFE**

Asymmetrie in einer statistischen Verteilung, bei der die Kurve entweder nach links oder nach rechts verzerrt ist oder schief erscheint

```
skewness (data$sDur [data$typeOfS == "nm"] )
```

```
## [1] 0.5884803
```

```
skewness (data$sDur [data$typeOfS == "pl"] )
```

```
## [1] 0.6259893
```

```
skewness (data$sDur [data$typeOfS == "is"] )
```

```
## [1] 0.8515867
```