

# 4

# Visualisierungen

Die relevanten Variablen eines Datensatzes werden häufig in einer Tabelle beschrieben

Für metrische Variablen werden häufig Anzahl, Spanne (Minimum und Maximum), Mittelwert mit Standardabweichung und der Median angegeben

Variable	n	Spanne	Mittelwert	Median
Schnabellänge	342 (2 NA)	32,1 - 59,6	43,9 $\pm$ 5,5	44,5
Schnabeltiefe	342 (2 NA)	13,1 - 21,5	17,2 $\pm$ 2,0	17,3
Flossenlänge	342 (2 NA)	172 - 231	200,9 $\pm$ 14,1	197
Körpergewicht	342 (2 NA)	2700 - 6300	4201,8 $\pm$ 802,0	4050

**Python:**

```
import pandas as pd
df.describe()
```

## Die relevanten Variablen eines Datensatzes werden häufig in einer Tabelle beschrieben

Für nominale Variablen können nur die Häufigkeiten angegeben werden. Bei ordinalen Variablen hängt es von der Anzahl Ausprägungen ab, ob eine Auflistung der Werte und/oder Minimum, Median und Maximum (und ggf. Quartile) angegeben werden.

Variable	n
<b>Art</b>	344 (0 NA)
Adelie	152 (44,2%)
Gentoo	124 (36,0%)
Chinstrap	68 (19,8%)

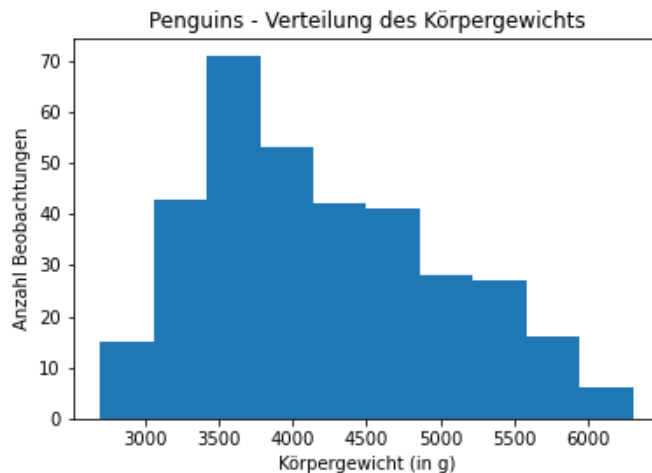
### Python:

```
import pandas as pd
df["spalte"].value_counts()
```

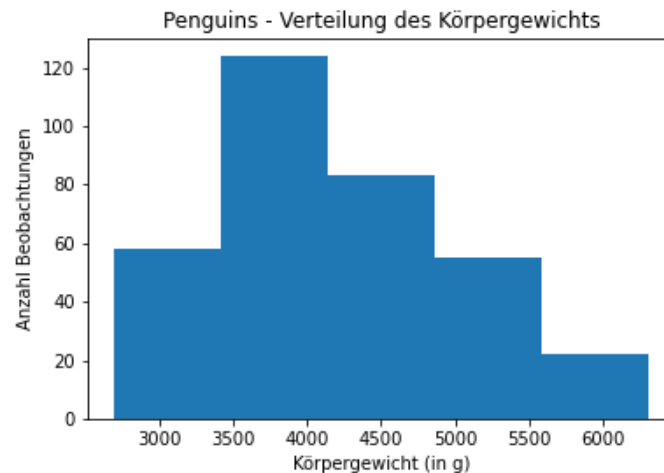
## Histogramm = Häufigkeitsverteilung metrischer Merkmale

- Das Intervall wird in Klassen eingeteilt. Die Anzahl Klassen (gleich große Intervalle) oder deren Grenzen beeinflussen, wie die Grafik aussieht
- Dann wird gezählt, wie viele Werte in einer der Klassen liegt

**10 Klassen**



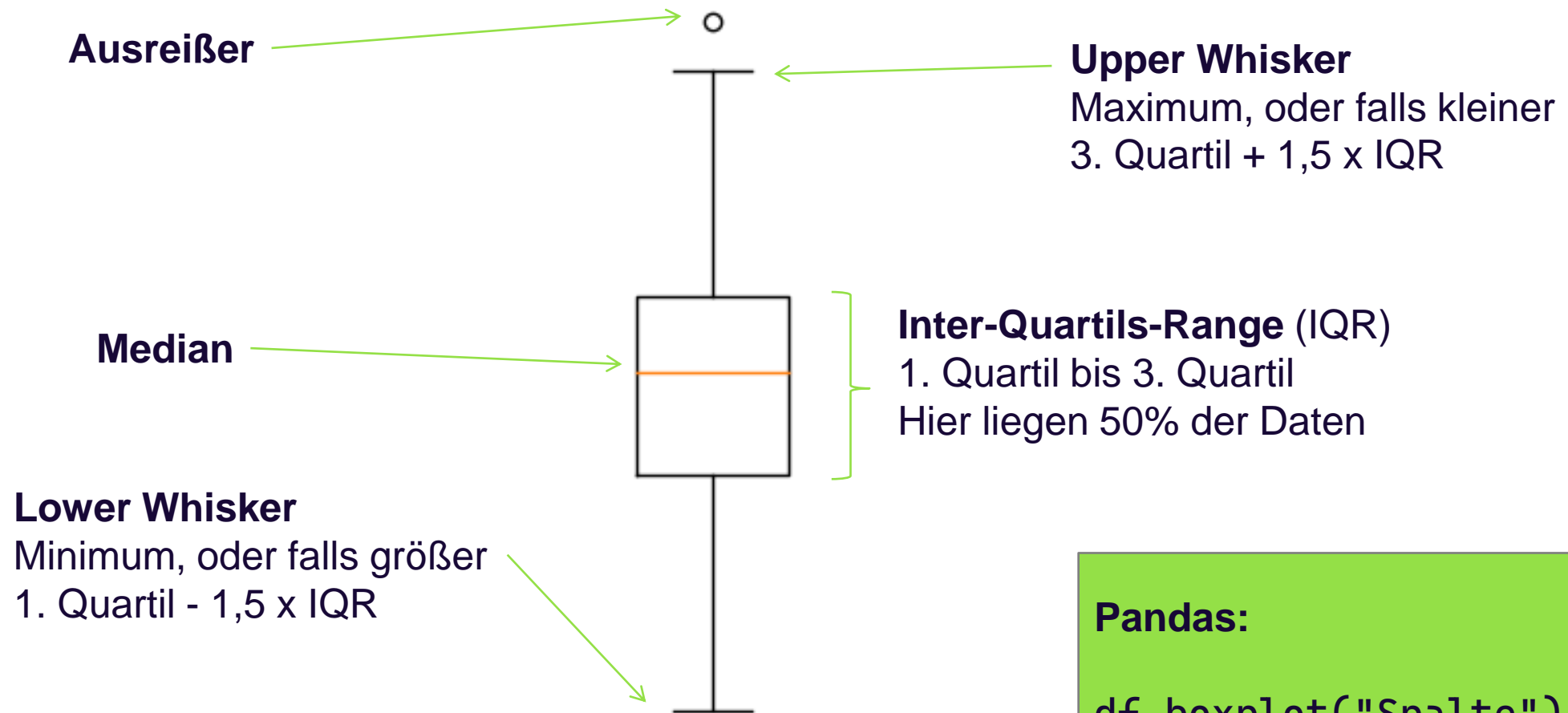
**5 Klassen**



### Pandas:

```
df["Spalte"].hist()  
df["Spalte"].hist(bins=5,  
                    grid=False)
```

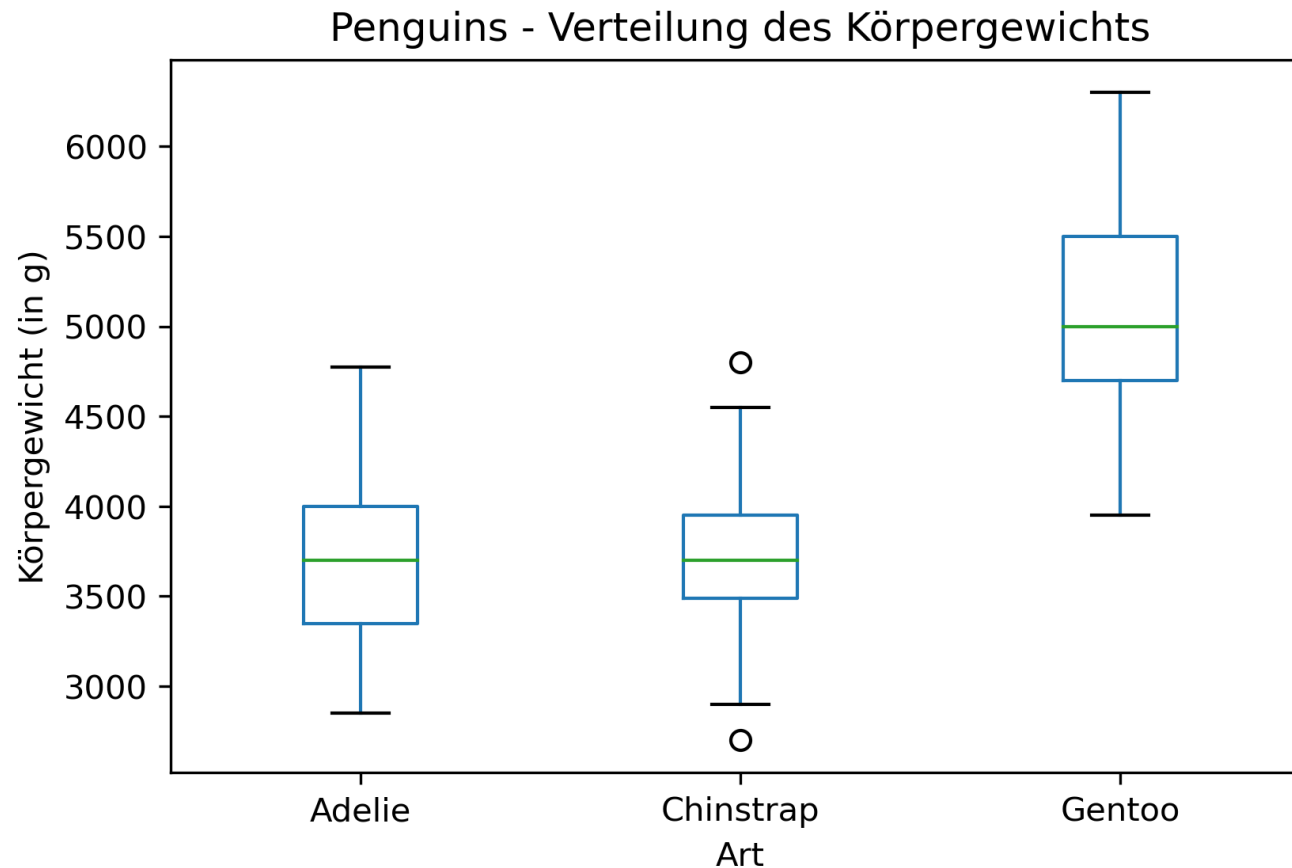
## Ein Boxplot visualisiert mehrere Kennzahlen für eine Variable



**Pandas:**

```
df.boxplot("Spalte")
```

**Boxplots sind gut geeignet, um die Verteilung einer Variable in mehreren Gruppen zu untersuchen**



**Pandas:**  
`df.boxplot(column="Spalte",  
by="Gruppe")`