

# Cheat Sheet – Deskriptive Statistik mit Python

Für Statistik I – Technische Hochschule Ingolstadt

## 0) Setup & Imports

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from scipy import stats
```

Tipps:

- `pd.set_option("display.max_columns", None)`
- `np.random.seed(42)` für Reproduzierbarkeit

## 1) Daten einlesen & speichern

```
1 df = pd.read_csv("file.csv")
2 # df = pd.read_excel("file.xlsx")
3 # df.to_csv("output.csv", index=False)
```

## 2) Überblick verschaffen

```
1 df.head()
2 df.info()
3 df.describe(include="all")
4 df.shape
5 df.columns
```

## 3) Auswahl, Filtern, Sortieren

```
1 df.loc[df["Age"] > 18, ["Sex", "Fare"]]
2 df.query("Age > 18 and Sex == 'female'")
3 df.sort_values(["Sex", "Age"], ascending=[True, False])
4 df.drop(columns=["Ticket"])
5 df.rename(columns={"old": "new"})
6 df.assign(BMI = df["Weight"] / (df["Height"]**2))
```

## 4) Datentypen, Strings, Kategorien

```
1 df["Age"] = pd.to_numeric(df["Age"], errors="coerce")
2 df["Date"] = pd.to_datetime(df["Date"])
3 df["Name"] = df["Name"].str.strip().str.lower()
4 df["CityFlag"] = df["City"].str.contains("Ingolstadt", case=False, na=False)
5 df["Class"] = pd.Categorical(df["Class"], ordered=True,
6                             categories=["low", "mid", "high"])
```

## 5) Fehlende Werte & Duplikate

```
1 df.isna().sum()
2 df.dropna(subset=["Age"])
3 df.fillna({"Age": df["Age"].median()})
4 df.duplicated().sum()
5 df.drop_duplicates()
```

## 6) Grundlegende Kennzahlen

```
1 df["x"].mean()
2 df["x"].median()
3 df["x"].std()
4 df["x"].quantile([0.25,0.5,0.75])
5 # Gruppiert:
6 df.groupby("Group")["x"].agg(["count","mean","median","std","min","max"])
```

## 7) Häufigkeiten & Kreuztabellen

```
1 df["cat"].value_counts(normalize=True).mul(100).round(1)
2 pd.crosstab(df["Sex"], df["Survived"], normalize="index", margins=True)
3 pd.crosstab(df["Class"], df["Sex"], normalize="columns", margins=True)
```

## 8) Binning (Klassen bilden)

```
1 pd.cut(df["x"], bins=[0,10,20,30], labels=["0-10","10-20","20-30"])
2 pd.qcut(df["x"], q=4) # Quartile
```

## 9) Ausreißer-Checks

```
1 Q1, Q3 = df["x"].quantile([0.25,0.75])
2 IQR = Q3 - Q1
3 lo, hi = Q1 - 1.5*IQR, Q3 + 1.5*IQR
4 mask_out = ~df["x"].between(lo, hi)
5
6 z = stats.zscore(df["x"], nan_policy="omit")
7 out_z = np.abs(z) > 3
```

## 10) Gewichtete Kennzahlen

```
1 wmean = np.average(df["x"], weights=df["w"])
2 # from statsmodels.stats.weightstats import DescrStatsW
3 # d = DescrStatsW(df["x"], weights=df["w"])
4 # d.mean, d.std
```

## 11) Visualisierung – Basics

```
1 plt.hist(df["x"].dropna()); plt.show()
2 plt.boxplot(df["x"].dropna()); plt.show()
3 plt.scatter(df["a"], df["b"]); plt.show()
```

## 12) Visualisierung – Seaborn

```
1 sns.histplot(data=df, x="x", bins=30)
2 sns.boxplot(data=df, x="cat", y="x")
3 sns.countplot(data=df, x="cat")
4 sns.scatterplot(data=df, x="a", y="b", hue="Group")
5 sns.heatmap(df.corr(numeric_only=True), annot=False)
```

## 13) Gruppierung + Visualisierung

```
1 g = df.groupby("cat")["x"].mean().reset_index()
2 plt.bar(g["cat"], g["x"]); plt.show()
3
4 shares = df["cat"].value_counts(normalize=True)
5 shares.plot(kind="bar"); plt.show()
```

## 14) Titanic-Beispiele

```
1 pd.crosstab(df["Pclass"], df["Survived"], normalize="index")
2 pd.crosstab(df["Sex"], df["Survived"], normalize="index")
3 df.groupby("Pclass")["Age"].agg(["count", "mean", "median", "std"])
4 sns.histplot(data=df, x="Fare", bins=40)
```

## 15) Export von Ergebnissen

```
1 summary = df.groupby("Group")["x"].agg(["count", "mean", "median"]).
   reset_index()
2 summary.to_csv("summary.csv", index=False)
3 plt.savefig("boxplot.png", dpi=150)
```

## 16) Häufige Stolperfallen

- Zahlen als Strings → `to_numeric()` verwenden
- `describe()` ignoriert NAs
- Immer `reset_index()` nach `groupby()` bei Bedarf
- Bei Visualisierungen Achsen prüfen und beschriften

## 17) Mini-Workflow

1. Laden (`read_csv`, `info`)
2. Typen & NAs bereinigen
3. Kennzahlen, Häufigkeiten, Kreuztabellen
4. Visualisierung (Histogramm, Box, Count, Scatter)
5. Ergebnisse exportieren und dokumentieren