

# Hypothesentest

## z-Test

## t-Test

Peter Büchel

HSLU I

Stat: Block 06

# Hypothesentest: Beispiele

- Hypothesentests: Wichtiges statistisches Mittel um zu entscheiden, ob eine Messreihe zu einer gewisse Grösse „passt“
- Brauerei bestellt neue Abfüllmaschine für 500 ml Büchsen
- Abfüllmaschine füllt *nie genau* 500 ml ab, sondern nur *ungefähr* 500 ml
- Brauerei ist daran interessiert, dass die Abfüllmaschine möglichst genau abfüllt
  - ▶ Füllt die Maschine zuviel ab, so ist dies schlecht für die Brauerei, da sie zuviel Bier für denselben Preis verkauft
  - ▶ Füllt sie zuwenig ab, sind die Kunden und der Konsumentenschutz unzufrieden, da sie für den entsprechenden Preis zuwenig Bier bekommen

- Herstellerfirma behauptet, dass Maschine die Büchsen normalverteilt mit  $\mu = 500$  ml und  $\sigma = 1$  ml abfüllt
- Die Brauerei macht 100 Stichproben
- Der Mittelwert dieser Stichproben ist 499.57 ml
- Weniger als 500 ml, aber liegt dies noch innerhalb der Angaben  $\mu = 500$  ml und  $\sigma = 1$  ml des Herstellers der Abfüllanlage?
- Wie können wir dies überprüfen?

# Beispiele

- Allgemeiner: Sie stellen eine Maschine her und müssen sich auf die Angaben der Spezifikationen der Hersteller für die Bestandteile verlassen können
- Wie können Sie feststellen, dass die Bestandteile die Spezifikationen auch erfüllen?
- Anfrage beim Bundesamt für Statistik: Durchschnittliche Körpergrösse der erwachsenen Frauen liegt in der Schweiz bei 180 cm mit einer Standardabweichung von 10 cm
- Angabe ist gefühlsmässig wohl falsch, da viel zu hoch
- Wie können wir dies aber mathematisch überprüfen und begründen, ohne uns auf unser Gefühl zu verlassen?

# Ziel

- Ziel: Standardisiertes, reproduzierbares Verfahren einzuführen, mit dem wir entscheiden können, ob der Mittelwert einer Messreihe zu einem bestimmten „wahren“ Mittelwert  $\mu$  passt oder nicht
- *Achtung*: Das kommende Verfahren liefert *niemals einen Beweis*, dass beispielsweise eine Grösse nicht zu einer Messreihe passt
- Wir können mit statistischen Mitteln nur zeigen, dass diese Grösse *mit grosser Wahrscheinlichkeit* nicht zu dieser Messreihe passt
- Lesen Sie in der Zeitung „... mit Statistik bewiesen...“, ist das ein Blödsinn!

# Problemstellung

- Vorgehen beim Hypothesentest durch folgendes Beispiel erklären
- *Datensatz*: Methode zur Bestimmung der Schmelzwärme von Eis
- Wiederholte Messungen der freigesetzten Wärme beim Übergang von Eis bei  $-0.7^{\circ}\text{C}$  zu Wasser bei  $0^{\circ}\text{C}$  ergaben die Werte (in cal/g):

Methode A	79.98	80.04	80.02	80.04	80.03	80.03	80.04
Methode A	79.97	80.05	80.03	80.02	80.00	80.02	

- Messungen als Realisierungen von unabhängigen, identisch verteilten Zufallsvariablen  $X_i$  betrachten
- Beispiel: 2. Messwert  $x_2 = 80.04$  Realisierung der Zufallsvariable  $X_2$

# Allgemein

- Messdaten  $x_1, \dots, x_n$  als Realisierungen von

$$X_1, \dots, X_n \text{ i.i.d. } \sim \mathcal{N}(\mu, \sigma_X^2)$$

- Zwei Kennzahlen der Zufallsvariablen  $X_i$  sind:

$$E(X_i) = \mu \quad \text{und} \quad \text{Var}(X_i) = \sigma_X^2$$

- Normalfall: Kennzahlen *unbekannt*
- Rückschlüsse darüber aus Daten wie in Beispiel Methode A machen
- (Punkt-) Schätzungen für den Erwartungswert und die Varianz sind:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i \quad \hat{\sigma}_X^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

- Schätzer hier als Funktionen der Zufallsvariablen  $X_1, \dots, X_n$
- $\hat{\mu}$  und  $\hat{\sigma}_X^2$  selbst wieder Zufallsvariablen mit Verteilungseigenschaften von  $\hat{\mu}$
- Hier: Schätzungen für das Beispiel der Methode A



## Beispiel: Methode A

- Schätzungen für Mittelwert  $\mu$  und die Varianz  $\sigma_X^2$ :

$$\hat{\mu} = 80.02 \quad \text{und} \quad \hat{\sigma}_X^2 = 0.024^2$$

- Berechnung mit **Python**:

```
from pandas import Series
import numpy as np

methodeA = Series([79.98, 80.04, 80.02, 80.04, 80.03, 80.03,
80.04, 79.97, 80.05, 80.03, 80.02, 80.00, 80.02])
print(methodeA.mean())
methodeA.mean()
methodeA.std()

## 80.02076923076923
## 0.023965787580611863
```

- Problem: *Andere* Messreihen haben andere Schätzwerte

# Simulation von neuen Messreihen

- Neue Messreihen simulieren, mit „ähnlichen“ Werten wie Methode A
- *Annahme*: Messwerte in Methode A normalverteilt mit

$$\mu = 80 \quad \text{und} \quad \sigma_X^2 = 0.02^2$$

- Generierung mit `norm.rvs()` aus **Scipy Stats** Zufallszahlen, die dieser Verteilung folgen
- Hier: Messreihen der Länge 6
- Rundung der meisten Resultate meist auf zwei Nachkommastellen:  
`np.round(...,2)`
- ~~Mit `np.random(...)` werden Zufallszahlen festgelegt~~

## ● Python-Code

```
from scipy.stats import norm
np.random.seed(1)

methodeA_sim1 = Series(np.round(norm.rvs(size=6, loc=80, scale=
0.02),2))

methodeA_sim1

methodeA_sim1.mean()
methodeA_sim1.std()

## 0      80.03
## 1      79.99
## 2      79.99
## 3      79.98
## 4      80.02
## 5      79.95
## dtype: float64
## 79.99333333333333
## 0.028751811537128993
```

## ● Geschätzten Werte $\hat{\mu}$ und $\hat{\sigma}^2$ (leicht) anders als vorher

- Fünf Simulationen mit Mittelwerten nahe bei  $\mu = 80$

```
np.random.seed(17)
for i in range(5):
    methodeA_sim1 = Series(np.round(norm.rvs(size=6, loc=80, scale=
0.02),2))
    print("Mittelwert:", np.round(methodeA_sim1.mean(), 3))
    print("Standardabw.:", np.round(methodeA_sim1.std(), 3))
    print()

## Mittelwert: 80.01
## Standardabw.: 0.027
##
## Mittelwert: 80.007
## Standardabw.: 0.02
##
## Mittelwert: 79.992
## Standardabw.: 0.028
##
## Mittelwert: 79.995
## Standardabw.: 0.016
##
## Mittelwert: 79.992
## Standardabw.: 0.013
```

- Mittelwert kann auch „weit“ von  $\mu = 80$  entfernt liegen:

```
np.random.seed(463137)

methodeA_sim2 = Series(np.round(norm.rvs(size=6, loc=80, scale=
0.02),2))

methodeA_sim2

methodeA_sim2.mean()
methodeA_sim2.std()

## 0      80.07
## 1      80.06
## 2      80.03
## 3      80.03
## 4      80.02
## 5      80.03
## dtype: float64
## 80.04
## 0.01999999999999998862
```

- Mittelwert dieser Messreihe ist 2 Standardabweichungen grösser als 80, was eben möglich ist, aber nicht sehr wahrscheinlich

- Aber was heisst hier „nicht sehr wahrscheinlich“?
- Verteilung des Mittelwertes dieser Messreihe:

$$\bar{X} \sim \mathcal{N}\left(80, \frac{0.02^2}{6}\right)$$

- W'keit, einen Mittelwert von *80.04 oder höher* erhalten ist  $5 \cdot 10^{-7}$ :

$$P(\bar{X}_6) \geq 80.04$$

```
from scipy.stats import norm
import numpy as np

1-norm.cdf(x=80.04, loc=80, scale=0.02/np.sqrt(6))

## 4.816785043049165e-07
```

- Punktw'keit  $P(\bar{X}_6 = 80.04)$  ist 0 und hier nichts
- Obwohl Daten zufällig aus der Verteilung  $\mathcal{N}(80, 0.02^2)$  stammen, ist diese Wahrscheinlichkeit sehr klein
- Wir können daran zweifeln, ob der wahre Mittelwert wirklich 80 ist

- Ein weiteres Beispiel:

```
np.random.seed(647)
methodeA_sim3 = Series(np.round(np.random.normal(size=6, loc=80,
scale= 0.02),2))

methodeA_sim3

methodeA_sim3.mean()
methodeA_sim3.std()

## 0      79.98
## 1      79.99
## 2      80.00
## 3      79.93
## 4      80.00
## 5      79.98
## dtype: float64
## 79.98
## 0.02607680962080759
```

- Mittelwert eine Standardabweichung unter 80
- W'keit,

$$P(\overline{X}_6) \leq 79.98$$

dass ein Messwert 79.98 *oder kleiner* ist, ist 0.007

```
from scipy.stats import norm
import numpy as np

norm.cdf(x=79.98, loc=80, scale=0.02/np.sqrt(6))

## 0.007152939217724509
```

- Möglich, aber doch wohl eher unwahrscheinlich, wenn wir  $\mu = 80$  annehmen
- Frage: Welche Mittelwerte sind noch tolerierbar oder bei welchen Werten beginnen wir am wahren (unbekannten) Wert  $\mu = 80$  zu zweifeln



# Fragestellungen:

- Ist eine Messreihe mit der Annahme  $\mu = 80$  noch kompatibel oder müssen an dieser Annahme zweifeln?
- Das heisst: Liegt der Mittelwert der Messreihe in der „Nähe“ des wahren Mittelwertes  $\mu = 80$  oder liegt er so „weit“ entfernt, dass wir an der Angabe des wahren  $\mu = 80$  zweifeln müssen?
- Frage, was heisst „nahe“ oder „weit“
- Beachte: Der wahre Mittelwert ist grundsätzlich *nicht* bekannt
- Verwenden Hypothesentest

# Vorgehen Hypothesentest

- Annahme: Daten normalverteilt sind mit  $\mu = 80.00$  und  $\sigma = 0.02$
- Wie können wir überprüfen, ob der Mittelwert  $\mu = 80$  auch stimmt?
- Grundidee: Mit einer Messreihe überprüfen, ob unter dieser Annahme  $\mu = 80$ , die Messreihe wahrscheinlich ist oder nicht
- Wählen dazu eine Messreihe der Länge 6 aus mit Modell

## Modell

Die 6 Messwerte sind Realisierungen der ZV  $X_1, X_2, \dots, X_6$ , wobei  $X_i$  eine kontinuierliche Messgrösse ist. Es soll gelten:

$$X_1, \dots, X_6 \text{ i.i.d. } \sim \mathcal{N}(80, 0.02^2)$$

- Wir wollen nun überprüfen, ob die *Annahme*  $\mu = 80$  auch gerechtfertigt ist
- Dazu führen wir folgende Begriffe ein

### **Nullhypothese**

$$H_0 : \mu = \mu_0 = 80$$

### **Alternativhypothese**

$$H_A : \mu \neq \mu_0 = 80 \quad (\text{oder „<“ oder „>“})$$

- Wählen Messreihe:

```
## 0    79.98
## 1    79.99
## 2    80.00
## 3    79.93
## 4    80.00
## 5    79.98
## dtype: float64
## Mittelwert: 79.98
```

- Der (geschätzte) Mittelwert ist hier  $\hat{\mu} = 79.98$
- Konkretisierung, was es heisst, dass dieser Mittelwert (un)wahrscheinlich ist
- W'keit

$$P(\overline{X}_6 = 79.98)$$

bringt uns hier nicht weiter, da diese 0 ist

- Da  $\hat{\mu} < 80$  ist, können wir aber folgende W'keit betrachten:

$$P(\overline{X}_6 \leq 79.98)$$

- Unter Annahmen  $\mu = 80$  und  $\sigma = 0.02$  ist  $\overline{X}_6$  wie folgt verteilt

$$\overline{X}_6 \sim \mathcal{N}\left(80, \frac{0.02^2}{6}\right)$$

- Testen mit dieser Verteilung, ob Annahme  $\mu = 80$  gerechtfertigt ist

### **Teststatistik**

Verteilung der Teststatistik  $T$  unter der Nullhypothese  $H_0$ :

$$T: \quad \overline{X}_6 \sim \mathcal{N}\left(80, \frac{0.02^2}{6}\right)$$

- W'keit:

$$P(\bar{X}_6 \leq 79.98) = 0.007$$

```
norm.cdf(x=79.98, loc=80, scale=0.02/np.sqrt(6))  
## 0.007152939217724509
```

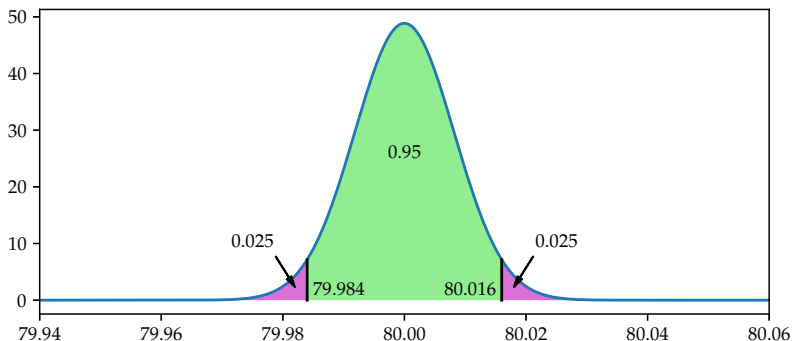
- Diese W'keit ist klein, 0.7 %
- Ist sie aber *zu* klein?
- *Abmachung*: Es hat sich als praktisch erwiesen, diese Grenze, was zu klein ist und was nicht bei 2.5 % festzulegen
- Gemäss dieser Abmachung ist

$$P(\bar{X}_6 \leq 79.98) < 0.025$$

- Betrachten diesen geschätzten Mittelwert  $\hat{\mu} = 79.98$  als *zu unwahrscheinlich*, als dieser zum Wert  $\mu = 80$  passen könnte
- *Wir gehen also davon aus, dass der angegebene Mittelwert von  $\mu = 80$  nicht stimmen kann!*
- Wir sagen: „Wir verwerfen die Nullhypothese und nehmen die Alternativhypothese an!“

# Graphische Darstellung

- Normalverteilungskurve in drei Teile aufteilen:
  - ▶ Symmetrischer Teil um Mittelwert  $\mu = 80$  soll 0.95, also 95 % betragen
  - ▶ Die beiden Teilen links und rechts müssen zusammen 0.05 ergeben
  - ▶ Also ergibt sich für jeden Teil 0.025
  - ▶ Abbildung:





- Begriff:

### **Signifikanzniveau $\alpha$**

Das Signifikanzniveau  $\alpha$ , gibt an, wie hoch das Risiko ist, das man bereit ist einzugehen, eine falsche Entscheidung zu treffen.

Für die meisten Tests wird ein  $\alpha$ -Wert von 0.05 bzw. 0.01 verwendet.

Wir verwenden hier

$$\alpha = 0.05$$

- Signifikanzniveau legt roten Bereich in Abbildung vorher fest
- Der rote Bereich in der Abbildung heisst *Verwerfungsbereich*

- Grenzen des Verwerfungsbereichs entsprechen den 0.025- und 0.975-Quantilen.

```
norm.ppf(q=0.025, loc=80, scale=0.02/np.sqrt(6))  
norm.ppf(q=0.975, loc=80, scale=0.02/np.sqrt(6))  
  
## 79.98399696107882  
## 80.01600303892118
```

- Liegt der gemessene Mittelwert im roten Bereich in Abbildung, so verwerfen wir die Nullhypothese, dass  $\mu = 80$
- Wir nennen diesen Bereich deshalb

### Verwerfungsbereich

$$K = (-\infty, 79.984] \cup [80.016, \infty)$$

- Wir gehen davon aus, dass ein Mittelwert einer Messreihe im Verwerfungsbereich so unwahrscheinlich ist, dass wir an der Richtigkeit von  $\mu = 80$  zweifeln
- Mit unserer Messreihe überprüfen, ob deren Mittelwert im Verwerfungsbereich liegt oder nicht
- Machen den sogenannten

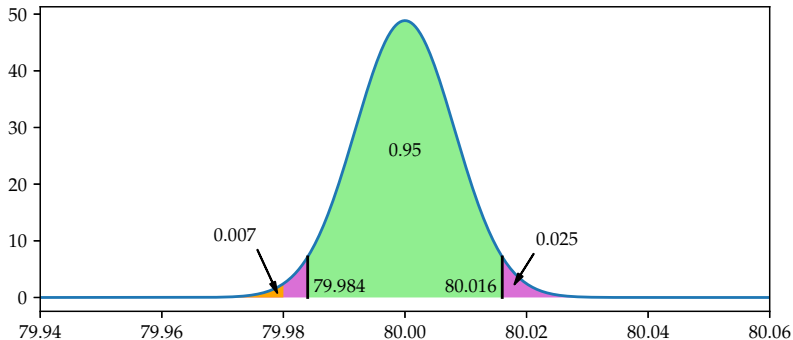
### Testentscheid

In unserem Beispiel hatten wir

$$\bar{X}_6 = 79.98 \in K$$

- Dieser Wert liegt im Verwerfungsbereich
- Wir verwerfen die Nullhypothese und nehmen Alternativhypothese an

- Graphisch:



# Bemerkungen

- Warum haben wir hier den Verwerfungsbereich nach oben und nach unten aufgeteilt, wenn wir schon wissen, dass der gemessene Mittelwert kleiner als  $\mu = 80$  ist?
- Nun, das wussten wir *vor* der Messung nicht
- Der gemessene Mittelwert hätte also auch grösser als  $\mu = 80$  sein können (siehe Beispiel nachher)
- Wir sprechen in diesem Fall von einem *zweiseitigen Test*
- Es gibt auch *einseitige Tests* (siehe Beispiel später)

- Wir haben hier eine *Annahme* gemacht, dass der gesamte Verwerfungsbereich 5 % (Signifikanzniveau 5 %) betragen soll
- Diese Annahme hat sich als praktisch erwiesen, aber wir hätten auch 1 % wählen können, was auch ab und zu gemacht wird
- In Beispiel oben folgt die zufällige Messreihe der Normalverteilung  $\mathcal{N}(80, 0.02^2/6)$  und doch wird der Parameter  $\mu = 80$  hier als unwahrscheinlich verworfen
- Dies heisst, wir haben hier einen *Fehler* gemacht
- Auf diese Problematik gehen wir später ein

- Wählen als *andere* Messreihe mit Modell, Nullhypothese, Alternativhypothese, Teststatistik, Signifikanzniveau und Verwerfungsbereich gleich wie im Beispiel vorher
- Wir müssen also nur noch den Testentscheid durchführen.

```
## 0      80.07
## 1      80.06
## 2      80.03
## 3      80.03
## 4      80.02
## 5      80.03
## dtype: float64
## Mittelwert: 80.04
```

- Der geschätzte Mittelwert ist im Verwerfungsbereich und somit wird auch hier die Nullhypothese verworfen.

- W'keit:

$$P(\overline{X}_6 > 80.04)$$

```
1-norm.cdf(x=80.04, loc=80, scale=0.02/np.sqrt(6))  
## 4.816785043049165e-07
```

- Bei weitem kleiner als 0.025 und damit so unwahrscheinlich, dass wir auch auf diese Weise  $\mu = 80$  als nicht richtig annehmen (müssen)
- Wir *verwerfen* die Nullhypothese
- Verwerfungsbereich: Nur Entscheidung möglich, ob der geschätzte Mittelwert im Verwerfungsbereich liegt oder nicht
- Wert von  $P(\overline{X}_6 > 80.04)$  macht noch Aussage über die Sicherheit des Verwerfen
- Hier:  $5 \cdot 10^{-7}$  *sehr viel kleiner* als 0.025 und damit können wir mit grosser Sicherheit davon ausgehen, dass  $\mu = 80$  *nicht* gilt



# Grössere Messreihen

- Beispiel vorher: Grösse  $\mu = 80$  verworfen, da die Messreihe einen zu tiefen und dann einen (viel) zu grossen Mittelwert lieferte
- Wie sieht es nun aber aus, wenn wir eine neue Messreihe bilden, die aus *beiden* Messreihen besteht?
- Oder anders gefragt: Welchen Einfluss hat die Anzahl der Messungen auf den Verwerfungsbereich?
- Wählen Messreihen verschiedener Länge  $n$ , die alle den geschätzten Mittelwert  $\hat{\mu} = 79.78$  haben

- Dann bestimmen wir für alle Messreihen den Wert

$$P(\overline{X}_n \leq 79.98)$$

mit

$$\overline{X}_6 \sim \mathcal{N}\left(80, \frac{0.02^2}{n}\right)$$

- Ist dieser Wert grösser als 0.025, dann wird die Nullhypothese nicht verworfen, ansonsten schon
- Für  $n = 2$  erhalten wir folgenden Wert für

$$P(\overline{X}_2 \leq 79.98) = 0.079 > 0.025$$

```
norm.cdf(x=79.98, loc=80, scale=0.02/np.sqrt(2))  
## 0.07864960352518385
```

- Nullhypothese wird also auf Signifikanzniveau 5 % nicht verworfen

- Für  $n = 4$  erhalten wir

$$P(\bar{X}_4 \leq 79.98) = 0.022 < 0.025$$

```
norm.cdf(x=79.98, loc=80, scale=0.02/np.sqrt(4))  
## 0.022750131948200674
```

- Hier wird die Nullhypothese (knapp) verworfen
- Für  $n = 6$  erhalten wir

$$P(\bar{X}_6 \leq 79.98) = 0.007 < 0.025$$

```
norm.cdf(x=79.98, loc=80, scale=0.02/np.sqrt(6))  
## 0.007152939217724509
```

- Die Nullhypothese wird klarer verworfen als für  $n = 4$

- Und schlussendlich noch für  $n = 8$ :

$$P(\bar{X}_6 \leq 79.98) = 0.002 < 0.025$$

```
norm.cdf(x=79.98, loc=80, scale=0.02/np.sqrt(8))  
## 0.0023388674905277422
```

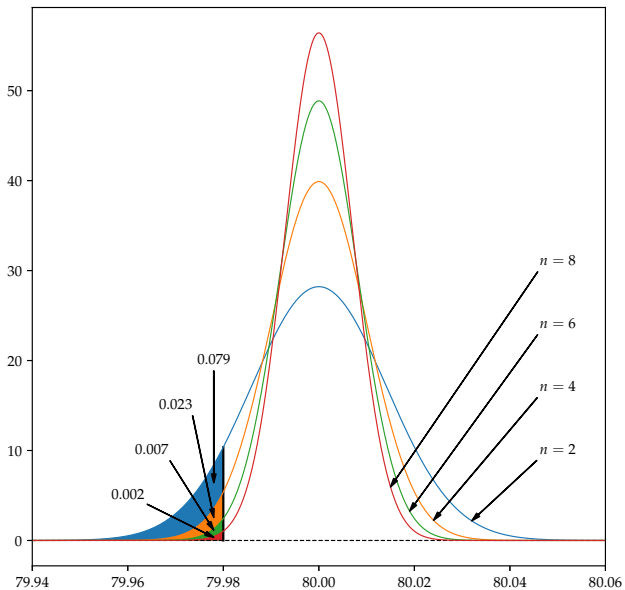
- Die Nullhypothese wird noch klarer verworfen, als bei  $n = 8$
- Mit zunehmendem  $n$  wird der Wert

$$P(\bar{X}_n \leq 79.98)$$

immer kleiner

- Dies liegt daran, dass die Standardabweichung mit grösser werdendem  $n$  kleiner wird und damit werden die Normalverteilungskurven schmaler (Abbildung nächste Folie).

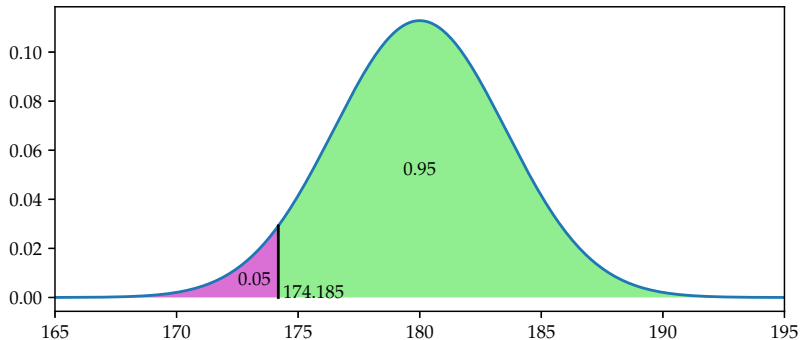
● Abbildung:



## Beispiel: Körpergrösse Frauen

- Bundesamt für Statistik behauptet, dass die durchschnittliche Körpergrösse der erwachsenen Frauen in der Schweiz bei 180 cm mit einer Standardabweichung von 10 cm liegt
- Vermutung: Dieser Mittelwert ist zu gross ist
- Hier macht ein zweiseitiger Test wenig Sinn, da wir „wissen“, dass dieser Mittelwert zu gross ist
- Das heisst, der wahre Wert liegt wohl eher tiefer
- Überlegung ähnlich vorher, aber Verwerfungsbereich nicht auf beide Seiten verteilen, sondern nur nach unten, da wir erwarten, dass der wahre Mittelwert tiefer als  $\mu = 180$  ist (Abbildung nächste Folie)
- Wir machen einen *einseitigen* Test

● Abbildung:



- *Modell:*

$$X_1, \dots, X_n \text{ i.i.d. } X_i \sim \mathcal{N}(180, 10^2)$$

- Annahme: Der wahre Mittelwert ist 180 cm

- *Nullhypothese:*

$$H_0 : \mu_0 = 180$$

- *Alternativhypothese:*

$$H_A : \mu < \mu_0 = 180$$

- Untersuchung unter  $n$  Personen und testen ob jetzt der Wert

$$P(\overline{X}_n < \overline{x}_n) < 0.05$$

ist oder nicht

- Der Verwerfungsbereich ist hier also einseitig nach unten



- In Abbildung vorher ist der Verwerfungsbereich für  $n = 8$  eingezeichnet pink eingezeichnet
- *Teststatistik unter der Nullhypothese  $H_0$ :*

$$\bar{X}_8 \sim \mathcal{N}\left(180, \frac{10^2}{8}\right)$$

- *Signifikanzniveau:*

$$\alpha = 0.05$$

- Grenze des Verwerfungsbereichs:

```
norm.ppf(q=0.05, loc=180, scale=10/np.sqrt(8))  
## 174.18456423161663
```

- *Verwerfungsbereich* (siehe Abbildung vorher):

$$K = (-\infty, 174.185)$$

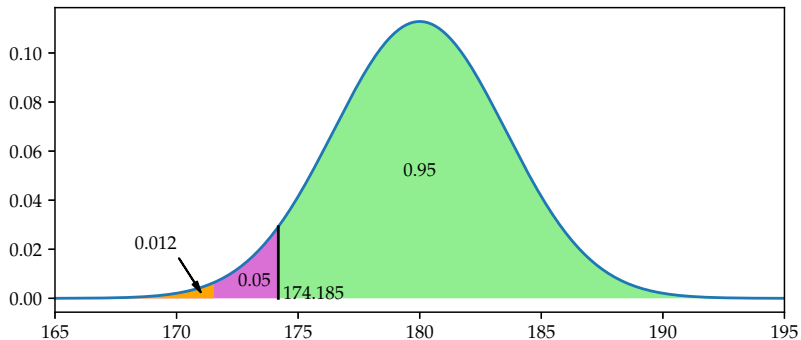
- Dieser Verwerfungsbereich ist natürlich viel zu gross, da wohl kaum Körpergrössen von erwachsenen Frauen unter 50 cm zu erwarten sind
- Wir arbeiten hier mit einem *Modell*, das eben nur in einem bestimmten Bereich Sinn macht
- Wählen wir nun zufällig acht erwachsene Frauen aus, messen deren Körpergrösse und bestimmen den Mittelwert, der bei 171.54 cm liegt
- *Testentscheid:*  
So ist Wert im Verwerfungsbereich und somit *verwerfen* wir die Nullhypothese, dass das wahre  $\mu = 180$  gilt
- Dieser Mittelwert der zufällig ausgewählten acht Frauen erscheint immer noch relativ hoch, aber er reicht schon, damit wir an der Annahme  $\mu = 180$  zweifeln müssen.

- Wert für  $P(\bar{X}_6 < 172)$  ist (siehe Abbildung unten).

$$P(\bar{X}_6 < 172) = 0.012$$

```
norm.cdf(x=172, loc=180, scale=10/np.sqrt(8))  
## 0.011825808327677984
```

- Abbildung:



- Dieser Wert heisst  $P$ -Wert und gibt die Sicherheit mit der wir den Testentscheid treffen
- Wird die Nullhypothese verworfen, so deutet ein sehr kleiner  $P$ -Wert darauf hin, dass die Nullhypothese sicherer verworfen wird, als wenn er in der Nähe des Signifikanzniveaus (hier  $\alpha = 0.05$ ) ist.

# P-Wert

- $P$ -Wert ist ein Wert zwischen 0 und 1, der angibt, wie gut *Nullhypothese* und *Daten* zusammenpassen (0: passt gar nicht; 1: passt sehr gut)
- Genauer: Definieren wir den  $P$ -Wert als die W'keit, unter Gültigkeit der Nullhypothese das erhaltene Ergebnis oder ein *extremes* zu erhalten
- Mit dem  $P$ -Wert wird also angedeutet, wie extrem das Ergebnis ist: Je kleiner der  $P$ -Wert, desto mehr spricht das Ergebnis gegen die Nullhypothese. Werte kleiner als eine im voraus festgesetzte Grenze, wie 5 %, 1 % oder 0.1 % sind Anlass, die Nullhypothese abzulehnen

## **$P$ -Wert**

Der  $P$ -Wert ist die Wahrscheinlichkeit, unter der Nullhypothese ein mindestens so extremes Ereignis (in Richtung der Alternative) zu beobachten wie das aktuell beobachtete.

- Testentscheid mit Hilfe des  $P$ -Wertes durchführen

### **$P$ -Wert und Statistischer Test**

Man kann anhand des  $P$ -Werts direkt den Testentscheid ablesen: Wenn der  $P$ -Wert kleiner als das Niveau ist, so verwirft man  $H_0$ , ansonsten nicht.

Verglichen mit dem reinen Testentscheid enthält der  $P$ -Wert aber mehr Information, da man direkt sieht, „wie stark“ die Nullhypothese verworfen wird.

Bei einem vorgegebenen Signifikanzniveau  $\alpha$  (z.B.  $\alpha = 0.05$ ) gilt aufgrund der Definition des  $P$ -Werts für einen einseitigen Test:

- 1 Verwerfe  $H_0$  falls  $P\text{-Wert} \leq \alpha$
- 2 Belasse  $H_0$  falls  $P\text{-Wert} > \alpha$

- Viele Computer-Pakete liefern den Testentscheid nur indirekt, indem der  $P$ -Wert angegeben wird
- Zusätzlich zu dieser Entscheidungsregel quantifiziert der  $P$ -Wert, wie *signifikant* eine Alternative ist (d.h. wie gross die Evidenz ist für das Verwerfen von  $H_0$ )
- Manchmal werden sprachliche Formeln oder Symbole anstelle der  $P$ -Werte angegeben:

$P\text{-Wert} \approx 0.05$  : schwach signifikant, “.”

$P\text{-Wert} \approx 0.01$  : signifikant, “\*”

$P\text{-Wert} \approx 0.001$  : stark signifikant, “\*\*\*”

$P\text{-Wert} \leq 10^{-4}$  : "ausserst signifikant, “\*\*\*\*”

- *Achtung:* Der  $P$ -Wert ist nicht die Wahrscheinlichkeit, dass die Nullhypothese stimmt. Darüber können wir hier gar keine Aussagen machen, da die Parameter fix und nicht zufällig sind.

## P-Wert für zweiseitigen Test

- Wir haben den  $P$ -Wert für einseitige Tests definiert
- Wie sieht nun aber der  $P$ -Wert für zweiseitige Tests aus?
- Beispiel von früher:

$$P(\bar{X}_6 \leq 79.98) = 0.007$$

der kleiner ist als 0.025

- Wir könnten dies als  $P$ -Wert betrachten
- Da aber das Signifikanzniveau auf  $\alpha = 0.05$  liegt, wird die W'keit oben auf 5 % umgerechnet, also verdoppelt:

$$P\text{-Wert} = 2 \cdot P(\bar{X}_6 \leq 79.98) = 0.014$$

- Dieser  $P$ -Wert dann mit dem Signifikanzniveau verglichen
- Computersoftware gibt den  $P$ -Wert *immer* auf Signifikanzniveau an.



# Beispiel: Statistischer Test für Durchschnittsgrösse

- Durchschnittsgrösse der Schweizer Frauen beträgt 1.64 m (Bundesamt für Statistik)
- Stimmt diese Aussage?
- Wie kann man sie überprüfen?
- *Lösung:* In Stadt Luzern zum Beispiel Grösse von 150 Frauen messen
- Gemessene Körpergrössen  $x_1, \dots, x_{150}$  als Realisierungen auffassen von

$$X_1, \dots, X_{150} \text{ i.i.d. } \sim \mathcal{N}(\mu, \sigma_X^2)$$

# Beispiel: Statistischer Test für Durchschnittsgrösse

- Bestimmen das arithmetische Mittel  $\bar{x}_{150}$  der 150 Messpunkte
- *Frage:* Wenn jeden Tag eine solche Messreihe erheben würden, was wird man feststellen?
- *Antwort:* Arithmetischen Mittel  $\bar{x}_{150}$  würden variieren
- Das gemessene arithmetische Mittel  $\bar{x}_{150}$  wird als Realisierung der Zufallsvariablen  $\bar{X}_{150}$  aufgefasst
- *Frage:* Welche Verteilung hat die Durchschnittsgrösse von 150 Frauen?

# Beispiel: Statistischer Test für Durchschnittsgrösse

- *Antwort:* Verteilung vom arithmetischen Mittel  $\bar{X}_{150}$ :

$$\bar{X}_{150} \sim \mathcal{N}(\mu, \sigma_{\bar{X}_{150}}^2)$$

- Wie entscheidet man, ob die gemessenen Durchschnittsgrössen zu dem vom Statistikamt angegebenen Wert von 1.64 m passt?
- *Lösung:* Berechnen W'keit, dass gemessener Wert  $\bar{X}_{150}$  oder einen extremeren Wert beobachtet wird, unter der Annahme, dass

$$\bar{X}_{150} \sim \mathcal{N}\left(\mu = 1.64, \frac{\sigma_X^2}{150}\right)$$

# Beispiel: Statistischer Test für Durchschnittsgrösse

- P-Wert bei einseitiger nach oben gerichteter Alternativhypothese:

$$P\left(\bar{x}_{150} < \bar{X}_{150}\right)$$

- P-Wert bei einseitiger nach unten gerichteter Alternativhypothese:

$$P\left(\bar{X}_{150} < \bar{x}_{150}\right)$$

- P-Wert bei zweiseitiger Alternativhypothese:

$$P\left(\bar{x}_{150} < |\bar{X}_{150}|\right)$$

# Beispiel: Statistischer Test für Durchschnittsgrösse

- Unterscheidung von zwei Fällen:
  - ▶  $\sigma_X$  ist bekannt (aus langjähriger Erfahrung)  $\rightarrow$  *z-Test*
  - ▶  $\sigma_X$  wurde aus den Daten geschätzt, d.h.,  $\hat{\sigma}_X$  bekannt (was in der Praxis normalerweise der Fall ist)  $\rightarrow$  *t-Test*

# z-Test: $\sigma_X$ bekannt

1. *Modell:*  $X_i$  ist eine kontinuierliche Messgrösse:

$$X_1, \dots, X_n \text{ iid } \mathcal{N}(\mu, \sigma_X^2), \sigma_X \text{ bekannt}$$

2. *Nullhypothese:*

$$H_0 : \mu = \mu_0$$

*Alternative:*

$$H_A : \mu \neq \mu_0 \quad (\text{oder } "<" \text{ oder } ">")$$

3. *Teststatistik:*

$$T = \bar{X}_n$$

*Verteilung der Teststatistik unter  $H_0$*

$$T \sim \mathcal{N}\left(\mu_0, \frac{\sigma_X^2}{n}\right)$$

## z-Test: $\sigma_X$ bekannt

oder durch Normalisierung

$$Z = \frac{(\bar{X}_n - \mu_0)}{\sigma_{\bar{X}_n}} = \frac{(\bar{X}_n - \mu_0)}{\sigma_X / \sqrt{n}} = \frac{\text{beobachtet} - \text{erwartet}}{\text{Standardfehler}}$$

Verteilung der Teststatistik unter  $H_0$

$$Z \sim \mathcal{N}(0, 1)$$

4. Signifikanzniveau:

$$\alpha$$

5. Verwerfungsbereich für die Teststatistik:

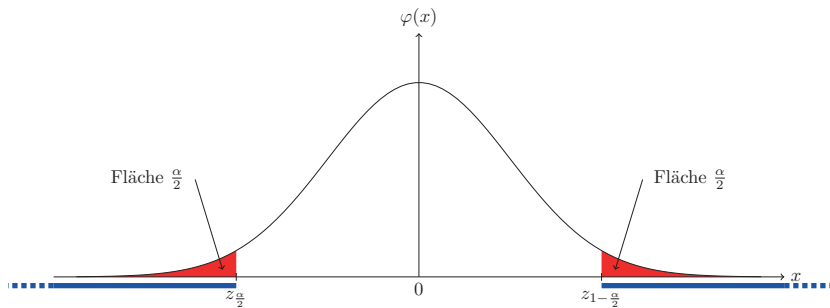
$$K = (-\infty, z_{\frac{\alpha}{2}}] \cup [z_{1-\frac{\alpha}{2}}, \infty) \quad \text{bei } H_A : \mu \neq \mu_0$$

$$K = (-\infty, z_{\alpha}] \quad \text{bei } H_A : \mu < \mu_0$$

$$K = [z_{1-\alpha}, \infty) \quad \text{bei } H_A : \mu > \mu_0$$

6. Testentscheid: Überprüfe, ob der beobachtete Wert der Teststatistik im Verwerfungsbereich liegt

# Zweiseitiger Verwerfungsbereich beim Z-Test





# Beispiel: z-Test

- Messreihe für die Körpergrösse von 150 Frauen in Luzern ergab:

$$\bar{x}_{150} = 168 \text{ cm}$$

- Vermutung: Durchschnittsgrösse der Schweizerinnen ist grösser als 164 cm (einseitiger nach oben gerichteter Test)
- $\sigma_X$  sei bekannt und beträgt 8 cm

# Beispiel: z-Test

- Berechnung vom P-Wert (mit **Python** )

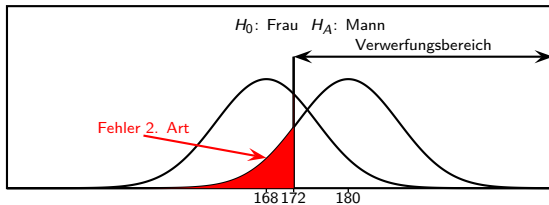
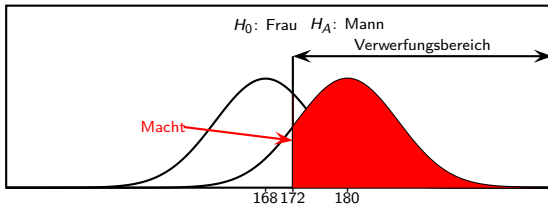
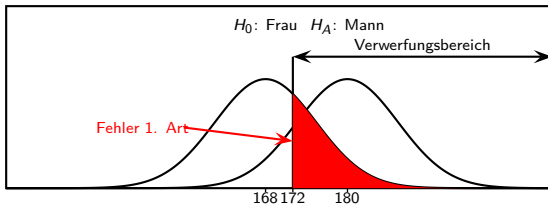
$$P[\bar{X}_{150} > 168] = 1 - P[\bar{X}_{150} \leq 168]$$

```
from scipy.stats import norm
import numpy as np

1-norm.cdf(x=168, loc=164, scale=8/np.sqrt(150))

## 4.5706494145036913e-10
```

- Nullhypothese auf dem Signifikanzniveau  $\alpha = 0.05$  verworfen werden
- Angabe vom Bundesamt stimmt nicht



## Problem in Praxis: $\sigma_X$ ist unbekannt!

- Falls  $\sigma_X$  unbekannt  $\rightarrow$  Varianz aus den Daten schätzen:

$$\hat{\sigma}_X^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

- Neue Teststatistik:

$$T = \frac{\bar{X}_n - \mu_0}{\frac{\hat{\sigma}_X}{\sqrt{n}}}$$

- Verteilung von  $T$ , falls  $H_0$  stimmt:

$$T \sim t_{n-1}$$

- $t_{n-1}$  ist die sogenannte *t-Verteilung mit  $n - 1$  Freiheitsgraden*

# „Student's" $t$ -Verteilung

- Annahme:

$$X_1, X_2, \dots, X_n \sim \mathcal{N}(\mu, \sigma_X^2) \quad \text{und unabhängig}$$

- Geschätzte Varianz

$$\hat{\sigma}_X^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

- Zufallsvariable

$$T = \frac{\bar{X}_n - \mu}{\frac{\hat{\sigma}_X}{\sqrt{n}}}$$

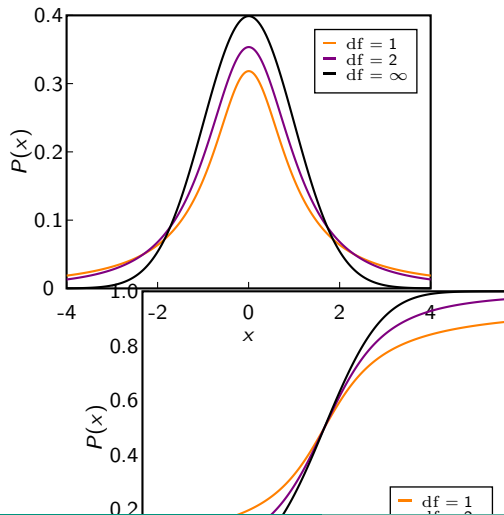
- $T$  folgt einer "t-Verteilung mit  $n - 1$  Freiheitsgraden"

$$T \sim t_{n-1}$$

# „Student's" $t$ -Verteilung

- Werte mit Computer ermittelbar

- Falls  $n = \infty$  :  $t_n = \mathcal{N}(0, 1)$



# $t$ -Verteilung: Eigenschaften

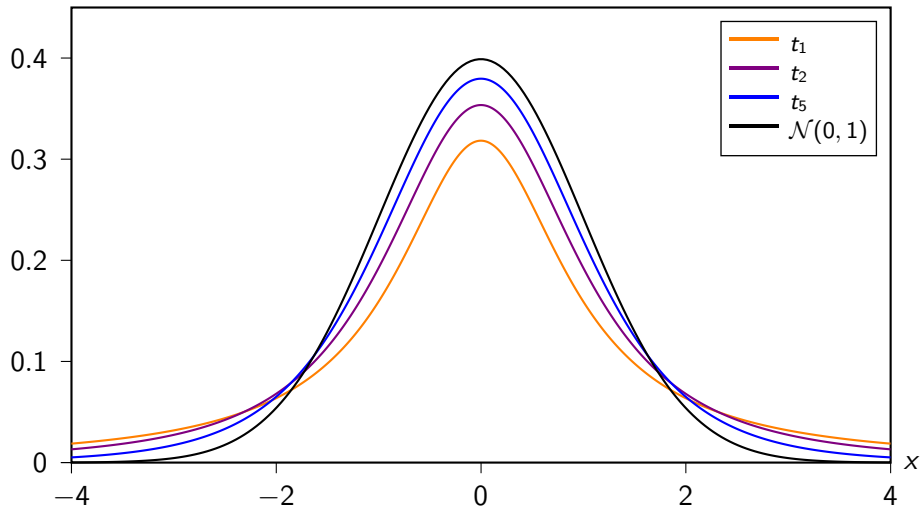
- Wie Standardnormalverteilung:  $t$ -Verteilung symmetrisch um 0
- Allerdings:  $t$ -Verteilung "langschwänziger"
- Peak in der Mitte ist weniger hoch und "weit aussen" ist die Dichte grösser (insbesondere falls die Anzahl Freiheitsgrade  $n$  klein ist)
- Verglichen mit Standardnormalverteilung liefert sie *eher grosse Werte*
- Es gilt:

$$t_n \rightarrow \mathcal{N}(0, 1) \quad \text{für} \quad n \rightarrow \infty$$

(siehe auch Plot mit Dichten)

# $t$ -Verteilung: Dichten

Dichte  $f(x)$





## t-Test: $\sigma_X$ unbekannt

1. Modell:  $X_i$  ist eine kontinuierliche Messgröße;

$X_1, \dots, X_n$  iid  $\mathcal{N}(\mu, \sigma_X^2)$ ,  $\sigma_X$  wird durch  $\hat{\sigma}_X$  geschätzt

2. Nullhypothese:

$$H_0 : \mu = \mu_0$$

Alternative:

$$H_A : \mu \neq \mu_0 \quad (\text{oder } "<" \text{ oder } ">")$$

3. Teststatistik:

$$T = \frac{(\bar{X}_n - \mu_0)}{\hat{\sigma}_{\bar{X}_n}} = \frac{(\bar{X}_n - \mu_0)}{\hat{\sigma}_X / \sqrt{n}} = \frac{\text{beobachtet} - \text{erwartet}}{\text{Standardfehler}}$$

Verteilung der Teststatistik unter  $H_0$  :

$$T \sim t_{n-1}$$

## $t$ -Test: $\sigma_X$ unbekannt

4. Signifikanzniveau:

$$\alpha$$

5. Verwerfungsbereich für die Teststatistik:

$$K = (-\infty, t_{n-1; \frac{\alpha}{2}}] \cup [t_{n-1; 1-\frac{\alpha}{2}}, \infty) \quad \text{bei } H_A : \mu \neq \mu_0$$

$$K = (-\infty, t_{n-1; \alpha}] \quad \text{bei } H_A : \mu < \mu_0$$

$$K = [t_{n-1; 1-\alpha}, \infty) \quad \text{bei } H_A : \mu > \mu_0$$

6. Testentscheid: Überprüfe, ob der beobachtete Wert der Teststatistik im Verwerfungsbereich liegt

## Beispiel: $t$ -Test

- Messreihe für Körpergrösse von 150 Frauen in Luzern:

$$\bar{x}_{150} = 168 \text{ cm}$$

- Vermutung: Durchschnittsgrösse der Schweizerinnen entweder grösser oder kleiner als 164 cm ist (zweiseitiger Test)
- $\sigma_X$  wurde nun geschätzt, und beträgt  $\hat{\sigma}_X = 10 \text{ cm}$
- Berechnen (mit **Python**) P-Wert für eine einseitige Alternativhypothese

$$P[\bar{X}_{150} \geq 168] = 1 - P[X \leq 168] = 1 - P\left[T \leq \frac{168 - 164}{10/\sqrt{150}}\right]$$

# Beispiel: t-Test

- *Achtung:* Bei der  $\bar{X}_{150}$  muss Freiheitsgrad 149 gewählt werden!

- Python

```
from scipy.stats import t
1-t.cdf(x=168, df=149, loc=164, scale=10/np.sqrt(150))
## 1.241988350275669e-06
```

oder standardisiert:

```
1-t.cdf(x=(168-164)/(10/np.sqrt(150)), df=149)
## 1.241988350275669e-06
```

- P-Wert:  $2 \cdot 1.241988 \cdot 10^{-6}$
- Nullhypothese auf dem Signifikanzniveau  $\alpha = 0.05$  verwerfen