

Dipl. Soz. Maximilian Sonnauer

# Methoden II

Mittelwertvergleiche





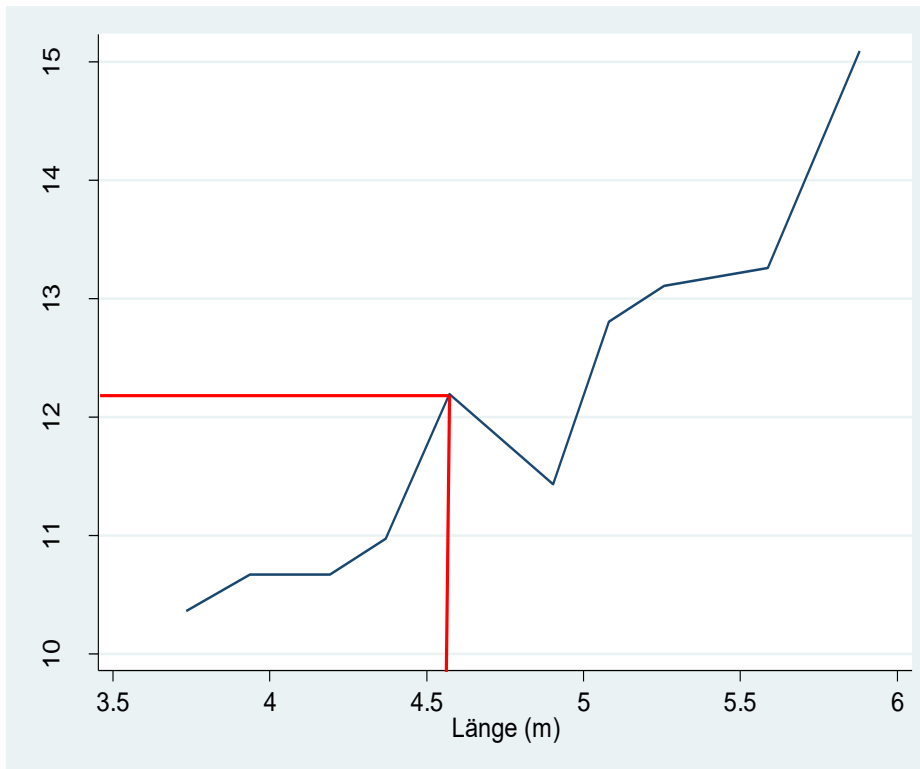
- **Termin:** 12.6.2017, 10.00 Uhr
- Geschrieben wird hier im Hörsaal (HGB-A240)
- Klausurmodus ist single-choice (nur eine Antwort ist richtig)
- Bearbeitungszeit: 45 Minuten
- **Stoff:** Inhalt sämtlicher Folien
- **Tutorium:** 07.06.2017 2 Termine von 16-18 & 18-20 Uhr c.t. im Raum 409 am Institut für Soziologie (Konradstraße 6)



1. Anmerkung letzte Woche
2. Wiederholungsfragen
3. Mittelwertvergleiche
  1. Einführung
  2. Ein-Stichproben t-test
  3. Zwei-Stichproben t-test
4. Übungsaufgabe



1. **Anmerkung letzte Woche**
2. Wiederholungsfragen
3. Mittelwertvergleiche
  1. Einführung
  2. Ein-Stichproben t-test
  3. Zwei-Stichproben t-test
4. Übungsaufgabe



## Median Band:

- Für Länge und Wendekreis werden  $k$  Abschnitte mit identischer Zahl von Beobachtungen gebildet
- Bei beiden Variablen werden lokale Mediane gebildet
- Grafik zeigt Schnittpunkt der Mediane im jeweiligen Abschnitt



1. Anmerkung letzte Woche
2. **Wiederholungsfragen**
3. Mittelwertvergleiche
  1. Einführung
  2. Ein-Stichproben t-test
  3. Zwei-Stichproben t-test
4. Übungsaufgabe



1. Anmerkung letzte Woche
2. Wiederholungsfragen
3. **Mittelwertvergleiche**
  1. **Einführung**
  2. Ein-Stichproben t-test
  3. Zwei-Stichproben t-test
4. Übungsaufgabe



## Bisher:

- Gibt es einen Zusammenhang zwischen zwei Variablen?

„Klassische“ Beispiele:

- Soziale Herkunft und erzielte Bildungsabschlüsse
- Einfluss von Bildung (Humankapital) auf Einkommen





## Vergleich mit bekanntem Wert aus Grundgesamtheit (Repräsentativität):

- Ist eine Stichprobe möglicherweise verzerrt, weil bestimmte Gruppen über-/ unterrepräsentiert sind?
  - Alter
  - Einkommensschichten
  - Herkunft



## Zusätzlich relevante soziologische Frage:

- Unterscheiden sich Zusammenhänge zwischen Gruppen?
- Unterscheidet sich der Einfluss sozialer Herkunft auf erzielte Bildungsabschlüsse **zwischen Einheimischen und Migrant\*innen**?
- **Geschlechterunterschiede** im Einfluss von Bildung (Humankapital) auf Einkommen



- **Zur Beantwortung dieser Fragen wendet man in der Sozialforschung Mittelwertvergleiche an**
  - Ist das Haushaltseinkommen in einer Stichprobe größer als in gesamt Deutschland?
  - Verdienen Männer mehr als Frauen?



- Mittelwertunterschiede in der Stichprobe können entweder auf einen **realen Unterschied in der Grundgesamtheit** zurück geführt werden, oder auf einen **zufällig auftretenden Unterschied**, der durch die Zufallsauswahl entstanden ist.
- Um zu prüfen, ob ein Mittelwertunterschied rein zufällig ist oder ein relevanter Unterschied in der Grundgesamtheit vorliegt, werden **statistische Testverfahren eingesetzt**.



- **Entscheidungskriterien bei der Auswahl des Verfahrens:**
  - Unabhängige oder abhängige Messung der Variable(n):  
Unabhängig: Bsp. Einkommen von Frauen und Männern;  
Abhängig: Einkommen der Eltern und eigenes Einkommen
  - Zahl der Ausprägungen der unabhängigen Variable: Zwei  
Ausprägungen: Mittelwertvergleich; Mehr als zwei  
Ausprägungen: Regressionsanalyse
  - Skalenniveau der abhängigen Variable(n): Metrisch



- **Metrische (und bei  $n \geq 30$  normalverteilte) abhängige Variable: t-Test**
  - t-test bei einer Stichprobe (Vergleich mit Referenzwert)
  - Zwei-Stichproben t-Test bei unabhängigen Stichproben
  - Zwei-Stichproben t-Test bei abhängigen Stichproben
- **Bei  $n < 30$ : t-Test kein sinnvolles Verfahren**



1. Anmerkung letzte Woche
2. Wiederholungsfragen
3. **Mittelwertvergleiche**
  1. Einführung
  2. **Ein-Stichproben t-test**
  3. Zwei-Stichproben t-test
4. Übungsaufgabe



1. Auswahl der Befragten mit einer Zufallsstichprobe (iid)
2. Mindestens intervallskalierte Daten
3. Normalverteilte Mittelwerte  
Bei  $n \geq 30$ :  
Asymptotisch normalverteilte Mittelwerte. D.h. bei einer Stichprobengröße von mehr als 30 Beobachtungen kann angenommen werden, dass die Mittelwerte normalverteilt sind.

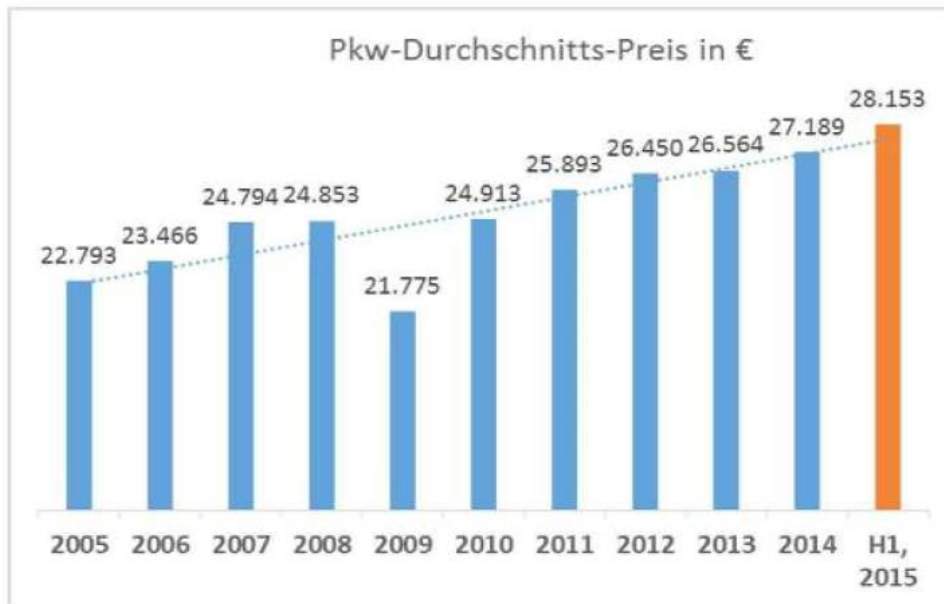
## Stata Syntax:

```
ttest abvar == #
```





- **Ein-Stichproben T-Test** prüft, ob ein Variablenmittelwert einer Stichprobe sich von einem Wert  $\mu_0$  in der Grundgesamtheit signifikant unterscheidet.



Quelle: CAR Universität Duisburg-Essen

### Beispiel:

Unterscheidet sich der mittlere Autopreis 2015 signifikant von Preisen aus dem Jahr 1978?



- Zu prüfende Hypothesen:

$$H_0: \bar{x} = \mu_0$$

$$H_1: \bar{x} \neq \mu_0$$

- $H_1$ : Gemessener Preis eines Autos aus 1978 unterscheidet sich signifikant von 28.153 € (Quelle: CAR 2015)
- $H_0$ : Gemessener Preis eines Autos aus 1978 unterscheidet sich nicht signifikant von 28.153 €



## Vorbereitung:

- Preise anpassen  $1\$_{1978} = 3,64\$_{2015}$  (*reine Inflationsbereinigung*)

gen `price_15=price*3.64`

- Umrechnung \$ in €:  $1\$_{2015} = 0,8780€_{2015}$

gen `price_15eu=price_15*0.8780`



Liegt der Preis eines Autos aus 1978 ( $\bar{x}$ ) bei 28153 € ( $\mu_0$ ) ?

### 1. Festlegung von $H_0$ und $H_1$

$$H_0: \bar{x} = 28.153\text{€}$$

$$H_1: \bar{x} \neq 28.153\text{€}$$

### 2. Wahl der Irrtumswahrscheinlichkeit ( $\alpha$ ) für $H_0$

$\alpha$  wird üblicherweise in den Sozialwissenschaften auf 0.05 gesetzt.



## 4. Berechnung der Prüfgröße und Entscheidung über die Verwerfung von $H_0$

```
. ttest price_15eu==28153
```

One-sample t test

Variable	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
price_~u	74	19703.67	1095.791	9426.353	17519.76	21887.58

```
mean = mean(price_15eu)          t = -7.7107
Ho: mean = 28153                degrees of freedom = 73
```

Ha: mean < 28153  
Pr(T < t) = 0.0000

Ha: mean != 28153  
Pr(|T| > |t|) = 0.0000

Ha: mean > 28153  
Pr(T > t) = 1.0000

- **Mittelwert:** Der Mittelwert der Stichprobe liegt bei 19703.67 €.
- **Signifikanztest:** Mit einer Irrtumswahrscheinlichkeit von weniger als 0,001% weicht der mittlere Preis von Autos aus 2015 von den Preisen aus 1978 ab.
- **$H_0$  wird verworfen:** Der mittlere Preis eines Autos aus 1978 ist signifikant kleiner als 28.153 €.



1. Anmerkung letzte Woche
2. Wiederholungsfragen
3. **Mittelwertvergleiche**
  1. Einführung
  2. Ein-Stichproben t-test
  3. **Zwei-Stichproben t-test**
4. Übungsaufgabe



- **Zwei-Stichproben T-Test für unabhängige Stichproben:**

Der Zwei-Stichproben T-Test prüft, ob in der Grundgesamtheit ein in der Stichprobe beobachteter Variablenmittelwert einer Gruppe (  $\bar{x}_1$  ), dem einer anderen Gruppe (  $\bar{x}_2$  ) entspricht.

- **Zu prüfende Hypothesen:**

$$H_0: \bar{x}_1 = \bar{x}_2$$

$$H_1: \bar{x}_1 \neq \bar{x}_2$$

- Beispiel für  $H_0$ : Einkommen von Männern (  $\bar{x}_1$  ) entspricht dem Einkommen von Frauen (  $\bar{x}_2$  ).



## Modellvoraussetzung für den Zwei-Stichproben T-Test bei unabhängigen Stichproben:

1. Auswahl der Befragten mit einer Zufallsstichprobe (iid)
2. Mindestens intervallskalierte Daten für beide Gruppen
3. Normalverteilte Mittelwerte für beide Gruppen
  - Bei  $n \geq 30$  in jeder Gruppe:  
Hat jede Gruppe mehr als 30 Datenpunkte, so können die Mittelwerte als asymptotisch normalverteilt angenommen werden.





## Prüfung auf Varianzgleichheit:

- Gleiche Varianzen der Gruppen erlauben effizienteren T-Test
- Überprüfung auf Varianzhomogenität mittels F-Test



## Levene F-Test:

- $H_0$ : Gruppen haben gleiche Varianz.
- $H_1$ : Gruppen haben ungleiche Varianz.

$p > 0,05$ :  $H_0$  beibehalten → T-Test für gleiche Varianzen

$p \leq 0,05$ :  $H_0$  wird abgelehnt → T-Test für ungleiche Varianzen

## Stata Syntax:

```
sdtest var, by (groupvar)
```

## Testentscheidung & Umsetzung in Stata

- Je nach Ergebnis des Levene F-Tests wird t-Test für verbundene Stichprobe oder unverbundene berechnet
- Stata berechnet zunächst immer t-Test für gleiche Varianzen
- Bei ungleichen Varianzen muss Option *unequal* ergänzt werden

### Stata Syntax:

```
ttest var , by(groupvar) unequal
```

Bei ungleicher  
Varianz zwischen den  
Gruppen  
(signifikanter F-Test)



**Beispiel: Wir vermuten, dass der Preis eines Autos nach Motorleistung variiert**

**H1:** Wenn Autos einen starken Motor haben, dann sind sie teurer.

### 1. Festlegung von $H_0$ und $H_1$

$$H_0: \bar{x}_1 = \bar{x}_2$$

$$H_1: \bar{x}_1 \neq \bar{x}_2$$

**2. Wahl der Irrtumswahrscheinlichkeit ( $\alpha$ ) für  $H_0$**   
 $\alpha$  wird auf 0.05 gesetzt.



## Operationalisierung der Fragestellung

Zielvariable:

price\_15eu (Preis der Autos in € auf 2015 inflationsbereinigt)

Gruppenvariable:

Hubraum (displacement in cu. in.) im Datensatz metrisch

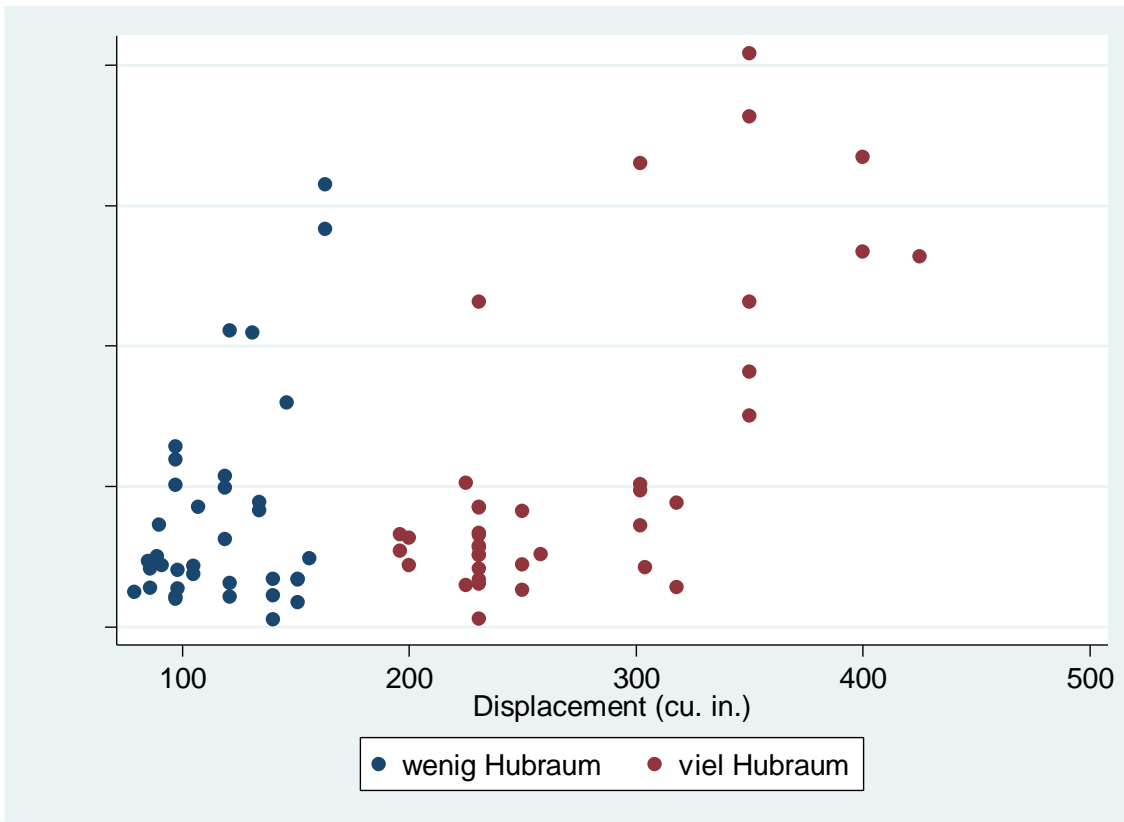
→ Kategorisieren mit Median-Split (Bildung von 2 Gruppen anhand Median)

```
sum displacement,d
gen disp_kat=.
replace disp_kat=1 if displacement<196
replace disp_kat=2 if displacement>=196
lab var disp_kat"Hubraum wenig/viel"
lab def hubk 1"wenig Hubraum" 2"viel Hubraum"
lab val disp_kat hubk
```



## Grafische Analyse

```
scatter price_15eu displacement if disp_kat==1 ///  
|| scatter price_15eu displacement if disp_kat==2
```



**Streudiagramm lässt  
Unterschied vermuten**

**Varianzhomogenität eher  
unwahrscheinlich**

## 4. Prüfung auf Varianzgleichheit:

```
. sdtest price_15eu, by(dispcat)
```

Variance ratio test

Group	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
wenig Hu	36	17765.32	1240.916	7445.494	15246.13	20284.51
viel Hub	38	21540	1745.317	10758.86	18003.65	25076.34
combined	74	19703.67	1095.791	9426.353	17519.76	21887.58

```
ratio = sd(wenig Hu) / sd(viel Hub)          f = 0.4789
Ho: ratio = 1                                degrees of freedom = 35, 37
```

```
Ha: ratio < 1
Pr(F < f) = 0.0154
```

```
Ha: ratio != 1
2*Pr(F < f) = 0.0309
```

```
Ha: ratio > 1
Pr(F > f) = 0.9846
```

Der F Test lehnt die  $H_0$  „Varianzen sind gleich ab.“

Es besteht ein signifikanter Varianzunterschied

Somit muss t-Test für ungleiche Varianzen verwendet werden.



```
. ttest price_15eu, by(dispcat) unequal
```

Option für  
ungleiche  
Varianz

Two-sample t test with unequal variances

Group	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
wenig Hu	36	17765.32	1240.916	7445.494	15246.13	20284.51
viel Hub	38	21540	1745.317	10758.86	18003.65	25076.34
combined	74	19703.67	1095.791	9426.353	17519.76	21887.58
diff		-3774.676	2141.495		-8050.277	500.9261

diff = mean(wenig Hu) - mean(viel Hub)

t = -1.7626

Ho: diff = 0

Satterthwaite's degrees of freedom = 66.0265

Ha: diff < 0

Pr(T < t) = 0.0413

Ha: diff != 0

Pr(|T| > |t|) = 0.0826

Ha: diff > 0

Pr(T > t) = 0.9587





- **Mittelwert:**

- Autos mit wenig Hubraum kosten im Mittel 17765 €
- Autos mit viel Hubraum kosten im Mittel 21540 €

- **Signifikanztest:**

Ha: diff < 0

$$\Pr(T < t) = 0.0413$$

Ha: diff != 0

$$\Pr(|T| > |t|) = 0.0826$$

Ha: diff > 0

$$\Pr(T > t) = 0.9587$$

- Autos mit wenig Hubraum sind signifikant günstiger zu einer Irrtumswahrscheinlichkeit von 4,13% (also signifikant bei  $\alpha=0.05$ )



Stichprobe	Variable(n)	Verteilung	Test
Ein-Stich- proben-Fall	Metrisch	NV oder $n > 30$	<b>t-Test (eine Stichprobe)</b> ttest var =Wert
Un- abhängige Stichprobe	Metrisch	NV, $n > 30$	<b>t-Test (unabhängige Stichproben)</b> ttest var, by(groupvar) ttest var, by(groupvar) unequal
	Ordinal	Keine NV, $n < 30$	<b>Mann-Whitney-U-Test</b> ranksum var , by(groupvar) <b>Mann-Whitney-Test</b> ranksum var , by(groupvar)
Abhängige Stichprobe	Metrisch	NV oder $n > 30$	<b>t-Test (abhängige Stichproben)</b> ttest var1 == var2
	Ordinal	Keine NV, $n < 30$	<b>Wilcoxon-Rangsummen-Test</b> ranksum var , by(groupvar) <b>Wilcoxon-Vorzeichen-Test</b> ranksum var , by(groupvar)



Die steigenden Preise für fossile Brennstoffe ändern möglicherweise die Präferenzen von Autokäufern. Demnach wären diese bereit, einen höheren Preis für Autos zu bezahlen, die einen niedrigeren Kraftstoffverbrauch aufweisen. Überprüfen Sie diese Überlegung.

- 1) Erstellen Sie eine Variable die den Kraftstoffverbrauch aus der amerikanischen Notierung (miles per gallon, Variable „mpg“ in auto.dta ) in das deutsche System (Liter pro 100 km) überträgt.

Formel zur Umrechnung: 
$$\frac{\text{Liter}}{100\text{km}} = \frac{235}{\text{mpg}}$$

- 2) Laut Umweltbundesamt lag der mittlere Kraftstoffverbrauch von Autos im Jahr 2015 bei 7.3 Liter/100 km. Überprüfen Sie, ob sich der mittlere Verbrauch von 1978 signifikant hiervon unterscheidet.



- 3) Überprüfen sie, ob Autos mit geringem Kraftstoffverbrauch teurer sind als Autos mit hohem Kraftstoffverbrauch
- Erstellen Sie eine Variable, die den Kraftstoffverbrauch am Median in zwei Kategorien einteilt
  - Überprüfen Sie den Zusammenhang zwischen Kraftstoffverbrauch und Preis zunächst grafisch
  - Wenden Sie den zweiseitigen t-Test an, um genauere Aussagen über den Zusammenhang treffen zu können. Überprüfen Sie zuvor welche Variante in diesem Fall geeignet ist.
  - Interpretieren Sie die Ergebnisse inhaltlich, lässt sich die Hypothese bestätigen, oder muss sie abgelehnt werden?