

### Professur Psychologie digitaler Lernmedien

Institut für Medienforschung

Philosophische Fakultät

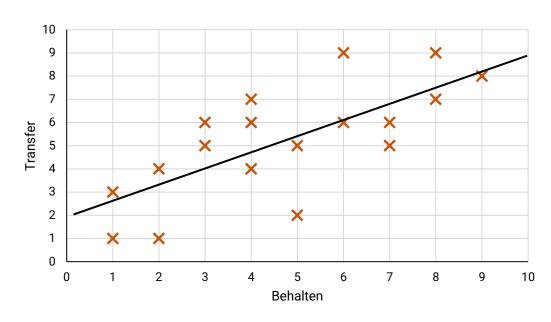


#### Überblick

- Lineare bivariate Regression
- Methode der kleinsten Quadrate
- Nichtlineare Zusammenhänge
- Multiple Regression
- Verfahren zur Auswahl unabhängiger Variablen
- Interaktionseffekte in der multiplen Regression
- Inkrement und Dekrement in der multiplen Regression
- Suppressorvariablen
- Prognosegüte und Kreuzvalidierung
- Indikatorcodierung
- Signifikanzprüfung zur multiplen Regression
- Konfidenzbänder
- Inferenzstatistische Voraussetzungen

## Lineare bivariate Regression (z. B. Rasch, Friese, Hofmann & Naumann, 2021)

- Lineare bivariate Regression: Statistisches Verfahren zur Vorhersage einer Kriteriumsvariable durch eine Prädiktorvariable mittels linearer Funktion (vgl. Sitzung "Regression" aus Statistik I)
- Fiktives Beispiel: Zusammenhang zwischen Behaltens- und Transferleistungen



- Regressionsgrade soll den Gesamttrend der Einzelwerte bestmöglich wiedergeben
- Regressionsgleichung zur Regressionsgraden:

$$\hat{y} = m \cdot x + b$$

```
ŷ = Vorhergesagte Kriteriumsvariable y
```

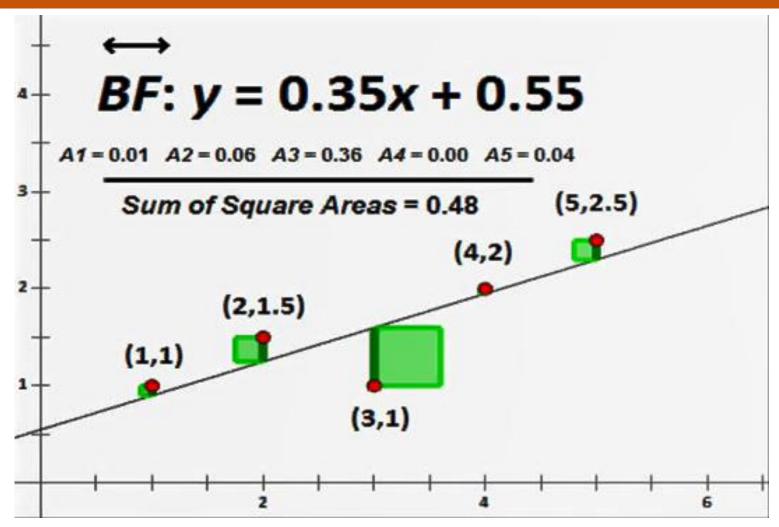
m = Steigung der Regressionsgraden

x = Prädiktorvariable x

b = Achsenabschnitt der Regressionsgraden

 Berechnung der Regressionsgewichte m und b mittels Methode der kleinsten Quadrate

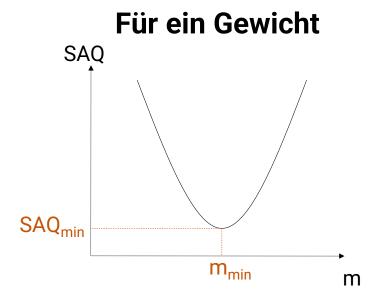
#### Methode der kleinsten Quadrate

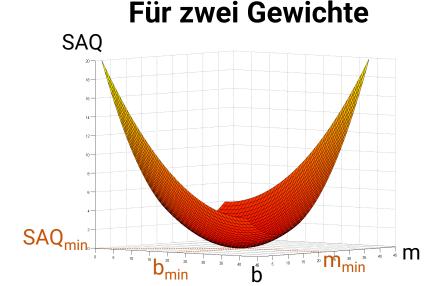


Quelle: <a href="http://www.youtube.com/watch?v=jEEJNz0RK4Q">http://www.youtube.com/watch?v=jEEJNz0RK4Q</a>

#### Methode der kleinsten Quadrate

- Summe der Abweichungsquadrate (SAQ) soll ein Minimum ergeben
- Ein Gewicht: Parabel (→ Regressionsgerade durch Achsenursprung)
- Zwei Gewichte: Paraboloid (→ Regressionsgrade)





#### Methode der kleinsten Quadrate

- Summe der Abweichungsquadrate (SAQ) soll ein Minimum ergeben
- Formel:

$$SAQ = \sum_{i=1}^{n} [y_i - \hat{y}_i]^2 = \sum_{i=1}^{n} [y_i - (m \cdot x_i + b)]^2 = \min$$

• Erste Ableitung bilden und auf Null setzen ergibt für m und b:

$$m_{yx} = \frac{\text{cov}(x, y)}{\sigma_x^2}$$

$$b_{yx} = \bar{y} - m_{yx} \cdot \bar{x}$$

y = Beobachtete Werte der Variablen y

ŷ = Vorhergesagte Kriteriumsvariable y

m = Steigung der Regressionsgraden

x = Prädiktorvariable x

b = Achsenabschnitt der Regressionsgraden

i = Person i

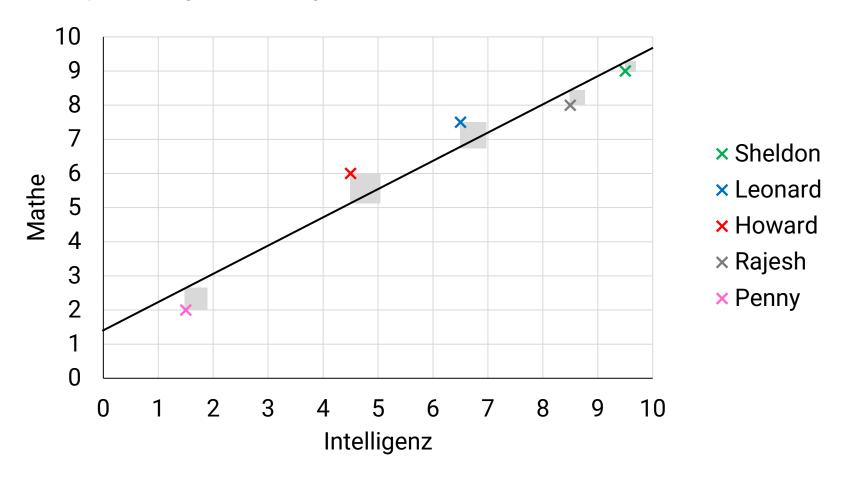
 Beispiel: Berechnung von b und m zu dem rechts dargestellten Datensatz:

$$m_{yx} = \frac{\text{cov}(x, y)}{\sigma_x^2} \approx \frac{8.5}{3.21^2} \approx 0.82$$

$$b_{yx} = \overline{y} - m_{yx} \cdot \overline{x} \approx 6.5 - 0.82 \cdot 6.1 = 1.47$$

VPN	IQ	Mathe	
Sheldon	9.5	9.0	
Leonard	6.5	7.5	
Howard	4.5	6.0	
Rajesh	8.5	8.0	
Penny	1.5	2.0	
М	6.1	6.5	
SD	3.21	2.74	

• Beispiel: Regressionsgrade mit b = 1.47 und m = 0.82:



Wie hoch ist die Intelligenz laut Regressionsgleichung für das Beispiel auf der vorherigen Folie bei einer Person mit einem Mathewert von 4?

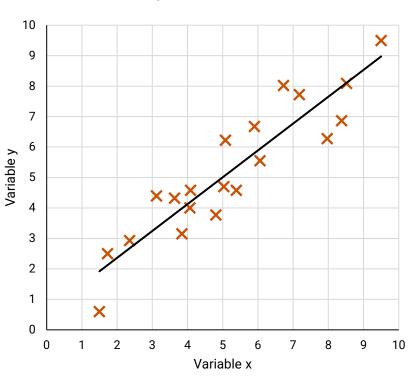
- □ A: 3
- ☐ B: 3.09 (gerundet)
- □ C: 4
- □ D: 4.75
- ☐ E: Wert kann nicht berechnet werden

# Nichtlineare Zusammenhänge (z.B. Rasch, Friese, Hofmann & Naumann, 2021)

Beispiele für lineare und nonlineare Zusammenhänge

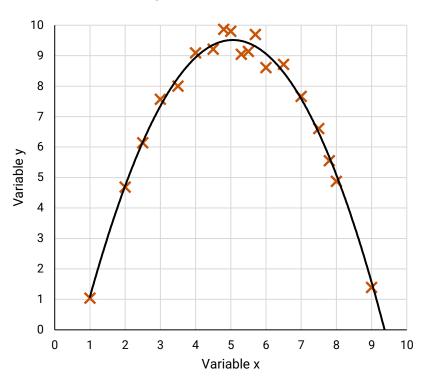
#### Linearer Zusammenhang

$$\hat{y} = m \cdot x + b$$



#### Nonlinearer Zusammenhang

$$\hat{y} = m \cdot x^2 + b$$



### Multiple lineare Regression

- Definition: Statistisches Verfahren zur Vorhersage einer Kriteriumsvariable durch mehrere Prädiktorvar. mittels Linearkombination
- Regressionsgleichung zur Regressions(hyper-)ebene:

$$\hat{y} = b_0 + x_1 \cdot b_1 + \dots + x_m \cdot b_m$$

- Bestimmung der Regressionsgewichte (Beta-Gewichte) wieder mittels Methode der kleinsten Quadrate
- Unterschied zur linearen bivariaten Regression: Berechnung mit Matrizen statt mit Zahlen

- ŷ VorhergesagteKriteriumsvariable y
- b<sub>0</sub> Achsenabschnittder Regressionsgraden
- x<sub>1</sub> Erste Prädiktorvariable
- b<sub>1</sub> Steigung zur ersten Prädiktorvariablen
- x<sub>m</sub> m-te Prädiktorvariable
- b<sub>m</sub> Steigung zur m-ten Prädiktorvariablen

# Verfahren zur Auswahl unabhängiger Variablen (Eid, Gollwitzer & Schmitt, 2017)

- Zwei generelle Strategien zur Auswahl unabhängiger Variablen (= Prädiktorvariablen) für ein Regressionsmodell
  - Auswahl aufgrund theoretischer Überlegungen
  - Datengesteuerte Auswahl zur Maximierung der Varianzaufklärung
- Strategien bei der datengesteuerten Auswahl
  - Vorwärtsselektion
  - Rückwärtselimination
  - Schrittweise Regression

# Verfahren zur Auswahl unabhängiger Variablen (Eid, Gollwitzer & Schmitt, 2017)

#### Theoriegeleitete Vorgehensweise

- Auswahl der Prädiktoren nach inhaltlicher Relevanz bzw. Hypothesen
- alle (oder in Blöcken) gleichzeitig ins Modell aufnehmen
- Interpretation erfolgt hypothesengeleitet, nicht datengesteuert

#### Hierarchische Regression

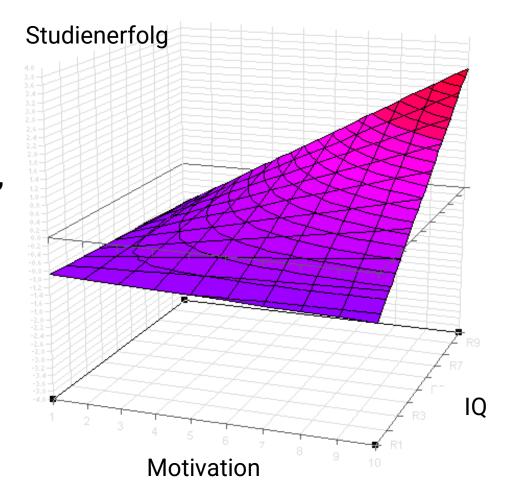
- Prädiktoren werden blockweise eingefügt (z. B. zuerst demografische, dann psychologische Variablen)
- Vergleich von ΔR² zwischen Blöcken → zeigt inkrementelle Varianzaufklärung (siehe Folie 16)

#### Empfehlung:

 theoretische bzw. hierarchische Verfahren sind interpretativ stabiler; datengesteuerte Verfahren nur explorativ einsetzen

### Interaktionseffekte in der multiplen Regression

- Interaktionseffekt (bzw. Moderatoreffekt bzw. Wechselwirkungseffekt)
- Fiktives Beispiel:
   Studienerfolg nur dann hoch, wenn IQ (x<sub>1</sub>) und Motivation (x<sub>2</sub>) hoch sind
- $\hat{y} = 1 \cdot b_0 + x_1 \cdot b_1 + x_1 \cdot x_2 \cdot b_3 + x_2 \cdot b_2$



# Inkrement und Dekrement in der multiplen Regression

- Beitrag zur Varianzaufklärung: Für jede einzelne Prädiktorvariable lässt sich ein solcher Beitrag bestimmen
- Unterscheidung zwischen Inkrement und Dekrement
  - Inkrement (R<sub>I</sub><sup>2</sup>): Zuwachs an aufgeklärter Varianz durch Hinzunahme weiterer Prädiktorvariablen
  - Dekrement (R<sub>D</sub><sup>2</sup>): Abnahme an aufgeklärter Varianz durch Verzicht auf bestimmte Prädiktorvariablen

## Inkrement und Dekrement in der multiplen Regression

- Orthogonaler Fall (sämtliche Prädiktorvariablen sind unkorreliert): Addition der Einzelkorrelationen zur Berechnung von  $R^2$ ;  $R_{l}^2$  (bzw.  $R_{D}^2$ ) =  $r_{x_i,y}^2$
- Kollinearer Fall (Prädiktoren sind korreliert)
  - R<sup>2</sup> kleiner als Summe der Einzelkorrelationen durch Informationsüberschneidungen (häufiger Fall)
  - R<sup>2</sup> größer als Summe der Einzelkorrelationen: Suppressoreffekte durch Informationspräzisierung (seltener Fall)

$$R^2 = \sum_{j=1}^{m} r_{x_j, y}^2$$

$$R^2 < \sum_{j=1}^m r_{x_j,y}^2$$

$$R^2 > \sum_{j=1}^{m} r_{x_j, y}^2$$

### Suppressorvariablen in der multiplen Regression

- Suppressorvariablen erhöhen die aufgeklärte Varianz durch Unterdrückung irrelevanter Varianzen anderer Variablen
- Bedingungen für eine Suppressorvariable
  - Keine oder geringe Korrelation mit der Kriteriumsvariable
  - Deutliche Korrelation mit mindestens einer Prädiktorvariable
  - Inkrement bzw. Dekrement der Variable ist (deutlich) größer als einfacher Determinationskoeffizient (R<sup>2</sup>) der Suppressorvariable
- Beispiel: Berufserfolg (AV) wird durch Abschlussnote im Studium (UV<sub>1</sub>) und Prüfungsangst (UV<sub>2</sub>) vorhergesagt
- Prüfungsangst könnte als mögliche Suppressorvariable irrelevante Varianz in der Abschlussnote unterdrücken

# Prognosegüte und Kreuzvalidierung (Eid, Gollwitzer & Schmitt, 2017)

- Prognosegüte: Wie gut kann das Regressionsmodell zukünftige Fälle vorhersagen?
- Überprüfung der Prognosegüte
  - Weitere Stichprobenziehung
  - Kreuzvalidierung
- Kreuzvalidierung (engl. cross-validation): Teilung des Datensatzes in Trainings-Substichprobe und Test-Substichprobe
- Ziel: Overfitting vermeiden

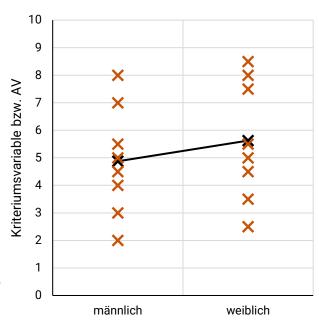
### Modellgüte und Residualdiagnostik (z. B. Field, 2018)

#### Residualanalyse:

- Streudiagramm der Residuen gegen vorhergesagte Werte zur Prüfung von Linearität und Homoskedastizität
- Q-Q-Plot der Residuen zur Prüfung der Normalverteilung
- Überprüfung des Einflusses einzelner Fälle auf Modellschätzung
  - Leverage/Hebelwert: sagt aus, wie stark ein Datenpunkt im Raum der Prädiktorvariablen vom Zentrum entfernt ist und damit, wie stark dieser Punkt die Form der Regressionsgerade "hebeln" kann
  - Mahalanobis-Distanz: misst wie der Hebelwert die Distanz eines Punktes vom Zentrum der Daten im Prädiktorraum, berücksichtigt aber die Korrelationen zwischen den unabhängigen Variablen
  - Cook-Distanz: misst als wichtigster Wert den tatsächlichen Einfluss eines Datenpunkts auf die geschätzten Regressionskoeffizienten; sie kombiniert die Informationen aus dem Hebelwert (Distanz im Prädiktorraum) und dem Residuum (Distanz im Kriteriumraum)

### Indikatorcodierung

- Regressionsanalyse mittels Indikatorcodierung auch bei fehlendem Intervallskalenniveau der Prädiktorvariable(n) möglich
- Indikatorcodierung: Umrechnung von nominaloder ordinalskalierten Prädiktorvariablen in künstliche, intervallskalierte Prädiktorvariablen
- Beispiel: Umrechnung der Variable Geschlecht in eine Indikatorvariable (z. B.  $\sigma = 0$  und 9 = 1)
- Äquidistanz: Diese Indikatorvariable enthält nur ein Intervall, welches zu sich selbst äquidistant ist und somit Intervallskalenniveau besitzt
- Wichtig: Durch Indikatorcodierung und das Allgemeine Lineare Modell gilt mathematisch: Varianzanalyse = Regressionsanalyse



# Signifikanzprüfung zur multiplen Regression (z. B. Eid, Gollwitzer & Schmitt, 2017)

- Signifikanztest f
  ür die multiple Regression mittels F-Test
- Berechnung des empirischen *F*-Wertes für die multiple Regression:

$$F = \frac{n-k-1}{k} \cdot \frac{R^2}{(1-R^2)}$$

n = Stichprobenumfang

k = Anzahl an Prädiktorvariablen

 $R^2$  = Determinationskoeffizient

- Formel für die Freiheitsgrade:
  - Zählerfreiheitsgrade:  $df_Z = k$
  - Nennerfreiheitsgrade:  $df_N = n k 1$
- Inferenzstatistische Entscheidung: Vergleich empirischer F-Wert mit kritischem F-Wert  $\rightarrow$  Ergebnis signifikant ( $F_{\rm emp} \ge F_{\rm krit}$ ) oder nicht signifikant ( $F_{\rm emp} < F_{\rm krit}$ )

# Signifikanzprüfung zur multiplen Regression (z. B. Eid, Gollwitzer & Schmitt, 2017)

- Beispiel: Berechnung des F-Wertes zu einer multiplen Regression mit zwei Prädiktorvariablen
- Weitere Angaben: R<sup>2</sup> = .338; N = 237
- Berechnung des F-Wertes:

$$F = \frac{n - k - 1}{k} \cdot \frac{R^2}{(1 - R^2)} = \frac{237 - 2 - 1}{2} \cdot \frac{0.338}{(1 - 0.338)} = \frac{234}{2} \cdot \frac{0.338}{0.662} \cong 59.737$$

 Inferenzstatistische Entscheidung: Vergleich empirischer F-Wert mit kritischem F-Wert

# Signifikanzprüfung zur multiplen Regression (z. B. Eid, Gollwitzer & Schmitt, 2017)

 "Varianzanalytische Darstellung": F-Wert für die multiple Regression auch als Quotient aus Quadratsummen darstellbar:

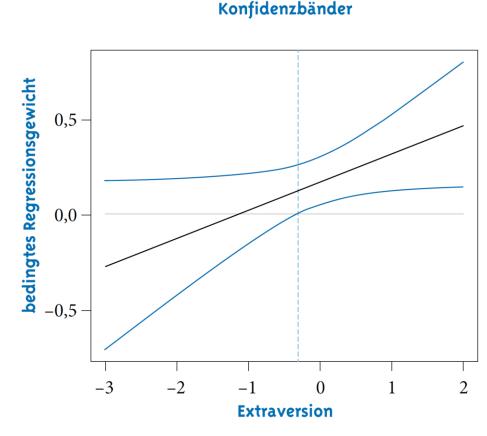
$$F = \frac{n - k - 1}{k} \cdot \frac{R^2}{(1 - R^2)}$$
 
$$F = \frac{MQSR}{MQSE}$$

MQSR = Mittlere Quadratsumme der RegressionMQSE = Mittlere Quadratsumme der Residuen

- Weiterer Beleg dafür, dass mathematisch gilt:
- Regressionsanalyse = Varianzanalyse

### Konfidenzbänder (Eid, Gollwitzer & Schmitt, 2017)

- Konfidenzintervall (= Vertrauensintervall): Beschreibt den Bereich um den geschätzten Populationsparameter, für den gilt, dass er mit einer Wahrscheinlichkeit von 1-α (z. B. 95%) den "wahren" Populationsparameter beinhaltet
- Konfidenzbänder: Vgl. Konfidenzintervalle
- Gekrümmte Form der Konfidenzbänder



Quelle: Eid, Gollwitzer und Schmitt (2017)

## Inferenzstatistische Voraussetzungen (z.B. Rasch, Friese, Hofmann & Naumann, 2021)

- Linearität zwischen den Prädiktoren und dem Kriterium
- Intervallskalenniveau der Kriteriumsvariable
- Normalverteilung der Residuen
- Unabhängigkeit der einzelnen Messwerte verschiedener Personen
- Homoskedastizität: Homogenität der Streuungen der zu einem x-Wert gehörenden y-Werte über den gesamten Wertebereich von x (vgl. inferenzstatistische Voraussetzungen der MANOVA ohne MW)
- keine Multikollinearität zwischen den Prädiktoren (keine zu hohe Korrelation)
- keine einflussreichen Ausreißer die das Modell unverhältnismäßig stark beeinflussen

## Beispiele für die multiple lineare Regression in Fachzeitschriften

We then used the attributes to predict perceived warmth and competence. Two separate linear multiple regressions were conducted with perceived warmth and competence as the dependent variables. We used multiple dummy codes for the gender displayed, developmental categories, and locomotion.

A robot was perceived with higher warmth if its mechanics were invisible ( $\beta = -0.13$ , SE = 0.05, and p < .05), it had more degrees of freedom for movement ( $\beta = 0.12$ , SE = 0.01, p < .05), younger age ( $\beta = -0.14$ , SE = 0.01, p < .05), and displayed as a child ( $\beta = -0.14$ , SE = 0.01, p < .05). The regression was significant F(25, 316) = 5.90, p < .001 with the four factors explaining 26.4% of the variance in perceived warmth. The standardized beta coefficients are presented in Table 6.

However, the three other cases (Fig. 1b–e) could not be covered with this multiple regression equation as in these cases both (suboptimal and optimized) conditions have non-linear progression lines. To solve this issue, it was decided to include two dummy variables in the regression equation, one indicating the optimized visual design (as dumOpt, see Eq. 1) and one indicating the suboptimal visual design (as dumSub, see Eq. 2):

$$perf = b_0 + b_1 abi + b_{21} dumOpt + b_{22} dumSub + b_{31} abi \times dumOpt + b_{32} abi \times dumSub + b_{41} (abi \times dumOpt)^2 + b_{42} (abi \times dumSub)^2$$
(3)

Quelle: Kühl, Fehringer und Münzer (2022)

Quelle: Reeves, Hancock und Liu (2020)

Table 3	Multiple regression to
predict	students' grade in the
first exa	nm

	В	SE B	β	p
Constant	-3.93	0.59		< 0.001
Cognitive strategy use	0.21	0.18	0.17	0.24
Metacognitive strategy use	0.00	0.14	0.00	0.99
Self-efficacy	0.23	0.11	0.24	0.04

Quelle: Gentner, Respondek und Seufert (2024)

### Zusammenfassung I

- Lineare bivariate Regression: Statistisches Verfahren zur Vorhersage einer Kriteriumsvariable durch eine Prädiktorvariable mittels linearer Funktion
- Methode der kleinsten Quadrate zur Berechnung der Regressionsgewichte
- Nichtlineare Zusammenhänge ebenfalls vorhersagbar
- Multiple Regression zur Vorhersage einer Kriteriumsvariable durch mehrere Prädiktorvariablen mittels Linearkombination
- Auswahl unabhängiger Variablen aufgrund theoretischer Überlegungen oder datengesteuert zur Maximierung der Varianzaufklärung
- Interaktionseffekte in der multiplen Regression
- Inkrement/Dekrement als Zuwachs/Abnahme an aufgeklärter Varianz in der multiplen Regression durch Hinzunahme/Verzicht von Prädiktorvariablen

### Zusammenfassung II

- Suppressorvariablen erhöhen die aufgeklärte Varianz durch Unterdrückung irrelevanter Varianzen anderer Variablen
- Überprüfung der Prognosegüte durch weitere Stichprobenziehung oder Kreuzvalidierung zur Overfitting-Vermeidung
- Indikatorcodierung als Umrechnung in künstliche, intervallskalierte Prädiktorvariablen
- Signifikanzprüfung zur multiplen Regression mittels F-Test (vgl. Varianzanalyse)
- Konfidenzbänder u. A. zur Abschätzung des Wertebereiches, in dem sich der "wahre" Wert mit hoher Wahrscheinlichkeit befindet
- Inferenzstatistische Voraussetzungen: Intervallskalenniveau & Normalverteilung der Kriteriumsvariable, Unabhängigkeit der einzelnen Messwerte & Homoskedastizität

### Prüfungsliteratur

- Eid, M., Gollwitzer, M., & Schmitt, M. (2017). Statistik und Forschungsmethoden (5. Aufl.). Weinheim: Beltz.
  - Multiple Regressionsanalyse (S. 629-726) → Ohne Rechenaufgaben

#### Weiterführende Literatur

- Rasch, B., Friese, M., Hofmann, W., & Naumann, E. (2021).
   Quantitative Methoden 1: Einführung in die Statistik für Psychologie, Sozial- & Erziehungswissenschaften (5. Aufl.). Heidelberg: Springer.
  - Merkmalszusammenhänge (S. 105–119)
- Bortz, J., & Schuster, C. (2010). Statistik für Human- und Sozialwissenschaftler (7. Aufl.). Berlin: Springer.
  - Partielle Korrelation und multiple lineare Regression (S. 339–361)
- Field, A. (2018). Discovering Statistics Using IBM SPSS Statistics.
   Sage.
  - The Lineare Model (Regression) (S. 369–436)
- Leonhart, R. (2022). Lehrbuch Statistik. Einstieg und Vertiefung (5. Auflage). Bern: Huber.
  - Multiple Korrelation und Multiple Regression (S. 355–400)