

Die Propensity Score-Methode - Eine anwendungsorientierte Einführung



Oliver Kuß¹, André Scherag²

¹ Institut für Biometrie und Epidemiologie, Deutsches Diabetes-Zentrum (DDZ),
Leibniz-Zentrum für Diabetes-Forschung an der Heinrich-Heine-Universität
Düsseldorf, Düsseldorf

² Integriertes Forschungs- und Behandlungszentrum (IFB) Sepsis und Sepsisfolgen
(CSCC), Universitätsklinikum Jena, Jena

Plan

- Einleitung
- Die Propensity Score-Methode
- Ein Beispiel aus der Herzchirurgie
- Propensity Score versus herkömmliche Regressionsmodelle
- Fazit

Einleitung

- Therapien sollen primär in randomisierten Studien geprüft werden.
- Nur die Randomisierung garantiert eine gleichmäßige Verteilung aller bekannten und unbekanntes (!) Störgrößen und Risikofaktoren auf die Therapiegruppen → kausale Aussagen möglich
- **Aber:** Randomisierte Studien sind in manchen Fällen „unnötig, ungeeignet, unmöglich oder ungenügend“ [Black, 1996]
- **Zusätzlich:** Geringe externe Validität [McKee 1999; Rothwell 2005]

Einleitung

- **Alternative:** Nichtrandomisierte Studien
- Bessere *externe* Validität, **aber:** fehlende *interne* Validität:
Die Therapiezuweisung erfolgt nichtrandomisiert und Interventions- und Kontrollgruppe können sich systematisch bzgl. bekannter und (schlimmer noch) unbekannter Störgrößen, so genannter „Confounder“, unterscheiden
- **Lösung:** Adjustierung für Confounder durch herkömmliche Regressionsmodelle (Standard) oder die Propensity Score-Methode [Rosenbaum/Rubin1983]

Die Propensity Score-Methode

- **Definition:** Der Propensity Score (Abk.: PS) ist die Wahrscheinlichkeit, die zu prüfende Therapie zu erhalten
- Der PS ist i.a. unbekannt und muss in einem **ersten Schritt** geschätzt werden (PS-Modell)
- Schätzung des PS-Modells durch (z.B.) logistische Regression:
 - Zielgröße (abhängige Variable): Therapie
 - Kovariablen (unabhängige Variablen): die zu Therapiebeginn bestehenden Patientenmerkmale (Risikofaktoren, Confounder)
- Aus den geschätzten Modellparametern dieses PS-Modells kann dann der Propensity Score für jeden einzelnen Patienten berechnet werden

Die Propensity Score-Methode

- **Frage:** Welche Kovariablen sollen ins PS-Modell aufgenommen werden?
- **Antwort:**
 1. Viele! (ruhig auch Interaktionen und nicht-lineare Terme)
 2. Vor allem diejenigen, die den späteren Therapieerfolg (und nicht etwa die Therapiezuweisung) vorhersagen, sonst droht Effizienzverlust ohne Biasgewinn [Brookhart2006]

Die Propensity Score-Methode

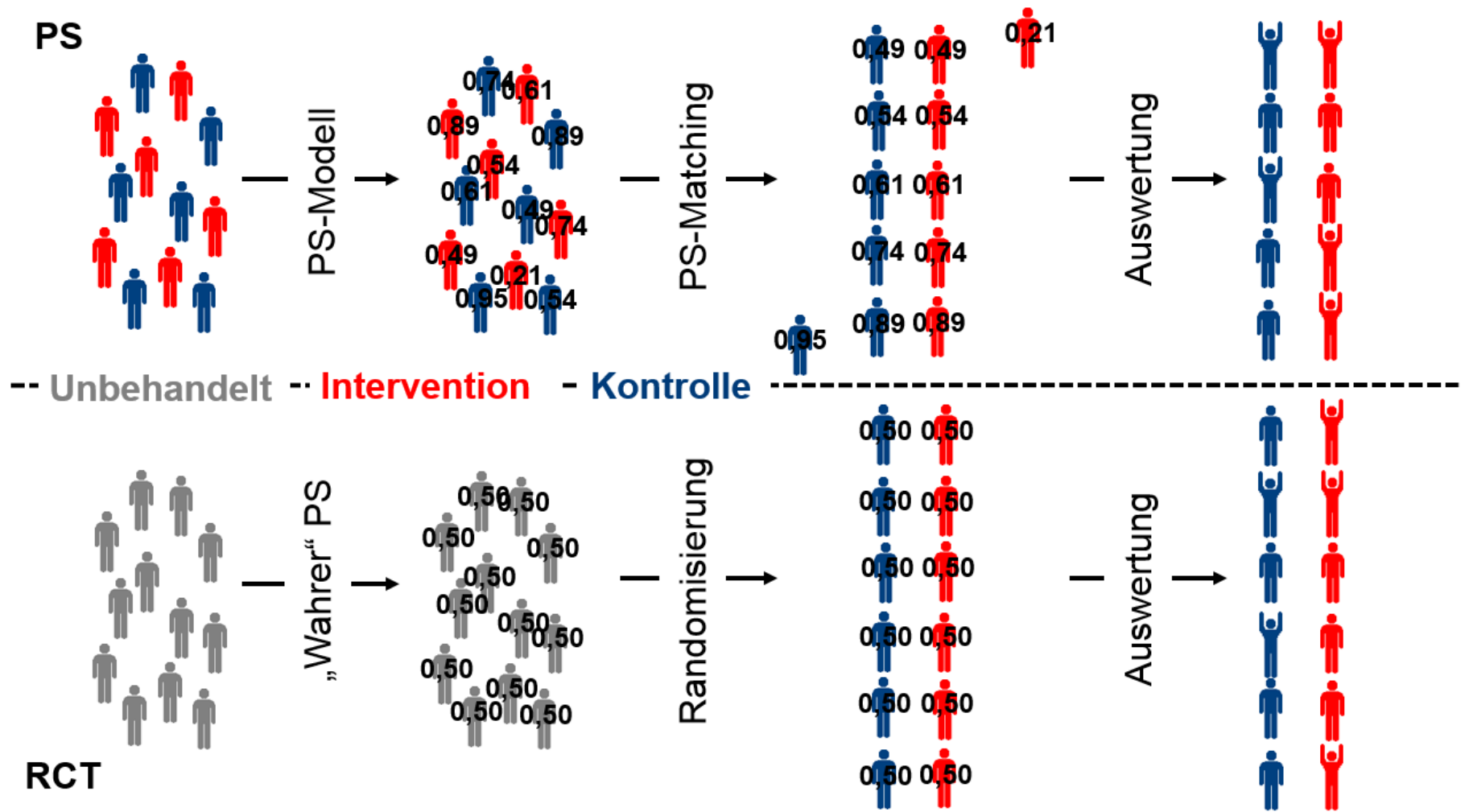
- **Zweiter Schritt:** Schätze den eigentlich interessierenden Therapieeffekt unter Zuhilfenahme des PS
- **Vier Methoden:**
 1. PS-Matching
 2. Regressionsadjustierung für den PS
 3. Stratifizierung
 4. IPTW(=Inverse probability of treatment weighting)-Methode

Die Propensity Score-Methode

- **(Nicht nur) unser Favorit: PS-Matching** [Austin2007, Morgan2006].

Beim PS-Matching wird jedem behandelten Patienten einer („1:1-Matching“) oder mehrere („1:n-Matching“, n kann sogar variieren) unbehandelte Patienten mit demselben (bzw. nur in kleinem Rahmen abweichender) PS zugeteilt

PS-Matching und RCT – Ein Vergleich



Die Propensity Score-Methode

Vorteile PS-Matching:

1. Explizite Darstellung sowohl der Eigenschaften von behandelten und unbehandelten Patienten („Table 1“ in einer randomisierten Studie) als auch der Balanciertheit der Confounder

	PS-gematchte Patienten (n = 788)		
	Less-OPCAB (n = 394)	cCABG (n = 394)	z-Differenz
Alter [Jahre]	69,3 ± 9,1	69,0 ± 8,9	0,46
Männlich [%]	78,2	77,9	0,09
Präoperativer Myokardinfarkt [%]	27,2	26,7	0,16

Die Propensity Score-Methode

Vorteile PS-Matching:

1. Explizite Darstellung sowohl der Eigenschaften von behandelten und unbehandelten Patienten („Table 1“ in einer randomisierten Studie) als auch der Balanciertheit der Confounder, als auch des Erfolges des Matchings

	Alle Patienten (n = 1.282)			PS-gematchte Patienten (n = 788)		
	Less-OPCAB (n = 395)	cCABG (n = 887)	z-Differenz	Less-OPCAB (n = 394)	cCABG (n = 394)	z-Differenz
Alter [Jahre]	69,3 ± 9,1	67,5 ± 9,4	3,24	69,3 ± 9,1	69,0 ± 8,9	0,46
Männlich [%]	78,2	77,9	0,13	78,2	77,9	0,09
Präoperativer Myokardinfarkt [%]	27,1	35,7	-3,14	27,2	26,7	0,16

Die Propensity Score-Methode

ACHTUNG!

	Alle Patienten (n = 1.282)			PS-gematchte Patienten (n = 788)		
	Less-OPCAB (n = 395)	cCABG (n = 887)	z-Differenz	Less-OPCAB (n = 394)	cCABG (n = 394)	z-Differenz
Alter [Jahre]	69,3 ± 9,1	67,5 ± 9,4	3,24	69,3 ± 9,1	69,0 ± 8,9	0,46
Männlich [%]	78,2	77,9	0,13	78,2	77,9	0,09
Präoperativer Myokardinfarkt [%]	27,1	35,7	-3,14	27,2	26,7	0,16
Unbekannter Confounder [%]	50,1	0	69,86	50,1	0	69,86

Die Propensity Score-Methode

ACHTUNG!

	Alle Patienten (n = 1.282)			PS-gematchte Patienten (n = 788)			
	Diese Gefahr besteht allerdings nur, wenn der unbekannte Confounder unabhängig von allen eingeschlossenen Merkmalen ist. Ansonsten wird er „mitbalanciert“.						erenz
Alter							46
Männer							09
Präoperativer Myokardinfarkt [%]	27,1	35,7	-3,14	27,2	26,7	0,16	
Unbekannter Confounder [%]	50,1	0	69,86	50,1	0	69,86	

Vorteile PS-Matching:

- Explizite Darstellung sowohl der Eigenschaften von behandelten und unbehandelten Patienten („Table 1“ in einer randomisierten Studie) als auch der Balanciertheit der Confounder, als auch des Erfolges des Matchings
- PS-Matching ist unter den vier genannten Methoden die beste, um Unbalanciertheiten zwischen behandelten und unbehandelten Patienten auszugleichen [Austin2009]

Vorteile PS-Matching:

- Im Vergleich zu Regressionsadjustierung für den PS: Spezifikation der korrekten funktionalen Form für den PS nicht nötig, Regressionsadjustierung für den PS hat keine kausale Interpretation [Williamson2012]
- Im Vergleich zu IPTW: Robustheit gegenüber extremen Beobachtungen
- Im Vergleich zu Regressionsadjustierung für den PS und IPTW: Der geschätzte PS wird nicht in der Auswertung verwendet und man erwartet dadurch eine größere Robustheit gegenüber Miss-Spezifikationen des PS-Modells [Deb et al.2015]

Die Propensity Score-Methode

Nachteil (?) PS-Matching:

Patienten, für die kein Matching-Partner gefunden wurde, werden ausgeschlossen.

→ Reduktion der Fallzahl, Verlust an statistischer Power

Aber:

- Es wird klar, über welche Subgruppe überhaupt Aussagen gemacht werden dürfen
- Theoretische Annahme [RosenbaumRubin1983] einer Überlappung der Verteilung des PS von behandelten und unbehandelten Patienten ist sichergestellt.
- Behandlungseffekt wird mit weniger Bias und nicht notwendigerweise erhöhter Varianz geschätzt [Rosenbaum2005]

Die Propensity Score-Methode

Konsequenz:

PS-Matching vs. die anderen PS-Methoden ist immer ein Trade-Off [Stuart2009] zwischen

- **Bias** (verzerrter Schätzung des Therapieeffektes) und
- **Varianz** (zu ungenaue Schätzung des Therapieeffektes)

Die Propensity Score-Methode

Nachteil PS-Matching:

Mit PS-matching kann nur der ATT (Average treatment effect of the treated) geschätzt werden, nicht jedoch der ATE (Average treatment effect).

Table 1. Propensity score methods and treatment effects

Methods	Average treatment effect of the treated	Average treatment effect
Matching	Yes	No
Stratification	Yes (weight adjusted)	Yes (equal weights)
Covariate adjustment	No	No
Inverse probability of treatment weighting	Yes (modified weighting)	Yes

“Yes” means that the selected method can estimate the indicated treatment effect. “No” means that the selected method cannot estimate the indicated treatment effect.

Die Propensity Score-Methode

Achtung:

Zwei Patienten mit identischem PS haben, sind nicht notwendigerweise identisch bzgl. ihrer Merkmale.

Die Balanciertheit gilt nur marginal, also beim Vergleich der beiden kompletten Gruppen.

Die Propensity Score-Methode

Die Güte eines PS-Modells sollte **alleine** dran gemessen werden, wie gut die Patientenmerkmale in den beiden Therapiegruppen balanciert ist.

Anpassungstests (z.B. Hosmer-Lemeshow-Test) oder die c-Statistik sind nicht geeignet, unbekannte Confounder zu entdecken [Weitzen2005].

Ein hoher Wert der c-Statistik ist weder notwendig noch hinreichend für eine gute Confounderadjustierung (vgl. RCT, dort $c=0.5$) [Westreich2011]

Die 10 Gebote der PS-Modellbildung

1. Du sollst für jede Kovariable mindestens 10 Beobachtungen haben.
2. Du sollst die Kovariablen auf Kollinearität prüfen.
3. Du sollst die Kovariablen auf statistische Signifikanz prüfen.
4. Du sollst sorgfältig die geschätzten Parameter interpretieren.
5. Du sollst Goodness-of-Fit-Tests berechnen und interpretieren.
6. Du sollst die c-Statistik berechnen und interpretieren.
7. Du sollst R^2 -Statistiken berechnen und interpretieren.
8. Du sollst die geschätzten Residuen auf Auffälligkeiten prüfen.
9. Du sollst die externe Validität des Modells an einer neuen Stichprobe prüfen.
10. ???

Das 10. Gebot:
Du sollst die Gebote 1-9 komplett ignorieren und nur sicherstellen, dass die Kovariablen in den beiden Therapiegruppen balanciert sind!

Ein Beispiel aus der Herzchirurgie



- Publierte PS-Analyse aus der koronaren Bypasschirurgie [Börgermann2012]
- 1282 Patienten, die zwischen Juli 2009 und November 2010 am Herz- und Diabeteszentrum NRW in Bad Oeynhausen isoliert koronarchirurgisch versorgt worden waren.
- Vergleich konventionelle Technik (cCABG, n=887, 69,2%) und Clampless-off-pump (Less-OPCAB, n=395, 30,8%)-Technik
- Entscheidung für Operationsmethode durch den jeweiligen Operateur

Ein Beispiel aus der Herzchirurgie

- Schätzung des PS-Modell durch logistisches Regressionsmodell
Alter, Geschlecht, BMI, Hauptstammstenose, LVEF, Präoperativer Myokardinfarkt, Hypertonie, Diabetes mellitus, COPD, Niereninsuffizienz, Schlaganfall, pAVK, Voroperationen, Dringlichkeit, Präoperative IABP
- 1:1-Matching mit logit-transformierten PS [Rubin2007], Optimal-Matching-Algorithmus, Caliper-Weite von 0,2 Standardabweichungen des logit-transformierten PS [Austin2011] .

Ein Beispiel aus der Herzchirurgie

Prüfen der Balanciertheit durch z-Differenz [Kuss2013]

- Maß für Balanciertheit, dass für metrische, binäre und ordinale Merkmale definiert und vergleichbar ist
- Prinzip: Standardisiere das jeweilige Unterschiedsmaß (Mittelwertsdifferenz, Risikodifferenz, Wilcoxon-Statistik) durch seinen Standardfehler (z-Standardisierung)

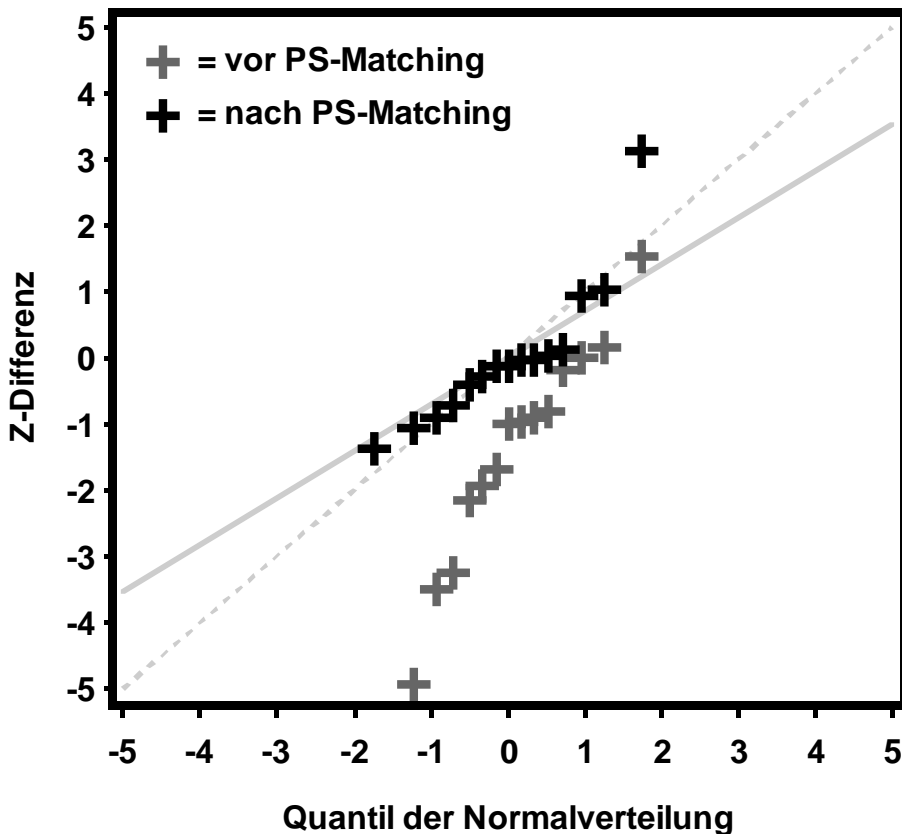
Ein Beispiel aus der Herzchirurgie

Prüfen der Balanciertheit durch z-Differenz [Kuss2013]

- In einem RCT sind die z-Differenzen standard-normalverteilt ($N(0,1)$)
- In einer perfekt gematchten Studie sind die z-Differenzen $N(0, \frac{1}{2})$ -verteilt [RubinThomas1996] (und damit besser als in einem RCT!)
- Berechne Mittelwert und Varianz der z-Differenzen vor und nach PS-Matching
- Besser: Zeichne Q-Q-Plot mit Referenzlinien für einen RCT und eine perfekt gematchte PS-Analyse vor und nach PS-Matching

Ein Beispiel aus der Herzchirurgie

Q-Q-Plot der z-Differenzen vor und nach PS-Matching



Vor PS-Matching:

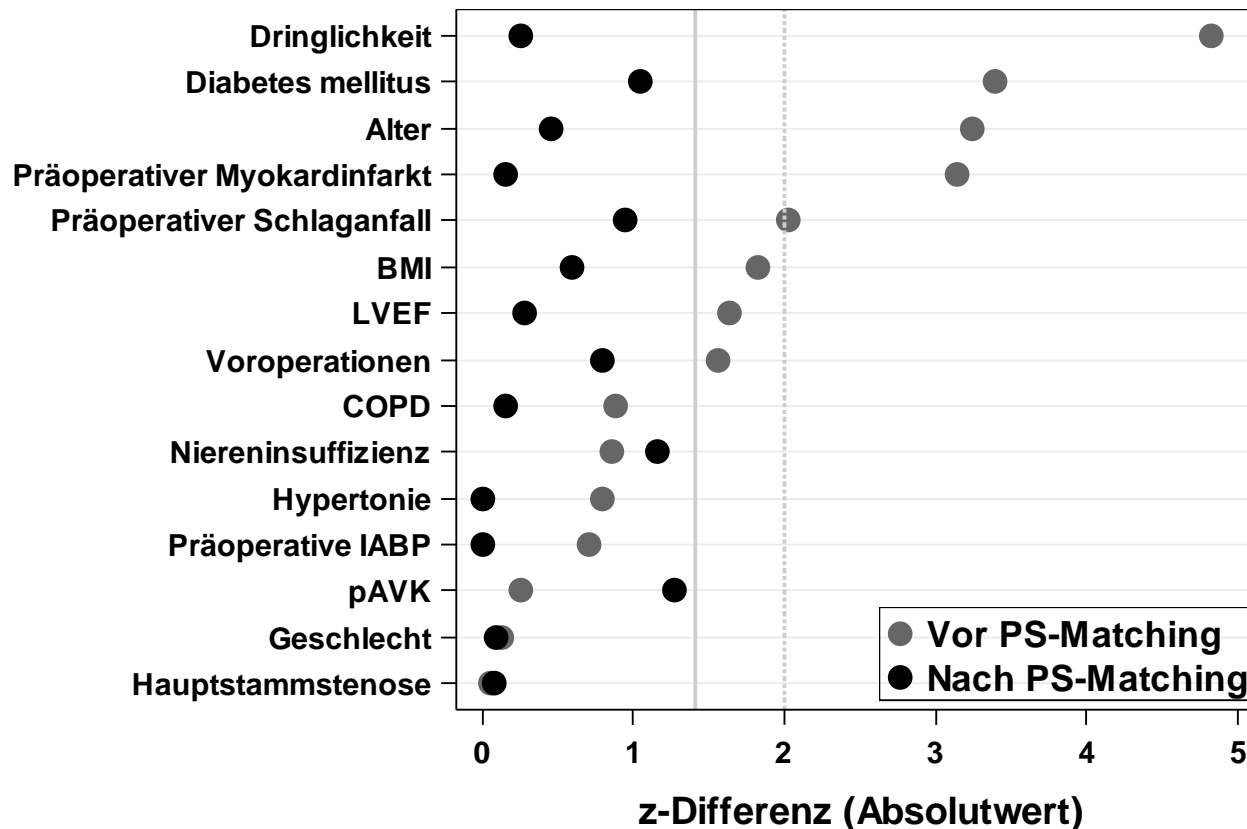
$$\bar{x}_{zDiff} = -0,99; s_{zDiff}^2 = 3,99$$

Nach PS-Matching

$$\bar{x}_{zDiff} = -0,05; s_{zDiff}^2 = 0,45$$

Ein Beispiel aus der Herzchirurgie

Absolute Werte der z-Differenzen vor und nach Matching



Ein Beispiel aus der Herzchirurgie

Eigentliche Auswertung des Therapieeffekt im PS-gematchten Kollektiv anhand dreier Zielgrößen:

- Tod oder Schlaganfall im Verlauf der Behandlung in der Klinik, ja/nein
- Operationsdauer
- Zeit bis Tod oder Schlaganfall in der Nachbeobachtung

Bei der Auswertung ist darauf zu achten, dass (z.B. durch eine konditionale Analyse) für das Matching-Stratum adjustiert wird [Austin2009]

Ein Beispiel aus der Herzchirurgie

	Less-OPCAB (n = 394)	cCABG (n = 394)	
Tod oder Schlaganfall im Verlauf der Behandlung i.d. Klinik [n (%)]	6 (1,5)	22 (5,6)	Odds ratio [95%-KI] 0,24 [0,09 – 0,63]
Operationsdauer in Minuten [Mittelwert (SD)]	175 (38)	180 (47)	MWD [95%-KI] -5 [-11; 1]
Zeit bis Tod oder Schlaganfall in der Nachbeobachtung [Ein-Jahres-Wahrscheinlichkeit für Ereignisfreiheit in %]	94,7	89,8	Hazard Ratio [95%-KI] 0,60 [0,35 – 1,03]

Propensity Score versus herkömmliche Regressionsmodelle

Vorteile PS-Methode

- PS-Analysen ähneln einem RCT
 - Zielgröße ist tabu/unbekannt bei der Schätzung des PS-Modells → geringere Manipulationsgefahr

“The most important flaw of regression adjustment for causal inference in observational studies is that study design is not separated from outcome analysis. How often does a researcher run only one regression model? It is tempting to fish for a certain result, fitting several models until the desired or expected answer appears”

[Pattanyak/Rubin/Zell2011]

Propensity Score versus herkömmliche Regressionsmodelle

Vorteile PS-Methode

- PS-Analysen ähneln einem RCT
 - Schätzung des PS-Modells gehört noch zum Design der Studie, nicht zur Auswertung [Rubin2007]
 - Beides Zwei-Schritt-Verfahren: Im ersten Schritt wird darauf Wert gelegt, dass beide Therapiegruppen ähnlich bzgl. der Patientenmerkmale sind (beim RCT durch Randomisierung, beim PS durch Berechnung des PS-Modells). Im zweiten Schritt wird dann im balancierten Sample der eigentlich interessierende Therapieeffekt geschätzt [Martens2008]

Propensity Score versus herkömmliche Regressionsmodelle

Vorteile PS-Methode

- Regressionsmodelle schätzen **immer** Therapieeffekte ...

Propensity Score versus herkömmliche Regressionsmodelle

Vorteile PS-Methode

- Regressionsmodelle schätzen **immer** Therapieeffekte ...
... selbst dann, wenn sich die beiden Gruppen von Behandelten und Nicht-Behandelten so extrem unterscheiden, dass eine solche Schätzung Unsinn ist

Regressionsmodelle machen Aussagen darüber was passiert wäre, wenn Behandelte nicht behandelt worden wären, nutzen dabei aber die Information von Nicht-Behandelten, die unter Umständen vollkommen anders sind als die Behandelten. Information über die Nicht-Behandelten wird dabei (durch Extrapolation) nur geschätzt, ist aber nicht wirklich beobachtet worden. [Stuart2009]

Propensity Score versus herkömmliche Regressionsmodelle

Vorteile PS-Methode

- In PS-Modelle können mehr Patientenmerkmale eingeschlossen werden können als in ein herkömmliches Regressionsmodell [Heinze2011].
In letzterem Fall würden zu viele Kovariablen schnell zu Overfitting und zu instabilen Schätzern für den Therapieeffekt führen.

Propensity Score versus herkömmliche Regressionsmodelle

Vorteile PS-Methode

- Für seltene Ereignisse ist die PS-Methode besonders überlegen [Cepeda2003].

Wenn die zu vergleichenden Therapie jeweils häufig angewandt werden, das eigentlich interessierende Zielereignis aber selten ist, dann wird es in der Regel so sein, dass nicht genug Information vorhanden ist, um den Zusammenhang zwischen Zielereignis und Patientenmerkmalen (einschließlich der Therapie) in einem herkömmlichen Regressionsmodell gut zu schätzen. Umgekehrt kann das PS-Modell, also den Zusammenhang zwischen Therapiezuweisung und Patientenmerkmalen gut geschätzt werden, weil dafür hinreichend Information vorhanden ist. [Braitman2002]

Fazit

- Die PS-Methode ist eine zwar nicht neue, aber doch innovative Methode zur Auswertung von nichtrandomisierten Therapiestudien, die sowohl statistische als auch erkenntnistheoretische Vorteile im Vergleich zur herkömmlichen Regression hat
- Unter den 4 Möglichkeiten zur Berücksichtigung des PS hat das PS-Matching eine Reihe von Vorteilen
- Aber: Der PS kann nur für die bekannten und tatsächlich gemessenen Confounder adjustieren!
Das gilt allerdings auch für die herkömmliche Regressionsmodelle für die Analyse von nichtrandomisierten Studien. [Austin2011]

Fazit

Die Frage, die über allem schwebt:

Können PS-Analysen eines Tages RCTs ersetzen?

Meine Antwort: Nicht ersetzen, aber PS-Analysen sollten viel häufiger auch für Therapieempfehlungen herangezogen werden!

Gründe:

- Auch RCTs haben Nachteile
- Immer bessere Evidenz, dass randomisierte und nicht-randomisierte in den meisten Fällen zu ähnlichen Ergebnissen führen [Anglemyer2014]
- Beispiele, wo sich Evidenz aus randomisierten und nicht-randomisierten Studien explizit widerspricht (z.B. WHI-Studie) sind wichtig aus historischen, pragmatischen oder pädagogischen Gründen [AbelKoch1999], bleiben aber Ausnahmen und können bei genauerer Analyse oft erklärt werden [Hernán2008].

Gründe:

- *Unbekannte* Confounder sind nur dann eine Gefahr, wenn diese *nicht* mit den *bekannt*en Confoundern assoziiert sind. Sind bekannte und unbekannte Confounder assoziiert, dann wird durch das Adjustieren für bekannte auch für die unbekannt e n Confounder mitadjustiert [Stuart2010]
- Es gibt schlicht zu viele Fragestellungen in der medizinischen Versorgung, als dass alle in randomisierten Studien beantwortet werden könnten. Die Gesellschaft wird sich weder die dazu nötigen Mittel noch die dazu nötige Zeit leisten können oder wollen [Borah2014]

Auf was alles nicht eingegangen wurde ...

- PS für mehr als zwei Behandlungen
- Schätzmethode, die mehrere PS-Methoden kombinieren
- Doubly-Robust-Methoden
- Theorie der kausalen Inferenz, Potential-Outcome-Modell, die konkreten Annahmen, die dem PS-Modell zugrunde liegen
- PS-Matching: Welches Matching-Verfahren sollte verwendet werden? (Optimal, Greedy usw.)

Auf was alles nicht eingegangen wurde ...

- PS-Matching: Machen auch 1:k-Matching-Verhältnisse Sinn? Antwort: Ja, aber es bringt nicht so viel wie erwartet
- Sensitivitätsanalysen für PS-Modelle (Rosenbaum: Design Sensitivity)
- Umgang mit fehlenden Werten [MitraReiter2012, QuLipkovich2009]
- Schätzung von marginalen Behandlungseffekten in RCTs und PS-Analysen vs. Schätzung von konditionalen Behandlungseffekten in herkömmlichen Regressionsanalysen [Martens2008]