

Bachelorarbeit

---

# Automatische Auswahl von maschinellen Lernverfahren für kausale Inferenz

von

Atanas Dimitrov

Pervasive Computing Systems / TECO

Institut für Telematik

Fakultät für Informatik

Abgabedatum: 02.09.2019

Verantwortlicher Betreuer:

Prof. Dr. Michael Beigl

Betreuerin:

Ployplearn Ravivanpong

## **Erklärung**

Hiermit erkläre ich, dass ich die vorliegende Bachelorarbeit selbstständig verfasst und keine anderen als die angegebenen Hilfsmittel und Quellen benutzt habe, die wörtlich oder inhaltlich übernommenen Stellen als solche kenntlich gemacht und weiterhin die Richtlinien des KIT zur Sicherung guter wissenschaftlicher Praxis beachtet habe.

Karlsruhe, 02.09.2019

---

## **Abstract**

Diese Bachelorarbeit beschäftigt sich mit dem Vergleich von Methoden für kausale Inferenz und mit der automatischen Auswahl von dem besten von denen abhängig von dem vorhandenen Datensatz. Dazu benutzen und erweitern wir Synth-Validation - ein Verfahren, mit dem von den echten Daten synthetische Daten mit einem gewünschtem durchschnittlichen Behandlungseffekt erstellt und dann ausgewertet werden. Dabei haben die Datensätze, auf denen wir unsere Experimente durchführen, unterschiedliche Natur - echten Rohdaten mit größeren oder kleineren Zahl von Kovariaten, von Rohdaten synthetisch generierten Daten und zufällig generierten Daten. Die kausale Inferenz Verfahren, von denen Synth-Validation auswählt, benutzen ausschließlich Algorithmen aus dem maschinellen Lernen. Es wird die Fähigkeit von Synth-Validation gemessen, den Verfahren zu wählen, der die nächste Schätzung von dem durchschnittlichen Behandlungseffekt hat. Das wird unter unterschiedlichen Konstellationen unterstellt - nach der Art der Daten, nach der Anzahl der Elementen in der Stichprobe usw.

# Inhaltsverzeichnis

<b>1</b>	<b>Einführung</b>	<b>1</b>
1.1	Kausale Inferenz . . . . .	1
1.2	Motivation . . . . .	3
1.3	Finden von ATE . . . . .	4
1.3.1	Experimentdurchführung . . . . .	4
1.3.2	Beobachtungsstudiedurchführung . . . . .	5
1.4	Ziele und Methodik . . . . .	7
<b>2</b>	<b>Methoden für kausale Inferenz</b>	<b>9</b>
2.1	Allgemeine Techniken . . . . .	9
2.1.1	Anpassung für Kovariaten . . . . .	9
2.1.2	Propensity Score Matching . . . . .	11
2.1.3	Inverse Probability Weighting . . . . .	13
2.2	Methoden in unserer Implementierung . . . . .	14
2.2.1	Quasi-Oracle Schätzung . . . . .	15
2.2.1.1	Lasso . . . . .	17
2.2.1.2	Gradient Boosting . . . . .	18
2.2.2	Kausale Wälder . . . . .	19
<b>3</b>	<b>Synth-Validation</b>	<b>22</b>
3.1	Generierung von synthetischen Daten . . . . .	22
3.1.1	Auswahl von synthetischen Effekten . . . . .	22
3.1.2	Schätzung von bedingten Erwartungswerten . . . . .	22
3.2	Methodenauswahl . . . . .	22
<b>4</b>	<b>Implementierung</b>	<b>23</b>
4.1	Fremde Bibliotheken . . . . .	23
4.2	Lesen/Schreiben von Daten . . . . .	23
4.3	Ziehen von Stichproben . . . . .	23
4.4	Methoden für kausale Inferenz . . . . .	23
4.5	Synth-Validation . . . . .	23
4.5.1	Datenstrukturen . . . . .	23
4.5.2	Schätzung . . . . .	23
4.5.3	Constrained Boosting . . . . .	23
4.5.4	Methodenauswahl . . . . .	23
4.6	Benchmark von Synth-Validation . . . . .	23
4.7	Erstellung von Abbildungen . . . . .	23
4.8	Experimenten . . . . .	23
4.9	Anderer Code . . . . .	23

<b>5</b>	<b>Ergebnisse und Evaluation</b>	<b>24</b>
5.1	Methodik und Daten . . . . .	24
5.2	TODO . . . . .	24
<b>6</b>	<b>Schlussfolgerung</b>	<b>25</b>
6.1	Zusammenfassung . . . . .	25
6.2	Diskussion . . . . .	25
	<b>Anhang</b>	<b>26</b>
	<b>Literatur</b>	<b>27</b>

# Abbildungsverzeichnis

1	Einfluss von Confounder auf Behandlung und Ergebnis . . . . .	3
---	---	---

# 1 Einführung

Um die Welt besser zu verstehen, haben die Menschen immer die Antwort auf der folgenden Frage gesucht: Was passiert, wenn eine bestimmte Tat oder Handlung durchgeführt wird? Manche Handlungen kann man als „einfach“ qualifizieren und bei denen ist diese Antwort leicht zu erreichen. Wenn man einen Apfel nach oben wirft, fällt er wieder nach unten. Es kostet (fast) nichts, diese Tatsache zu prüfen. Fast jeder kann dieses Ergebnis vorhersagen, weil man irgendwann mal ein fallendes Objekt beobachtet hat, unabhängig davon, ob man weiß, warum die Objekte fallen. Handlungen können aber deutlich komplexer sein wie z.B. Einnahme einer neuen Steuerpolitik oder die medizinische Behandlung mit einem Medikamenten. Diese Handlungen sind schwerer durchzuführen, sind von mehreren Faktoren oder Naturgesetzen betroffen, sind „abhängiger“. Auf Basis früherer Erfahrung im Bereich können Experten Rahmen davon setzen, was passieren wird. Sicher kann man nur dann sein, wenn die Handlung durchgeführt ist und man beobachtet und bemisst, was passiert hat. Die letzte Schlussfolgerung gilt sowohl für die „einfacheren“, als auch für die „komplexeren“ Handlungen. Also um eine sichere Kenntnis zu erschaffen, brauchen wir Erfahrung - in dem Fall mit dem Apfel entweder werfen wir den alleine oder beobachten wir jemanden, der das macht oder der das irgendwann mal gemacht hat. Bei den „komplexeren“ Handlungen ist meistens die Antwort auf der Frage „Was?“ nicht ausreichend - man braucht „Wie viel?“. Das führt uns langsam zu dem Sachverhalt der kausalen Inferenz.

## 1.1 Kausale Inferenz

Wenn wir über die Kausalität in ihrem wissenschaftlichen Sinn sprechen, muss es klar sein, dass „**A** verursacht **B**“ nicht bedeutet, dass **A** der Hauptgrund ist, warum **B** passiert hat. Es bedeutet, dass **A** den Unterschied oder die Differenz gemacht hat, die zu **B** geführt haben. Es existieren möglicherweise noch weitere Gründe z.B. **C** und **D**, ohne deren Existenz und ohne deren verursachten Differenz **B** nicht möglich war oder nicht den Wert halten würde, den es hält. **A** soll man nicht als eine Zuschreibung von **B** anschauen, sondern als einen Beitrag[Mac]. Dabei sollen wir unter **kausale Inferenz** ähnlich wie bei der statistischen Inferenz den Prozess und die Methodiken zum Finden von Kenntnissen über die Kausalität von irgendeiner Handlung, insbesondere das Finden von der Differenz, die sie verursacht.

Die kausale Inferenz hat aber ein fundamentales Problem. Wenn wir den Effekt von einer Handlung wissen wollen, führen wir diese Handlung aus und bemessen den Wert von irgendeine Variable, die uns interessiert. Der echte Effekt ist aber die Differenz zwischen dieser Variable bedingt, dass wir die Handlung durchgeführt haben und die gleiche Variable, bedingt, dass wir die Handlung nicht durchgeführt haben[Hol86]. Natürlich ist nur eine der beiden Optionen direkt beobachtbar. Formal kann man es so dargestellt:

$$\mathbf{TE} = (\mathbf{Y}|\mathbf{W} = 1) - (\mathbf{Y}|\mathbf{W} = 0) \quad (1.1)$$

Hier bezeichnet  $\mathbf{TE}$  das Behandlungseffekt,  $\mathbf{Y}$  ist die beobachtete Variable und  $\mathbf{W}$  ist eine binäre Variable, die bezeichnet, ob die Handlung durchgeführt wurde oder nicht.

Obwohl wir den Behandlungseffekt von einer Handlung individuell nicht beobachten können, können wir schätzen, wie hoch durchschnittlich diesen Effekt ist, wenn wir die Handlung auf mehreren Objekten durchführen beziehungsweise auch auf mehreren nicht durchführen und die Differenz zwischen diesen beiden Ergebnissen als **durchschnittliches Behandlungseffekt** oder auf Englisch **average treatment effect (ATE)** bezeichnen. Meistens können wir die Behandlung nicht auf allen (oder eigentlich auf der Hälfte von allen) Objekten durchführen, sondern müssen wir eine Stichprobe ziehen. Formal stellen wir das so dar:

$$\tau = E_Y[\mathbf{Y}|\mathbf{W} = 1] - E_Y[\mathbf{Y}|\mathbf{W} = 0] \quad (1.2)$$

In dieser Modelldarstellung sind  $\mathbf{W}$  und  $\mathbf{Y}$  Zufallsvariablen.  $\mathbf{W}$  ist binär, wo  $\mathbf{1}$  für eine Behandlung und  $\mathbf{0}$  für keine Behandlung steht.  $\mathbf{Y}$  ist die Zufallsvariable von einem gewünschten Wert nach der potenziellen Behandlung.  $E_Y$  steht für den Erwartungswert von  $\mathbf{Y}$ . Schon in dieser Darstellung treffen wir allen Problemen, die man sonst in der statistischen Inferenz trifft - damit unsere Schätzung über die Stichprobe näher zu der echten ATE der Gesamtpopulation liegt, müssen bestimmte Regeln eingehalten werden.

*Gleichung 1.2* ist immer noch kein gutes Model der Wirklichkeit, weil es oft der Fall ist, dass die Objekten, die wir behandeln, sich durch einigen Merkmalen unterscheiden. Diese andere Differenzen können Einfluss auf  $\mathbf{Y}$  sowohl bedingt  $\mathbf{W}$ , als auch unbedingt  $\mathbf{W}$  haben und diese Tatsache soll in der Gleichung einbezogen werden:

$$\tau = E_{Y,X}[\mathbf{Y}|\mathbf{X}, \mathbf{W} = 1] - E_{Y,X}[\mathbf{Y}|\mathbf{X}, \mathbf{W} = 0] \quad (1.3)$$

$\mathbf{X}$  ist eine Zufallsvariable für die so genannten **Kovariaten**.  $\mathbf{X}$  ist mehrdimensional und stellt den Zustand, in dem sich ein Objekt oder allgemeiner gesagt eine Welt befindet. Die Kovariaten sollen nicht, aber können einen Einfluss auf das Ergebnis  $\mathbf{Y}$  haben. Diejenigen mit Einfluss nennen wir **Confounders** oder **confounding Variablen**. Auf Deutsch kann man den Begriff am nächsten mit dem Wort Verwirrungsvariablen übersetzen. Die Existenz von Confounders erschwert unsere Aufgabe, den Effekt zu finden, den die Behandlung verursacht.[VS13].

In der Tat können Confounders auch indirekten Einfluss auf die Behandlung haben. Beispiel: Sei  $\mathbf{X}$  das Alter sein,  $\mathbf{W}$  - ob man raucht oder nicht und  $\mathbf{Y}$  - eine Gesundheitsmaß. Es ist klar, dass das Alter an sich einen Effekt auf der Gesundheit hat. Die empirische Daten zeigen



aber auch, dass das Rauchen unter älteren Menschen verbreiteter ist als unter jüngeren, d.h das Alter bewirkt sowohl das Rauchen, als auch die Gesundheit, die an sich vom Rauchen bewirkt wird. Diese Abhängigkeiten werden in *Abbildung 1* veranschaulicht.

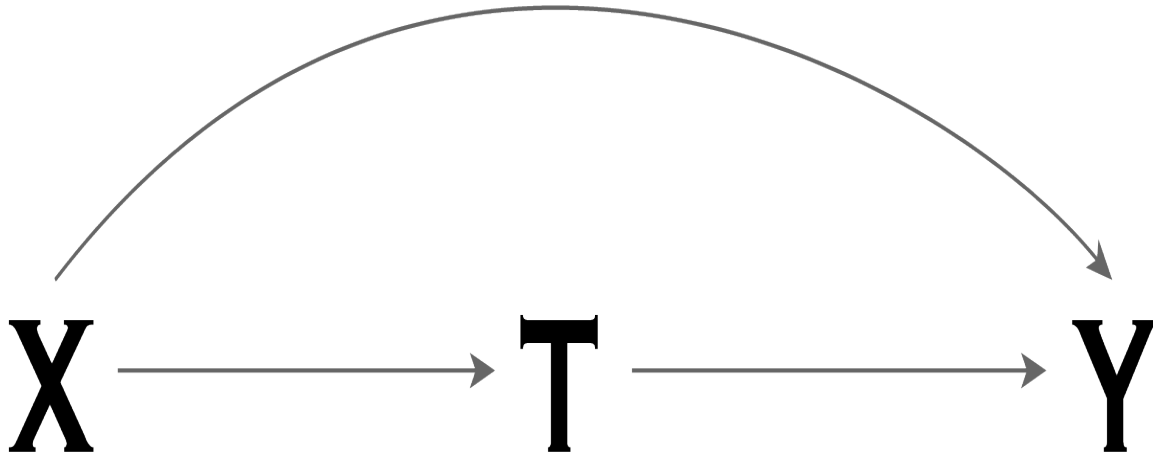


Abbildung 1: Einfluss von Confounder auf Behandlung und Ergebnis. Hier ist ***X*** der Confounder, ***T*** - die Behandlung und ***Y*** - das Ergebnis. Es wird gezeigt, dass das Ergebnis nicht nur von der Behandlung abhängt, sondern auch vom Confounder. Der Confounder hat auch eine Auswirkung auf die Behandlung.[Reb]

Wir schließen diesen Abschnitt mit der Gedanke von dem Philosophen David Lewis über die Kausalität, der sie folgendermaßen beschreibt: „etwas, was einen Unterschied verursacht und der Unterschied, den dieses etwas verursacht, soll der Unterschied davon sein, was ohne dieses etwas passieren würde“ [Lew74].

## 1.2 Motivation

Es ist wichtig zu wissen, was genau eine Handlung verursacht oder welche Differenz sie macht. Mit diesen Kenntnissen können wir besser die passenden Handlungen und deren genauen Maßen zum Erreichen unseres Ziels einschätzen und dann auswählen. Diese Behauptungen sind allgemeingültig und wir betrachten konkret drei Beispiele von unterschiedlichen Bereichen, wo das Kennen von dem durchschnittlichen Behandlungseffekt von großer Bedeutung ist.

**Beispiel 1** In der Pharmazie spielt die Kausalität eine besonders große Rolle. Wenn man einen Medikamenten entwickelt, soll man ganz genau wissen, welchen durchschnittlichen Effekt dieses Medikament auf unterschiedenen variablen körperlichen Werten hat - sowohl für wohltuende, als auch für schädliche Effekte. Hier ist das Nutzen, das die kausale Inferenz bringt, die Gesundheit der Menschen.

**Beispiel 2** Weitere Bedeutung hat die kausale Inferenz im Marketing und in der Bewertung von Werbungen. Eine Werbung kann eine Kette von unterschiedlichen Zielen haben, aber am Ende steht die Erhöhung von tatsächlichen Einkäufen von dem geworbenen Produkt oder Dienstleistung. Natürlich ist die Werbung nicht der einzige Treiber des Umsatzes. Also man muss nur den Effekt der Werbung einschätzen. Dann abhängig von diesem durchschnittlichen Effekt und dem Preis der Werbung stellt man fest, ob eine Werbung Nutzen bringt bzw. welche Werbung am nützlichsten ist.

**Beispiel 3** Die kausale Inferenz ist bedeutend für die Auswahl vom eigenen Verhalten sein. Zum Beispiel kann man den Einfluss vom regelmäßigen Rauchen von Rauschgift auf die mentalen Fähigkeiten untersuchen. Wenn man diesen durchschnittlichen Effekt kennt, sollte es wahrscheinlicher sein, dass man auf diese schlechte Sucht verzichtet oder am besten gar nicht damit anfängt. Diese deutliche wissenschaftliche Daten können durch den Gesetzgeber als Rechtfertigung für einen Verbot von diesem Stoff benutzt werden.

### 1.3 Finden von ATE

*Wir gehen davon aus, dass unser Problem schon klar definiert ist: Welchen durchschnittlichen Effekt(ATE) verursacht eine Behandlung? In diesem Abschnitt erläutern wir, wie man genau mit diesem Problem in den unterschiedlichen Situationen umgeht.*

Grundsätzlich muss man *Gleichung 1.3* lösen, die wir zur Bequemlichkeit hier noch mal hinbringen.

$$\tau = E_{Y,X}[Y|X, W = 1] - E_{Y,X}[Y|X, W = 0] \quad (1.4)$$

Dazu gibt es zwei Vorgehensweisen, die sich durch die Art von Datensammlung unterscheiden - Durchführung von einem Experimenten und Durchführung von einer Beobachtungsstudie. Die Daten aus den unterschiedlichen Studien wertet man dann anders aus.

#### 1.3.1 Experimentdurchführung

Es gilt, dass der durch Daten von Experimenten geschätzte Behandlungseffekt weniger Bias hat und deswegen näher zu dem echten durchschnittlichen Effekt der Population liegt. Der Grund dafür ist, dass die Objekten für die beiden Gruppen zufällig ausgewählt sind und somit wie die Population verteilt sind.[Rub74] Man soll immer diese Art von Datensammlung für kausale Inferenz wählen, wenn sie möglich ist. Von den im *Abschnitt 1.2* erwähnten Beispielen ist es bei Beispielen 1 und 2 möglich.

Das Experiment läuft folgendermaßen durch: zwei Gruppen von Objekten werden gebildet. Die erste Gruppe heißt Behandlungsgruppe und die zweite - Kontrollgruppe. Die Objekten von der ersten Gruppe werden mit irgendetwas z.B. einem Medikamenten behandelt, diese in der zweite jedoch nicht. Alle anderen Konditionen, auf die die Objekten in der Kontrollgruppe unterlegt sind, sollen sich nicht unterscheiden. Natürlich bleiben aber die einzelnen Objekten unterschiedlich, was zu unterschiedlichen Ergebnissen bei den individuell gemessenen Werten führt. Uns interessiert aber der durchschnittliche Wert von jeder Gruppe. Deswegen muss sichergestellt werden, dass die Objekte in jeder Gruppe gleichmäßig verteilt sind. Es gibt unterschiedliche Strategien das zu erreichen. Eine der besten, wenn sie richtig durchgeführt wird, ist die zufällige Auswahl von Objekten aus der Population für jede Gruppe. Am Ende wird für jedes Objekt der Wert von Bedeutung gemessen. Die Mittelwerten von diesen Werten werden für jede Gruppe berechnet und die Differenz dazwischen ist der erwartete durchschnittliche Behandlungseffekt.

Man kann auch Daten von den beiden Gruppen benutzen, um signifikanten Aussagen über den durchschnittlichen Behandlungseffekt mit einer bestimmten Konfidenz zu treffen. Dazu führt man meistens einen statistischen  $t$ -test durch. Man erstellt erstmal eine Nullhypothese und prüft, ob diese abzulehnen ist. Dafür berechnet man erstmal die Teststatistik  $t$ , die hier dargestellt ist:

$$t = \sqrt{\frac{n_0 n_1}{n_0 + n_1}} \frac{\bar{Y}_1 - \bar{Y}_0 - \omega}{s} \quad (1.5)$$

$$s = \sqrt{\frac{(n_0 - 1)s_0^2 + (n_1 - 1)s_1^2}{n_0 + n_1 - 2}} \quad (1.6)$$

$n_0$  und  $n_1$  sind die Anzahl von Objekten in den beiden Gruppen,  $\bar{Y}_1$  und  $\bar{Y}_0$  - die Mittelwerte von den Ergebnissen,  $s_0^2$  und  $s_1^2$  - die Varianzen,  $s$  - die gewichtete Abweichung,  $\omega$  - den Wert von dem durchschnittlichen Behandlungseffekt, den man für die Nullhypothese genommen hat. Dann abhängig davon, was man testen will und mit welcher Konfidenz die Behauptung gültig sein soll, vergleicht man den Wert von  $t$  mit dem Wert von einem Quantil von der studentischen  $t$ -Verteilung. Am Ende hat man eine Behauptung, die mit irgendeiner Signifikanz wahr ist.[Wikb]

### 1.3.2 Beobachtungsstudiedurchführung

Wir wiederholen noch mal, dass man lieber ein Experiment anstatt einer Beobachtungsstudie durchführen muss, um die Daten zu sammeln, wenn es möglich ist. Und natürlich ist es wegen unterschiedlichen Gründen nicht immer so.

Grundsätzlich kann ein Experiment viel mehr als eine Beobachtungsstudie kosten, weil man

die Objekten (die Leuten) in experimentellen Bedingungen haben möchte. Man hat (höhere) Kosten für:

- Vergütung der Leute
- die Behandlung
- Mitarbeiter, die das Experiment durchführen
- usw.

Bei manchen Studien kann der Prozess der Behandlung sogar länger dauern, was die oben genannten Kosten multipliziert.

Es gibt Behandlungen, wo es klar ist, dass sie einen negativen Effekt haben; man weiß doch nicht genau, wie hoch er ist. In diesem Fall ist es unethisch, Leuten auf diese negative Behandlung in Rahmen eines Experiments zu unterlegen. Vielleicht ist die Behandlung auch durchs Gesetz verboten. Genauso das ist die Situation im Beispiel 3 vom *Abschnitt 1.2*. Man darf nicht die Gesundheit der Leuten opfern, um etwas zu erfahren. In diesem Fall führt man eine Beobachtungsstudie unter Menschen durch, die auf die Behandlung unterlegt sind und zur Kontrolle auch unter solchen, die nicht unterlegt sind.

Und hier kommt ein großes Problem. Die Objekten, die behandelt sind, können anders verteilt sein als die ganze Population. Beispielsweise gibt es unter den Behandelten viel mehr Männer als Frauen. Die Behandlung auf Männer könnte aber einen anderen Effekt haben als diesen auf Frauen. Also wenn man seine Stichprobe durch zufälliges Ziehen bildet, bekommt man einen geschätzten durchschnittlichen Effekt, der von der höheren Anwesenheit von Männern gestreut ist und nicht an dem durchschnittlichen Effekt der Population entspricht. Man kann sich die *Abbildung 1* noch mal anschauen und feststellen, dass die Behandlung  $T$  von den Confounders  $X$  abhängt. In einem Experiment ist das ausgeschlossen, da es bedingt ist, dass am Anfang niemand behandelt war.

Weil die Daten aus Beobachtungsstudien viel Bias haben, würde eine bloße Berechnung von den Mittelwerten von den beiden Gruppen nicht ausreichend gut sein. Deswegen hat man unterschiedliche Methoden entwickelt, die wir ab hier **Methoden für kausale Inferenz** nennen, die grundsätzlich das Ziel haben, die Behandlungs- und die Kontrollgruppe zu normalisieren, damit die vergleichbar sind und damit das Ergebnis näher am Ergebnis der Population liegt. In dieser Bachelorarbeit haben wir über einige dieser Methoden im *Abschnitt 2* erzählt. Es ist natürlich zu erwarten, dass diese Methoden unterschiedliche Ergebnisse liefern. Dabei gibt es keine allgemein Beste - in unterschiedlichen Situationen funktioniert die eine besser als die anderen. Deswegen wurde ein Verfahren entwickelt, den die Methode für einen Datensatz automatisch wählt und auf diesen Verfahren basieren wir diese Bachelorarbeit.

## 1.4 Ziele und Methodik

*Nachdem wir in das Thema „Kausale Inferenz“ eingestiegen sind, werden wir im folgenden Abschnitt die Arbeitsmethodik dieser Bachelorarbeit erläutern. Dabei nennen wir auch die Ziele und Aufgaben, die wir uns setzen.*

Wir basieren unsere Arbeit auf den Synth-Validation Verfahren[SJT<sup>+</sup>17]. Synth-Validation ist ein Verfahren, der für einen bestimmten Datensatz, zu dem wir den durchschnittlichen Behandlungseffekt schätzen wollen, versucht, von einer Menge von Methoden für kausale Inferenz, die Beste auszuwählen und setzt diese ein. Im *Abschnitt 3* erzählen wir ausführlicher darüber. Wir arbeiten grundsätzlich an den Aufgaben, die im Abschnitt **Future work** im Artikel über Synth-Validation[SJT<sup>+</sup>17] genannt sind.

Das erste Ziel, die wir uns setzen, ist Synth-Validation auf R zu implementieren. Dazu verfügen wir über einen großen Teil des Codes. Der ist auf Julia geschrieben und wir haben den persönlich von einem der Autoren von Synth-Validation bekommen. Außer dem Übersetzen von Julia auf R müssen einige Teile von Synth-Validation, die im Artikel beschrieben sind, aber im vorgegebenen Julia Code jedoch fehlen, neu geschrieben werden. Darunter sind das Lesen von echten Daten, die Wahl von synthetischen Effekten, die finale Auswahl von einer Methode für kausale Inferenz von den synthetischen Daten und den ganzen Benchmarkprozess von Synth-Validation. Wir erzählen mehr darüber im *Abschnitt 4*.

Wir wählen R als Implementierungssprache, weil sie sehr verbreitet und ein Standard für statistikbasierte Ausarbeitungen ist. Es gibt eine große Gemeinschaft (auch unter Wissenschaftlern), die sie benutzt und deswegen sind auch viele Bibliotheken verfügbar. Dabei bietet die Scriptsprache einen einfachen Syntax, der erlaubt schnell einzusteigen und ein Basis für kürzere Implementierungszeit ist. Wir kennen natürlich auch die potentiellen Nachteilen, die die Sprache hat und die dazu geführt haben, dass man am Anfang entschieden hat, Synth-Validation auf Julia anstatt auf R zu implementieren. Diese besprechen wir im *Abschnitt 5*.

Eine weitere Aufgabe ist zu testen, wie Synth-Validation funktioniert, wenn sie auf echten Daten eingesetzt ist. Im Artikel von Synth-Validation wird beschrieben, dass sie nur mit zufällig generierten Daten getestet wurde. Das Finden von solchen Daten ist keine triviale Aufgabe, weil der echte Behandlungseffekt mit einer großer Sicherheit bekannt sein muss, damit wir Synth-Validation benchmarken können. Bei den automatisch generierten Daten kann man natürlich immer einen beliebigen Effekt haben, aber bei den echten ist das Finden von diesem Effekt eigentlich das, was wir von Anfang an erzielen wollen.

Die Daten, über die wir verfügen, kommen aus Wettbewerben für kausale Inferenz, deswegen haben sie einen bekannten durchschnittlichen Behandlungseffekt. Dabei unterscheiden wir zwischen Daten, die roh oder gar nicht bearbeitet wurden und solche, die durch echten Daten generiert wurden. Die Datensätze haben auch eine unterschiedliche Anzahl von Beobachtungen und eine unterschiedliche Anzahl von Kovariaten. Neben diesen echten Daten testen wir auch

mit solchen, die zufällig generiert sind und vergleichen unsere Ergebnisse mit diesen von den Autoren von Synth-Validation, soweit es möglich ist. Mehr über die Ergebnisse, ihre Auswertung und den Sachverhalt der Daten erzählen wir im *Abschnitt 5*.

Eine weitere Aufgabe ist Synth-Validation mit neuen Methoden zu testen. Seitdem Synth-Validation entwickelt und der Artikel darüber geschrieben wurde, sind einige neuen Methoden für kausale Inferenz entstanden, die bessere Ergebnisse im Allgemeinen aufweisen. Diese Methoden benutzen Algorithmen aus dem maschinellen Lernen. Es wird nützlich zu wissen, ob Synth-Validation immer noch erfolgreich unter denen wählt oder ob eine von denen sich besser als Synth-Validation vorstellt.

## 2 Methoden für kausale Inferenz

Man kann den durchschnittlichen Behandlungseffekt einer Handlung am einfachsten schätzen, indem man von dem Mittelwert der Ergebnisvariable der Behandlungsgruppe den Mittelwert der Ergebnisvariable der Kontrollgruppe abzieht. Diese einfache Schätzung ist sehr gut, wenn wir davon ausgehen, dass die beiden Gruppen homogen aneinander sind. Also wenn in der ersten Gruppe es keine Untergruppe von Objekten mit spezifischen Kovariaten gibt, die in der zweiten nicht vorgestellt sind. In diesem Fall sind die Daten im Gefallen von dieser Untergruppe gestreut und dieses Bias wird auf die Schätzung übertragen. Um das auszuschließen, wählt man zufällig Objekten aus, die nicht behandelt sind und führt ein Experiment mit diesen Objekten durch, in dem man die Hälfte behandelt und erst dann die Ergebnisvariable bemesst. Es gibt aber zu teure, unmögliche oder unethische Behandlungen, was dazu führt, dass die Daten darüber nur aus Beobachtungsstudien kommen. Diese Daten sind oft gestreut, was nämlich zu einer gestreuten Schätzung von der allereinfachste Methode führt[CM82a].

Um dieses Bias zu bekämpfen, hat man unterschiedliche Methoden und Techniken entwickelt. Im ersten *Unterabschnitt 2.1* werden wir über einige allgemeine, konzeptuelle Techniken erzählen, die zum Schätzen von kausalen Effekten aus Beobachtungsstudien erarbeitet wurden. Es gibt kaum einen konkreten Verfahren, an dem irgendeine davon nicht teilnimmt. Im zweiten *Unterabschnitt 2.2* stellen wir die Methoden vor, die wir in unserer Implementierung von Synth-Validation integriert haben, damit ihre Fähigkeit getestet wird, für einen bestimmten Datensatz die beste Methode zu finden.

### 2.1 Allgemeine Techniken

*In diesem Abschnitt erzählen wir über einige allgemeine Techniken, die für die bessere Modellierung unseres Problems beigetragen haben. Man benutzt sie in vielen konkreten Verfahren zur Schätzung von kausalen Effekten unter der Konstellation, dass die Daten aus Beobachtungsstudien kommen und die allereinfachste Schätzung ungenau ist. Diese Techniken wurden originell für lineare Schätzer erarbeitet und aus diesem Grund werden wir bei deren Darstellung hier einen linearen Schätzer benutzen.*

#### 2.1.1 Anpassung für Kovariaten

Wir schauen uns erstmal die einfachste Methode zum Schätzen von dem durchschnittlichen Behandlungseffekt genauer an.

$$\tilde{\tau} = \overline{Y_1} - \overline{Y_0} \tag{2.1}$$

Die Schätzung bekommen wir, indem wir die Differenz zwischen dem Mittelwert der Ergebnissen der Behandlungsgruppe und dem Mittelwert der Ergebnissen der Kontrollgruppe finden. Man stellt nicht so schwierig fest, dass die Kovariaten  $\mathbf{X}$  an der Rechnung gar nicht teilnehmen. Wenn die Annahme trifft, dass sie homogen in den beiden Gruppen sind, sollte das nicht problematisch sein. Wenn aber nicht, wird die Schätzung gestreut sein.

In diesem Abschnitt erzählen wir über einen Ansatz, der sich darauf basiert, wie  $\mathbf{Y}$  gebildet ist und wovon es abhängt. Nämlich wissen wir, dass  $\mathbf{Y}$  sowohl von den Kovariaten  $\mathbf{X}$ , als auch von der Behandlung  $\mathbf{W}$  abhängen kann. Wenn wir annehmen, dass diese Abhängigkeit linear ist, kann  $\mathbf{Y}$  durch das folgende lineare Model dargestellt werden:

$$\mathbf{Y}_i = \alpha + \beta \mathbf{W}_i + \gamma \mathbf{X}_i + \epsilon_i \quad (2.2)$$

In dieser Gleichung sind  $\mathbf{Y}_i$ ,  $\mathbf{W}_i$  und  $\mathbf{X}_i$  wie üblich das Ergebnis, die binäre Behandlungsvariable und die Kovariaten.  $\alpha$ ,  $\beta$  und  $\gamma$  sind die Koeffizienten, die den Beitrag von  $\mathbf{W}_i$  und  $\mathbf{X}_i$  in  $\mathbf{Y}_i$  zeigen.  $\gamma$  ist natürlich ein Vektor, da  $\mathbf{X}_i$  ein Vektor ist.  $\epsilon_i$  ist ein Störterm. Wir können diese Parameter durch lineare Regression schätzen. Unser durchschnittlicher Behandlungseffekt ist in dem Fall der geschätzte Wert von  $\beta$ [CM82a].

Damit wir aber kein gestreutes Ergebnis bekommen, müssen wir einige Bedingungen einhalten:

- Nur solche Kovariaten ins Modell nehmen, die  $\mathbf{Y}$  tatsächlich beeinflussen
- Keine Kovariaten nehmen, die von der Behandlung beeinflusst sind
- Die Anzahl der Objekten in der Stichprobe soll viel größer als die Anzahl gewählten Kovariaten sein.

Wir müssen nur solche Kovariaten in unseren Modell nehmen, die einen Einfluss auf das Ergebnis haben. Ansonsten wird die lineare Regression die Koeffizienten mit Bias anpassen, unter denen auch sich die Schätzung vom unseren durchschnittlichen Behandlungseffekten  $\beta$  befindet. Und das wollen wir natürlich nicht. Ein Problem ist aber, dass man nicht immer weiß, welche Variablen Einfluss haben. Eine hohe Korrelation zwischen dem Ergebnis und einem Kovariaten kann ein positives Zeichen dafür sein. Als Faustregel kann angenommen werden, dass diejenigen mit einer absoluten Korrelation kleiner als 0.2 nicht ins Modell genommen werden sollen. Es gilt, dass man die Kovariaten für das Modell bestimmen muss, bevor das Modell mit Daten angepasst ist. Ansonsten könnte man zu dem Ergebnis kommen, zu dem man kommen möchte, ohne dass es in der Grundgesamtheit repräsentativ ist[CM82a][Win].

Soweit man einen Kovariaten ins Modell nimmt, der von der Behandlung beeinflusst wurde, bekommt man Bias in der Schätzung von dem durchschnittlichen Effekt. Dieses Problem kann bei Experimenten nicht vorkommen, da man die Kovariaten immer vor der Behandlung bemisst. Bei den Beobachtungsstudien muss aufgepasst werden, dass ein Kovariat nicht mit der



Behandlung korreliert ist.[Lin]

### 2.1.2 Propensity Score Matching

Wenn wir Daten aus einer Beobachtungsstudie zum Finden von kausalen Effekten benutzen, ist es möglich, dass die Objekten, die behandelt sind, sich im Durchschnitt wesentlich oder zumindest in einiger Maßen von den Objekten, die nicht behandelt sind, unterscheiden. Das erklärt sich mit den Confoundingvariablen, die einen Effekt auf den Zustand der Behandlung haben. Wenn wir versuchen, den Behandlungseffekt von diesen qualitativ unterschiedenen Gruppen zu schätzen, bekommen wir ein Ergebnis mit Bias. Aus diesem Grund sind die beiden Gruppen nicht direkt vergleichbar.

Bevor man den Effekt schätzt, ist es sinnvoll, die beiden Gruppen irgendwie in Gleichgewicht zu bringen. Dieses Gleichgewicht werden wir durch zwei Eigenschaften definieren, die unseren beiden Gruppen erfüllen sollen. Diese Eigenschaften sind bei den Gruppen in Experimenten immer erfüllt und sind ausreichend für eine nicht gestreute Schätzung des durchschnittlichen Behandlungseffekts.

Die erste Eigenschaft ist folgende: bedingt von den Kovariaten  $\mathbf{X}$  sind die geschätzte Werte  $\mathbf{Y}_0$  und  $\mathbf{Y}_1$  unabhängig von  $\mathbf{W}$ . Also gegeben  $\mathbf{X}$  tragen  $\mathbf{Y}_0$  und  $\mathbf{Y}_1$  die gleichen Werte sowohl für  $\mathbf{W} = 1$ , als auch für  $\mathbf{W} = 0$  [RR83]. Formal wird diese Eigenschaft so dargestellt:

$$\mathbf{Y}_0, \mathbf{Y}_1 \perp\!\!\!\perp \mathbf{W} | \mathbf{X} \quad (2.3)$$

Die zweite Eigenschaft besagt, dass jedes Objekt von der Grundgesamtheit die Chance haben muss, sowohl behandelt, als auch unbehandelt zu sein. Anders gesagt ist die Wahrscheinlichkeit für ein Objekt mit bestimmten Kovariaten, behandelt zu sein, streng größer 0 und streng kleiner 1 [RR83]. Formal haben wir:

$$0 < pr(\mathbf{W} = 1 | \mathbf{X}) < 1 \quad (2.4)$$

Wir brauchen also irgendeiner Funktion, die unseren Kovariaten normalisiert und nach deren Einsatz die oberen Bedingungen erfüllt sind. Deswegen führen wir die Balancing Score oder die Gleichgewichtsmaßfunktion  $\mathbf{b}(\mathbf{X})$  ein, die die Kovariaten als Argument nimmt und die folgende Eigenschaft hat: die bedingte Verteilung von  $\mathbf{X}$  gegeben  $\mathbf{b}(\mathbf{X})$  ist gleich für  $\mathbf{W} = 1$  und  $\mathbf{W} = 0$  oder anders gesagt ist unabhängig von  $\mathbf{W}$  [CM82b]. Formal wird das so dargestellt:

$$\mathbf{X} \perp\!\!\!\perp \mathbf{W} | \mathbf{b}(\mathbf{X}) \quad (2.5)$$

Die einfachste und die feinste mögliche Balancing Score ist  $\mathbf{b}(\mathbf{X}) = \mathbf{X}$ . Eine Balancing Score  $\mathbf{b}(\mathbf{X})$  ist feiner als eine andere Balancing Score  $\mathbf{e}(\mathbf{X})$ , wenn es eine Funktion  $\mathbf{f}$  gibt, für die nicht die Identität ist und für die  $\mathbf{e}(\mathbf{X}) = \mathbf{f}(\mathbf{b}(\mathbf{X}))$  gilt. Je feiner eine Balancing Score ist, desto besser, soweit sie natürlich die Eigenschaften von einer Balancing Score trägt. Wenn wir aber viele Kovariaten haben, ist das Finden von der feinsten Balancing score keine leichte Aufgabe. Dazu ist die Schätzung von der grobsten Balancing Score relativ einfach. Die nennen wir Propensity Score oder die Neigungsmaßfunktion  $p(\mathbf{X})$ . Die Propensity Score ist die Wahrscheinlichkeit für einen Objekt mit bestimmten Kovariaten  $\mathbf{X}$ , dass es behandelt wird. Formal:

$$p(\mathbf{X}) = pr(W = 1|\mathbf{X}) \quad (2.6)$$

Dementsprechend kann es mehr unterschiedlichen Objekten geben, die gleiche Wahrscheinlichkeit zur Behandlung haben. Rosenbaum und Rubin haben bewiesen, dass Gruppen, die von Objekten mit einander gleichen Propensity Scores erstellt sind, die obengenannten Eigenschaften erfüllen. Aus diesem Grund folgt, dass der geschätzte durchschnittliche Behandlungseffekt von diesen Gruppen kein Bias hat[RR83].

Die Propensity Score können wir durch logistische Regression schätzen. Wir haben das Modell:

$$\frac{pr(W = 1|\mathbf{X})}{1 - pr(W = 1|\mathbf{X})} = e^{\alpha + \beta \mathbf{X}} \quad (2.7)$$

Wir vereinfachen den Ausdruck, so dass wir die Propensity Score Funktion sauber haben:

$$p(\mathbf{X}) = pr(W = 1|\mathbf{X}) = \frac{1}{1 + e^{-(\alpha + \beta \mathbf{X})}} \quad (2.8)$$

Hier sind  $\mathbf{X}$  die Kovariaten,  $\alpha$  und  $\beta$  sind die Koeffizienten, die wir durch die logistische Regression schätzen, indem  $\beta$  ein Vektor ist, da  $\mathbf{X}$  auch Vektor ist. Hier ist es wichtig, dass wir für die Schätzung nur Confounder benutzen, also Kovariaten, für die wir vermuten, dass sie den Zustand der Behandlung beeinflussen. Ansonsten riskieren wir, dass unsere Propensity Scores von unwesentlichen Kovariaten gestreut sind.

Nachdem wir die Propensity Score Funktion geschätzt haben, können wir die Propensity Scores von allen beobachteten Objekten berechnen. Dann bilden wir unsere Behandlungsgruppe, indem wir eine Stichprobe von behandelten Objekten ziehen. Jetzt müssen wir die zweite Stichprobe bilden, indem wir jedes Objekt aus der Ersten mit einem oder mehr Kontrollobjekten nach Propensity Score verbinden. Für diese Verbindung können wir unterschiedliche Strategien benutzen. Man kann nach exakt gleichen Propensity Scores verbinden, was schwierig ist, da wir eine sehr große Stichprobe und/oder Kovariaten mit wenigen optionalen Werten haben sollen, damit wir genug gleiche Propensity Scores haben. Andere Option wäre, unterschiedliche Intervallen für die Propensity Score zu bilden und für jedes behandeltes Objekt ein Kontrollobjek

zu wählen, dessen Propensity Score im gleichen Intervall liegt wie die Propensity Score von dem behandelten Objekt. Man kann auch die Methode der nächsten Nachbarn nutzen, in dem man für jedes behandelte Objekt ein Kontrollobjekt wählt, deren Propensity Score am nächsten zu der Propensity Score des behandelten Objekts liegt. Oder von allen Kontrollobjekten, deren Propensity Score in einem vordefinierten Abstand von der Propensity Score von jedem behandelten Objekt liegt, einen wählen[Stu10].

Ein Nachteil von den oben vorgestellten direkten Verbindungenstechniken ist die Tatsache, dass einige beobachtete Kontrollobjekte von der Analyse ausgeschlossen werden. Wenn man das nicht will, kann man allen Objekten nach Propensity Score gewichten. Darüber erzählen wir im *Abschnitt 2.1.3*.

Nachdem die beiden Gruppen gebildet sind, bleibt nur die Schätzung von dem durchschnittlichen Behandlungseffekt, die wie üblich stattfindet: die Differenz zwischen den Mittelwerten von der Behandlung - und Kontrollgruppe.

### 2.1.3 Inverse Probability Weighting

*In diesem Abschnitt werden wir eine Technik zur Schätzung von dem durchschnittlichen Behandlungseffekt vorstellen, die man teilweise als Erweiterung des oberen Abschnitts ansehen kann, weil sie die Propensity Scores benutzt.*

Inverse probability of treatment weighting (IPTW) oder auf Deutsch Gewichtung mit der inversen Wahrscheinlichkeit der Behandlung ist ein Verfahren, mit dem man ein Gewicht für jedes Objekt von der Beobachtungsstudie schätzt und dieses Gewicht zur Schätzung des durchschnittlichen Behandlungseffekts nutzt. Wie der Name der Methode sagt, nehmen wir in der Gewichtsmaße die Wahrscheinlichkeit von einer Behandlung für dieses Objekt oder die Propensity Score, die wir im oberen *Abschnitt 2.1.2* erläutert haben. Das Gewicht eines Objekts  $i$  berechnen wir folgendermaßen:

$$w_i = \frac{T_i}{p(X_i)} + \frac{1 - T_i}{1 - p(X_i)} \quad (2.9)$$

Hier ist  $T_i$  die Behandlung, die wir in dieser Bachelorarbeit sonst mit  $W_i$  darstellen.  $p(X_i)$  ist die Propensity Score vom Objekt  $i$ . Wir beobachten, dass ein Gewicht  $w_i$  Werte von 1 bis  $\infty$  nimmt, da die Propensity Score eine Wahrscheinlichkeitsmaß ist, die zwischen 0 und 1 liegt. Der Einsatz der Gewichten kann man in diesem Fall als Erstellung von einer Pseudopopulation ansehen, in der jedes Objekt von unseren Beobachtungen  $w_i$ -mal daran teilnimmt. In dieser Pseudopopulation sind die Objekten, für die eine Behandlung wahrscheinlicher ist, deutlich weniger als diese, für die sie unwahrscheinlicher ist. Deswegen gibt es in dieser neuen Population keine Confounder mehr, also keine Kovariaten, die Einfluss auf der Behandlung haben, und wir

können den durchschnittlichen Behandlungseffekt ungestreut schätzen[RHB00].

Es gibt eine alternative Weise, die Gewichte berechnen. Da bleibt das Gewicht der behandelten Objekten konstant 1 und das Gewicht der Kontrollobjekten wird durch seine Propensity Score bestimmt. Formal ist für ein Objekt  $i$ :

$$w_i = T_i + (1 - T_i) \frac{p(\mathbf{X}_i)}{1 - p(\mathbf{X}_i)} \quad (2.10)$$

Die Signatur der Variablen ist wie bei *Gleichung 2.9*. Jetzt haben wir eine Pseudopopulation, in der jedes behandelte Objekt einmal vorkommt. Je wahrscheinlicher ist es, dass für ein Kontrollobjekt die Behandlung **1** wäre, desto öfter kommt es vor[HIR03].

Nachdem wir die Gewichte geschätzt haben, können wir unseren durchschnittlichen Behandlungseffekt schätzen. Im Unterschied zu dem üblichen Fall, wo wir die Mittelwerte der Ergebnisvariable  $\mathbf{Y}$  für den beiden Gruppen ausrechnen, nehmen wir jetzt die gewichteten Mittelwerte mit dem individuellen Gewicht von jedem Objekt und ziehen diese ab.

Der Vorteil von einer Gewichtung ist, dass diese Methode alle Objekten für die Schätzung von dem durchschnittlichen Behandlungseffekt nimmt, was bei den üblichen Verbindungsstrategien, die wir kurz am Ende des *Abschnitt 2.1.2* erwähnt haben, nicht der Fall ist. Der Nachteil ist, dass ein Gewicht bei einer zu hohen oder zu niedrigen Propensity Score sehr groß werden kann, was Varianz in der finalen Schätzung bringt. Wenn die Propensity Scores richtig geschätzt sind, sollte dieser Nachteil keine Rolle spielen. Wenn sie aber gestreut sind, wird dieses Bias noch stärker in der finalen Schätzung übertragen. Deswegen kann man einen Maximumwert einführen, den genommen wird, wenn ein Gewicht es überschreitet [Stu10].

## 2.2 Methoden in unserer Implementierung

*In diesem Abschnitt werden wir die Methoden für kausale Inferenz vorstellen, die wir in unserer Implementierung von Synth-Validation integriert haben. Sie sind nach der Entwicklung von Synth-Validation entstanden und zeigen bessere Ergebnisse im Allgemeinen als die älteren Methoden, deswegen haben wir sie genommen. Da die drei getesteten Methoden heterogene Behandlungseffekte schätzen, erklären wir zuerst, wie sich der heterogene Behandlungseffekt von dem durchschnittlichen Behandlungseffekt unterscheidet. Dann deuten wir an, wie wir aus den heterogenen Schätzungen den durchschnittlichen Effekt berechnen. Im Abschnitt 2.2.1 stellen wir die Modellbedingungen der R-learner Methode vor, die sich für die beiden Schätzer (das Lasso und das Gradient Boosting) nicht unterscheiden. Dann stellen wir individuell die beiden Schätzer dar. Im Abschnitt 2.2.2 erzählen wir über eine andere Schätzmethode - die kausale Wälder, die Random Forests für die Schätzung benutzt.*

Die Methoden, die wir in unserer Implementierung integriert haben, schätzen den heterogenen Behandlungseffekt. Wir wiederholen, dass der durchschnittliche Effekt die Differenz zwischen die durchschnittliche Ergebnisvariable aller behandelten Objekten und die durchschnittliche Ergebnisvariable aller Kontrollobjekten ist, wenn wir annehmen, dass keine Confounder gibt oder ihr Effekt entfernt ist. Wir bringen noch mal *Gleichung 1.3* zur Bequemlichkeit:

$$\tau = E_{Y,X}[Y|X, W = 1] - E_{Y,X}[Y|X, W = 0] \quad (2.11)$$

Der heterogene Behandlungseffekt ist nicht der durchschnittliche, sondern der individuelle Effekt, den ein Objekt oder eine Ergebnisvariable dieses Objekts nach einer Behandlung bekommt. Der heterogene Effekt ist die Differenz zwischen die potentiellen Werten dieser Variable mit und ohne Behandlung für ein spezifisches Objekt. Die Schätzung ist nicht mehr einen einzigen Wert  $\tau$ , sondern eine Funktion  $\tau(\mathbf{x})$ , die von den Werten der Kovariaten  $\mathbf{x}$  eines Objekt abhängt. Formal stellen wir die Funktion folgendermaßen dar:

$$\tau(\mathbf{x}) = E_X \left[ E_Y[Y|W = 1] - E_Y[Y|W = 0] \middle| X = \mathbf{x} \right] \quad (2.12)$$

Da wir für unsere Implementierung den durchschnittlichen Effekt brauchen, wenden wir die geschätzte Funktion  $\tau(\mathbf{x})$  für alle unserer Objekten an und finden dann den Mittelwert von diesen individuellen Behandlungseffekten. Formal:

$$\tau = \frac{1}{N} \sum_{i=1}^N \tau(\mathbf{x}_i) \quad (2.13)$$

### 2.2.1 Quasi-Oracle Schätzung

*In diesem Abschnitt werden wir über die Modellrahmen von der Quasi-Oracle Schätzung von heterogenen Behandlungseffekten[NW17] erzählen, die von Xinkun Nie und Stefan Wager entwickelt wurde. Intern hat die Methode entweder das Lasso oder Gradient Boosting als Schätzer und wir haben die beiden Varianten der Methode in unserer Implementierung von Synth-Validation integriert. Über die beiden Schätzer erzählen wir in den Unterabschnitten.*

Die Methode der Quasi-Oracle Schätzung von heterogenen Behandlungseffekte oder der R-learner, der von Nie und Wager entwickelt wurde, definiert das allgemeine Problem zum Schätzen vom heterogenen Behandlungseffekt als Optimierungsproblem, das man mit einer Menge von unterschiedliche Verfahren lösen kann. Davor muss man zwei Funktionen schätzen, die im Problem eingesetzt sind. Wenn man zum Lösen des Optimierungsproblems das Lasso benutzt, verhält sich die Methode als ein Quasi-Orakel Schätzer - sogar im Fall von einer ungenauen Schätzung von den beiden Funktionen im ersten Schritt ist der Schätzwert der Methode

so gut wie dieser von einem Quasi-Orakel. Der Quasi-Orakel kennt die beiden Funktionen und unterscheidet sich mit dem normalen Orakel damit, dass er nur den echten Behandlungseffekt nicht kennt. Wir werden jetzt das von der Methode benutzte Modell formal ausführlich anschauen. Alle mit  $*$  hochgestellten Signaturen in diesem Abschnitt sind unbekannt, können aber geschätzt werden. Zuerst stellen wir die potenziellen Ergebnissen  $\mathbf{Y}_i(\mathbf{1})$  und  $\mathbf{Y}_i(\mathbf{0})$  vor, die die Variable  $\mathbf{Y}$  nimmt, wenn man das Objekt  $i$  behandelt oder nicht behandelt. Weiter benutzt das Modell die Propensity Score  $e^*(\mathbf{x})$ , über die wir im *Abschnitt 2.1.2* erzählt haben:

$$e^*(\mathbf{x}) = P[\mathbf{W} = 1 | \mathbf{X} = \mathbf{x}] \quad (2.14)$$

Danach stellen wir der bedingte Oberfläche vom Ergebnis  $\mu_{(w)}^*(\mathbf{x})$  vor, den wir für die Behandlung  $\mathbf{w} = \{0, 1\}$  haben:

$$\mu_{(w)}^*(\mathbf{x}) = E[Y(\mathbf{w}) | \mathbf{X} = \mathbf{x}] \quad (2.15)$$

Noch haben wir den von der Behandlung bedingten Restwert  $\epsilon_i(\mathbf{w})$  von dem potenziellen Ergebnis  $\mathbf{Y}_i(\mathbf{w})$ :

$$\epsilon_i(\mathbf{w}) = \mathbf{Y}_i(\mathbf{w}) - (\mu_{(0)}^*(\mathbf{X}_i) + \mathbf{w}\tau^*(\mathbf{X}_i)) \quad (2.16)$$

Weiter stellen wir den bedingten Erwartungswert des Ergebnisses  $\mathbf{m}^*(\mathbf{x})$  dar:

$$\mathbf{m}^*(\mathbf{x}) = E[\mathbf{Y} | \mathbf{X} = \mathbf{x}] = \mu_{(0)}^*(\mathbf{X}_i) + e^*(\mathbf{x})\tau^*(\mathbf{X}_i) \quad (2.17)$$

Man kann jetzt die Funktion des heterogenen Behandlungseffekts  $\tau^*(\mathbf{x})$  anders darstellen, indem man die folgende Transformation benutzt:

$$\mathbf{Y}_i - \mathbf{m}^*(\mathbf{X}_i) = (\mathbf{W}_i - e^*(\mathbf{X}_i))\tau^*(\mathbf{X}_i) + \epsilon_i(\mathbf{W}_i) \quad (2.18)$$

Diese Darstellung nennt man die Transformation von Robinson[Rob88] und um ihn zu beehren, haben die Autoren von dieser Methode für kausale Inferenz sie nach ihm R-learner benannt. Der heterogene Behandlungseffekt nimmt dann die Form:

$$\tau^*(.) = \underset{r}{\operatorname{argmin}} \left\{ E \left[ \left( (\mathbf{Y}_i - \mathbf{m}^*(\mathbf{X}_i)) - (\mathbf{W}_i - e^*(\mathbf{X}_i))\tau^*(\mathbf{X}_i) \right)^2 \right] \right\} \quad (2.19)$$

Man bemerkt, dass der Term von  $\epsilon_i(\mathbf{W}_i)$  im oberen Ausdruck 2.19 im Vergleich zum Ausdruck 2.18 ausgefallen hat. Das hat ganz gezielt passiert, weil wir beobachten, dass die Confounder genauso dann keinen Einfluss auf den finalen Wert von  $\mathbf{Y}_i$  habe, wenn gilt:

$$E[\epsilon_i(\mathbf{W}_i)|\mathbf{X}_i, \mathbf{W}_i] = 0 \quad (2.20)$$

Wir sind jetzt bereit, unser finalen Problem in zwei Schritte zu zerlegen:

1. Schätzung von  $\mathbf{m}^*(\mathbf{x})$  und  $\mathbf{e}^*(\mathbf{x})$
2. Lösung von dem folgenden Optimierungsproblem:

$$\begin{aligned} \tau^*(.) &= \operatorname{argmin}_r \{L_n(r(.)) + \Lambda_n(r(.))\}, \\ L_n(r(.)) &= \frac{1}{n} \sum_{i=1}^n \left( (Y_i - m(X_i)) - (W_i - e(X_i))\tau^*(X_i) \right)^2 \end{aligned} \quad (2.21)$$

$L_n(r(.))$  ist die Verlustfunktion, für die wir  $\mathbf{r}(.)$  finden wollen, für die die Funktion minimal ist.  $\Lambda_n(r(.))$  ist ein Regulierungsterm, der durch Cross-Validation bestimmt werden kann. Bei einer solchen Definition des Problems hat man die Möglichkeit und die Flexibilität, es mit unterschiedlichen Verfahren zu lösen - lineare Modellen, Boosting, neuronale Netze usw. Die Autoren vom R-learner benutzen das Lasso oder Gradient Boosting für den ersten Schritt und das Lasso für den zweiten. In den folgenden zwei Unterabschnitten werden wir über diese zwei Schätzmethoden erzählen.

### 2.2.1.1 Lasso

Das Lasso oder LASSO steht für Least Absolute Shrinkage and Selection Operator oder auf Deutsch Selektion Operator mit absoluter Schwindung. Es ist eine Erweiterung von der üblichen linearen Regression und hat viele Varianten. Wir stellen hier die allereinfachste Form vor, die aber schon die wichtigsten Charakteristiken der Methode beinhaltet.

Die lineare Regression oder genauer gesagt die Methode der kleinsten Quadraten ist im Grunde genommen für eine Datenmenge von  $\mathbf{y}_i$  und  $\mathbf{x}_i$  eine Lösung des folgenden Optimierungsproblems:

$$\min_{\alpha, \beta} \left\{ \frac{1}{N} \sum_{i=1}^n (y_i - \alpha - \mathbf{x}_i^T \beta)^2 \right\} \quad (2.22)$$

Dabei sind die Ergebnisse  $\alpha$  und  $\beta$ , die die Methode angepasst hat, die Koeffizienten von einem linearen Model. Durch Anwendung dieser Koeffizienten kann man weitere  $\mathbf{y}_i$  für gegebene  $\mathbf{x}_i$

vorhersagen. Das Fehler dieser Vorhersagen bemisst man mit der Größe des Bias und der Varianz. Und soweit die lineare Regression mit wenig Bias schätzt, hat sie oft viel Varianz. Das liegt daran, dass manche von den Inputvariablen in  $\mathbf{x}$  in der Tat keinen Einfluss auf das Ergebnis  $\mathbf{y}$  haben. Diese falsche Modellierung führt dazu, dass die  $\alpha$  und  $\beta$  nicht im allgemein, sondern nur für den Trainingsdatensatz stimmen.

Um dieses Problem zu vermindern, kann man die Betas begrenzen oder regulieren, so dass sie nicht beliebig viel wachsen und damit Werte nehmen, die das Optimierungsproblem eigentlich besser lösen, aber in der Tat ungültig sind. Man bestimmt einen Parameter  $t$ , den die absolute Summe von  $\beta$  nicht überschreiten soll[Tib96]. Wir fügen diese neue Bedingung in das Optimierungsproblem ein:

$$\min_{\alpha, \beta} \left\{ \frac{1}{N} \sum_{i=1}^n (y_i - \alpha - \mathbf{x}_i^T \beta)^2 \right\} \text{ beschränkt von } \sum_{j=1}^p |\beta_j| \leq t \quad (2.23)$$

Wir geben das Problem auch in seiner Lagrangeform:

$$\min_{\alpha, \beta} \left\{ \frac{1}{N} \sum_{i=1}^n (y_i - \alpha - \mathbf{x}_i^T \beta)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\} \quad (2.24)$$

Durch diese Regulierung erzielen wir folgendes: unwichtige Faktoren bekommen kleinere Werte in ihre Betas. Manche Betas werden zu  $\mathbf{0}$  und sind somit als unwichtig von dem Modell ausgeschlossen. Das bedeutet, dass wir das Lasso benutzen können, wenn wir unsicher sind, welche  $\mathbf{x}$  den Wert von  $\mathbf{y}$  bestimmen[Tib96]. Wir müssen nur den passenden Wert von  $\lambda$  finden. Je höher  $\lambda$  ist, desto mehr Faktoren reduzieren wir. Je kleiner ( $\mathbf{0}$  als Minimum), desto ähnlich ist das Lasso an der linearen Regression. Zu niedrige  $\lambda$  bedeutet potenzielle hohe Varianz, zu hohe - ein gestreutes Ergebnis. Für ihre Bestimmung können wir Cross-Validation benutzen[Tib96].

### 2.2.1.2 Gradient Boosting

Gradient Boosting ist eine Technik aus dem maschinellen Lernen, die als Ziel hat, für gegebenen Daten  $\mathbf{y}_i$  und  $\mathbf{x}_i$  eine Funktion  $\mathbf{y} = \mathbf{F}(\mathbf{x})$  zu finden, die ihre Beziehung am besten beschreibt. Dazu benutzt sie viele „schwache“ Lerner, die individuell nicht so gut sind, aber zusammen eine insgesamt gute Schätzung produzieren. Jeder neue Lerner trägt iterativ bei. Wie üblich im maschinellen Lernen haben wir auch hier eine Verlustfunktion  $L(\mathbf{y}, \hat{\mathbf{y}})$ , in der  $\mathbf{y}$  für die echten Werten von den Daten und  $\hat{\mathbf{y}}$  für die von dem Modell geschätzten Werten steht. Wenn die Summe der Verlustfunktion eingesetzt mit jedem  $\mathbf{y}_i$  minimal ist, haben wir die beste Anpassung erreicht. Beim Gradient Boosting passiert diese Anpassung iterativ: wir fangen mit einem schwachen Schätzer  $\mathbf{F}_0(\mathbf{x})$  an, der oft den Mittelwert  $\bar{\mathbf{y}}$  ist[Wika]. Jeder folgende Schätzer  $\mathbf{F}_{m+1}$  stellen wir dann so dar:



$$\mathbf{F}_{m+1} = \mathbf{F}_m(\mathbf{x}) + \mathbf{h}(\mathbf{x}) \quad (2.25)$$

Hier ist  $\mathbf{h}(\mathbf{x})$  die Funktion, die den vorherigen Schätzer am besten ergänzt, so dass die neue Schätzung am wenigsten von den echten  $\mathbf{y}$  abweicht. Diese Abweichung bemessen wir durch die Verlustfunktion. Unser finaler Schätzer  $\hat{\mathbf{F}}(\mathbf{x})$  hat dann die Form:

$$\hat{\mathbf{F}}(\mathbf{x}) = \sum_{i=1}^M \gamma_i \mathbf{h}_i(\mathbf{x}) + \text{const} \quad (2.26)$$

Hier ist  $\mathbf{h}_i(\mathbf{x})$  die Ergänzung von jedem weiteren schwachen Schätzer und  $\gamma_i$  ist sein Gewicht. **const** ist der Schätzwert von dem allerersten Schätzer. In jedem Schritt schätzen wir dann die  $\mathbf{h}_m(\mathbf{x})$  mit einem schwachen Schätzer  $\mathbf{m}$ , der als Eingabe alle  $\mathbf{x}_i$  und die sogenannten Pseudorestwerte  $\mathbf{r}_{im}$  nimmt. Danach wird  $\gamma_m$  gefunden, für die die Verlustfunktion  $\mathbf{L}(\mathbf{y}_i, \mathbf{F}_{m-1}(\mathbf{x}_i) + \gamma_m \mathbf{h}_m(\mathbf{x}_i))$  mit der geschätzten  $\mathbf{h}_m(\mathbf{x})$  minimal ist. Die Pseudorestwerte  $\mathbf{r}_{im}$  sind der negierte Gradient von der Verlustfunktion  $\mathbf{L}(\mathbf{y}_i, \mathbf{F}_{m-1}(\mathbf{x}_i))$  mit den letzten Schätzwerten. Eine oft benutzte Verlustfunktion ist  $\frac{1}{2}(\mathbf{y} - \mathbf{F}(\mathbf{x}))^2$ , da ihre Pseudorestwerte an den üblichen Restwerte entsprechen[Wika]. Wir geben den Algorithmus von Gradient Boosting an:

1. Konstantenwert für den ersten Schätzer setzen:

$$\mathbf{F}_0(\mathbf{x}) = \text{argmin}_{\gamma} \sum_{i=1}^n \mathbf{L}(\mathbf{y}_i, \gamma)$$

2. Für  $\mathbf{m} = 1$  bis  $\mathbf{M}$ :

- a) Pseudorestwerte ausrechnen:

$$\mathbf{r}_{im} = - \left[ \frac{\partial \mathbf{L}(\mathbf{y}_i, \mathbf{F}(\mathbf{x}_i))}{\partial \mathbf{F}(\mathbf{x}_i)} \right]_{\mathbf{F}(\mathbf{x}) = \mathbf{F}_{m-1}(\mathbf{x})} \text{ für } i = 1, \dots, n$$

- b) Modell  $\mathbf{m}$  zur Schätzung von  $\mathbf{h}_m(\mathbf{x})$  durch  $\mathbf{r}_{im}$  anpassen

- c)  $\gamma_m$  durch das Lösen vom folgenden Optimierungsproblem finden:

$$\gamma_m = \text{argmin}_{\gamma} \sum_{i=1}^n \mathbf{L}(\mathbf{y}_i, \mathbf{F}_{m-1}(\mathbf{x}_i) + \gamma \mathbf{h}_m(\mathbf{x}_i))$$

- d) Model aktualisieren:

$$\mathbf{F}_m(\mathbf{x}) = \mathbf{F}_{m-1}(\mathbf{x}) + \gamma_m \mathbf{h}_m(\mathbf{x})$$

3.  $\mathbf{F}_M(\mathbf{x})$  zurückgeben

### 2.2.2 Kausale Wälder

Die kausale Wälder oder Causal Forests sind eine Methode zum Schätzen von dem heterogenen Behandlungseffekt  $\tau(\mathbf{x})$  von einer Behandlung. Dabei sind sie von den Random

Forests inspiriert und von der Struktur her unterscheiden sich mit ihnen nicht. Das „Kausale“ bei den kausalen Wäldern ist die Tatsache, dass wir zwei Arten von Kovariaten in den Blättern haben - von behandelten Objekten und von Kontrollobjekten. Das verlangt von uns, dass wir eine andere finale Schätzung von den Blattelementen nutzen, den Baum etwa anders aufspalten und noch einige Bedingungen einhalten[WA18].

Random Forest ist eine Technik aus dem maschinellen Lernen, die zur Klassifikation von Objekten oder zur Regression von Funktionen benutzt wird. Der Wald besteht aus  $N$  Bäume. Jeder Baum ist ein einzelner Schätzer und der Schätzwert des Waldes ist der durchschnittliche Schätzwert von allen Bäumen. So vermindert man die Varianz, die sonst eine Schätzung eines einzigen Baums hat. Ein Baum besteht aus Knoten und Blättern. Bei einer Modellvorhersage läuft man für ein  $\mathbf{x}$  den Baum durch und sucht, in welchem Blatt dieses  $\mathbf{x}$  liegen würde. Dann wird der Durchschnitt von allen Elementen in diesem Blatt als Schätzwert für  $\mathbf{y}$  herausgegeben. Man muss aber zuerst durch den gegebenen Datensatz  $(\mathbf{x}_i, \mathbf{y}_i)$  den Baum anpassen. Dazu wird der Raum von  $\mathbf{x}$  abhängig von den Daten in Unterräumen aufgespalten und die Knoten des Baumes merken diese Aufspaltungsregel. Wenn ein Unterraum eine minimale Anzahl  $k$  von Trainingdaten in sich hat oder eine weitere Aufspaltung keine Verbesserung bei der Schätzung verspricht, werden die Elemente in diesem Unterraum in ein Blatt gesteckt. Die Daten für jeden Baum werden durch zufälliges Stichprobenziehen ohne Ersatz von den gemeinsamen Datensatz erstellt [Bre01].

Die Daten, die die Kausalen Wälder anpassen müssen, haben außer den üblichen  $(\mathbf{X}_i, \mathbf{Y}_i)$  auch die Variable der Behandlung  $\mathbf{W}_i$ . Damit die Schätzung von dem heterogenen Behandlungseffekt nicht gestreut ist, muss jedes Blatt eine minimale Anzahl  $k$  von Objekten in sich haben, die sowohl behandelt als auch unbehandelt sind. Dann hat die finale Blattschätzung eine andere Form als bei dem normalen Regressionbaum - wir suchen nach der Differenz zwischen den Mittelwerten von den behandelten und unbehandelten Objekten[WA18]. Wir stellen formal die beiden Varianten dar:

$$\hat{\mu}(\mathbf{x}) = \frac{1}{|\{i : \mathbf{X}_i \in L(\mathbf{x})\}|} \sum_{\{i : \mathbf{X}_i \in L(\mathbf{x})\}} (\mathbf{Y}_i) \quad (2.27)$$

$$\hat{\tau}(\mathbf{x}) = \frac{1}{|\{i : \mathbf{W}_i = 1, \mathbf{X}_i \in L(\mathbf{x})\}|} \sum_{\{i : \mathbf{W}_i = 1, \mathbf{X}_i \in L\}} (\mathbf{Y}_i) - \frac{1}{|\{i : \mathbf{W}_i = 0, \mathbf{X}_i \in L(\mathbf{x})\}|} \sum_{\{i : \mathbf{W}_i = 0, \mathbf{X}_i \in L\}} (\mathbf{Y}_i) \quad (2.28)$$

*Gleichung 2.27* ist die finale Blattschätzung von dem Regressionsbaum, *Gleichung 2.28* - von dem kausalen Baum.  $L(\mathbf{x})$  ist die Menge von allen  $\mathbf{X}_i$  aus den Trainingdaten, die in diesem Blatt liegen. Damit die Blattschätzung nicht gestreut ist, muss jeder Baum in unserem Wald

„ehrlich“ sein. Wagner und Athey[WA18] haben die Bedingung der Ehrlichkeit so erklärt: für jedes Objekt  $i$  aus den Trainingsdaten wird entweder sein  $\mathbf{Y}_i$  für die Blattschätzung, oder sein  $\mathbf{X}_i$  für die Aufspaltung bei dem Erstellen des Baums benutzt, aber nicht beides. Um das zu erzielen, teilt man die Daten für jeden Baum auf 2 Untergruppen  $\mathbf{A}$  und  $\mathbf{B}$ . Dann darf man für die Aufspaltung des Baums alle Daten aus  $\mathbf{A}$  und die  $\mathbf{X}_i$  und  $\mathbf{W}_i$  aus  $\mathbf{B}$  benutzen. In den Blättern werden aber nur die Objekten aus  $\mathbf{B}$  gesteckt, sodass nur seine  $\mathbf{Y}_i$  für die Blattschätzung benutzt werden. Noch eine sehr wichtige Bedingung ist, dass jedes Blatt mindesten  $k$  behandelten und  $k$  unbehandelten Objekten enthalten muss. Je kleiner  $k$  ist, desto weniger gestreut wird die Schätzung sein. Desto länger wird aber die Anpassung des Baums dauern[WA18].

## **3 Synth-Validation**

### **3.1 Generierung von synthetischen Daten**

#### **3.1.1 Auswahl von synthetischen Effekten**

#### **3.1.2 Schätzung von bedingten Erwartungswerten**

### **3.2 Methodenauswahl**

## **4 Implementierung**

### **4.1 Fremde Bibliotheken**

### **4.2 Lesen/Schreiben von Daten**

### **4.3 Ziehen von Stichproben**

### **4.4 Methoden für kausale Inferenz**

### **4.5 Synth-Validation**

#### **4.5.1 Datenstrukturen**

#### **4.5.2 Schätzung**

#### **4.5.3 Constrained Boosting**

#### **4.5.4 Methodenauswahl**

### **4.6 Benchmark von Synth-Validation**

### **4.7 Erstellung von Abbildungen**

### **4.8 Experimenten**

### **4.9 Anderer Code**

## **5 Ergebnisse und Evaluation**

### **5.1 Methodik und Daten**

### **5.2 TODO**

## **6 Schlussfolgerung**

### **6.1 Zusammenfassung**

### **6.2 Diskussion**

# Anhang



# Literatur

- [Bre01] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001. (Seite 20).
- [CM82a] David R Cox and Peter McCullagh. A biometrics invited paper with discussion. some aspects of analysis of covariance. *Biometrics*, pages 541–561, 1982. (Seiten 9, 10).
- [CM82b] David R Cox and Peter McCullagh. A biometrics invited paper with discussion. some aspects of analysis of covariance. *Biometrics*, pages 541–561, 1982. (Seite 11).
- [HIR03] Keisuke Hirano, Guido W Imbens, and Geert Ridder. Efficient estimation of average treatment effects using the estimated propensity score. *Econometrica*, 71(4):1161–1189, 2003. (Seite 14).
- [Hol86] Paul W Holland. Statistics and causal inference. *Journal of the American statistical Association*, 81(396):945–960, 1986. (Seite 1).
- [Lew74] David Lewis. Causation. *The journal of philosophy*, 70(17):556–567, 1974. (Seite 3).
- [Lin] Lindsay Dolan. 10 things to know about covariate adjustment. <https://egap.org/methods-guides/10-things-know-about-covariate-adjustment>. [Online; accessed 18-August-2019]. (Seite 11).
- [Mac] Macartan Humphreys. 10 things to know about causal inference. <https://egap.org/methods-guides/10-things-you-need-know-about-causal-inference>. [Online; accessed 14-August-2019]. (Seite 1).
- [NW17] Xinkun Nie and Stefan Wager. Quasi-oracle estimation of heterogeneous treatment effects. *arXiv preprint arXiv:1712.04912*, 2017. (Seite 15).
- [Reb] Rebecca Barter. Confounding in causal inference: what is it, and what to do about it? <http://www.rebeccabarter.com/blog/2017-07-05-confounding/>. [Online; accessed 16-August-2019]. (Seite 3).
- [RHB00] James M Robins, Miguel Angel Hernan, and Babette Brumback. Marginal structural models and causal inference in epidemiology. *Epidemiology*, 2000. (Seite 14).
- [Rob88] Peter M Robinson. Root-n-consistent semiparametric regression. *Econometrica: Journal of the Econometric Society*, pages 931–954, 1988. (Seite 16).
- [RR83] Paul R Rosenbaum and Donald B Rubin. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55, 1983. (Seiten 11, 12).
- [Rub74] Donald B Rubin. Estimating causal effects of treatments in randomized and nonran-

domized studies. *Journal of educational Psychology*, 66(5):688, 1974. (Seite 4).

- [SJT<sup>+</sup>17] Alejandro Schuler, Ken Jung, Robert Tibshirani, Trevor Hastie, and Nigam Shah. Synth-validation: Selecting the best causal inference method for a given dataset. *arXiv preprint arXiv:1711.00083*, 2017. (Seite 7).
- [Stu10] Elizabeth A Stuart. Matching methods for causal inference: A review and a look forward. *Statistical science: a review journal of the Institute of Mathematical Statistics*, 25(1):1, 2010. (Seiten 13, 14).
- [Tib96] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996. (Seite 18).
- [VS13] Tyler J VanderWeele and Ilya Shpitser. On the definition of a confounder. *Annals of statistics*, 41(1):196, 2013. (Seite 2).
- [WA18] Stefan Wager and Susan Athey. Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113(523):1228–1242, 2018. (Seiten 20, 21).
- [Wika] Wikipedia Autoren. Gradient boosting. [https://en.wikipedia.org/wiki/Gradient\\_boosting](https://en.wikipedia.org/wiki/Gradient_boosting). [Online; accessed 24-August-2019]. (Seiten 18, 19).
- [Wikb] Wikipedia Autoren. Zweistichproben-t-test. <https://de.wikipedia.org/wiki/Zweistichproben-t-Test>. [Online; accessed 16-August-2019]. (Seite 5).
- [Win] Winston Lin. Regression adjustment in randomized experiments: Is the cure really worse than the disease? <https://web.archive.org/web/20151024055802/http://blogs.worldbank.org/impactevaluations/node/847>. [Online; accessed 18-August-2019]. (Seite 10).