

Technische Universität Dortmund  
Fakultät für Mathematik  
Sommersemester 2021

Vorlesungsskript:  
**Konzentrationsungleichungen**

Dozent: Prof. Dr. Ivan Veselić  
erstellt von: Dennis Andreas Malcherczyk

Dieses Skript ist aus einer vierstündigen Vorlesung entstanden, dass ich 2017/2018 an der TU Dortmund hielt. Es orientierte sich im Wesentlichen an dem Buch: *Concentration Inequalities: A Nonasymptotic Theory of Independence* von Stéphane Boucheron, Gábor Lugosi und Pascal Massart.

Herr Malcherczyk war einer der Hörer und hat im Anschluss ein detailliert ausgearbeitetes Skript erstellt. Der vorliegende Text stellt einer Bearbeitung meinerseits dar und verwendet einige von Christoph Schumacher erzeugte Graphiken.

Dortmund, März-Juli 2021

Ivan Veselić

## INHALTSVERZEICHNIS

1. Motivation	5
2. Grundlegende Ungleichungen	8
2.1. Markov-Ungleichung und Co.	8
2.2. Cramér-Chernoff-Methode	10
2.2.1. Die kumulantenerzeugenden Funktion und die Cramér-Transformierte	12
2.2.2. Bestimmung des Supremums mittels der Ableitung	14
2.2.3. Cramér-Transformierte verschiedener Verteilungsklassen	15
2.3. Sub-Gaußsche Zufallsvariablen	19
2.4. Sub-Gamma-Zufallsvariablen	26
2.5. Eine Maximal-Ungleichung	28
2.5.1. Ähnliche Resultate gelten auch für sub- $\Gamma$ -Zufallsvariablen	30
2.6. Hoeffding-Ungleichung	32
2.7. Bennett-Ungleichung	34
2.8. Bernstein-Ungleichung	36
2.9. Johnson-Lindenstrauss-Lemma	39
2.10. Assoziations- und Korrelationsungleichungen	42
2.11. Anwendung der Harris-Ungleichung: Janson-Ungleichung	45
2.12. Anwendung der Harris-Ungleichung: Perkolation	49
Frage: Wir kann ich für einen passenden Wert $p \in [0, 1]$ zeigen, dass f.s. ein unendlicher Cluster existiert?	50
2.13. Negativ assoziierte ZV	52
2.14. Minkowski-Ungleichung	53
3. Schranken an die Varianz	55

3.1.	Efron-Stein-Ungleichung	55
3.2.	Funktionen mit beschränkter Differenz	60
3.3.	Selbstbeschränkende Funktionen	63
	Anwendung bei verschiedenen Modellen	65
3.4.	Exkurs: Ursprung der VC-Theorie:	69
3.5.	Weitere Anwendungen von self-bounding	75
3.6.	Eine konvexe Poincaré-Ungleichung	79
3.7.	Anwendung der Efron-Stein-Ungleichung auf Tail-Events	82
3.8.	Gaußsche Poincaré-Ungleichung	88
3.9.	Beweis für die ES-Ungleichung mittels Dualität	90
4.	Entropie	92
4.1.	Shannon-Entropie und relative Entropie	92
4.2.	Entropie von Produkten und Kettenregel	94
4.3.	Han-Ungleichung	96
4.4.	Isoperimetrische Ungleichung auf $BHC$	96
4.5.	Kombinatorische Entropien	100
4.6.	Han-Ungleichung für relative Entropien	103
4.7.	Sub-Additivität der Entropie	104
4.8.	Entropie für allgemeine Zufallsvariablen	107
4.9.	Dualität und Variationsformel	110
4.10.	Transportkosten-Abschätzung	113
4.11.	Pinsker-Ungleichung	113
4.12.	Birgé-Ungleichung	113
4.13.	Subadditivität der Entropie für allgemeine Zufallsvariablen	113
4.14.	Brunn-Minkowski-Ungleichung	113
5.	Logarithmische Sobolev-Ungleichung (LSU)	113
5.1.	LSU für symmetrische Bernoulli-Verteilungen	114
5.2.	Herbst-Argument	116
5.3.	Gaußsche LSU	118
5.4.	TSI-Konzentrationsungleichungen für Gauß-Zufallsvariablen	118
5.5.	Konzentrationsungleichung für Suprema von Gauß-Prozessen	118
5.6.	Zufällige Gaußsche Projektionen	118
5.7.	Hyperkontraktivität	118
	Literatur	119



## 1. MOTIVATION

Zunächst sollen klassische Situationen als motivierendes Fundament vorgestellt werden, in denen man Konzentrationsungleichungen begegnen kann.

### (A) Gesetz der großen Zahlen:

Für unabhängig, identisch verteilte Zufallsvariablen  $X_1, \dots, X_n$  in  $\mathcal{L}^1(\Omega, P)$ , äquivalent durch  $E|X_1| < \infty$  ausgedrückt, gilt das *Gesetz der großen Zahlen*:

$$(1.1) \quad \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow[n \rightarrow \infty]{} E(X_1).$$

Eine andere nützliche Formulierung ist folgende

$$\frac{1}{n} \sum_{i=1}^n (X_i - E(X_i)) \xrightarrow[n \rightarrow \infty]{} 0.$$

Letztere Schreibweise kann vor allem in Fällen nützlich sein, in denen wir keine identisch verteilten Zufallsvariablen vorliegen haben. Man beachte aber, dass in solchen Fällen nicht allgemein die Konvergenz (1.1) gelten muss.

Bisher haben wir offen gelassen, welche Art von Konvergenz hier vorliegt.

Typischerweise formuliert man das Gesetz der großen Zahlen z.B. in

- fast sicherer Konvergenz (starkes Gesetz der großen Zahlen),
- stochastischer Konvergenz (schwaches Gesetz der großen Zahlen),
- der  $\mathcal{L}^2$ -Norm.

Ein wichtiger Aspekt bei der Anwendung von Grenzwertsätzen in der Stochastik ist die Konvergenzgeschwindigkeit. Typischerweise fragt man sich, wie groß der Approximationsfehler für endliche  $n$  ist. Solche *nicht asymptotische* Fragestellungen sind in der Anwendung wichtig, um die Güte einer Approximation für den Erwartungswert von echten gegebenen Daten  $x_1, \dots, x_n$  abzuschätzen.

Dabei werden Abschätzungen der folgenden Bauart angestrebt:

$$\left\| \frac{1}{n} \sum_{i=1}^n (X_i - E(X_i)) \right\| \leq f(n, P_{X_1}),$$

wobei  $f(n, P_{X_1})$  eine obere Schranke ist, die von dem Stichprobenumfang  $n$  und von der Verteilung  $P_{X_1}$  selbst abhängt. Wünschenswert wäre es, wenn  $f(n, P_{X_1})$  wenige Informationen über die Verteilung benötigte, um möglichst allgemeine Aussagen zu erhalten. In der  $\mathcal{L}^2$ -Norm für unabhängige, aber nicht notwendigerweise identisch verteilte, quadrat-integrierbare Zufallsvariablen

können wir folgende Abschätzung angeben

$$\left\| \frac{1}{n} \sum_{i=1}^n (X_i - E(X_i)) \right\|_{\mathcal{L}^2(P)} \leq \frac{\max_{i=1, \dots, n} \{\sigma_i\}}{\sqrt{n}}.$$

Dabei sind  $\sigma_i$  die Standardabweichungen der ZVen  $X_i$  für  $i = 1, \dots, n$ .

Die hier vorliegende Abschätzung bildet im Fall von identisch verteilten Zufallsvariablen sogar Gleichheit.

In Abschnitt 2 wird u.a. die Frage untersucht, welche Schranken sich durch Verwendung von höheren Momenten finden lassen.

### (B) Rekonstruktion von Verteilungen in statistischer Lerntheorie

Wir betrachten den Datensatz  $x_1, \dots, x_n \in \mathbb{R}$  als Realisationen von unabhängigen, identisch verteilten ZVen  $X_1, \dots, X_n$  mit unbekannter Verteilung  $P_X = \mu$ .

Frage: Wie lässt sich aus geg. Daten die wahre Verteilung rekonstruieren?

Wir verwenden das sogenannte *empirische Maß*

$$\mu_n = \mu_n^{x_1, \dots, x_n} = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$$

mit  $\delta_{x_i}$  als das Punktmaß in  $x_i$ . Das liefert eine intuitive Möglichkeit für eine Schätzung von  $\mu$ . Die empirische Verteilungsfunktion aus den Datensatz  $x_1, \dots, x_n$  ist dabei zum empirischen Maß assoziiert.

Auch hier stellt man sich die Frage nach der Güte der Approximation von  $\mu_n$  zu  $\mu$ . In welchem Sinne können wir hier überhaupt eine Konvergenz formulieren? Ein elementarer Ansatz ist durch den *Fundamentalsatz der Statistik von Glivenko-Cantelli* gegeben, durch den die Konvergenz der empirischen Verteilungsfunktion gegen die wahren Verteilungsfunktion in der Supremumsnorm geliefert wird.

Auf Ebene der Maßtheorie liegt eine schwache Konvergenz des empirischen Maßes gegen das wahre Wahrscheinlichkeitsmaß für Mengen der Bauart

$$A = (-\infty, x] \text{ für } x \in \mathbb{R}$$

vor. (Hier evtl. noch Arten der schwachen Konvergenz in der Sprache Funktionalanalysis diskutieren.)

Lässt sich die schwache Konvergenz für andere Klassen von Mengen formulieren? Diese Frage wird in einem Exkurs zur statistischen Lerntheorie in Abschnitt 3 untersucht.

### (C) Irrfahrten auf $\mathbb{Z}^2$

Wir betrachten unabhängig, identisch verteilte ZVen der Form

$$X_1, \dots, X_n : \Omega \rightarrow \left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ -1 \end{pmatrix} \right\}.$$

Die Werte der ZVen beschreiben dabei Bewegungsrichtungen im  $\mathbb{Z}^2$ -Gitter. Die Wahrscheinlichkeit für alle Richtungen soll  $\frac{1}{4}$  betragen. Nun summieren wir die ersten  $n$  Bewegungsschritte auf und erhalten eine neue ZVe

$$Z_n := \sum_{i=1}^n X_n.$$

Sie beschreibt für verschiedene Zeiten  $i = 1, \dots, n$  einen Pfad auf  $\mathbb{Z}^2$ . (Hier evtl. noch eine Abbildung einfügen.) Man nennt solche Prozesse *Irrfahrten*.

Wir fragen uns in diesem Kontext beispielsweise, wie sich der euklidische Abstand  $\|Z_n\|$  im Verlauf eines Pfades typischerweise verhält. Es kann z.B. nach einer oberen Schranke für  $\max_{i=1, \dots, n} \|Z_i\|$  gefragt werden. Eine sehr einfache Antwort wäre

$$\max_{i=1, \dots, n} \|Z_i\| \leq n,$$

da jeder Schritt maximal Länge 1 hat. Das ist allerdings eine grobe Abschätzung.

Mit dem *Zentralen Grenzwertsatz* können wir oft Aussagen der folgenden Bauart

$$\max_{i=1, \dots, n} \|Z_i\| \leq C \cdot \sqrt{n}$$

gewinnen. Die Frage, wann solche Abschätzungen gelten und wovon die Konstante  $C$  abhängt, bleibt noch offen. Unklar bleibt auch, wie die Verteilungen  $P_{X_i}$  in die Abschätzung eingehen.

### Optional (D) Brownsche Bewegung

Hier interessieren wir uns für Suprema von stochastischen Prozessen. Ein prominentes Beispiel dafür ist die Brownsche Bewegung mit Zeithorizont  $T$

$$B : \Omega \times [0, T] \rightarrow \mathbb{R}.$$

Mit einer funktionalen Version des zentralen Grenzwertsatzes gilt:

$$\text{Irrfahrt} \xrightarrow{\text{Donsker}} \text{Brownsche Bewegung}$$

$t \mapsto B_\omega(t)$  ist f.s. nicht differenzierbare Trajektorie.

Ferner: Für jede Zeit  $t > 0$  gilt  $\sup_{\omega \in \Omega} \|B_t(\omega)\| = \infty$ . Aber es gilt wiederum:

$$\sup_{t \in [0, T]} \|B_t\| \leq C\sqrt{T}, \quad C > 0.$$

Ziel: Präzisiere die Konstante  $C$  und die Wahrscheinlichkeit!

## 2. GRUNDLEGENDE UNGLEICHUNGEN

In diesem Kapitel wollen wir erste Beispiele von Konzentrationsungleichungen kennenlernen.

### 2.1. Markov-Ungleichung und Co.

Sei  $X$   $\mathbb{R}$ -wertige ZVe auf Wahrscheinlichkeitsraum  $(\Omega, P)$  mit endlichem Erwartungswert  $E(X)$ .

**Frage:** Um wie viel weicht  $X$  von seinem Erwartungswert  $E(X)$  ab?

Wir wollen obere Schranken für  $t > 0$  der folgenden Form finden:

$$P(X - E(X) \geq t) \leq \dots$$

$$P(X - E(X) \leq -t) \leq \dots$$

einfachste Antwort: *Markov-Ungleichung*

Sei  $Y \geq 0$  eine ZVe mit  $E(Y) < \infty$ . Dann gilt für alle  $t \geq 0$ :

$$Y \cdot \mathbf{1}_{\{Y \geq t\}} \geq t \cdot \mathbf{1}_{\{Y \geq t\}} \text{ auf } \Omega$$

Integration liefert:

$$E(Y \cdot \mathbf{1}_{\{Y \geq t\}}) \geq t \cdot E(\mathbf{1}_{\{Y \geq t\}}) = t \cdot P(Y \geq t).$$

Weitere Abschätzung der linken Seite:

$$E(Y \cdot \mathbf{1}_{\{Y \geq t\}}) \leq E(Y).$$

Falls die Dichtefunktion  $f$  von  $Y$  gegeben ist, ist es natürlich, den Wertebereich von  $Y$  auf die horizontale Achse aufzutragen.

Zusammenfassend: Für jede ZVe  $Y: \Omega \rightarrow [0, \infty)$  in  $\mathcal{L}^1(\Omega, P)$  gilt nach obigen Ausführungen für  $t > 0$  die Markov-Ungleichung:

$$(2.1) \quad P(Y \geq t) \leq \frac{1}{t} E(Y \cdot \mathbf{1}_{\{Y \geq t\}}) \leq \frac{1}{t} E(Y).$$

Setze  $Y = |E(X) - X|$ . Das liefert eine erste Antwort:

$$P(X - E(X) \geq t) + P(X - E(X) \leq -t) = P(Y \geq t) \leq \frac{E(Y)}{t} = \frac{E(|X - E(X)|)}{t}.$$

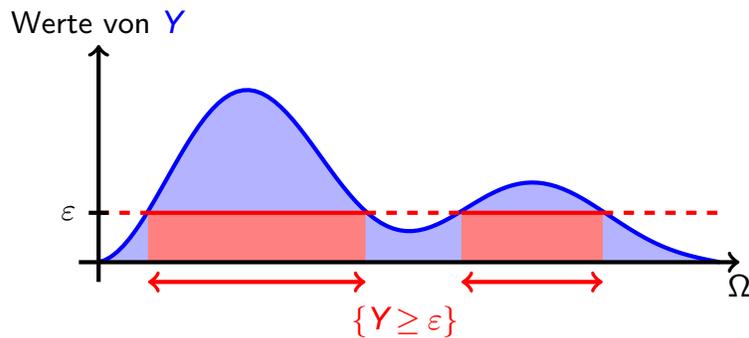


ABBILDUNG 1.  $\Omega$  auf der horizontalen Achse.

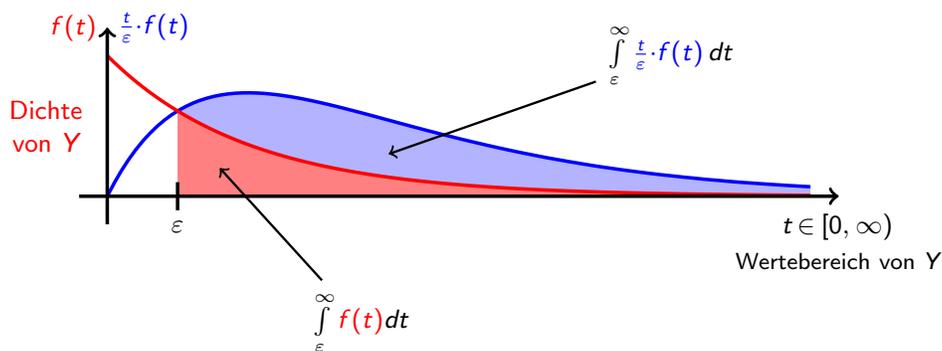


ABBILDUNG 2. Wertebereich auf der horizontalen Achse.

Frage: Gibt es eine bessere Wahl von  $Y$  in (2.1)?

Hat  $X$  z.B. eine endliche Varianz  $\text{Var}(X)$ , so gilt für  $\Phi(Y)$  mit  $\Phi(y) = y^2$

$$E(\Phi(Y)) = E(Y^2) = E(|X - E(X)|^2) = \text{Var}(X) < \infty \Rightarrow \Phi(Y) \in \mathcal{L}^1(\Omega, P).$$

Wende nun für  $t \geq 0$  die Markov-Ungleichung an:

$$(2.2) \quad P(|X - E(X)| \geq t) = P(\Phi(Y) \geq \Phi(t)) \leq \frac{E(\Phi(Y))}{\Phi(t)} = \frac{\text{Var}(X)}{t^2}$$

und erhalte damit die *Markov-Čebyšev-Ungleichung*. Die Ungleichung (2.2) gilt für sämtliche isotone (monoton wachsende) Funktionen  $\Phi: I \rightarrow [0, \infty)$  auf einem Intervall  $I \subset \mathbb{R}$ , sodass  $\Phi(t) > 0$  ist.<sup>1</sup> Für die Anwendbarkeit der

<sup>1</sup>In dieser Vorlesung benutzen wir der Kürze halber den Begriff *isoton* für monoton wachsend und *antiton* für monoton fallend.

Methode braucht man

$$\Phi(Y) = \Phi(|X - E(X)|) \in \mathcal{L}^1(\Omega, P).$$

Gilt für eine ZVe  $X$ :  $E|X^q| < \infty \forall q \in \mathbb{N}$ , erhalten wir  $\forall q, t \geq 0$ :

$$P(|X - E(X)| \geq t) \leq \frac{E(|X - E(X)|^q)}{t^q}.$$

→ schöne Form, da linke Seite unabhängig von  $q$  ist (Optimierungsaspekt)

$$(2.3) \quad \Rightarrow \quad P(|X - E(X)| \geq t) \leq \inf_{q>0} \frac{E(|X - E(X)|^q)}{t^q}.$$

Erstes Fazit: Je mehr Informationen über eine ZVe vorliegen, desto breiter das Spektrum an potentiellen Abschätzungen und Methoden.

### Warum spielt der Fall $q = 2$ eine zentrale Rolle?

Seien  $X_1, \dots, X_n$  unabhängige ZVen mit endlichen Varianzen. Für  $Z = \sum_{i=1}^n X_i$  gilt nach dem Additionssatz (Bienaymé):

$$\text{Var}(Z) = \sum_{i=1}^n \text{Var}(X_i).$$

→ Formulierung einer Konzentrationsungleichung für gemittelte ZVen mit Hilfe von Ungleichung (2.2):

$$\begin{aligned} P\left(\left|\frac{1}{n} \sum_{i=1}^n (X_i - E(X_i))\right| > t\right) &= P\left(\left|\sum_{i=1}^n (X_i - E(X_i))\right| > t \cdot n\right) \\ &\leq \frac{\text{Var}(Z)}{t^2 \cdot n^2} = \frac{\sigma^2}{t^2 \cdot n}, \end{aligned}$$

wobei  $\sigma^2 := \frac{1}{n} \sum_{i=1}^n \text{Var}(X_i)$  die gemittelte Varianz ist.

Bemerkung: In dieser Vorlesung spielt die Eigenschaft der *Identischen Verteilung* eine weniger zentrale Rolle als die Unabhängigkeit, da uns nicht der explizite Wert des Limes interessiert, sondern gute Abschätzungen für *endliches*  $n$  (bzw. die Konvergenzordnung). Natürlich vereinfachen sich viele Aussagen, falls die ZVen identisch verteilt sind.

### 2.2. Cramér-Chernoff-Methode.

Die Methoden dieses Kapitels werden in den nachfolgenden Kapiteln 2.3 - 2.9 benötigt. Zusammenhänge zur Entropie werden in Kapitel 4.9 dargestellt. In Kapitel 5.2, 5.4, 5.5 taucht sie ebenso auf.

Idee: Wähle statt  $\Phi(t) = t^2$  nun  $\Phi(t) = e^{\lambda t}$  für  $\lambda > 0$ .

Analog zu (2.2) mit  $\Phi(y) = e^{\lambda y}$  nach Anwendung der Markov-Ungleichung:

$$(2.4) \quad P(X \geq t) = P\left(e^{\lambda X} \geq e^{\lambda t}\right) \leq \frac{E\left(e^{\lambda X}\right)}{e^{\lambda t}}.$$

→ Schranke mit exponentiellem Abfall gewonnen.

Studiere nun die Schranken genauer. Die folgende Abbildung

$$M: \mathbb{R} \rightarrow [0, \infty], M(\lambda) := E\left(e^{\lambda X}\right)$$

nennen wir die *momentenerzeugenden Funktion* von  $X$  (kurz: MEF von  $X$ ).

Sie ist gegebenenfalls unendlich. Betrachte  $Z = \sum_{i=1}^n X_i$  mit unabhängigen ZVen  $X_1, \dots, X_n$ . Für die MEF von  $Z - E(Z)$  gilt dann

$$E\left(e^{\lambda \sum_{i=1}^n (X_i - E(X_i))}\right) = \prod_{i=1}^n E\left(e^{\lambda (X_i - E(X_i))}\right),$$

wegen der Unabhängigkeit der  $X_1, \dots, X_n$ . Sind  $(X_i - E(X_i))$  zusätzlich identisch verteilt mit MEF  $M(\lambda) := M_{X_1 - E(X_1)}(\lambda)$ , so gilt mit (2.4):

$$P\left(\frac{1}{n}(Z - E(Z)) \geq t\right) \leq \frac{\prod_{i=1}^n M(\lambda)}{e^{\lambda t n}} = \frac{(M(\lambda))^n}{e^{\lambda t n}}.$$

Vorgehensweise bisher:

- Gewinne Klassen von Ungleichungen für verschiedene  $\lambda$
- Schranke über den Parameter  $\lambda$  optimieren (Minimierungsaufgabe)
- Lösen der Optimierungsaufgabe liefert eine (hoffentlich) gute Abschätzung

Bemerkungen:

- (a) Im Allg. liefern polynomielle Transformationen  $\Phi(t) = t^q$  aus (2.3) bessere Schranken als exponentielle Transformationen  $\Phi(t) = e^{\lambda t}$  aus (2.4). D.h.: für beliebige  $t > 0$  und ZVe  $X \geq 0$ :

$$\inf_{q>0} \frac{E(X^q)}{t^q} \leq \inf_{\lambda>0} \frac{E(e^{\lambda X})}{e^{\lambda t}}.$$

Beweisidee: Taylorentwicklung der Exponentialfunktion (Übung).

- (b) Wir interessieren uns für die Wahrscheinlichkeit der Abweichungen vom Mittelwert mit folgender Bauart:

$$P(|Z - E(Z)| \geq t) = P(Z - E(Z) \geq t) + P(Z - E(Z) \leq -t) \quad \text{für } t > 0.$$

Wegen (b) genügt es o.B.d.A. die zentrierte Version der ZVe  $\tilde{Z} := (Z - E(Z))$  zu betrachten und Abschätzungen für  $P(|\tilde{Z}| \geq t)$  herzuleiten.

2.2.1. *Die kumulantenerzeugenden Funktion und die Cramér-Transformierte.*

Sei  $Z$  ZVe mit MEF  $M_Z(\lambda)$ , d.h. (äquivalent zu (2.4)):

$$(2.5) \quad P(Z \geq t) \leq e^{-\lambda t} M(\lambda) = e^{-\lambda t + \ln(M(\lambda))}.$$

Fasse die obere Schranke als Klasse von Funktionen auf, um sie dann zu minimieren. Dazu definieren wir die *kumulantenerzeugende Funktion*  $\psi_Z(\lambda)$  (kurz: KEF):

$$\psi_Z(\lambda) := \ln(M(\lambda)) = \ln(E(e^{\lambda Z})).$$

Betrachte dazu die sogenannte *Cramér-Transformierte*:

$$\psi_Z^*(t) := \sup_{\lambda \geq 0} (\lambda t - \psi_Z(\lambda)) \quad \text{für } \lambda \geq 0.$$

Die Transformierte  $\psi^*$  einer konvexen Funktion  $\psi$  wird abstrakt in Lemma 2.14 untersucht. Bedeutung der Cramér-Transformierte:

Minimiere die Schranke von (2.5) über  $\lambda$  und betrachte nur noch den Exponenten. Beachte: Vorzeichenwechsel liefert Maximierungsproblem:

$$P(Z \geq t) \leq \inf_{\lambda \geq 0} e^{-(\lambda t - \ln(M(\lambda)))} = e^{-\psi_Z^*(t)}.$$

Untersuche nun den Definitions- und Wertebereich der Cramér-Trafo  $\psi_Z^*$ . Zunächst bemerkt man eine Eigenschaft der KEF für beliebige ZVen  $Z$ :

$$\psi_Z(0) = \ln(1) = 0.$$

Dieser Zusammenhang ist nützlich für Randwertuntersuchungen. Auf die Cramér-Trafo überträgt sich dies, wenn man für  $\lambda = 0$  setzt, wie folgt:

$$\psi_Z^*(t) = \sup_{\lambda \geq 0} (\lambda t - \psi_Z(\lambda)) \geq 0 - 0 = 0.$$

Insbesondere wissen wir nun, dass der Wertebereich von  $\psi_Z^*$  nichtnegativ ist.

Warum betrachten wir nicht  $\lambda \in \mathbb{R}$ , sondern nur  $\lambda \geq 0$  im Supremum?

Falls für die ZV  $Z \in \mathcal{L}^1$  gilt, so gilt nach der Jensen-Ungleichung

$$e^{\lambda E(Z)} \leq E(e^{\lambda Z}) = M(\lambda),$$

wobei der letzte Ausdruck auch unendlich sein könnte. Logarithmieren dieser Ungleichung ergibt

$$\lambda \cdot E(Z) \leq \ln(M(\lambda)) = \psi_Z(\lambda).$$

Betrachte  $\lambda < 0$  und  $t \geq E(Z)$  und schätze die linke Seite ab

$$\lambda \cdot E(Z) \geq \lambda \cdot t$$

Insgesamt folgt aus beiden Ungleichungen für  $\lambda < 0$  und  $t \geq E(Z)$ :

$$\lambda \cdot t - \psi_Z(\lambda) \leq 0$$

Die Annahme  $t \geq E(Z)$  lässt sich allgemein dadurch motivieren, dass wir zentrierte ZVe  $Z$  mit  $E(Z) = 0$  betrachten wollen.

Fazit: Obige Randbetrachtung für  $\lambda = 0$  liefert, dass das Supremum über  $\lambda$  nicht im negativen Bereich angenommen wird, d.h.

$$(2.6) \quad \tilde{\psi}_Z(t) := \sup_{\lambda \in \mathbb{R}} (\lambda t - \psi_Z(\lambda)) = \sup_{\lambda \geq 0} (\lambda t - \psi_Z(\lambda)) = \psi_Z^*(t) \quad \text{für } t \geq E(Z).$$

Wir nennen die Funktion  $\tilde{\psi}_Z$  in (2.6) die *Fenchel-Legendre-Transformierte* oder auch *Fenchel-Legendre-Duale* von  $\psi_Z$ .

Nicht jedes  $t \geq 0$  liefert brauchbare Chernoff-Schranken. Falls  $\psi_Z^*(t) = 0$  ist, ergibt sich eine triviale Schranke  $e^{-\psi_Z^*(t)} = 1$ . In welchen Fällen tritt dies noch ein?

Einerseits ist der Fall  $\psi_Z(\lambda) \equiv \infty$  für  $\lambda > 0$  problematisch. Dann folgt

$$\psi_Z^*(t) = \sup_{\lambda \geq 0} (\lambda t - \psi_Z(\lambda)) = 0.$$

Andererseits ist der Fall  $t \leq E(Z)$  problematisch wegen

$$\begin{aligned} \lambda t &\leq \lambda E(Z) \leq \psi_Z(\lambda) \quad \text{für } \lambda \geq 0 \\ \Rightarrow \lambda t - \psi_Z(\lambda) &\leq 0 \quad \text{und} \quad = 0 \quad \text{für } \lambda = 0. \end{aligned}$$

Um diese Situation zu vermeiden, nehmen wir in diesem Kapitel an, dass ein  $\lambda_0 > 0$  existiert, sodass  $E(e^{\lambda_0 Z}) < \infty$  gilt. Mit der Hölder-Ungleichung lässt sich zeigen, dass dann auch das exponentielle Moment für  $\lambda \leq \lambda_0$  existiert (Übung). Es gilt dann also:

$$\text{für alle } \lambda \in [0, \lambda_0] : E(e^{\lambda Z}) < \infty.$$

Setze dazu die Zahl  $b := \sup\{\lambda \geq 0 \mid E(e^{\lambda Z}) < \infty\} \in [0, \infty]$ . Da die Voraussetzung  $E(e^{\lambda Z}) < \infty$  für  $\lambda = 0$  immer erfüllt ist, genügt es auch hier statt  $\lambda \in \mathbb{R}$  nur den Bereich  $\lambda \geq 0$  zu untersuchen. In den meisten Fällen ist  $b \in \{0, \infty\}$ . Dagegen ist für exponentialverteilte ZVen der Wert  $b$  gerade der Parameter der Exponentialverteilung.

*Bemerkung 2.1* (Eigenschaften der kumulantenerzeugenden Funktion). Folgende Eigenschaften werden sich für Optimierungsaufgaben als nützlich erweisen:

- (a)  $\psi = \psi_Z$  ist konvex auf  $I := (0, b)$  (auch gültig, falls  $b = \infty$ )
- (b)  $\psi$  ist strikt konvex auf  $I$ , falls  $Z$  nicht fast sicher konstant
- (c)  $\psi: I \rightarrow \mathbb{R}$  ist  $\mathcal{C}^\infty$

Für zentrierte Zufallsvariablen  $Z$  gilt darüber hinaus:

- (d)  $\psi_Z: [0, b) \rightarrow \mathbb{R}$  ist  $\mathcal{C}^1$ . Beachte: 0 ist *zusätzlich* im Definitionsbereich.
- (e)  $\psi'_Z(0) = 0$  (zusätzlich zu  $\psi_Z(0) = 0$ )
- (f) Es genügt die Cramér-Trafo auf dem Intervall  $I$  zu bestimmen:

$$\psi_Z^*(t) = \sup_{\lambda \geq 0} (\lambda t - \psi_Z(\lambda)) = \sup_{\lambda \in I} (\lambda t - \psi_Z(\lambda))$$

Die Beweise von (a) - (f) werden als Übungsaufgaben gestellt.

### 2.2.2. Bestimmung des Supremums mittels der Ableitung. Ansatz:

- $\psi_Z$  ist  $\mathcal{C}^1 \rightarrow$  mittels Ableitung stationäre Punkte berechnen
- strikte Konvexität liefert Eindeutigkeit des Optimums auf  $I$

$$\begin{aligned} 0 &= \frac{d}{d\lambda} (\lambda t - \psi(\lambda)) = t - \psi'(\lambda) \\ \Leftrightarrow t &= \psi'(\lambda) \end{aligned}$$

Sei  $\lambda_t$  eine Lösung dieser Gleichung. Falls man den trivialen Fall einer fast sicher konstanten Zufallsvariable ausschließt, ist  $\psi$  strikt konvex.

$\Rightarrow \lambda_t$  ist eindeutig.

Definition: Sei  $B := \psi'_Z(b)$ . Dann ist  $\psi'_Z: I \rightarrow (0, B)$  wegen strikter Monotonie bijektiv mit strikt monotoner Inversen  $(\psi'_Z)^{-1}$ .

Daher gilt für alle  $t \in (0, B)$ :  $\lambda_t = (\psi'_Z)^{-1}(t)$ .

Diese Formel können wir nun nutzen, um für konkrete Verteilungen die Cramér-Trafo auszurechnen.

### 2.2.3. Cramér-Transformierte verschiedener Verteilungsklassen.

#### (a) Cramér-Transformierte für zentrierte Normalverteilungen:

Die Kenntnis der Cramér-Transformierten der zentrierten Normalverteilung ist beim Verständnis von sub-gaußschen Zufallsvariablen in Kapitel 2.2 relevant. Sei  $Z \sim \mathcal{N}(0, \sigma^2)$  mit Varianz  $\sigma^2$ . Beachte, dass wir Zentriertheit brauchen.

- zunächst berechne MEF:  $M(\lambda) = e^{\frac{\lambda^2 \sigma^2}{2}}$  (Übung)
- MEF ergibt sofort die KEF  $\psi_Z(\lambda) = \frac{\lambda^2 \sigma^2}{2}$ .
- erste Ableitung ist  $\psi'_Z(\lambda) = \lambda \sigma^2$
- obige Formel ergibt  $\lambda_t = (\psi')^{-1}(t) = \frac{t}{\sigma^2}$  als Lösung der Optimierungsaufgabe

$\forall t > 0$  besitzt die Cramér-Trafo also die folgende Gestalt

$$\begin{aligned}\psi_Z^*(t) &= \lambda_t t - \psi_Z(\lambda_t) \\ &= \frac{t^2}{\sigma^2} - \frac{\sigma^2 t^2}{2\sigma^4} \\ &= \frac{t^2}{2\sigma^2}.\end{aligned}$$

Das liefert die Chernoff-Schranke für  $\forall t > 0 : P(Z \geq t) \leq e^{-\frac{t^2}{2\sigma^2}}$ .

Wie gut ist diese Schranke? Kann man sie noch verbessern? Zur Beantwortung dieser Fragen formuliere die Chernoff-Abschätzung um

$$P(Z \geq t) \cdot e^{\frac{t^2}{2\sigma^2}} \leq 1.$$

Diese Abschätzung lässt sich aber verbessern.

$$\text{Für alle } t > 0 \text{ gilt } P(Z \geq t) \cdot e^{\frac{t^2}{2\sigma^2}} \leq \frac{1}{2} \text{ (Übung).}$$

→ globale Vorfaktor  $\frac{1}{2}$  wird verschenkt, Größenordnung passt zumindest.

Zudem kann man zeigen, dass letztere Abschätzung scharf ist:

$$\sup_{t>0} \left( P(Z \geq t) \cdot e^{\frac{t^2}{2\sigma^2}} \right) = \frac{1}{2} \text{ (ebenfalls Übung).}$$

Unter Normalität sind Gewinnungen von Abschätzungen mittels anderer Techniken noch einfach, wodurch wir einen Vergleich zwischen ihnen und der Chernoff-Methode erhalten. In anderen Fällen ist dies nicht mehr so einfach möglich.

#### (b) Cramér-Transformierte für zentrierte Poisson-Verteilungen:

Die Kenntnis der KEF der zentrierten Poisson-Verteilung wird im Beweis

der Bennett-Ungleichung 2.20 in Kapitel 2.7 benutzt. Sei  $Y \sim Poi(\nu)$  für ein  $\nu > 0$ . Nach Definition der Poisson-Verteilung gilt

$$\forall k \in \mathbb{N}_0 : P(Y = k) = \frac{\nu^k}{k!} \cdot e^{-\nu}.$$

$\Rightarrow E(Y) = \nu$ . Wir arbeiten mit zentrierten ZVen und definieren daher  $Z := Y - E(Y)$  mit  $E(Z) = 0$ . Berechne im ersten Schritt die MEF:

$$\begin{aligned} M_Z(\lambda) &= E(e^{\lambda Z}) = e^{-\lambda\nu} \sum_{k \in \mathbb{N}_0} \left( e^{\lambda k} \cdot \frac{\nu^k}{k!} \right) \cdot e^{-\nu} \\ &= e^{-\lambda\nu - \nu} \sum_{k \in \mathbb{N}_0} \frac{(\nu e^\lambda)^k}{k!} = e^{-(\lambda+1)\nu} \cdot e^{\nu e^\lambda}. \end{aligned}$$

Logarithmieren liefert für  $\lambda > 0$  die KEF  $\psi_Z$ . Berechne außerdem  $\psi'_Z$

$$\psi_Z(\lambda) = \nu \left( e^\lambda - \lambda - 1 \right), \quad \psi'_Z(\lambda) = \nu \left( e^\lambda - 1 \right).$$

Ansatz:  $t = \psi'_Z(\lambda)$ , um  $\lambda_t$  als Lösung des Optimierungsproblems herzuleiten:

$$\begin{aligned} t &= \psi'_Z(\lambda) = \nu \left( e^{\lambda_t} - 1 \right) \\ \Leftrightarrow e^{\lambda_t} &= \frac{t}{\nu} + 1 \\ \Leftrightarrow \lambda_t &= \ln \left( \frac{t}{\nu} + 1 \right). \end{aligned}$$

Der optimierende Parameter  $\lambda_t$  liefert nun die Gestalt der Cramér-Trafo:

$$\begin{aligned} \psi_Z^*(t) &= t\lambda_t - \psi_Z(\lambda_t) \\ &= t \ln \left( \frac{t}{\nu} + 1 \right) - \nu \left( \frac{t}{\nu} + 1 - \ln \left( \frac{t}{\nu} + 1 \right) - 1 \right) \\ &= (t + \nu) \cdot \ln \left( \frac{t}{\nu} + 1 \right) - t \\ &= \nu \cdot h \left( \frac{t}{\nu} \right), \end{aligned}$$

wobei wir für  $x \geq -1$  die Funktion  $h$  einführen:

$$h(x) := (1 + x) \cdot \ln(1 + x) - x.$$

Die ZVe  $-Z$  ist ebenfalls zentriert und analog folgt:

$$\psi_{-Z}^*(t) = \nu \cdot h \left( -\frac{t}{\nu} \right), \text{ sofern } t \leq \nu \text{ gilt.}$$

(c) Cramér-Transformierte für zentrierte Bernoulli-Verteilungen:

Eigentliches Ziel: Cramér-Trafo von binomialverteilten ZVen.

Betrachte zunächst in (c) Bernoulli-ZVen und in (d) Summen von unabhängigen ZVen, um in Teil (e) die Binomialverteilung zu untersuchen.

Sei  $Y : \Omega \rightarrow \{0, 1\}$  eine ZVe mit  $P(Y = 1) = p = 1 - P(Y = 0)$  für ein  $p \in [0, 1]$ . Der Erwartungswert ist dann  $E(Y) = p$ . Die ZVe  $Z := Y - p$  ist dann die zentrierte Version. Die MEF von  $Z$  ergibt sich aus der Definition

$$M_Z(\lambda) = (p \cdot e^\lambda + (1 - p)) e^{-\lambda p},$$

mit KEF

$$\psi_Z(\lambda) = -\lambda p + \ln(p e^\lambda + 1 - p).$$

Schließlich folgt mit gleicher Strategie (Übung) für  $t \in (0, 1 - p)$ :

$$\begin{aligned} \psi_Z^*(t) &= (1 - p - t) \cdot \ln\left(\frac{1 - p - t}{1 - p}\right) + (p + t) \cdot \ln\left(\frac{p + t}{p}\right) \\ &= (1 - a) \cdot \ln\left(\frac{1 - a}{1 - p}\right) + a \cdot \ln\left(\frac{a}{p}\right) \\ &=: h_p(a). \end{aligned}$$

Hierbei haben wir  $a := t + p$  gesetzt, so dass  $a \in (p, 1)$ .

**Definition 2.2.** Die Funktion  $h_p : (0, 1) \rightarrow \mathbb{R}$ ,  $h_p(a) = (1 - a) \cdot \ln\left(\frac{1 - a}{1 - p}\right) + a \cdot \ln\left(\frac{a}{p}\right)$  nennt man die *Kulback-Leibler-Divergenz*  $D(P_a \| P_p)$  zwischen zwei Bernoulli-Verteilungen mit Parameter  $a$  bzw.  $p$ .

Wir werden im Kapitel 4 den Begriff der *Entropie* einführen und auch die Funktion  $h_p$  diesem Begriff zuordnen können (vgl. insbesondere die *Dualitätsformel 4.14 aus Kapitel 4.9.*). Die Schranken einiger Konzentrationsungleichungen entsprechen solchen Entropien.

(d) Cramér-Transformierte für Summen unabhängiger Zufallsvariablen:

Die Cramér-Chernoff-Methode erlaubt einfachen Umgang mit Summen von unabhängig, identisch verteilten ZVen  $X_1, \dots, X_n$ . Dazu betrachten wir die ZVe  $Z := \sum_{i=1}^n X_i$  und die KEF  $\psi_{X_1}$ , sowie die Cramér-Trafo  $\psi_{X_1}^*$  für  $X_1$ .

Ziel: Darstellung der KEF von  $Z$  in Abhängigkeit der KEF der  $X_1, \dots, X_n$ .

Für  $\lambda$  mit  $\psi_{X_1}(\lambda) < \infty$  bestimmen wir die KEF von  $Z$ :

$$\begin{aligned}
 \psi_Z(\lambda) &= \ln \left( E \left( e^{\lambda Z} \right) \right) = \ln \left( E \left( e^{\lambda \sum_{i=1}^n X_i} \right) \right) \\
 &= \ln \left( \prod_{i=1}^n \left( E \left( e^{\lambda X_i} \right) \right) \right) \quad (\text{wegen Unabhängigkeit}) \\
 (2.7) \quad &= \ln \left( E \left( e^{\lambda X_1} \right)^n \right) \quad (\text{wegen identischer Verteilung}) \\
 &= n \cdot \ln \left( E \left( e^{\lambda X_1} \right) \right) \\
 &= n \cdot \psi_{X_1}(\lambda).
 \end{aligned}$$

Falls die  $X_i$  nur unabhängig sind, erhalten wir mit denselben Schritten zumindest

$$(2.8) \quad \psi_Z(\lambda) = \sum_{i=1}^n \psi_{X_i}(\lambda).$$

*Bemerkung 2.3.* Wir nehmen in diesem Abschnitt durchgängig an, dass die betrachtete ZV  $X$  integrierbar ist und ein  $\lambda > 0$  existiert mit  $M(\lambda) < \infty$ , was sich dann auch auf die KEF überträgt. Insbesondere ist  $b > 0$  in der Definition des Intervalls  $I = (0, b)$ . Im Spezialfall einer f.s. konstanten ZV gilt

$$\psi^*(t) = \begin{cases} 0, & \text{falls } t = EX \\ +\infty, & \text{falls } t > EX, \end{cases}$$

anderenfalls

$$\psi^*(t) = \begin{cases} 0, & \text{falls } t = EX \\ \in (0, \infty), & \text{falls } t > EX. \end{cases}$$

Auch die Cramér-Trafo hat ein einfaches Verhalten bei einer i.i.d.-Summe:

**Lemma 2.4.** *Seien  $X, X_1, \dots, X_n$  unabhängige, identisch verteilten ZVen und  $Z := \sum_{i=1}^n X_i$ . Dann gilt:*

$$\psi_Z^*(t) = n \cdot \psi_{X_1}^* \left( \frac{t}{n} \right).$$

*Beweis.* Es sei an die Gestalt des optimierenden Parameters  $\lambda_t = (\psi')^{-1}(t)$  der Cramér-Trafo erinnert. Also untersuchen wir:

$$\begin{aligned}
 \psi'_Z(\lambda) &\stackrel{(2.7)}{=} n \cdot \psi'_X(\lambda) \\
 &= \mathcal{M}_n(\psi'_X(\lambda)) = (\mathcal{M}_n \circ \psi'_X)(\lambda)
 \end{aligned}$$

mit  $\mathcal{M}_n(x) = n \cdot x$ . Außerdem ist  $(\mathcal{M}_n)^{-1}(x) = \frac{x}{n}$ .

$$\Rightarrow \lambda_t = (\psi'_Z)^{-1}(t) = (\psi'_X)^{-1}((\mathcal{M}_n)^{-1}(t)) = (\psi'_X)^{-1}\left(\frac{t}{n}\right).$$

Für die Cramér-Trafo von  $Z$  folgt schließlich

$$\begin{aligned} \psi_Z^*(t) &= t\lambda_t - \psi_Z(\lambda_t) \\ &= t(\psi'_Z)^{-1}(t) - \psi_Z((\psi'_Z)^{-1}(t)) \\ &= t(\psi'_X)^{-1}\left(\frac{t}{n}\right) - n\psi_X\left((\psi'_X)^{-1}\left(\frac{t}{n}\right)\right) \\ &= n\psi_X^*\left(\frac{t}{n}\right) \end{aligned}$$

□

(e) Cramér-Transformierte für zentrierte Binomialverteilungen:

Jetzt sind alle Vorbereitungen getroffen, die Cramér-Trafo der Binomialverteilung zu berechnen. Seien  $X_1, \dots, X_n \sim \text{Ber}(p)$  unabhängig.

$\Rightarrow Y := \sum_{i=1}^n X_i \sim \text{Bin}(n, p)$ . Betrachte zentrierte Version  $Z := Y - np = \sum_{i=1}^n (X_i - p)$ , da ja  $E(X_i) = p$ . Für  $t \in (0, n(1-p))$  gilt nach (d) und (c) mit  $a = p + \frac{t}{n}$ :

$$\psi_Z^*(t) = n \cdot \psi_{X-p}^*\left(\frac{t}{n}\right) = n \cdot h_p\left(\frac{t}{n} + p\right),$$

wobei  $h_p$  die in (c) bereits eingeführte Kulback-Leibler-Divergenz zwischen zwei Bernoulli-Verteilungen ist.

Wir erhalten insbesondere die Chernoff-Schranke:

$$P(Z \geq t) \leq \exp\left(-n \cdot h_p\left(\frac{t}{n} + p\right)\right).$$

$\rightarrow$  Guter Vorfaktor  $n$  für  $n \rightarrow \infty$ , aber wie schnell und wohin konvergiert  $h_p$  für  $n \rightarrow \infty$ ?  $\rightarrow$  Noch zu untersuchen!

### 2.3. Sub-Gaußsche Zufallsvariablen.

Gauß-ZVen  $\hat{=}$  gut umgängliche Klasse von ZVen

Klasse von Verteilungen, die von Gauß-ZVen dominiert werden, sind ähnlich gut umgänglich.

**Definition 2.5.** Eine reellwertige, zentrierte ZVe  $X$  heißt *sub-gaußsch* mit Varianzfaktor  $\nu$ , in Zeichen:  $P_X \in \mathcal{G}(\nu)$ , falls

$$(2.9) \quad \forall \lambda \in \mathbb{R} : \psi_X(\lambda) \leq \frac{\lambda^2 \nu}{2}.$$

Die Klasse aller solchen ZVen, bezeichnen wir (auch) mit  $\mathcal{G}(\nu)$ . Die Schranke entspricht der KEF einer zentrierten Normalverteilung mit Varianz  $\nu$ . Wir betrachten also die Klasse von Verteilungen, die durch gaußsche ZVen im Sinne der KEF dominiert werden.

Die Konzept von sub-gaußsch findet in den nachfolgenden Kapiteln 2.5 und 2.9 Anwendung.

*Bemerkung 2.6* (Eigenschaften sub-gaußscher Zufallsvariablen). (a) Varianz:

Es folgt  $\text{Var}(X) \leq \nu$ , im Allg. gilt aber  $\text{Var}(X) \neq \nu$ . Die Ungleichung kann man durch eine Taylorentwicklung zeigen (Übung). Sie ist insbesondere scharf. Betrachte dazu Rademacher-verteilte ZVen.

(b) Für normalverteilte ZVen  $X \sim \mathcal{N}(m, \nu)$  ist die momentenerzeugende Funktion  $M_X(\lambda) = \exp\left(m\lambda + \frac{\lambda^2}{2}\nu\right)$ , wie man durch direkte Rechnung sieht. Also gilt

$$X \sim \mathcal{N}(0, \nu) \implies X \in \mathcal{G}(\nu)$$

(c) In der Tat, Bedingung (2.9) kann sogar als Vergleich mit Gauß-ZV verstanden werden:

$$X = Y - E(Y) \in \mathcal{G}(\nu) \iff \psi_X \leq \psi_{\mathcal{N}(0, \nu)} \text{ auf } \mathbb{R}$$

(d) Verträglichkeit mit der Faltung:

Sind  $X_1, \dots, X_n$  unabhängige ZVen mit  $X_i \in \mathcal{G}(\nu_i)$  für  $i = 1, \dots, n$ , so gilt folgende Stabilitätseigenschaft für  $Z = \sum_{i=1}^n X_i \in \mathcal{G}(\sum_{i=1}^n \nu_i)$ . Nachweis erfolgt z.B. mit dem Additionssatz der Varianz und der Stabilität der KEF:

$$\psi_Z = \sum_{i=1}^n \psi_{X_i} \leq \sum_{i=1}^n \psi_{\mathcal{N}(0, \nu_i)} = \psi_{\mathcal{N}(0, \sum_{i=1}^n \nu_i)}.$$

(e) Inklusionen

$$\lambda \leq \tilde{\lambda} \implies \mathcal{G}(\lambda) \subset \mathcal{G}(\tilde{\lambda})$$

**Definition 2.7.** Ein messbares  $X: \Omega \rightarrow \{-1, 1\}$  mit  $P(X = 1) = P(X = -1) = \frac{1}{2}$  heißt *Rademacher-Zufallsvariable*.

Statt über Ungleichungen für KEF  $\psi_X$  kann man die Eigenschaft sub-gaußsch auch über Ungleichungen für Momente oder über Ungleichungen von Abfall bei  $\infty$  (engl. *tail-probability*) charakterisiert werden.

*Bemerkung 2.8* (Äquivalente *tail*-Charakterisierungen von sub-gaußsch). Sei  $X \in \mathcal{G}(\nu)$ . Chernoff liefert für  $t > 0$

$$\max \{P(X > t), P(X < -t)\} \leq e^{-\frac{t^2}{2\nu}}.$$

Für  $X \in \mathcal{G}(\nu)$  gilt nämlich  $\psi_X^*(t) \geq \psi_{\mathcal{N}(0,\nu)}^*(t)$  (Übung) und damit

$$P(X > t) \leq e^{-\psi_X^*(t)} \leq e^{-\psi_{\mathcal{N}(0,\nu)}^*(t)} = e^{-t^2/2\nu}$$

Analoge Rechnung für  $P(-X > t)$  wegen  $-X \in \mathcal{G}(\nu)$  liefert Behauptung.

**Satz 2.9.**

Sei  $X \in \mathcal{L}^1$  eine zentrierte ZVe.

(a) Gibt es ein  $\nu > 0$ , sodass  $\forall s > 0$

$$\max \{P(X > s), P(X < -s)\} \leq e^{-\frac{s^2}{2\nu}}$$

gilt, so folgt  $\forall q \in \mathbb{N}$ :

$$(2.10) \quad E(X^{2q}) \leq 2q!(2\nu)^q \leq q!(4\nu)^q.$$

(b) Existiert ein  $C \in (0, \infty)$  mit

$$(2.11) \quad \forall q \in \mathbb{N} : E(X^{2q}) \leq q!C^q,$$

so gilt  $X \in \mathcal{G}(4C)$ . Insbesondere folgt daraus

$$(2.12) \quad \max \{P(X > s), P(X < -s)\} \leq e^{-\frac{s^2}{8C}}.$$

Theorem 2.9 kann als Hilfsmittel dienen, um die Voraussetzung der Bernstein-Ungleichung 2.22 nachzurechnen, wie es im Beweis vom Johnson-Lindenstrauß-Lemma 2.26 in 2.9 getan wird.

*Beweis.* zu (a): Gilt  $X \in \mathcal{G}(\nu)$ , so ist  $Y = \frac{1}{\sqrt{\nu}}X \in \mathcal{G}(1)$ , denn

$$\psi_Y(\lambda) = \ln \left( E \left( e^{\lambda Y} \right) \right) = \ln \left( E \left( e^{\frac{\lambda X}{\sqrt{\nu}}} \right) \right) = \psi_X \left( \frac{\lambda}{\sqrt{\nu}} \right) \leq \frac{\frac{\lambda^2}{\nu} \cdot \nu}{2} = \frac{\lambda^2}{2}.$$

Betrachte zunächst den Fall  $\nu = 1$  und beginne auf der linken Seite von (2.10):

$$E(Y^{2q}) = \int_0^\infty P(|Y|^{2q} > x) \, dx.$$

Erste Substitution mit  $y = x^{\frac{1}{2q}}$  ergibt

$$= 2q \int_0^\infty P(|Y| > y) \cdot y^{2q-1} \, dy.$$

Anwendung der Voraussetzung ermöglicht folgende Abschätzung:

$$\leq 4q \int_0^\infty e^{-\frac{y^2}{2}} \cdot y^{2q-1} dy.$$

Mit zweiter Substitution  $t = \frac{y^2}{2}$  erhalten wir

$$= 4q \int_0^\infty e^{-t} \cdot (2t)^{q-\frac{1}{2}} \cdot (2t)^{-\frac{1}{2}} dt.$$

Diese Substitution ist nützlich, da das Integral explizite Darstellung hat (vgl. Bronstein):

$$E(Y^{2q}) \leq 2^{q+1} q!.$$

Für allgemeinere Varianzen  $\nu$  folgt:

$$E(X^{2q}) = E\left((\sqrt{\nu}Y)^{2q}\right) = \nu^q E(Y^{2q}) \leq 2^{q+1} \nu^q q! \leq 2 \cdot 2^q \nu^q q! = (4\nu)^q q!.$$

zu (b): Es gelte  $E(X^{2q}) \leq q! C^q$ .

Sei  $\tilde{X}$  eine unabhängig, identisch verteilte Kopie von  $X$ .

$\Rightarrow X - \tilde{X}$  ist symmetrisch verteilt, d.h.  $P(X - \tilde{X} > s) = P(\tilde{X} - X > s) \forall s \in \mathbb{R}$ . Wir nutzen den Satz zur monotone Konvergenz und dass die ungeraden Momente verschwinden:

$$\begin{aligned} & E\left(e^{\lambda X}\right) \cdot E\left(e^{-\lambda X}\right) \stackrel{\text{id. vert.}}{=} E\left(e^{\lambda X}\right) \cdot E\left(e^{-\lambda \tilde{X}}\right) \stackrel{\text{unabh.}}{=} E\left(e^{-\lambda(X-\tilde{X})}\right) \\ &= E\left(\sum_{q \in \mathbb{N}_0} \left[ \frac{\lambda^{2q}(X-\tilde{X})^{2q}}{(2q)!} + \frac{\lambda^{2q+1}(X-\tilde{X})^{2q+1}}{(2q+1)!} \right]\right) \\ &= \sum_{q \in \mathbb{N}_0} \left[ E\left(\frac{\lambda^{2q}(X-\tilde{X})^{2q}}{(2q)!}\right) + \underbrace{E\left(\frac{\lambda^{2q+1}(X-\tilde{X})^{2q+1}}{(2q+1)!}\right)}_{=0 \text{ wegen Symmetrie}} \right] \\ &= \sum_{q \in \mathbb{N}_0} \frac{\lambda^{2q} E\left(\left(X-\tilde{X}\right)^{2q}\right)}{(2q)!} \quad \forall \lambda \in \mathbb{R} \end{aligned}$$

Frage/Übung: Existieren überhaupt die ungeraden Momente? Sind sie summierbar?

Wegen Konvexität von  $x \mapsto x^{2q} = x^m$  für  $m \in 2\mathbb{N}$  folgt (vgl. auch Abb. ??):

$$\begin{aligned} & (a-b)^m \leq 2^{m-1}(a^m - b^m) \leq 2^{m-1}(a^m + b^m) \\ \Rightarrow & E\left(\left(X-\tilde{X}\right)^{2q}\right) \leq 2^{2q-1} \left(E(X^{2q}) + E(\tilde{X}^{2q})\right) \stackrel{\text{id. vert.}}{=} 2^{2q} E(X^{2q}) \end{aligned}$$

$$\begin{aligned}
\Rightarrow E(e^{\lambda X}) \cdot E(e^{-\lambda X}) &\stackrel{\text{s.o.}}{=} \sum_{q \in \mathbb{N}_0} \frac{\lambda^{2q} E((X - \tilde{X})^{2q})}{(2q)!} \\
&\leq \sum_{q \in \mathbb{N}_0} \frac{\lambda^{2q}}{(2q)!} 2^{2q} \underbrace{E(X^{2q})}_{\leq q! C^q \text{ nach Vor.}} \\
&\leq \sum_{q \in \mathbb{N}_0} 2^{2q} \lambda^{2q} C^q \frac{q!}{(2q)!}
\end{aligned}$$

Nebenrechnung:  $\frac{(2q)!}{q!} = \prod_{j=1}^q (q+j) \geq \prod_{j=1}^q 2j = 2^q \cdot q!$

$$\begin{aligned}
\Rightarrow E(e^{\lambda X}) \cdot E(e^{-\lambda X}) &\stackrel{\text{s.o.}}{\leq} \sum_{q \in \mathbb{N}_0} 2^{2q} \lambda^{2q} C^q \frac{q!}{(2q)!} \\
&\stackrel{\text{N.R.}}{\leq} \sum_{q \in \mathbb{N}_0} 2^{2q} \lambda^{2q} C^q \frac{1}{2^q q!} \\
&= \sum_{q \in \mathbb{N}_0} \frac{2^q \lambda^{2q} C^q}{q!} = e^{2C\lambda^2}.
\end{aligned}$$

Da  $X$  zentriert, folgt mit der Jensen-Ungleichung:

$$\begin{aligned}
\Rightarrow M_X(\lambda) &\leq E(e^{\lambda X}) \cdot \underbrace{E(e^{-\lambda X})}_{\geq 1} \leq e^{2C\lambda^2} \\
\Rightarrow \text{Also: } \psi_X(\lambda) &\leq \frac{4C}{2} \lambda^2 \Rightarrow X \in \mathcal{G}(4C)
\end{aligned}$$

□

Werfen wir noch einen zweiten Blick auf das *tail*-Verhalten:

**Lemma 2.10.** *Die Momentenbedingung (2.11) ist äquivalent zu folgender Bedingung :*

$$(2.13) \quad \exists \alpha > 0, \text{ sodass } E(e^{\alpha X^2}) \leq 2.$$

Interpretation: Da  $\exp(\cdot)$  schnell wächst, muss  $\alpha X^2$  schnell abfallen, falls (2.13) gilt.

*Beweis.* Nach Voraussetzung und Konvergenzsatz

$$2 \geq \sum_{k \in \mathbb{N}_0} \frac{\alpha^k E(X^{2k})}{k!}$$

$$\Leftrightarrow 1 \geq \sum_{k \in \mathbb{N}} \frac{\alpha^k E(X^{2k})}{k!}$$

Alle Summanden nichtnegativ  $\Rightarrow \forall k \in \mathbb{N}$  gilt:  $E(X^{2k}) \leq \alpha^{-k} k!$ .

Satz 2.9 Teil b)  $\Rightarrow X \in \mathcal{G}\left(\frac{4}{\alpha}\right)$ .

Gegenrichtung: mit Satz 2.9 Teil a)

$$X \in \mathcal{G}(\nu) \stackrel{\text{Satz 2.9}}{\Rightarrow} E(X^{2q}) \leq C^q q! \text{ mit } C = 4\nu.$$

Setze:  $\alpha = \frac{1}{2C} = \frac{1}{8\nu}$

$$\Rightarrow E\left(e^{\alpha X^2}\right) = \sum_{q \in \mathbb{N}_0} \frac{\alpha^q E(X^{2q})}{q!}$$

$$\leq \sum_{q \in \mathbb{N}_0} \left(\frac{1}{2C}\right)^q \frac{C^q q!}{q!} = \sum_{q \in \mathbb{N}_0} \left(\frac{1}{2}\right)^q = 2.$$

□

Quantitative Variante der Charakterisierung (2.13)? Es sei  $\alpha > 0$  und  $X$  zentriert. Dann:

$$X \in \mathcal{G}\left(\frac{1}{8\alpha}\right) \Rightarrow E\left(e^{\alpha X^2}\right) \leq 2 \Rightarrow X \in \mathcal{G}\left(\frac{4}{\alpha}\right)$$

Beschränkte zentrierte ZVen sind sub-gaußsch:

**Lemma 2.11** (Hoeffding-Lemma).

Sei  $Y$  eine  $[a, b]$ -wertige (schreibe zukünftig  $Y \in [a, b]$ ) zentrierte ZVe. Sei  $\psi_Y(\lambda) = \ln(E(e^{\lambda Y}))$ .

$$\Rightarrow \psi_Y''(\lambda) \leq \frac{(b-a)^2}{4} \text{ und}$$

$$Y \in \mathcal{G}\left(\frac{(b-a)^2}{4}\right)$$

Das Hoeffding-Lemma 2.11 wird beim Beweis der Hoeffding-Ungleichung 2.18 in 2.6 benötigt.

*Beweis.*

$$\begin{aligned} \left| Y - \frac{b+a}{2} \right| \leq \frac{b-a}{2} &\Rightarrow \text{Var}(Y) = \text{Var}\left(Y - \frac{b+a}{2}\right) \\ &\leq \left(\frac{b-a}{2}\right)^2 = \frac{(b-a)^2}{4}. \end{aligned}$$

Sei  $P_Y$  Verteilung von  $Y$  und  $P_\lambda(dx) = e^{-\psi_Y(\lambda)} e^{\lambda x} P_Y(dx)$  modifiziertes absolutstetiges Maß. Sei  $Z \in \mathbb{R}$  mit Verteilung  $P_Z = P_\lambda$ .

Frage/Übung: Ist überhaupt  $P_\lambda$  ein W-Mass?

Da  $Z \in [a, b]$  gilt ebenso:

$$\text{Var}(Z) \leq \left(\frac{b-a}{2}\right)^2.$$

Direkte Rechnung ergibt:

$$\begin{aligned} \psi_Y''(\lambda) &= \frac{d^2}{d\lambda^2} \ln(M_Y(\lambda)) \\ &= e^{-\psi_Y(\lambda)} E\left(Y^2 e^{\lambda Y}\right) - e^{-2\psi_Y(\lambda)} \left(E\left(Y e^{\lambda Y}\right)\right)^2 \\ &= e^{-\psi_Y(\lambda)} \int_{\mathbb{R}} y^2 e^{\lambda y} dP_Y(y) - \left(e^{-\psi_Y(\lambda)} \int_{\mathbb{R}} y e^{\lambda y} dP_Y(y)\right)^2 \\ &= \int_{\mathbb{R}} y^2 dP_\lambda(y) - \left(\int_{\mathbb{R}} y dP_\lambda(y)\right)^2 \\ &= \text{Var}(Z) \leq \frac{(b-a)^2}{4} \quad \forall \lambda \in \mathbb{R}. \end{aligned}$$

Da  $Y$  zentriert:  $\psi_Y(0) = \psi_Y'(0) = 0$  und  $\psi \in \mathcal{C}^2(0, \infty) \cap \mathcal{C}^1[0, \infty)$ . Taylorentwicklung mit Lagrange-Restglied liefert:

$$\begin{aligned} \exists \theta \in [0, \lambda] : \psi_Y(\lambda) &= \psi_Y(0) + \lambda \psi_Y'(0) + \frac{\lambda^2}{2} \psi_Y''(\theta) \\ &\leq 0 + 0 + \frac{\lambda^2(b-a)^2}{8} \Rightarrow Y \in \mathcal{G}\left(\frac{(b-a)^2}{4}\right) \end{aligned}$$

□

Beispiel: Ungleichung des Lemmas ist scharf: Führe sogenannte *Rademacher-Zufallsvariable*  $X: \Omega \rightarrow \{-1, 1\}$  ein mit  $P(X=1) = P(X=-1) = \frac{1}{2}$ . Es gilt  $X \in [-1, 1]$  mit  $a = -1$  und  $b = 1$ . Wie in Lemma 2.11 zeigen wir:

$$\Rightarrow \psi_X''(\lambda) \leq \frac{(b-a)^2}{4} = 1 \quad \forall \lambda \geq 0$$

Nutzen nun die charakteristische Eigenschaft der KEF.  $\lambda = 0$  einsetzen liefert:  $\text{Var}(X) = \psi_X''(0) \leq 1$ .

Nun berechnen wir die Varianz exakt:

$$\text{Var}(X) = E(X^2) = \frac{1}{2}(-1)^2 + \frac{1}{2}(-1)^2 = 1.$$

→ Ungleichung lässt sich nicht mehr verbessern.

#### 2.4. Sub-Gamma-Zufallsvariablen.

Einige wichtige Verteilungen haben Dichten, die schnell abfallen bei  $\pm\infty$ , jedoch nicht ganz so schnell wie  $\mathcal{G}(\nu)$ . Daher führen wir ein:

**Definition 2.12.** Sei  $X \in \mathcal{L}^1$  mit  $E(X) = 0$ .  $X$  heißt *sub- $\Gamma$ -verteilt von rechts* mit Varianzfaktor  $\nu > 0$  und Skalenparameter  $c \geq 0$ , in Symbolen

$$X \in \Gamma_+(\nu, c) :\Leftrightarrow \begin{cases} \forall \lambda \in (0, \frac{1}{c}) \text{ gilt } \psi_X(\lambda) \leq \frac{\lambda^2 \nu}{2(1 - c\lambda)}, & \text{falls } c > 0 \\ \forall \lambda \in (0, \infty) \text{ gilt } \psi_X(\lambda) \leq \frac{\lambda^2 \nu}{2}, & \text{falls } c = 0. \end{cases}$$

Die Klasse  $\Gamma_-(\nu, c)$  führen wir ein, indem wir setzen:

$$X \in \Gamma_-(\nu, c) \Leftrightarrow -X \in \Gamma_+(\nu, c),$$

d.h.:  $X$  ist sub- $\Gamma$  von links mit Varianzfaktor  $\nu$  und Skalenparameter  $c :\Leftrightarrow -X \in \Gamma_+(\nu, c) \Leftrightarrow X \in \Gamma_-(\nu, c)$ .

$X$  heißt sub- $\Gamma$ -verteilt mit Varianzfaktor  $\nu$  und Skalenparameter  $c \Leftrightarrow X \in \Gamma(\nu, c) := \Gamma_+(\nu, c) \cap \Gamma_-(\nu, c)$ .

Das Konzept von sub-gamma findet im Kapitel 2.5 Verwendung, kann (dort) aber auch weggelassen werden. Es soll ein alternatives, flexibleres Kriterium zur sub-gauß-Eigenschaft bereitstellen.

Einige Bemerkungen:

- (a) Insbesondere  $\Gamma(\nu, 0) = \mathcal{G}(\nu)$

(b) Sei ZV  $Y$  gammaverteilt mit Parameter  $a, b \geq 0$ ,  $X := Y - E(Y)$  zentrierte Version.  $Y$  hat die Dichte

$$\begin{aligned} \forall x \geq 0 : f(x) &= \frac{x^{a-1} e^{-\frac{x}{b}}}{\Gamma(a) b^a} \\ \Rightarrow E(Y) &= ab, \text{Var}(Y) = ab^2 \text{ und} \\ E(e^{\lambda X}) &= \int_0^\infty e^{\lambda(y-ab)} f(y) dy = e^{-\lambda ab - a \ln(1-\lambda b)} \\ \Rightarrow \forall \lambda \in \left(0, \frac{1}{c}\right) : \psi_X(\lambda) &= a(-\lambda b - \ln(1-\lambda b)) \\ &\stackrel{\text{NR.}}{\leq} \frac{\lambda^2 \nu}{2(1-c\lambda)}, \text{ mit } \nu = ab^2, c = b \end{aligned}$$

Nebenrechnung: Für  $x := \lambda b \in (0, 1)$  gilt mit der Reihenentwicklung des Logarithmus:

$$\begin{aligned} -\ln(1-x) - x &= \frac{x^2}{2} + \frac{x^3}{3} + \frac{x^4}{4} + \dots = \frac{x^2}{2} \left(1 + \frac{2}{3}x + \frac{2}{4}x^2 + \dots\right) \\ &\leq \frac{x^2}{2} (1 + x + x^2 + \dots) \leq \frac{x^2}{2} \frac{1}{1-x} \end{aligned}$$

Also sind  $\Gamma$ -ZVen  $\text{sub-}\Gamma(ab^2, b)$ . Allerdings ist  $X$  nicht symmetrisch verteilt. Die Verteilung von  $-X$  fällt sogar schneller ab, da ja

$$\psi_{-X}(\lambda) = \ln E(e^{-\lambda X}) = \ln E(e^{(-\lambda)X}) = a(\lambda b - \ln(1 + \lambda b)) \leq \frac{\lambda^2}{2} ab^2 = \frac{\lambda^2}{2} \nu$$

wegen

$$\begin{aligned} y - \ln(1-y) &= y - y + \frac{y^2}{2} - \frac{y^3}{3} + \frac{y^4}{4} - \dots \\ &= \frac{y^2}{2} - \left(\frac{y^3}{3} - \frac{y^4}{4}\right) - \left(\frac{y^5}{5} - \frac{y^6}{6}\right) - \dots \leq \frac{y^2}{2} \text{ für } y \in (0, 1). \end{aligned}$$

Also:  $X \in \Gamma_-(ab^2, 0) \subset \Gamma_-(ab^2, b)$  und damit  $X \in \Gamma(ab^2, b)$ .

Um das tail-Verhalten von sub-gamma ZV zu verstehen, untersuchen wir die Fenchel-Legendre-Duale von  $\psi(\lambda) = \frac{\lambda^2 \nu}{2(1-c\lambda)}$ . Setze dazu  $h_1(u) = 1 + u - \sqrt{1+2u}$  für  $u > 0$ . In der Übung wird gezeigt, dass

$$(2.14) \quad \psi^*(t) = \sup_{\lambda \in (0, \frac{1}{c})} \left( t\lambda - \frac{\lambda^2 \nu}{2(1-c\lambda)} \right) = \frac{\nu}{c^2} h_1\left(\frac{ct}{\nu}\right),$$

dass die Funktion  $h_1 : (0, \infty) \rightarrow (0, \infty)$  strikt monoton wachsend, und bijektiv mit  $h_1^{-1}(u) = u + \sqrt{2u}$  ist. Damit ergibt sich:

$$(2.15) \quad \begin{aligned} (\psi^*)^{-1}(u) &= (\mathcal{M}_{c/\nu})^{-1} \circ h_1^{-1} \circ (\mathcal{M}_{\nu/c^2})^{-1}(u) \\ &= \frac{\nu}{c} h_1^{-1} \left( \frac{c^2 u}{\nu} \right) = \frac{\nu}{c} \left( \frac{c^2 u}{\nu} + \sqrt{2c^2 u / \nu} \right) = \sqrt{2\nu u} + cu. \end{aligned}$$

Die Chernoff-Schranken lauten:

$$(i) \quad X \in \Gamma_+(\nu, c) \Rightarrow \forall t > 0 \text{ gilt } P(X > t) \leq \exp \left( -\frac{\nu}{c^2} h_1 \left( \frac{ct}{\nu} \right) \right).$$

Mit der Substitution  $s := \psi^*(t) = \frac{\nu}{c^2} h_1 \left( \frac{ct}{\nu} \right)$  und der berechtigten Formel für die Inverse erhalten wir die oftmals praktischere äquivalente Darstellung

$$\forall s > 0 \text{ gilt } P(X > \sqrt{2\nu s} + cs) \leq e^{-s}.$$

$$(ii) \quad X \in \Gamma(\nu, c) \Rightarrow \forall t > 0 :$$

$$\max\{P(X > \sqrt{2\nu t} + ct), P(-X > \sqrt{2\nu t} + ct)\} \leq e^{-t}.$$

Es gilt wieder (fast) eine Äquivalenz zu einer Momentenbedingung:

**Satz 2.13.**

Sei  $X \in \mathcal{L}^1$  eine zentrierte Zufallsvariable.

(a) Gilt für ein  $\nu > 0$  für jedes  $t > 0$ :

$$(2.16) \quad \max \left\{ P(X > \sqrt{2\nu t} + ct), P(X < -(\sqrt{2\nu t} + ct)) \right\} \leq e^{-t},$$

so folgt für jedes  $q \in \mathbb{N}$

$$E(X^{2q}) \leq q!(8\nu)^q + (2q)!(4C)^{2q}.$$

(b) Gilt umgekehrt für zwei Parameter  $A, B \geq 0$  und jedes  $q \in \mathbb{N}$

$$(2.17) \quad E(X^{2q}) \leq q!A^q + (2q)!B^{2q},$$

so ist bereits  $X \in \Gamma(4(A + B^2), 2B)$ .

Man beachte die Analogien zu Kapitel 2.3, u.a. Satz 2.9.

**2.5. Eine Maximal-Ungleichung.**

In Teil (C) der Motivation wollten wir das Supremum einer Familie von ZVen abschätzen. Konkrete Schranke wird in diesem Kapitel für den Erwartungswert des Supremums hergeleitet.

Seien ZVen  $Z_1, \dots, Z_N \in \mathbb{R}$  mit  $Z_i \in \mathcal{G}(\nu)$  für  $i = 1, \dots, N$  und  $\nu > 0$ .

Ziel: Schätze  $E\left(\max_{i=1,\dots,N} Z_i\right)$  nach oben ab.

Idee: Nutze sub-gaußsche Eigenschaft der ZVen wie folgt aus.

$$\begin{aligned}
 \exp\left(\lambda E\left(\max_{i=1,\dots,N} Z_i\right)\right) &\stackrel{\text{Jensen}}{\leq} E\left(\exp\left(\lambda \max_{i=1,\dots,N} Z_i\right)\right) \\
 &\leq E\left(\max_{i=1,\dots,N} \exp(\lambda Z_i)\right) \\
 (2.18) \qquad &\leq E\left(\sum_{i=1}^N \exp(\lambda Z_i)\right) \\
 &= \sum_{i=1}^N M_{Z_i}(\lambda) \leq N \cdot \exp\left(\frac{\lambda^2 \nu}{2}\right).
 \end{aligned}$$

Durch Logarithmieren erhalten wir eine Abschätzung in gewünschter Form:

$$E\left(\max_{i=1,\dots,N} Z_i\right) \leq \frac{\ln(N) + \frac{\lambda^2 \nu}{2}}{\lambda}.$$

Wähle nun

$$(2.19) \qquad \lambda = \sqrt{2(\ln N)/\nu}.$$

Wir erhalten:

$$(2.20) \qquad E\left(\max_{i=1,\dots,N} Z_i\right) \leq \frac{\sqrt{\ln(N)}}{\sqrt{2/\nu}} + \frac{\sqrt{\nu \ln(N)}}{\sqrt{2}} = \sqrt{2\nu \ln(N)}.$$

Fragen:

- Kann man  $\lambda$  geschickter wählen als in (2.19)?
- Ist die Schranke in (2.20) optimal?

Zur zweiten Frage: Sind  $Z_1, \dots, Z_N \sim \mathcal{N}(0, 1)$  unabhängig, folgt aus dem obigen:

$$\frac{E\left(\max_{i=1,\dots,N} Z_i\right)}{\sqrt{2\nu \ln(N)}} \leq 1.$$

Optimalität würde bedeuten, dass sogar:

$$\frac{E\left(\max_{i=1,\dots,N} Z_i\right)}{\sqrt{2\nu \ln(N)}} = 1.$$

Für dieses Beispiel gilt zumindest (Übung):

$$\lim_{N \rightarrow \infty} \frac{E \left( \max_{i=1, \dots, N} Z_i \right)}{\sqrt{2\nu \ln(N)}} = 1.$$

2.5.1. *Ähnliche Resultate gelten auch für sub- $\Gamma$ -Zufallsvariablen.*

Dazu holen wir etwas aus und entwickeln allgemeinere Resultate. Die folgenden Resultate dieses Kapitels sind anwendbar auf Klassen von Verteilungen, die geeignet dominiert werden. Die sub-gauß oder sub- $\Gamma$ -ZVen sind lediglich Beispiele für solche Klassen, sodass auch ohne Kapitel 2.4 nachfolgende Ausführungen (bis auf Korollar 2.16) verstanden werden können.

Untersuchung von Eigenschaftern der abstrakten Legendre-Fenchel-Dualen.

**Lemma 2.14.**

Seien  $b \in [0, \infty)$  und  $\psi: [0, b) \rightarrow \mathbb{R}$  konvex und  $\mathcal{C}^1([0, \infty))$  mit  $\psi(0) = \psi'(0) = 0$ . Für  $t \geq 0$  setzen wir:

$$\psi^*(t) := \sup_{\lambda \in [0, b)} (\lambda t - \psi(\lambda)).$$

Dann ist  $\psi^*$  auf  $[0, \infty)$  nichtnegativ, monoton wachsend, konvex und unbeschränkt. Die verallgemeinerte Inverse (Pseudo-Inverse):

$$(\psi^*)^{-1}(y) := \inf \{t \geq 0 \mid \psi^*(t) > y\}$$

erfüllt die Gleichung:  $(\psi^*)^{-1}(y) = \inf_{\lambda \in (0, b)} \left( \frac{y + \psi(\lambda)}{\lambda} \right)$ .

Die Voraussetzungen sind für die KEF  $\psi_X$  einer zentrierten ZV  $X$  erfüllt.

Man beachte inhaltliche Analogien zum Kapitel 2.2 .

*Beweis.*  $t \mapsto \lambda t - \psi(\lambda)$  ist affin-linear, konvex und isoton für  $\lambda \geq 0$ .

Übung  $\Rightarrow \psi^*$  auch konvex und isoton als Supremum solcher Funktionen.

$$\psi^*(0) = \sup_{0 < \lambda < b} (0 - \psi(\lambda)) = - \inf_{\lambda \in (0, b)} \psi(\lambda) = 0,$$

**Denn:**  $\psi$  konvex  $\Rightarrow \psi'$  isoton, mit  $\psi'(0) = 0$ . Also ist  $\psi'$  nichtnegativ.

$\Rightarrow \psi$  isoton, mit  $\psi(0) = 0$  und somit ist  $\psi$  nichtnegativ.

Damit ist auch  $\psi^*$  nichtnegativ.

Sei nun  $\lambda_0 \in (0, b)$  fix. Da für jedes  $t \geq 0$

$$\psi^*(t) \geq \sup_{\lambda \in (0, b)} (\lambda t - \psi(\lambda)) \geq \lambda_0 t - \psi(\lambda_0)$$

vorliegt, ist  $\psi^*$  unbeschränkt und  $\{t \geq 0 \mid \psi^*(t) > y\} \neq \emptyset, \forall y \geq 0$ .

Setze  $u := \inf_{\lambda \in (0, b)} \left( \frac{y + \psi(\lambda)}{\lambda} \right)$ , so gilt  $\forall t \geq 0$ :

$$\begin{aligned} u \geq t &\Leftrightarrow \forall \lambda \in (0, b) : \frac{y + \psi(\lambda)}{\lambda} \geq t \\ &\Leftrightarrow \forall \lambda \in (0, b) : y \geq \lambda t - \psi(\lambda) \\ &\Leftrightarrow y \geq \sup_{\lambda \in (0, b)} (\lambda t - \psi(\lambda)) = \psi^*(t) \end{aligned}$$

Komplementbildung liefert äquivalente Aussage dazu:

$$u < t \Leftrightarrow \psi^*(t) > y.$$

Also:  $(u, \infty) = \{t \geq 0 \mid \psi^*(t) > y\}$

$$\Rightarrow u = \inf((u, \infty)) = \inf\{t \geq 0 \mid \psi^*(t) > y\} = \underbrace{(\psi^*)^{-1}(t)}_{\text{verallg. Inverse}}. \quad \square$$

Vergleichen Sie die obige Vorgehensweise mit der aus der Stochastik bekannten Definition/Konstruktion der Quantil-Transformierten.

Anwendung von Lemma 2.14:

**Satz 2.15** (Maximal-Ungleichung).

Seien  $Z_1, \dots, Z_N \in \mathbb{R}$  ZVen,  $b \in (0, \infty)$ ,  $\psi \in C^1([0, b])$  konvex mit  $\psi(0) = \psi'(0) = 0$ , sodass

$$\psi_{Z_i}(\lambda) = \ln(M_{Z_i}(\lambda)) \leq \psi(\lambda) \quad \text{für } \lambda \in [0, b] \text{ und } i = 1, \dots, N.$$

Dann folgt

$$\text{für } i = 1, \dots, N : Z_i \in \mathcal{L}^1 \text{ und } E \left( \max_{i=1, \dots, N} Z_i \right) \leq (\psi^*)^{-1}(\ln(N)).$$

Es wird keine Unabhängigkeit verlangt!

*Beweis.* Mit Jensen-Ungleichung ist bekannt für  $\lambda \in (0, b)$ :

$$\begin{aligned} \exp \left( \lambda E \left( \max_{i=1, \dots, N} Z_i \right) \right) &\leq E \left( \exp \left( \lambda \max_{i=1, \dots, N} Z_i \right) \right) \\ &\leq E \left( \left( \max_{i=1, \dots, N} e^{\lambda Z_i} \right) \right) \leq \sum_{i=1}^N E \left( e^{\lambda Z_i} \right) \stackrel{\text{Vorr.}}{\leq} N e^{\psi(\lambda)} \end{aligned}$$

$$\begin{aligned}
&\Rightarrow \forall \lambda \in (0, b) : \lambda E \left( \max_{i=1, \dots, N} Z_i \right) \leq \psi(\lambda) + \ln(N) \\
&\Rightarrow \forall \lambda \in (0, b) : E \left( \max_{i=1, \dots, N} Z_i \right) \leq \frac{\psi(\lambda) + \ln(N)}{\lambda} \\
&\Leftrightarrow E \left( \max_{i=1, \dots, N} Z_i \right) \leq \inf_{\lambda \in (0, b)} \frac{\psi(\lambda) + \ln(N)}{\lambda} \stackrel{2.14}{=} (\psi^*)^{-1}(\ln(N))
\end{aligned}$$

□

**Korollar 2.16.**

Seien  $Z_1, \dots, Z_n \in \mathcal{L}^1$  zentrierte ZVen. Es gelten:

(a) Falls  $Z_i \in \mathcal{G}(\nu)$  für  $i = 1, \dots, N$ , so:

$$E \left( \max_{i=1, \dots, N} Z_i \right) \leq \sqrt{2\nu \ln(N)}.$$

(b) Falls  $Z_i \in \Gamma_+(\nu, c)$  für  $i = 1, \dots, N$ , so:

$$E \left( \max_{i=1, \dots, N} Z_i \right) \leq \sqrt{2\nu \ln(N)} + c \ln(N).$$

*Beweis.* zu (a): siehe Idee/Rechnung (2.18) am Anfang von Kapitel 2.5.

zu (b): Einsetzen von (2.15) in  $(\psi^*)^{-1}(\ln(N))$  liefert die angegebene obere Schranke. □

*Beispiel 2.17* ( $\chi^2$ -Verteilung).

Sind  $Y_1, \dots, Y_k \sim \mathcal{N}(0, 1)$  unabhängige ZVen, so ist  $X := \sum_{i=1}^k Y_i^2$   $\chi^2$ -verteilt mit  $k$  Freiheitsgraden. Die Dichte von  $X$  ist

$$f(x) = \frac{x^{k/2-1} e^{-x/2}}{\Gamma(k/2) \sqrt[2]{2}},$$

d.h. Dichte der  $\Gamma$ -Verteilung mit  $a = k/2, b = 2$ . Insbesondere  $X - E(X) = X - k \in \Gamma_+(2k, 2) \cap \Gamma_-(2k, 0)$ . Für  $X_1, \dots, X_N \sim \chi^2$  mit  $k$  Freiheitsgraden impliziert Korollar 2.16

$$E \left( \max_{i=1, \dots, N} X_i - k \right) \leq 2\sqrt{k \ln(N)} + 2 \ln(N) \square$$

**2.6. Hoeffding-Ungleichung.**

Nächstes Ziel: Schranken für Wahrscheinlichkeiten für große Werte von Summen von unabhängigen ZVen. Seien  $X_1, \dots, X_n \in \mathbb{R}$  unabhängige  $\mathcal{L}^1$ -ZVen, sodass ein echtes Intervall  $I \subset \mathbb{R}$  existiert mit  $M_{X_i}(\lambda) = E(e^{\lambda X_i}) < \infty$  für  $i = 1, \dots, n, \lambda \in I$ . Für  $S := \sum_{i=1}^n (X_i - E(X_i))$  gilt:

$$\forall \lambda \in I : \psi(\lambda) = \sum_{i=1}^n \ln(E(e^{\lambda(X_i - E(X_i))})).$$

Sind die  $X_i$  sogar beschränkt, genauer:  $X_i \in [a_i, b_i]$  für  $i = 1, \dots, n$ , so gilt nach Hoeffding-Lemma 2.11:

$$\begin{aligned} \psi''_{X_i - E(X_i)}(\lambda) &\leq \frac{(b_i - a_i)^2}{4} \Rightarrow \psi_{X_i - E(X_i)}(\lambda) \leq \frac{\lambda^2 (b_i - a_i)^2}{8} \\ (2.8) \Rightarrow \psi_S(\lambda) &\leq \frac{\lambda^2}{8} \sum_{i=1}^n (b_i - a_i)^2 \\ \Rightarrow S &\in \mathcal{G} \left( \frac{1}{4} \sum_{i=1}^n (b_i - a_i)^2 \right) \end{aligned}$$

und somit

**Satz 2.18** (Hoeffding-Ungleichung).

Seien  $X_1 \in [a_1, b_1], \dots, X_n \in [a_n, b_n]$  unabhängige ZVen mit zentrierter Summe  $S$ . Dann gilt  $\forall t \geq 0$ :

$$P(S \geq t) \leq \exp \left( - \frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2} \right)$$

*Beweis.* s.o. und Bemerkung 2.8. □

Bemerkung/Beispiel: Wähle  $X_i = \alpha_i \varepsilon_i$ , wobei  $\varepsilon_1, \dots, \varepsilon_n$  unabh. Rademacher-ZVen:

$$\stackrel{\text{Satz 2.18}}{\Rightarrow} P(S \geq t) \leq \exp \left( - \frac{2t^2}{4 \sum_{i=1}^n \alpha_i^2} \right)$$

Da  $\text{Var}(S) = \text{Var}(\alpha_1 \varepsilon_1 + \dots + \alpha_n \varepsilon_n) = \sum_i \alpha_i^2 \text{Var}(\varepsilon_i) = \sum_i \alpha_i^2$  identifizieren wir

$$\exp \left( - \frac{2t^2}{4 \sum_{i=1}^n \alpha_i^2} \right) = \exp \left( - \frac{t^2}{2 \text{Var}(S)} \right)$$

Im Allgemeinen gilt aber:  $\text{Var}(S) < \frac{1}{4} \sum_{i=1}^n (b_i - a_i)^2$ .

In diesen Fällen gibt es feinere Ungleichungen: die *Bennett-* und *Bernstei- nungleichung*.

*Bemerkung 2.19* (Extremalitätseigenschaft der Rademacher-ZVen). Wie lässt sich die Varianz beschränkter ZVen maximieren?

→ schiebe Werte so weit wie möglich nach außen, an die Ränder des Wertebereiches

→ Rademacher-ZVen sind gerade diejenigen ZVen, die auf einem Intervall  $[-\alpha_i, \alpha_i]$  die Varianz maximieren (zum Wert  $\sum_{i=1}^n \alpha_i^2$ ).

## 2.7. Bennett-Ungleichung.

Wie zuvor:  $X_1, \dots, X_n$  unabhängige ZVen.  $S := \sum_{i=1}^n (X_i - E(X_i))$ .

$$(2.21) \quad \psi_S(\lambda) = \sum_{i=1}^n \left( \ln(E(e^{\lambda X_i})) - \lambda E(X_i) \right)$$

$$(2.22) \quad \leq \sum_{i=1}^n E(e^{\lambda X_i} - \lambda X_i - 1)$$

da ja  $\ln u \leq u - 1$  für  $u > 0$ .

(2.21) und (2.22) sind Startpunkte für Bennett- bzw. Bernstein-Ungleichung.

**Satz 2.20** (Bennett-Ungleichung).

Sei  $b > 0$ . Seien  $X_1, \dots, X_n \in \mathcal{L}^2$  unabhängige ZVen, sodass  $X_i \leq b$  fast sicher für  $i = 1, \dots, n$  (einseitige Beschränktheit).

Sei  $\nu = \sum_{i=1}^n E(X_i^2)$ . Dann gilt:

$$(2.23) \quad \psi_S(\lambda) \leq n \ln \left( 1 + \frac{\nu}{nb^2} \phi(b\lambda) \right) \leq \frac{\nu}{b^2} \phi(b\lambda)$$

mit  $\phi(u) := e^u - u - 1$  für  $u \in \mathbb{R}$ . Ferner gilt:

$$P(S \geq t) \leq \exp \left( -\frac{\nu}{b^2} h \left( \frac{bt}{\nu} \right) \right),$$

wobei  $h(u) := (1 + u) \ln(1 + u) - u$  für  $u \geq 0$ .

Die Bennett-Ungleichung wird mit der Bernstein-Ungleichung in Kapitel 2.8 verglichen. Ansonsten wird die Bennett-Ungleichung im weiteren Verlauf des Skripts nicht verwendet.

*Beweis.* O.E. sei  $b = 1$  (skaliere sonst nach). Die Abbildung

$\mathbb{R} \setminus \{0\} \ni u \mapsto \phi(u)/u^2$  ist isoton und stetig auf  $\mathbb{R}$  fortsetzbar (Übung).

Sei  $\lambda > 0$ . Nach Voraussetzung ist  $\lambda X_i \leq \lambda b = \lambda$ . Also gilt  $\forall i = 1, \dots, n$

$$\begin{aligned} e^{\lambda X_i} - \lambda X_i - 1 &\leq X_i^2 (e^\lambda - \lambda - 1), \\ \Rightarrow E(e^{\lambda X_i}) &\leq 1 + \lambda E(X_i) + E(X_i^2) \phi(\lambda). \end{aligned}$$

Mit Aufsummieren und (2.21) folgt:

$$\begin{aligned}
 \psi_S(\lambda) &\stackrel{(2.21)}{=} \sum_{i=1}^n (\ln(E(e^{\lambda X_i})) - \lambda E(X_i)) \\
 &\leq \frac{n}{n} \sum_{i=1}^n (\ln(1 + \lambda E(X_i) + E(X_i^2)\phi(\lambda)) - \lambda E(X_i)) \\
 (2.24) \quad &\stackrel{\text{konkav}}{\leq} n \left( \ln \left( 1 + \lambda \frac{\sum_{i=1}^n E(X_i)}{n} + \underbrace{\frac{\sum_{i=1}^n E(X_i^2)}{n}}_{=\nu/n} \phi(\lambda) \right) - \lambda \frac{\sum_{i=1}^n E(X_i)}{n} \right) \\
 &\stackrel{\ln(x+1) \leq x}{\leq} \nu \phi(\lambda).
 \end{aligned}$$

Um die andere in (2.23) behauptete Schranke zu erhalten, schätzen wir (2.24) ab mit Hilfe von:

$$\begin{aligned}
 &\text{Für } x \geq 1, a > 0 : 1 + \frac{a}{x} \leq 1 + a \leq e^a \\
 &\Rightarrow \text{Für } a, b > 0 : \ln(1 + a + b) - a \leq \ln(1 + b).
 \end{aligned}$$

In Kapitel 2.2 wurde bereits gesehen:

$$\nu \phi(\lambda) \text{ ist KEF von } Y = X - E(X) \text{ für } X \sim Poi(\nu)$$

Aus der Ungleichung zwischen zwei KEFs ergibt sich eine zwischen den entsprechenden Cramer-Transformierten:

$$\psi_S^*(t) \geq \psi_Y^*(t) = \nu h\left(\frac{t}{\nu}\right) \stackrel{\text{Chernoff}}{\Rightarrow} P(S \geq t) \leq \exp(-\psi_S^*(t)) \leq \exp\left(-\nu h\left(\frac{t}{\nu}\right)\right)$$

□

*Bemerkung 2.21.* In einer Übungsaufgabe wird gezeigt

$$\begin{aligned}
 h(u) &= (1 + u) \ln(1 + u) - u \geq \frac{u^2}{2\left(1 + \frac{u}{3}\right)} \\
 &\stackrel{\text{Bennett}}{\Rightarrow} \underbrace{P(S \geq t) \leq \exp\left(-\frac{t^2}{2\left(\nu + \frac{bt}{3}\right)}\right)}_{\text{Bernstein-Ungleichung}}
 \end{aligned}$$

Das ist die *Bernstein-Ungleichung*. Für  $\nu \gg tb$  sind Bennett und Bernstein im Wesentlichen äquivalent. Für  $t \gg \frac{\nu}{b}$  verliert Bernstein gegenüber Bennett einen  $\ln(t)$ -Faktor im Exponenten.

→ Bernstein lässt sich aber allgemeiner beweisen (für schwächere Annahmen an ZVen)!

## 2.8. Bernstein-Ungleichung.

**Satz 2.22** (Bernstein-Ungleichung). *Seien  $X_1, \dots, X_n$  unabhängige ZVen mit  $\nu, c > 0$ , sodass  $\sum_{i=1}^n E(X_i^2) \leq \nu$  und  $\sum_{i=1}^n E((X_i)_+)^q \leq \frac{q!}{2} \nu c^{q-2}$  für alle  $q \in \mathbb{N}_{\geq 3}$ . Dann gilt  $\forall \lambda \in (0, \frac{1}{c}), \forall t > 0$ :*

$$\psi_S(\lambda) \leq \frac{\nu \lambda^2}{2(1 - c\lambda)}, \quad \psi_S^*(t) \geq \frac{\nu}{c^2} h_1\left(\frac{ct}{\nu}\right),$$

wobei  $h_1(u) = 1 + u - \sqrt{1 + 2u}$ ,  $u > 0$ . Insbesondere gilt  $\forall t > 0$ :

$$P(S \geq \sqrt{2\nu t} + ct) \leq e^{-t}.$$

Die Bernstein-Ungleichung wird in Kapitel 2.9 beim Johnson-Lindenstrauf-Lemma verwendet.

*Beweis.* Im Beweis der Bennett-Ungleichung wurde  $\phi(u) := e^u - u - 1$  definiert. Es gilt  $\phi(u) \leq \frac{u^2}{2}$  für  $u \leq 0$  (Übung). Also folgt für  $\lambda > 0, i = 1, \dots, n$ :

$$\begin{aligned} \phi(\lambda X_i) &\leq \begin{cases} \sum_{q=3}^{\infty} \frac{\lambda^q (X_i)^q}{q!}, & \text{falls } X \geq 0 \\ \frac{\lambda^2 (X_i)^2}{2}, & \text{falls } X \leq 0. \end{cases} \\ &\leq \frac{\lambda^2 (X_i)^2}{2} + \sum_{q=3}^{\infty} \frac{\lambda^q (X_i)_+^q}{q!} \end{aligned}$$

$$\Rightarrow E(\phi(\lambda X_i)) \leq \frac{\lambda^2 E(X_i^2)}{2} + \sum_{q=3}^{\infty} \frac{\lambda^q E((X_i)_+^q)}{q!}$$

$$\Rightarrow \sum_{i=1}^n E(\phi(\lambda X_i)) \leq \frac{\nu}{2} \sum_{q=2}^{\infty} \lambda^q c^{q-2}. \quad \text{nach Vorauss.}$$

Letzte Summe ist endlich, falls  $\lambda \in (0, \frac{1}{c})$ . Mit (2.22) ergibt sich

$$\psi_S(\lambda) \leq \sum_{i=1}^n E(\phi(\lambda X_i)) \leq \frac{\nu}{2} \sum_{q=2}^{\infty} \lambda^q c^{q-2} = \frac{\nu \lambda^2}{2(1 - c\lambda)}.$$

(Der letzte Schritt beleuchtet, warum der Bruch ein Bezugswert ist.) Damit folgt mit (2.15) insgesamt

$$\psi_S^*(t) \geq \sup_{\lambda \in (0, \frac{1}{c})} \left( t\lambda - \frac{\lambda^2 \nu}{2(1 - c\lambda)} \right) = \frac{\nu}{c^2} h_1\left(\frac{ct}{\nu}\right)$$

Die Abschätzung an  $P(S \geq t)$  folgt mit der Bemerkung über  $h_1^*$  vor Satz 2.13 in Kapitel 2.4, vgl. Übung.  $\square$

**Korollar 2.23.**

Seien  $X_1, \dots, X_n$  unabhängige ZVen wie in Theorem 2.22. Dann gilt  $\forall t > 0$ :

$$P(S \geq t) \leq \exp\left(-\frac{t^2}{2(\nu + ct)}\right)$$

*Beweis.* Direkte Folgerung aus  $h_1(u) \geq \frac{u^2}{2(1+u)}$  (Übungsaufgabe).  $\square$

Man kann zeigen, dass dies für  $X_i \leq b$  die gleiche Schranke liefert, wie am Ende von Kapitel 2.7 aus Bennett hergeleitet wurde (Übungsaufgabe).

*Beispiel 2.24* (Gaußsches Chaos der Ordnung Zwei).

Sei  $X = (X_1, \dots, X_n)$  mit unabhängigen Komponenten  $X_j \sim \mathcal{N}(0, 1)$ , also  $X \sim \mathcal{N}(0, Id_n)$ . Sei  $A = (a_{ij})_{i,j=1}^n \in \mathbb{R}^{n \times n}$  symmetrisch mit  $a_{jj} = 0$  für  $j = 1, \dots, n$ , insbesondere  $\text{Spur}(A) = \sum_{j=1}^n a_{jj} = 0$ .

Definiere ZVe  $Z$  als quadratische Form:  $Z = X^\top A X = \sum_{i,j=1}^n X_i a_{ij} X_j$ .

$$\Rightarrow E(Z) = \sum_{i,j=1, i \neq j}^n a_{ij} \underbrace{E(X_i X_j)}_{E(X_i)E(X_j)=0} + \sum_{i=1}^n \underbrace{a_{ii}}_{=0} E(X_i X_i) = 0.$$

Frage:

Wie stark schwankt die ZVe  $Z$  um den Mittelwert 0? ( $\hat{=}$  Konzentration).

Da  $A$  symmetrisch, existiert  $B \in \mathbb{R}^{n \times n}$  orthogonal mit  $A = B^\top D B$ , wobei

$$D = \begin{pmatrix} \mu_1 & 0 & \dots & 0 \\ 0 & \mu_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \mu_n \end{pmatrix} \text{ Diagonalmatrix mit Eigenwerten } \mu_1, \dots, \mu_n \text{ von } A.$$

Sei  $B = (b_{ij})_{i,j=1}^n$ . Setze  $Y_i = \sum_{j=1}^n b_{ij} X_j$  für  $i = 1, \dots, n$ .

$$\begin{aligned} \Rightarrow \sum_{i=1}^n \mu_i Y_i^2 &= \sum_{i=1}^n \mu_i \left( \sum_{j=1}^n b_{ij} X_j \right)^2 = \sum_{i=1}^n \mu_i (B X)_i^2 \\ &= \langle X, B^\top D B X \rangle = X^\top A X = Z. \end{aligned}$$

Ist diese andere Darstellung von  $Z$  über  $Y_i$  besser?

Da  $X \sim \mathcal{N}(0, Id_n)$  und  $B$  orthogonal, ist auch  $Y = (Y_1, \dots, Y_n) \sim \mathcal{N}(0, Id_n)$  wegen Rotationsinvarianz.

$$\begin{aligned} \Rightarrow P_X &= P_Y, P_{(X_1^2, \dots, X_n^2)} = P_{(Y_1^2, \dots, Y_n^2)}. \\ \Rightarrow P_Z &= P_{\sum_{i=1}^n \mu_i Y_i^2} = P_{\sum_{i=1}^n \mu_i X_i^2}. \end{aligned}$$

Da  $0 = \text{Spur}(A) = \sum_{i=1}^n \mu_i$  gilt:

$$\sum_{i=1}^n \mu_i X_i^2 = \sum_{i=1}^n \mu_i (X_i^2 - 1).$$

Eigenschaft von  $X_i^2$ :  $E(X_i^2) = \text{Var}(X_i) = 1 \Rightarrow X_i^2 - 1$  zentriert  
 $X_i^2 \sim \chi^2$ -verteilt mit einem Freiheitsgrad, also insbesondere  $\Gamma$ -verteilt mit  
 Parameter  $a = \frac{1}{2}$  und  $b = 2$ . Nach Bsp. 2.17 in Kapitel 2.4 gilt:

$$\begin{aligned} \psi_{X_i^2-1}(\lambda) &= \frac{1}{2} [-\ln(1 - 2\lambda) - 2\lambda] \\ &\leq 2 \frac{\lambda^2}{2(1 - 2\lambda)} = \frac{\lambda^2}{1 - 2\lambda}. \end{aligned}$$

Also folgt für die KEF  $\psi_Z$  von  $Z$  ähnlich wie (2.8) (in Kapitel 2.6) wegen  
 Unabhängigkeit:

$$\begin{aligned} \psi_Z(\lambda) &= \psi_{\sum_{i=1}^n \mu_i (X_i^2 - 1)}(\lambda) = \sum_{i=1}^n \psi_{\mu_i (X_i^2 - 1)}(\lambda) \\ &= \sum_{i=1}^n \psi_{X_i^2 - 1}(\mu_i \lambda) = \frac{1}{2} \sum_{i=1}^n (-\ln(1 - 2\mu_i \lambda) - 2\mu_i \lambda) \\ &\leq \sum_{i=1}^n \frac{\mu_i^2 \lambda^2}{1 - 2\mu_i \lambda}, \text{ sofern } \lambda \in (0, (2 \max_{i=1, \dots, n} \mu_i)^{-1}) =: J \end{aligned}$$

Nun gilt  $\forall j = 1, \dots, n$

$$\mu_i \leq |\mu_i| \leq \|A\| = \sup_{\|x\|=1} \|Ax\|.$$

Sei  $\|A\|_{\text{HS}} := \sqrt{\sum_{i=1}^n \mu_i^2}$  die *Hilbert-Schmidt-* oder *Frobenius-Norm* von  $A$ .  
 Dann folgt für  $\lambda \in J$ :

$$\psi_Z(\lambda) \leq \sum_{i=1}^n \frac{\mu_i^2 \lambda^2}{1 - 2\lambda \|A\|} = \frac{\lambda^2 \|A\|_{\text{HS}}^2}{1 - 2\lambda \|A\|}.$$

In der Notation der sub- $\Gamma$ -Verteilungen ist  $\nu = 2\|A\|_{\text{HS}}^2$ ,  $c = 2\|A\|$  und nach  
 Kapitel 2.4 gilt:

$$\begin{aligned} P(Z > 2\|A\|_{\text{HS}} \sqrt{t} + 2\|A\|t) &\leq e^{-t} \text{ (vgl. Satz 2.13), oder} \\ P(Z > t) &< \exp\left(-\frac{\|A\|_{\text{HS}}^2}{2\|A\|} h_1\left(\frac{\|A\|t}{\|A\|_{\text{HS}}^2}\right)\right) \text{ oder} \\ P(Z > t) &\leq \exp\left(-\frac{t^2}{4(\|A\|_{\text{HS}}^2 + t\|A\|)}\right) \text{ Übung wie in Kor. 2.23,} \end{aligned}$$

wobei  $h_1$  die schon bekannte Entropie-Funktion ist.

## 2.9. Johnson-Lindenstrauss-Lemma.

Folgende Ausführungen werden in einem viel allgemeineren Rahmen in Kapitel 5.6 ausgeführt.

**Definition 2.25.** Sei  $\mathcal{H}$  ein separabler Hilbertraum (z.B.  $\mathcal{H} = \mathbb{R}^D$  mit Dimension  $D \in \mathbb{N}$ ). Für gegebenes  $\varepsilon \in (0, 1)$  und  $A \subset \mathcal{H}$  heißt  $f: \mathcal{H} \rightarrow \mathbb{R}^d$   $\varepsilon$ -Isometrie auf  $A$ , falls für alle  $a, b \in A$  gilt:

$$(2.25) \quad (1 - \varepsilon)\|a - b\|_{\mathcal{H}} \leq \|f(a) - f(b)\|_{\mathbb{R}^d} \leq (1 + \varepsilon)\|a - b\|_{\mathcal{H}}$$

Falls  $\mathcal{H}$  hochdimensional,  $A$  große Kardinalität hat und  $\varepsilon \in (0, 1)$  und  $d \in \mathbb{N}$  klein sind, ist völlig unklar, ob ein solches  $f$  existiert.

Das Johnson-Lindenstrauss-Lemma besagt, dass es eine universelle Konstante  $\kappa$  gibt, sodass für jedes  $A$  mit Kardinalität  $\text{card}(A) = n < \infty$  und

$$d \geq \frac{\kappa}{\varepsilon^2} \ln(n)$$

tatsächlich ein  $f: \mathcal{H} \rightarrow \mathbb{R}^d$  mit Eigenschaft (2.25) existiert.

Allerdings kann man dieses  $f$  nicht explizit angeben, da es mit einem *Zufallsmechanismus* konstruiert wird. (Vergleiche die *probabilistic method* von Paul Erdős in der Kombinatorik und Graphentheorie.) Dafür ist es möglich  $f$ , linear zu wählen. Wir zeigen sogar, dass aus einem vorgegebenen Ensemble linearer Funktionen  $f: \mathcal{H} \rightarrow \mathbb{R}^d$ , die meisten (2.25) erfüllen.

Einfachheitshalber nehmen wir  $\mathcal{H} = \mathbb{R}^D$  und typischerweise  $D \gg d$  an.

Ansatz: Sei  $W: \mathbb{R}^D \rightarrow \mathbb{R}^d$  linear und zufällig gewählt mit der Eigenschaft

$$\forall \alpha \in \mathbb{R}^D : E(\|W\alpha\|^2) = \|\alpha\|^2,$$

d.h. im Mittel haben wir eine exakte Isometrie. Gewünscht ist die Eigenschaft, dass die Zufallsvariable  $\|W\alpha\|^2$  wenig streut.

### Konstruktion von $W$

Seien  $X_{ij}, i = 1, \dots, d$  und  $j = 1, \dots, D$  unabhängig identisch verteilte Zufallsvariablen. Ferner seien  $E(X_{ij}) = 0$  und  $\text{Var}(X_{ij}) = 1$ . Für  $\alpha = (\alpha_1, \dots, \alpha_D) \in \mathbb{R}^D$  und  $i = 1, \dots, d$  setze

$$(2.26) \quad \tilde{W}_i(\alpha) = \sum_{j=1}^D \alpha_j X_{ij} \quad \text{und} \quad W_\alpha = \left( \frac{1}{\sqrt{d}} \tilde{W}_i(\alpha) \right)_{i=1}^d.$$

Wegen der Unabhängigkeit gilt für jedes  $i = 1, \dots, d$ :

$$E\left(\tilde{W}_i(\alpha)^2\right) = E\left(\left(\sum_{j=1}^D \alpha_j X_{ij}\right)^2\right) = \sum_{i,j=1}^D \alpha_i \alpha_j E(X_{ij}^2) \delta_{ij},$$

wobei  $\delta_{ij}$  das Kronecker-Delta ist. Daraus folgt für alle  $\alpha \in \mathbb{R}^D$

$$E(\|W_\alpha\|^2) = E\left(\frac{1}{d} \sum_{i=1}^d (\tilde{W}_i(\alpha))^2\right) = \frac{1}{d} \sum_{i=1}^d \|\alpha\|^2 = \|\alpha\|^2.$$

Wir wollen nun solche Zufallsvariablen  $X_{ij}$  wählen, die sich typischerweise so verhalten, wie ihr Mittelwert. Das ist eine typische Konzentrationsbedingung, wie sie z.B. für sub-gaußsche Zufallsvariablen erfüllt ist.

**Satz 2.26** (Johnson-Lindenstrauss-Lemma).

Seien  $A \subset \mathbb{R}^D$  mit  $\text{card}(A) = n \in \mathbb{N}$ , sowie  $\varepsilon, \delta \in (0, 1)$ . Für  $\nu \geq 1$  seien  $X_{ij} \in \mathcal{G}(\nu)$ ,  $i = 1, \dots, d$  und  $j = 1, \dots, D$  unabhängig identisch verteilt und  $W$  wie in (2.26).

Sobald  $d \geq 100 \cdot \frac{\nu^2}{\varepsilon^2} \ln\left(\frac{n}{\sqrt{\delta}}\right)$  gilt:

$$P(W \text{ ist eine } \varepsilon\text{-Isometrie auf } A) \geq 1 - 2\delta.$$

*Beweis.* Sei  $S \subset \mathbb{R}^d$  die Einheitssphäre und  $T \subset S$  definiert durch

$$T = \left\{ \frac{a-b}{\|a-b\|} \mid a, b \in A, a \neq b \right\}$$

Wir wollen beweisen, dass mit hoher Wahrscheinlichkeit gilt:

$$\max_{\alpha \in T} \left| \|W\alpha\|^2 - 1 \right| \leq \varepsilon.$$

Denn dann folgt wegen Linearität  $\forall a, b \in A$ :

$$\begin{aligned} & \left| \|Wa - Wb\|^2 - \|a - b\|^2 \right| \leq \varepsilon \|a - b\|^2 \\ \Rightarrow & (1 - \varepsilon) \|a - b\|^2 \leq \|Wa - Wb\|^2 \leq (1 + \varepsilon) \|a - b\|^2. \end{aligned}$$

Für jedes  $\alpha \in S$  und  $i \leq d$  gilt:

$$\begin{aligned}
E(\exp(\lambda \tilde{W}_i(\alpha))) &= E(\exp(\lambda \sum_{j=1}^D \alpha_j X_{ij})) \\
&= \prod_{j=1}^D E(\exp(\lambda \alpha_j X_{ij})) \leq \prod_{j=1}^D \exp\left(\frac{\lambda^2 \alpha_j^2 \nu}{2}\right) \\
&= \exp\left(\frac{\lambda^2 \nu}{2} \sum_{j=1}^D \alpha_j^2\right) = \exp\left(\frac{\lambda^2 \nu \|\alpha\|^2}{2}\right) = \exp\left(\frac{\lambda^2 \nu}{2}\right),
\end{aligned}$$

also  $\tilde{W}_i(\alpha) \in \mathcal{G}(\nu)$ .

Insbesondere impliziert Satz 2.9  $\forall k \in \mathbb{N}_{\geq 2} : E(\tilde{W}_i(\alpha)^{2k}) \leq \frac{k!}{2} (4\nu)^k$ .

Weiterhin sind alle Komponenten  $\tilde{W}_1(\alpha), \dots, \tilde{W}_d(\alpha)$  unabhängig  $\forall \alpha \in S$ .

Damit können wir die Bernstein-Ungleichung 2.22 anwenden:

$$\Rightarrow \forall \alpha \in T, t > 0 : P\left(\left|\sum_{i=1}^d (\tilde{W}_i(\alpha)^2 - 1)\right| \geq 4\nu\sqrt{2dt} + 4\nu t\right) \leq 2e^{-t}.$$

$$\begin{aligned}
&\Rightarrow P\left(\max_{\alpha \in T} \left|\sum_{i=1}^d (\tilde{W}_i(\alpha)^2 - 1)\right| \geq 4\nu\sqrt{2dt} + 4\nu t\right) \\
&= P\left(\bigcup_{\alpha \in T} \left\{\left|\sum_{i=1}^d (\tilde{W}_i(\alpha)^2 - 1)\right| \geq 4\nu\sqrt{2dt} + 4\nu t\right\}\right) \\
&\leq \sum_{\alpha \in T} P\left(\left|\sum_{i=1}^d (\tilde{W}_i(\alpha)^2 - 1)\right| \geq 4\nu\sqrt{2dt} + 4\nu t\right) \\
&\leq 2e^{-t}|T| \leq 2e^{-t}n^2
\end{aligned}$$

Wähle  $t = \ln\left(\frac{n^2}{\delta}\right) = 2 \ln\left(\frac{n}{\sqrt{\delta}}\right)$ , dann:

$$\begin{aligned}
&P\left(\max_{\alpha \in T} \left|\sum_{i=1}^d (\tilde{W}_i(\alpha)^2 - 1)\right| \geq 8\nu\sqrt{d \ln\left(\frac{n}{\sqrt{\delta}}\right)} + 8\nu \ln\left(\frac{n}{\sqrt{\delta}}\right)\right) \leq 2\delta \\
&\Leftrightarrow P\left(\max_{\alpha \in T} \left|\|W\alpha\|^2 - 1\right| \geq 8\nu\left(\sqrt{\frac{\ln(n/\sqrt{\delta})}{d}} + \frac{\ln(n/\delta)}{d}\right)\right) \leq 2\delta
\end{aligned}$$

Falls  $d \geq 100 \frac{\nu^2}{\varepsilon^2} \ln \left( \frac{n}{\sqrt{\delta}} \right)$  gewählt wird, folgt

$$8\nu \left[ \sqrt{\frac{\ln(n/\sqrt{\delta})}{d}} + \frac{\ln(n/\sqrt{\delta})}{d} \right] \leq \frac{4\varepsilon}{5} + \frac{2\varepsilon^2}{25\nu} \leq \frac{20\varepsilon}{25} + \frac{2\varepsilon}{25} \leq \varepsilon,$$

da nach Annahme  $\nu \geq 1$ . Also haben wir tatsächlich gezeigt:

$$P \left( \sup_{\alpha \in T} \left| \|W\alpha\|^2 - 1 \right| \leq \varepsilon \right) \geq 1 - 2\delta$$

□

Übung: Welche Konstanten verbessern sich, falls  $X_{ij} \sim \mathcal{N}(0, 1)$ ?

### 2.10. Assoziations- und Korrelationsungleichungen.

Dieser Abschnitt wird bei der Janson-Ungleichung und der Perkolationstheorie (optional) angewendet. Für das sonstige weitere Verständnis kann dieses Kapitel übersprungen werden.

Sind  $X, Y$  unabh. ZV, gelten praktische Rechenregeln für Produkte von  $X$  und  $Y$  oder von Funktionen davon. Unter welchen Annahmen kann man etwas von diesen Rechenregeln ‚retten‘, falls  $X$  und  $Y$  nicht unabh. sind?

**Satz 2.27** (Čebyšev-Assoziationsungleichung).

Seien  $f, g: \mathbb{R} \rightarrow \mathbb{R}$  isotone Funktionen, ZVen  $X \in \mathbb{R}$  und  $Y \in \mathbb{R}_+$ , so dass  $f, g \in \mathcal{L}^2(P_X)$ . Dann gilt

$$E(Y)E(Yf(X)g(X)) \geq E(Yf(X))E(Yg(X))$$

Bemerkung: Sei  $Y \equiv 1$ . Dann impliziert Satz 2.27:

$$E(f(X)g(X)) \geq E(f(X))E(g(X)) \Leftrightarrow E(f(X)g(X)) - E(f(X))E(g(X)) \geq 0$$

$$\stackrel{f, g \in \mathcal{L}^2}{\Leftrightarrow} \text{Cov}(f(X), g(X)) \geq 0.$$

Ist  $f$  antiton und  $g$  weiterhin isoton, gilt

$$E(Y)E(Yf(X)g(X)) \leq E(Yf(X))E(Yg(X)).$$

*Beweis.* Sei der Vektor  $(X', Y') : \Omega \rightarrow \mathbb{R} \times \mathbb{R}_+$  unabhängige Kopie von  $(X, Y)$  mit  $P_{(X, Y)} = P_{(X', Y')}$ . Sind  $f, g$  monoton wachsend, so gilt:

$$\begin{aligned}
& (f(X) - f(X'))(g(X) - g(X')) \geq 0 \\
& \Rightarrow YY'(f(X) - f(X'))(g(X) - g(X')) \geq 0 \\
& \Rightarrow E(YY'(f(X) - f(X'))(g(X) - g(X'))) \geq 0 \\
& \stackrel{\text{Linearität}}{\Leftrightarrow} E(YY'f(X)g(X) + YY'f(X')g(X')) \\
& \qquad \geq E(YY'f(X)g(X')) + E(YY'f(X')g(X)) \\
& \stackrel{\text{unabh.}}{\Leftrightarrow} E(Y')E(Yf(X)g(X)) + E(Y)E(Y'f(X')g(X')) \\
& \qquad \geq E(Yf(X))E(Y'g(X')) + E(Yg(X))E(Y'f(X')) \\
& \stackrel{\text{id. vert.}}{\Leftrightarrow} 2E(Y)E(Yf(X)g(X)) \geq 2E(Yf(X))E(Yg(X)).
\end{aligned}$$

□

Der Beweis illustriert die Methode der *unabhängigen Kopie*. Führe dazu unabhängig, identisch verteilte ZVen ein und im zweiten Schritt durch Erwartungswerte auf die ursprünglichen ZVen zurück.

Es gibt auch eine multivariate Variante dieser Ungleichung. Dazu:

**Definition 2.28.**  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  heißt *isoton* (oder *n-isoton*)  $:\Leftrightarrow$

$$\forall i \in \{1, \dots, n\} \text{ und } (y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_n) \in \mathbb{R}^{n-1} \text{ ist}$$

$$\mathbb{R} \ni x \mapsto f(y_1, \dots, y_{i-1}, x, y_{i+1}, \dots, y_n) \text{ monoton wachsend bzw. isoton,}$$

d.h.  $f$  isoton in jeder Koordinate.

$f$  heißt *antiton* (oder *n-antiton*), falls  $-f$  isoton ist.

**Satz 2.29** (Harris-Ungleichung).

Seien  $f, g: \mathbb{R}^n \rightarrow \mathbb{R}$  *n-isotone Funktionen* und *unabhängige ZVen*  $X_1, \dots, X_n \in \mathbb{R}$ . Setze  $X = (X_1, \dots, X_n) \in \mathbb{R}^n$ . Dann gilt

$$(2.27) \qquad E(f(X)g(X)) \geq E(f(X))E(g(X)).$$

*Bemerkung 2.30* (FKG-Ungleichung). Die Aussage (2.27) gilt auch im Fall  $n = \infty$ , d.h. falls  $X_k, k \in \mathbb{N}$ , eine Folge unabhängiger ZV,  $f, g: \mathbb{R}^{\mathbb{N}} \rightarrow \mathbb{R}$  isoton in jedem Argument sind und  $f(X), g(X)$  endliche Varianz besitzen. Diese Aussage wird mit Hilfe eines Martingal-Konvergenzsatzes bewiesen.

*Bemerkung 2.31.* Da die  $X_i$  unabhängig sind, integrieren wir bezüglich eines Produktmaßes. Dennoch ist es in solchen Situationen sinnvoll, bedingte Erwartungen zu benutzen, um Schreibarbeit zu sparen.

Illustration für  $n = 3$ : Sei  $f(X_1, X_2, X_3) = f(X)$ . Nach dem Trafo-Satz gilt:

$$\begin{aligned} E(f(X) | X_1) &= \int_{\mathbb{R}} \int_{\mathbb{R}} f(X_1, t, s) dP_{X_2}(t) dP_{X_3}(s), \\ E(f(X) | X_1, X_2) &= \int_{\mathbb{R}} f(X_1, X_2, s) dP_{X_3}(s) \text{ und} \\ E(E(f(X) | X_1, X_2)) &= \int_{\mathbb{R}} \int_{\mathbb{R}} \left( \int_{\mathbb{R}} f(r, t, s) dP_{X_3}(s) \right) dP_{X_2}(t) dP_{X_1}(r) \\ &= E(f(X)), \end{aligned}$$

was sich auch aus der Turmeigenschaft für bedingte Erwartungen herleiten lässt.

*Beweis.* Zeige (2.27) in Satz 2.29 durch vollständige Induktion.

Induktionsanker: Fall  $n = 1$  folgt aus Bemerkung zu Satz 2.27.

Induktionsschritt: Nehme an, dass (2.27) bekannt ist für alle  $n < k$ . Für das 1-dim. Integral bezüglich  $P_{X_k}$  gilt wieder wegen Satz 2.27:

$$E(f(X)g(X) | X_1, \dots, X_{k-1}) \geq E(f(X) | X_1, \dots, X_{k-1})E(g(X) | X_1, \dots, X_{k-1}),$$

denn bei eingefrorenen  $X_1, \dots, X_{k-1}$  sind  $f$  und  $g$  isotone Funktionen im  $k$ ten Argument. Turmeigenschaft und Monotonie des Integrals ergeben:

$$\begin{aligned} E(f(X)g(X)) &= E(E(f(X)g(X) | X_1, \dots, X_{k-1})) \\ &\geq E[E(f(X) | X_1, \dots, X_{k-1})E(g(X) | X_1, \dots, X_{k-1})] \end{aligned}$$

Nun ist

$$\begin{aligned} f_{k-1}: (X_1, \dots, X_{k-1}) &\mapsto E(f(X) | X_1, \dots, X_{k-1}) \\ &\stackrel{\text{unabh.}}{=} \int_{\mathbb{R}} f(X_1, \dots, X_{k-1}, t) dP_{X_k}(t) \quad (k-1)\text{-isoton,} \end{aligned}$$

ebenso wie

$$(X_1, \dots, X_k) \mapsto E(g(X) | X_1, \dots, X_{k-1}).$$

Wende Induktionsvoraussetzung (IV) an und erhalte:

$$\begin{aligned} &E \left[ \underbrace{E(f(X) | X_1, \dots, X_{k-1})}_{\text{unabh.}} E(g(X) | X_1, \dots, X_{k-1}) \right] \\ &\stackrel{\text{(IV)}}{\geq} E[E(f(X) | X_1, \dots, X_{k-1})] E[E(g(X) | X_1, \dots, X_{k-1})] \\ &\stackrel{\text{Turmeig.}}{=} E(f(X)) E(g(X)). \end{aligned}$$

□

*Bemerkung 2.32* (Frage/Überlegung:). Funktionieren solche Aussage auch für andere Klassen von Funktionen? Z.B.:

- sphärisch-symmetrische,
- radial-abfallende

Funktionen? Betrachte dazu folgenden Satz.

**Satz 2.33** (Gaußsche-Korrelationsungleichung).

Sei  $P = \mathcal{N}(0, C)$ , mit  $C$  positiv definit, ein zentriertes Gaußmaß auf  $\mathbb{R}^d$  und  $A, B$  zwei (bzgl. Spiegelung am Ursprung) symmetrische<sup>2</sup>, konvexe Mengen.

$$\Rightarrow P(A \cap B) \geq P(A)P(B).$$

Falls  $P(A) > 0$ , gilt das besser interpretierbare

$$P(B | A) \geq P(B)$$

Diesen Satz gibt es als Vermutung seit den 60ern motiviert durch zwei Arbeiten in 1955 und 1959. Beweis und Veröffentlichung von Thomas Royen im Jahr 2014. Seit 2017 gibt es weiteren, strukturierteren Beweis von Latala et. al..

**2.11. Anwendung der Harris-Ungleichung: Janson-Ungleichung.** Wir betrachten folgende kombinatorische Situation, die man in der Theoretischen Informatik, der probabilistischen Methode in der Kombinatorik und dem Studium von zufälligen Graphen antrifft.

Seien  $n \in \mathbb{N}$ ,  $[n] := \{1, \dots, n\}$ ,  $X_1, \dots, X_n \in \{0, 1\}$  unabh. ZV mit

$$p_k := P(X_k = 1) = 1 - P(X_k = 0) \text{ für } k \in [n] \text{ und } p_1, \dots, p_n \in [0, 1]$$

Für  $A \subset [n]$  setze

$$Y_A = \prod_{i \in A} X_i$$

und  $\mathcal{I} \subset \mathcal{P}([n])$

$$Z := \sum_{A \in \mathcal{I}} Y_A$$

was ein Polynom in den 0/1-Variablen  $X_1, \dots, X_n$  ist.

---

<sup>2</sup>d.h.  $x \in A \Rightarrow -x \in A$

*Bemerkung 2.34* (Übung). Seien  $f, g: \mathbb{R}^n \rightarrow \mathbb{R}$  zwei  $n$ -antitone Funktionen. Leiten Sie aus der Harris-Ungleichung

$$E(f(X)g(X)|Y_A = 1) \geq E(f(X)|Y_A = 1)E(g(X)|Y_A = 1)$$

her. Hierbei ist  $X = (X_1, \dots, X_n)$ .

Offensichtlich gilt für  $A, B \in \mathcal{I}$  mit  $A \cap B = \emptyset$

$$E(Y_A Y_B) = E\left(\prod_{i \in A} X_i \prod_{k \in B} X_k\right) = \prod_{i \in A} \prod_{k \in B} E(X_i)E(X_k) = E(Y_A)E(Y_B)$$

Daher reduziert sich die

$$\begin{aligned} \text{Var}(Z) &= E(Z^2) - E(Z)^2 = \sum_{A, B \in \mathcal{I}} E(Y_A Y_B) - \sum_{A, B \in \mathcal{I}} E(Y_A)E(Y_B) \\ (2.28) \quad &= \sum_{A, B \in \mathcal{I}, A \cap B \neq \emptyset} [E(Y_A Y_B) - E(Y_A)E(Y_B)] \leq \sum_{A, B \in \mathcal{I}, A \cap B \neq \emptyset} E(Y_A Y_B) =: \Delta \end{aligned}$$

Čebyšev liefert die symmetrische Schranke

$$P(|Z - EZ| > t) \leq \frac{\Delta}{t}$$

Auch wenn die Summanden von  $Z$  nicht unabhängig sind, gibt es eine (zumindest einseitige) exponentielle *tail*-Schranke.

$$Y_A = \prod_{i \in A} X_i, \quad Z = \sum_{A \in \mathcal{I}} Y_A, \quad \Delta = \sum_{A, B \in \mathcal{I}, A \cap B \neq \emptyset} E(Y_A Y_B)$$

**Satz 2.35.** Seien  $\mathcal{I} \subset \mathcal{P}([n])$  sowie  $Z$  und  $\Delta$  wie soeben definiert. Dann gilt für alle  $\lambda \leq 0$

$$\psi_{Z-EZ}(\lambda) \leq \phi\left(\frac{\lambda\Delta}{EZ}\right) \frac{(EZ)^2}{\Delta}, \quad \text{mit } \phi(x) = e^x - x - 1$$

Insbesondere gilt für  $t \in [0, EZ]$

$$P(Z - EZ < -t) \leq \exp\left(\frac{-t^2}{2\Delta}\right)$$

*Beweis.* Für die KEF von  $Z - EZ$  gilt

$$\psi'(\lambda) = \frac{E(Ze^{\lambda Z})}{E(e^{\lambda Z})} - E(Z) = \sum_{A \in \mathcal{I}} \frac{E(Y_A e^{\lambda Z})}{E(e^{\lambda Z})} - E(Z)$$

Wie wollen jeden  $A$ -Summanden auf der rechten Seite einzeln abschätzen.  
Zu jedem  $A \in \mathcal{I}$  setze

$$U_A = \sum_{B \in \mathcal{I}, A \cap B \neq \emptyset} Y_B$$

$$Z_A = \sum_{B \in \mathcal{I}, A \cap B = \emptyset} Y_B,$$

so dass  $Z = U_A + Z_A \geq Z_A$  für jedes  $A \in \mathcal{I}$ . Wegen der Fallunterscheidungsformel

$$E(Y_A e^{\lambda Z}) = E(Y_A e^{\lambda Z} | Y_A = 1) P(Y_A = 1) + 0 = E(e^{\lambda Z} | Y_A = 1) E(Y_A)$$

reicht es, die bedingte Erwartung von oben abzuschätzen. Da  $\lambda$  negativ ist, sind  $X \mapsto e^{\lambda U_A}$  und  $X \mapsto e^{\lambda Z_A}$  antitone Funktionen. Es folgt

$$(2.29) \quad \begin{aligned} E(e^{\lambda Z} | Y_A = 1) &= E(e^{\lambda U_A} e^{\lambda Z_A} | Y_A = 1) \\ &\stackrel{\text{(Harris)}}{\geq} E(e^{\lambda U_A} | Y_A = 1) E(e^{\lambda Z_A} | Y_A = 1) \\ &\stackrel{(Z_A, Y_A \text{ unabhängig})}{=} E(e^{\lambda U_A} | Y_A = 1) E(e^{\lambda Z_A}) \\ &\stackrel{(Z \geq Z_A)}{\geq} E(e^{\lambda U_A} | Y_A = 1) E(e^{\lambda Z}) \\ &\stackrel{\text{(Jensen)}}{\geq} e^{\lambda E(U_A | Y_A = 1)} E(e^{\lambda Z}) \end{aligned}$$

Damit erhalten wir

$$\begin{aligned} \frac{E(Z e^{\lambda Z})}{E(Z)} &= \sum_{A \in \mathcal{I}} \frac{E(Y_A e^{\lambda Z})}{E(Z)} = \sum_{A \in \mathcal{I}} \frac{E(e^{\lambda Z} | Y_A = 1) E(Y_A)}{E(Z)} \\ &\geq E(e^{\lambda Z}) \underbrace{\sum_{A \in \mathcal{I}} \frac{E(Y_A)}{E(Z)}}_{e^{\lambda E(U_A | Y_A = 1)}} \end{aligned}$$

$$\stackrel{\text{(Jensen)}}{\geq} E(e^{\lambda Z}) \exp \left[ \lambda \sum_{A \in \mathcal{I}} \frac{E(Y_A)}{E(Z)} E(U_A | Y_A = 1) \right]$$

Nun wollen wir die bedingte Erwartung loswerden.

$$\begin{aligned} \Delta &= \sum_{A, B \in \mathcal{I}, A \cap B \neq \emptyset} E(Y_A Y_B) = \sum_{A \in \mathcal{I}} E(Y_A U_A) \\ &= \sum_{A \in \mathcal{I}} E(Y_A U_A | Y_A = 1) E(Y_A) + E(Y_A U_A | Y_A = 0) P(Y_A = 0) \\ &= \sum_{A \in \mathcal{I}} E(U_A | Y_A = 1) E(Y_A) \end{aligned}$$

Umstellen ergibt

$$\frac{E(Ze^{\lambda Z})}{E(e^{\lambda Z})} \geq E(Z) \exp \left[ \lambda \frac{\Delta}{E(Z)} \right]$$

Damit folgt für die Ableitung der KEF

$$\frac{E(Ze^{\lambda Z})}{E(e^{\lambda Z})} - E(Z) \geq E(Z) \left[ \exp \left( \lambda \frac{\Delta}{E(Z)} \right) - 1 \right]$$

und wegen  $\psi(0) = 0$  für die KEF selbst :

$$\begin{aligned} \psi(\lambda) &= \psi(0) - \int_{\lambda}^0 \psi'(t) dt \leq -E(Z) \int_{\lambda}^0 (e^{t\Delta/E(Z)} - 1) dt \\ &= -E(Z) \int_{\lambda\Delta/E(Z)}^0 (e^s - 1) \frac{E(Z)}{\Delta} ds = -\frac{(EZ)^2}{\Delta} [e^s - s]_{\lambda\Delta/E(Z)}^0 \\ &= \frac{(EZ)^2}{\Delta} \phi \left( \frac{\lambda\Delta}{E(Z)} \right) \\ &= \frac{\lambda^2}{2} \Delta \end{aligned}$$

Somit folgt auch  $P(Z - EZ \leq -t) \leq e^{-t^2/(2\nu)}$  mit  $\nu = \Delta$ .  $\square$

*Bemerkung 2.36.* In Anwendungen will man oft zeigen, dass es Vektoren  $X = (X_1, \dots, X_n)$  gibt, die gewisse Klauseln erfüllen. Ja, man will sogar zeigen, dass dies für ein zufällig ausgewähltes  $X$  mit hoher W'keit zutrifft.

Hier gilt:

- $Y_A = 1 \Leftrightarrow A$ -te Bedingung ist erfüllt und
- $Z = \sum_{A \in \mathcal{I}} Y_A > 0 \Leftrightarrow$  mindestens eine der Klauseln in  $\mathcal{I}$  ist erfüllt.

Wir schätzen die W'keit des Komplements  $\{Z = 0\}$  von  $\{Z > 0\}$  ab (da ja  $Z$  nichtnegativ). Janson liefert:

$$P(Z = 0) = P(Z \leq EZ - EZ) \leq \exp \left( -\frac{(EZ)^2}{\Delta} \right).$$

Ob diese Abschätzung gut ist, ergibt sich aus der Abhängigkeit der Größen  $(EZ)^2$  und  $\Delta$  von den Modellparametern.

*Beispiel 2.37* (Dreiecke in zufälligen Graphen). Betrachte  $n$  Vertices und verbinde sie alle mit Kanten: Das sind dann  $m = \binom{n}{2}$  Stück. Der entstandene Graph heißt *vollständiger Graph auf/mit  $n$  Ecken*.

Seien  $X_1, \dots, X_m$  iid Bernoulli-ZV mit  $P(X_i = 1) = p \in [0, 1]$ . Falls  $X_i = 0$ , entfernen wir die entsprechende Kante aus dem Graphen. Das sich ergebende

Ensemble von Teilgraphen wir mit  $G(n, p)$  bezeichnet. Ein Element davon wird als *Erdős-Rényi-Graph* bezeichnet.

Er *enthält ein Dreieck*, falls es drei Ecken  $u, v, w$  im vollständigen Graphen gibt, so dass die Kanten-ZV  $X_{(u,v)}, X_{(v,w)}, X_{(w,u)}$  der entsprechenden Kanten  $(u, v), (v, w), (w, u)$  alle den Wert Eins annehmen, d.h. die drei Kanten ‚überleben‘ den Zufallsprozess.

Sei  $A = \{(u, v), (v, w), (w, u)\} \in \mathcal{P}([m])$ . Dann bedeutet  $Y_A = 1$ , dass das entsprechende Dreieck in dem Erdős-Rényi-Graphen vorhanden ist. Sei  $\mathcal{I} \subset \mathcal{P}([m])$  die Menge aller solchen  $A$ -s. Dann ist  $Z = \sum_{A \in \mathcal{I}} Y_A$  die Anzahl der Dreiecke im Erdős-Rényi-Graphen.

Man rechnet als Übung nach:

$$(2.30) \quad EZ = \binom{n}{3} p^3$$

$$(2.31) \quad \text{Var}(Z) = \binom{n}{3} (p^3 - p^6) + 2 \binom{n}{4} \binom{4}{2} (p^5 - p^6)$$

$$(2.32) \quad \Delta = \binom{n}{3} p^3 + 2 \binom{n}{4} \binom{4}{2} p^5$$

und mit Janson:

$$(2.33) \quad P(Z = 0) \leq \exp\left(-\frac{\binom{n}{3} p^3}{2(1 + 3np^2)}\right).$$

Beachte, dass  $\binom{n}{3} \sim n^3 p^2$  für großes  $n$ . Also wächst der Exponent bei festem  $p$  wie  $n^2$ .

**2.12. Anwendung der Harris-Ungleichung: Perkolation.** Gegeben:

- Gitter  $\mathbb{Z}^d$
- $p \in [0, 1]$
- Kante  $e$  wird mit Wahrscheinlichkeit  $p$  beibehalten bzw. mit Wahrscheinlichkeit  $(1 - p)$  entfernt, unabh. von allen anderen Kanten.
- Modell entspricht Produktmaß von Bernoulli-Verteilungen mit Parameter  $p$  bzgl.  $\mathbb{Z}^d$ , in Zeichen:  $\mathcal{B}_p^{\otimes \mathbb{Z}^d}$

Es entsteht ein zufälliger Untergraph auf  $\mathbb{Z}^d$  (vgl. Abbildung ?? im Fall  $d = 2$ ).

Man interessiert sich für die erzeugten Graphenstrukturen ohne die entfernten Kanten: Welche Eigenschaften besitzen die Cluster, d.h. die Zusammenhangskomponenten?

Das *0-1-Gesetz von Kolmogorov* liefert: Entweder

ABBILDUNG 3. **Graph.**

- (1) es existiert ein  $\infty$ -Cluster fast sicher  
oder
- (2) es existiert ein  $\infty$ -Cluster fast sicher nicht.

<sup>Monotonie</sup>  
 $\Rightarrow$  Es existiert ein kritischer Wert  $p_c \in (0, 1)$ , sodass

$$p > p_c \Rightarrow (1)$$

$$p < p_c \Rightarrow (2)$$

Bei  $p = p_c$  unklar, was passiert!

*Frage: Wir kann ich für einen passenden Wert  $p \in [0, 1]$  zeigen, dass f.s. ein unendlicher Cluster existiert?*

→ Ein möglicher Zugang über (immer größere) Würfel bzw. Rechtecke bzw. Quader  $\subset \mathbb{Z}^d$ . Unendliche Cluster sind bei endlichen Unterstrukturen natürlich nie möglich, also brauchen wir eine modifizierte Sichtweise.

Strategie in  $\mathbb{Z}^2$ : Mit welcher W'keit gibt es einen durchgängigen Pfad zwischen dem linken und rechten Rändern eines Rechtecks?

→ Studiere dieses Verhalten in Abhängigkeit der Länge und Breite des Rechtecks.

→ Lasse die Größe der Box dann gegen unendlich laufen, um asymptotisches Perkulationsverhalten zu verstehen.

**Was hat das mit der Harris-Ungleichung zu tun?**

Bei der Frage, ob ‚Perkolation‘ bereits auf einem endlichem großen Rechteck sichtbar ist (vgl. Abbildung ??) spielen Ereignisse der folgenden Bauart

$A = \{\text{Es gibt einen durchgängigen Weg von der linken zur rechten Kante.}\}$   
und

$B = \{\text{Es gibt einen durchgängigen Weg von der unteren zur oberen Kante.}\}$

eine Rolle (vgl. Abbildung ??). Wie stehen die Wahrscheinlichkeiten der Ereignisse  $A, B$  zueinander?

Klar ist: Man kann nicht mit Unabhängigkeit folgern:  $P(A \cap B) = P(A)P(B)$ ,

Wegen Monotonie liefert Anwendung der Harris-Ungleichung:

$$P(A \cap B) \geq P(A)P(B).$$

### Was hat das mit Konzentrationsungleichungen zu tun?

Russo-Seymour-Welsh Theorie: „Auf größeren Skalen sind Durchquerungswahrscheinlichkeiten größer“.

Mittels einer induktiven Zerlegung von grossen Rechtecken in kleinere kann man folgende Konzentrationsungleichung für unabhängige Kantenperkolati-on auf  $\mathbb{Z}^2$  zeigen.

#### Definition 2.38.

$$\Lambda_{L,n} := \mathbb{Z}^2 \cap ([-nL, nL] \times [-L, L])$$

$$\Omega_{L,n} := \{\exists \text{ aktiver Pfad, der linken und rechten Rand von } \Lambda_{L,n} \text{ verbindet}\},$$

$$R_{L,n}(p) := P_p(\Omega_{L,n}).$$

Ereigniss  $\Omega_{L,n}$  bzw. ZV  $\mathbb{1}_{\Omega_{L,n}} : \{0, 1\}^{\Lambda_{L,n}} \rightarrow \{0, 1\}$  ist isoton im Sinne von Definition (2.28) und dessen Wahrscheinlichkeit ist Polynom in  $p$ .

#### Lemma 2.39 (Reskalierung).

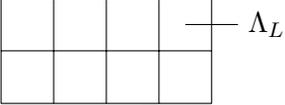
Sei  $p \in [0, 1]$ . Angenommen es existiert ein  $L \in \mathbb{N}$  und  $c \leq \frac{1}{16}$ , s.d.  $R_{L,2}(p) \geq 1 - c \cdot e^{-1}$ . Dann gilt  $\forall k \in \mathbb{N} : R_{2^k L, 2}(p) \geq 1 - c \cdot e^{-2^k}$

*Beweis.* Zerlege Rechteck  $\Lambda_{L,4}$  in vier Quadrate vom Typ  $\Lambda_L := \Lambda_{L,1} \Rightarrow$

$$\begin{aligned} R_{L,4}(p) &= \mathbb{P}_p \left( \text{Diagramm 1} \right) \geq \mathbb{P}_p \left( \text{Diagramm 2} \right) \\ &\stackrel{\text{Harris}}{\geq} \left( \mathbb{P}_p \left( \text{Diagramm 3} \right) \right)^3 \left( \mathbb{P}_p \left( \text{Diagramm 4} \right) \right)^2 \\ &= R_{L,2}(p)^3 R_L(p)^2 \end{aligned}$$

Wegen  $R_{L,2} = \mathbb{P}_p(\text{Diagramm 5}) \leq \mathbb{P}_p(\text{Diagramm 6}) = R_L^2$  gilt:

$$R_{L,4} \geq (R_{L,2})^4 \geq \left(1 - \frac{c}{e}\right)^4 \geq 1 - \frac{4c}{e}$$

Zerlege  $\Lambda_{2L,2}$  in zwei parallele Streifen oben und unten   $\Lambda_L$   
 Streifen disjunkt  $\Rightarrow$  Ereignisse im oberen Streifen unabhängig von Ereignissen im unteren.

$$R_{2L,2} \geq \mathbb{P}_p \left( \left( \begin{array}{|c|c|c|c|} \hline \text{wavy line} \\ \hline \end{array} \right) \text{ oder } \left( \begin{array}{|c|c|c|c|} \hline \text{wavy line} \\ \hline \end{array} \right) \right)$$

$\Rightarrow$  gehe zu Komplementen über, damit " $\cap$ " erscheint.

$$\Rightarrow 1 - R_{2L,2} \leq (1 - R_{L,4})^2 \leq \left(\frac{4c}{e}\right)^2 = \frac{16c^2}{e^2} \leq \frac{c}{e^2}$$

da  $c \leq \frac{1}{16}$ .

Fertig für  $k = 1$ . Weiter induktiv. □

2.13. **Negativ assoziierte ZV.** Betrachte reelwertige ZV  $X_i, i \in [n]$ .

**Definition 2.40** (Negative Assoziation). Die ZV  $X_i, i \in [n]$  heißen *negativ assoziiert*, falls für beliebige disjunkte  $I, J \subset [n]$  und für beliebige isotone  $f: \mathbb{R}^I \rightarrow \mathbb{R}, g: \mathbb{R}^J \rightarrow \mathbb{R}$

$$E(f(X_I)g(X_J)) \leq E(f(X_I))E(g(X_J))$$

gilt. Hierbei haben wir die Abkürzung  $X_I = (X_i)_{i \in I}$  benutzt.

Geben Sie ein Beispiel für neg. assoz. ZV an!

Welche Eigenschaften folgen aus dieser Definition?

*Bemerkung 2.41.* Seien die ZV  $X_i, i \in [n]$  negativ assoziiert. Dann:

- (1) Für jedes Paar  $i \neq j \in [n]$  gilt:  $E(X_i X_j) \leq E(X_i)E(X_j)$
- (2) Sind  $f_i, i \in [n]$ , isotone, nichtnegative Funktionen so gilt

$$E \left( \prod_{i \in [n]} f_i(X_i) \right) \leq \prod_{i \in [n]} E(f_i(X_i)),$$

- (3) insbesondere

$$E \left( \exp\left(\lambda \sum_{i \in [n]} X_i\right) \right) \leq \prod_{i \in [n]} E \left( e^{\lambda X_i} \right) \quad \text{für } \lambda \geq 0.$$

<sup>2</sup>See: Concentration of Measure for the Analysis of Randomised Algorithms, by: Devdatt P. Dubhashi and Alessandro Panconesi

Nun können wir Ideen der Chernoff und der Hoeffding-Schranke auf diese Situation übertragen.

**Lemma 2.42.** Seien  $a_i \leq b_i \in \mathbb{R}$ ,  $X_i \in [a_i, b_i]$ ,  $i \in [n]$ , zentrierte, negativ assoziierte ZV und  $S = \sum_{i=1}^n X_i$ . Dann gilt für  $\lambda \geq 0$

$$\psi_S(\lambda) \leq \frac{\lambda^2}{2} \sum_{i=1}^n \frac{(b_i - a_i)^2}{4}$$

$$P(S > t) \leq e^{-\frac{t^2}{2\nu}} \quad \text{für } t \geq 0.$$

*Beweis.*

$$(2.34) \quad \psi_S(\lambda) = \ln E \left( \exp \left( \lambda \sum_{i=1}^n X_i \right) \right)$$

$$(2.35) \quad \leq \ln \left( \prod_{i=1}^n E \left( e^{\lambda X_i} \right) \right)$$

$$(2.36) \quad = \sum_{i=1}^n \psi_{X_i}(\lambda)$$

$$(2.37) \quad \text{(Hoeffding)} = \frac{\lambda^2}{2} \sum_{i=1}^n \frac{(b_i - a_i)^2}{4}$$

also  $S \in \mathcal{G}(\nu)$  mit  $\nu = \sum_{i=1}^n \frac{(b_i - a_i)^2}{4}$ .

Mit der Chernoff-Methode und Argumenten wie in Bemerkung 2.8 folgt

$$P(S > t) \leq e^{-\frac{t^2}{2\nu}}$$

□

Fragen:

- Kann man die Annahmen in dem Lemma abschwächen?
- Falls  $X_1, \dots, X_n$  negativ assoz., sind dann auch die zentrierten Versionen  $X_1 - E(X_1), \dots, X_n - E(X_n)$  negativ assoz.?
- Falls  $X_1, \dots, X_n$  negativ assoz., sind dann auch die zentrierten Versionen  $-X_1, \dots, -X_n$  negativ assoz.?

#### 2.14. Minkowski-Ungleichung.

Die Minkowski-Ungleichung wird im Beweis der Bonami-Beckner-Ungleichung über Hyperkontraktivität verwendet. Bekannte Minkowski-Ungleichung:  $X, Y \in \mathcal{L}^q$ , dann gilt

$$E(|X + Y|^q)^{\frac{1}{q}} \leq E(|X|^q)^{\frac{1}{q}} + E(|Y|^q)^{\frac{1}{q}}.$$

**Satz 2.43** (Minkowski-Ungleichung).

Seien  $X : \Omega \rightarrow E, Y : \Omega \rightarrow F$  unabhängige ZVen und

$f : E \times F \rightarrow \mathbb{R}$ , messbar und  $Z = f(X, Y) : \Omega \rightarrow \mathbb{R}$   $P_X$ -f.s.  $P_Y$ -intbar

Für  $q \geq 1$  gilt (auch für  $q = \infty$ ):

$E_X(|E_Y(Z)|^q)^{\frac{1}{q}} \leq E_Y((E_X(|Z|^q))^{\frac{1}{q}})$ , wobei

$$E_X(Z) = E(Z | Y) = \int_{\mathbb{R}} f(t, Y) dP_X(t) \text{ die Integration bzgl. } P_X.$$

Also ist

$$\left( \int_{\mathbb{R}} \int_{\mathbb{R}} |f(t, s)|^q dP_X(t) dP_Y(s) \right)^{\frac{1}{q}} \leq \int_{\mathbb{R}} \left( \int_{\mathbb{R}} |f(t, s)|^q dP_X(t) \right)^{\frac{1}{q}} dP_Y(s).$$

Für den Fall  $q = 1$  werden lediglich Betragsstriche reingezogen.

Frage: Gibt es einen Zusammenhang zur klassischen Minkowski-Ungleichung?

Seien dazu

- $F = \{1, 2\}, P_Y(\{1\}) = P_Y(\{2\}) = \frac{1}{2}$
- $X = (X_1, X_2), f(X, 1) = X_1, f(X, 2) = X_2$

$$\begin{aligned} \Rightarrow \left( E_X \left( \left| \frac{1}{2}(X_1 + X_2) \right|^q \right) \right)^{\frac{1}{q}} &\leq \frac{1}{2} \left( (E_X(|X_1|^q))^{\frac{1}{q}} + (E_X(|X_2|^q))^{\frac{1}{q}} \right) \\ &\leq \frac{1}{2} \left( \int_{\mathbb{R}} |X_1(e)|^q dP_X(e) \right)^{\frac{1}{q}} + \frac{1}{2} \left( \int_{\mathbb{R}} |X_2(e)|^q dP_X(e) \right)^{\frac{1}{q}} \end{aligned}$$

*Beweis.* Für  $q = 1$ : Wende Fubini und Dreiecksungleichung an.

Für  $q = \infty$ :  $E_X(\text{ess-sup}(|E_Y(Z)|)) \leq E_Y(E_X(\text{ess-sup}|Z|))$ .

Nun  $q = (1, \infty)$ : o.E. sei  $Z \geq 0$ , daher  $|Z| = Z$ . Sei  $U$  unabhängige Kopie von  $Y$  und unabhängig von  $X$ .

$$\begin{aligned} E_X(E_Y(Z)^{q-1+1}) &= E_X [E_U(f(X, U))^{q-1} E_Y(f(X, Y))] \\ &= E_Y [E_X((E_U(f(X, U))^{q-1} f(X, Y)))] \\ &\leq E_Y \left[ (E_X(E_U(f(X, U)))^q)^{\frac{q-1}{q}} (E_X(f(X, Y)^q))^{\frac{1}{q}} \right] \\ &= [E_X(E_U(f(X, U)))^q]^{\frac{q-1}{q}} E_Y \left[ (E_X(f(X, Y)^q))^{\frac{1}{q}} \right] \\ &= [E_X((E_Y(Z))^q)]^{1-\frac{1}{q}} E_Y \left[ (E_X(Z^q))^{\frac{1}{q}} \right] \end{aligned}$$

Dividiere durch  $(E_X(E_Y(Z)^q))^{1-\frac{1}{q}}$ , dann folgt

$$(E_X(E_Y(Z)^q))^{\frac{1}{q}} \leq E_Y((E_X(Z^q))^{\frac{1}{q}})$$

□

### 3. SCHRANKEN AN DIE VARIANZ

Wir interessieren uns nun für ZVen der Form  $Z = f(X_1, \dots, X_n)$ , wobei  $X_1, \dots, X_n : \Omega \rightarrow \mathcal{X}$  unabhängige ZVen seien (nicht notwendig identisch verteilt) und  $\mathcal{X}$  ein messbarer Raum und  $f : \mathcal{X}^n \rightarrow \mathbb{R}$  eine messbare Abbildung. Ferner lautet die Generalannahme in diesem Kapitel  $Z \in \mathcal{L}^2$ , da wir an Abschätzungen für die Varianz interessiert sind.

Beispiel: Sei  $Z = \sum_{i=1}^n X_i$ .

In diesem Fall erzielt der Beweis des Gesetzes der großen Zahlen bereits eine Konzentrationsungleichung mithilfe der Čebyšev-Ungleichung. Jetzt möchten wir allgemeinere  $f$  betrachten. **Theorem 3.1 in Kapitel 3.1 findet in allen Unterkapitel von 3.2 bis 3.8 Verwendung. In Kapitel 3.9 wird ein alternativer Beweis für Theorem 3.1 angegeben. Die Kapitel 3.2 bis 3.8 lassen sich fast unabhängig voneinander lesen. Lediglich unter den Kapiteln 3.2, 3.3 und Bsp. 3.14 gibt es kleine Bezüge.**

#### 3.1. Efron-Stein-Ungleichung.

Zur Motivation betrachte:

$$Z = \sum_{i=1}^n X_i$$

Es gilt:  $X_1, \dots, X_n \in \mathcal{L}^2$  unabhängig  $\Rightarrow X_1, \dots, X_n$  unkorreliert  $\Rightarrow$

$$\text{Var}(Z) = \sum_{i=1}^n \text{Var}(X_i).$$

Idee um zu Verallgemeinern: Drücke  $Z - E(Z)$  für allgemeines  $Z = f(X_1, \dots, X_n) \in \mathcal{L}^2$  als Summe von *Martingaldifferenzen* bzgl. der Doob-Filtrierung aus. Notation dazu: Sei  $Y$  intbare ZVe und

$$E_i(Y) := E(Y \mid X_1, \dots, X_i) \text{ für } i = 1, \dots, n,$$

wobei  $E_0(Y) = E(Y)$ .

Setze  $\Delta_i := E_i(Z) - E_{i-1}(Z)$  für  $i = 1, \dots, n$ , sodass

$$Z - E(Z) = \sum_{i=1}^n \Delta_i \text{ (Teleskopsumme).}$$

$$\begin{aligned} \Rightarrow \text{Var}(Z) &= E((Z - E(Z))^2) = E\left(\left(\sum_{i=1}^n \Delta_i\right)^2\right) \\ &= \sum_{i=1}^n E(\Delta_i^2) + 2 \sum_{i=1}^n \sum_{j>i}^n E(\Delta_i \Delta_j). \end{aligned}$$

Für  $i \leq k$  gilt wegen der Turmeigenschaft:

$$\begin{aligned} E_i(E_k(Z)) &= E(E_k(Z) \mid X_1, \dots, X_i) = E[E(Z \mid X_1, \dots, X_k) \mid X_1, \dots, X_i] E(Z \mid X_1, \dots, X_i) \\ &= E_i(Z) \end{aligned}$$

Für  $i < j$  folgt damit:

$$\begin{aligned} E_i(\Delta_j) &= E_i(E_j(Z) - E_{j-1}(Z)) = 0 \\ \Rightarrow E_i(\Delta_i \Delta_j) &= E_i(\Delta_i(E_j Z - E_{j-1} Z)) = \Delta_i E_i(\Delta_j) = 0 \\ \Rightarrow E(E_i(\Delta_i \Delta_j)) &= E(\Delta_i E_i(\Delta_j)) = 0. \end{aligned}$$

Also gilt auch *ohne* jegliche Unabhängigkeit/Unkorreliertheit:

$$(3.1) \quad \text{Var}(Z) = E\left(\left(\sum_{i=1}^n \Delta_i\right)^2\right) = \sum_{i=1}^n E(\Delta_i^2).$$

Nimmt man die Unabhängigkeit der  $X_1, \dots, X_n$  an, folgt wie in Kapitel 2.9 bereits einmal:

$$E_i(Z) = \int_{\mathcal{X}^{n-i}} f(X_1, \dots, X_i, t_{i+1}, \dots, t_n) dP_{X_{i+1}}(t_{i+1}) \dots dP_{X_n}(t_n).$$

Ist  $X^{(i)} := (X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$ , so gilt analog

$$E^{(i)}(Z) := E(Z \mid X^{(i)}) = \int_{\mathcal{X}} f(X_1, \dots, X_{i-1}, x, X_{i+1}, \dots, X_n) dP_{X_i}(x).$$

Mit Fubini folgt

$$(3.2) \quad E_i(E^{(i)}(Z)) = E_{i-1}(Z).$$

Dies wird im folgenden Beweis benutzt

**Satz 3.1** (Efron-Stein-Ungleichung).

Seien  $X_1, \dots, X_n$  unabhängige ZVen,  $X = (X_1, \dots, X_n)$  und  $Z = f(X) \in \mathcal{L}^2$ .

$$\Rightarrow \text{Var}(Z) \leq \nu := \sum_{i=1}^n E \left( (Z - E^{(i)}(Z))^2 \right).$$

Ist  $X' = (X'_1, \dots, X'_n)$  unabhängige, identische Kopie zu  $X$  und  $Z'_i = f(X_1, \dots, X_{i-1}, X'_i, X_{i+1}, \dots, X_n)$ , so gilt

$$\begin{aligned} \nu &= \frac{1}{2} \sum_{i=1}^n E \left( (Z - Z'_i)^2 \right) = \sum_{i=1}^n E \left( (Z - Z'_i)_+^2 \right) \\ &= \sum_{i=1}^n E \left( (Z - Z'_i)_-^2 \right) = \inf_{Z_1, \dots, Z_n} \sum_{i=1}^n E \left( (Z - Z_i)^2 \right). \end{aligned}$$

wobei das Infimum über alle  $n$ -Tupel  $Z_1, \dots, Z_n$  läuft, wobei jedes  $Z_i$   $X^{(i)}$ -messbar und quadratintegrierbar ist.

*Beweis.* Mit (3.2) folgt direkt  $\Delta_i = E_i(Z - E^{(i)}(Z))$ . Verwende Jensen-Ungleichung für die bedingte Erwartung:

$$\begin{aligned} \Delta_i^2 &= (E_i(Z - E^{(i)}(Z)))^2 \leq E_i \left( (Z - E^{(i)}(Z))^2 \right) \\ \Rightarrow \text{Var}(Z) &= \sum_{i=1}^n E(\Delta_i^2) \leq \sum_{i=1}^n E(E_i(Z - E^{(i)}(Z))^2) = \sum_{i=1}^n E((Z - E^{(i)}(Z))^2). \end{aligned}$$

Nun sind die Gleichheiten für  $\nu$  noch zu zeigen:

Für eine ZVe  $Y$  setze

$$\text{Var}^{(i)}(Y) = E((Y - E(Y | X^{(i)}))^2 | X^{(i)}) = E^{(i)}((Y - E^{(i)}(Y))^2)$$

Turmeigenschaft

$$\Rightarrow \nu = \sum_{i=1}^n E(E^{(i)}((Z - E^{(i)}(Z))^2)) = \sum_{i=1}^n E(\text{Var}^{(i)}(Z)).$$

Für zwei unabhängig, identisch verteilte ZVen  $W, Y \in \mathcal{L}^2$  gilt:

$$\begin{aligned} \text{Var}(W) &= \frac{1}{2} \text{Var}(W) + \frac{1}{2} \text{Var}(Y) \\ &= \frac{1}{2} (E(W^2) - \underbrace{E(W)^2}_{=E(W)E(Y)} + E(Y^2) - \underbrace{E(Y)^2}_{=E(W)E(Y)}) \\ &= \frac{1}{2} (E(W^2) + E(Y^2) - 2E(Y)E(W)) \\ &= \frac{1}{2} E((W - Y)^2). \end{aligned}$$

Analog gilt:

$$\text{Var}^{(i)}(Z) = \frac{1}{2} E^{(i)}((Z - Z'_i)^2),$$

falls nur die 1-dim. Integration  $E^{(i)}(\cdot)$  ausgeführt und die Unabhängigkeit benutzt wird.

**Vorüberlegung::** Die Verteilung von  $Y - W$  ist symmetrisch, da für unabhängig, identisch verteilte ZVen  $W, Y \in \mathbb{R}$  gilt

$$P_{(W,Y)} = P_W \otimes P_Y \stackrel{\text{id. vert.}}{=} P_Y \otimes P_W = P_{(Y,W)}$$

Sei nun  $g(a, b) := a - b$ , dann sind

$$\begin{aligned} W - Y &= g(W, Y) \text{ und } Y - W = g(Y, W). \\ \Rightarrow P_{W-Y} &= P_{(W,Y)} \circ g^{-1} = P_{(Y,W)} \circ g^{-1} = P_{Y-W} \\ \Rightarrow \forall t \in \mathbb{R}: P(W - Y \geq t) &= P(Y - W \geq t). \end{aligned}$$

Unter der Annahme  $W, Y \in \mathcal{L}^2$  folgt:

$$\frac{1}{2} E((Y - W)^2) = \frac{1}{2} \int_{-\infty}^{\infty} t^2 dP_{Y-W}(t) = \int_0^{\infty} t^2 dP_{Y-W}(t) = \int_{-\infty}^0 t^2 dP_{Y-W}(t).$$

Bezüglich der 1-dim.  $E^{(i)}$ -Integration sind die ZVen  $Z$  und  $Z'_i$  unabhängig und identisch verteilt. Also:

$$\frac{1}{2} E^{(i)}((Z - Z'_i)^2) = E^{(i)}((Z - Z'_i)_+^2) = E^{(i)}((Z - Z'_i)_-^2).$$

Zur letzten Gleichheit: Zeige ähnlich wie oben die Aussage für eine reellwertige ZVe und folgere sie im multivariaten Fall.

$$\text{Es gilt für } Y \in \mathcal{L}^2 : \text{Var}(X) = E((X - E(X))^2) = \inf_{q \in \mathbb{R}} E((X - q)^2).$$

Das Infimum wird bei  $q = E(X)$  angenommen. Daher gilt für jedes  $i = 1, \dots, n$

$$\text{Var}^{(i)}(Z) = \inf_{q \in \mathbb{R}} E^{(i)}((Z - q)^2),$$

falls wir nur die 1-dim.  $E^{(i)}$ -Integration ausführen. Der Minimierer  $q_{\min}$  ist eine Funktion der  $(n-1)$  Variablen  $X^{(i)}$ . Wie oben ist der Minimierer gerade  $q_{\min} = E^{(i)}(Z)$ . Dieser ist  $X^{(i)}$ -messbar und in  $\mathcal{L}^2$ . Also darf dies bei der Klasse von Funktionen über das Infimum verwendet werden.  $\square$

Im Fall  $Z = \sum_{i=1}^n X_i$  gilt:

$$Z - E^{(i)}(Z) = \sum_{j=1}^n X_j - \left( \sum_{j \neq i} X_j + E(X_i) \right) = X_i - E(X_i).$$

$$\Rightarrow \nu = \sum_{i=1}^n E((X_i - E(X_i))^2) = \sum_{i=1}^n \text{Var}(X_i) \stackrel{\text{unabh.}}{=} \text{Var}(Z).$$

Beispiel zeigt: Efron-Stein-Ungleichung ist scharfe Abschätzung.

*Bemerkung 3.2* (Jackknife Resampling).

In der Statistik benutzt man Subpopulationen von Datensätzen, um Parameter besser zu schätzen. *Jackknife Resampling* ist eine Vorgängerversion des *Bootstraps*, bei dem aus einer Population von  $N$  Individuen (samples)  $N - 1$  Subpopulationen durch Weglassen eines Individuums erzeugt werden.

Seien  $X_1, \dots, X_n$  unabhängig, identisch verteilte ZVen mit Verteilung  $P_X$ . Der zu schätzende Parameter  $\theta$  ist eine Funktion von  $P_X$ , den wir durch eine Funktion  $Z = f_n(X_1, \dots, X_n)$  der Daten schätzen wollen. Um die Qualität des Schätzers zu beschreiben, verwendet man u.a.

Bias:  $E(Z) - \theta$

Mittlere quadratische Abweichung:  $\text{Var}(Z) + (E(Z) - \theta)^2$

Allerdings können wir  $E(Z)$  und  $\text{Var}(Z)$  nicht explizit ausrechnen, da  $P_X$  unbekannt ist. Also benutzt man eine empirische Mittlung anhand der Daten.

Jackknife-Schätzer für den Bias:

$$(n-1) \left( \frac{1}{n} \sum_{i=1}^n (Z_i - Z) \right), \text{ wobei } Z_i := f_{n-1}(x^{(i)}).$$

$X^{(i)}$  heißt auch  $i$ -tes Jackknife-Sample und  $Z_i$  die  $i$ -te Jackknife-Replikation.

Jackknife-Schätzer für die Varianz:

$$\frac{n-1}{n} \sum_{i=1}^n (Z - Z_i)^2$$

Die Efron-Stein-Ungleichung besagt, dass dieser Schätzer immer einen nicht-negativen Bias besitzt. D.h. wir schätzen die Varianz eher zu hoch als zu tief ein. Bei statistischen Verfahren sind wir damit eher auf der sicheren Seite.

Prinzip: Pessimismus > Optimismus (bei Varianzen)

Es werden in wachsender Allgemeinheit verschiedene Anwendungen der Efron-Stein Ungleichung vorgestellt.

### 3.2. Funktionen mit beschränkter Differenz.

**Definition 3.3.**  $f: \mathcal{X}^n \rightarrow \mathbb{R}$  hat *beschränkte Differenzen* (oder die beschränkte Differenzeigenschaft), falls es  $c_1, \dots, c_n \geq 0$  gibt, sodass

$$\sup_{x_1, \dots, x_n, x'_i \in \mathcal{X}} |f(x_1, \dots, x_n) - f(x_1, \dots, x'_i, \dots, x_n)| \leq c_i \forall i.$$

#### Korollar 3.4.

Hat  $f$  beschränkte Differenzen mit Konstanten  $c_1, \dots, c_n$  und sind  $X_1, \dots, X_n$  unabhängige ZVen mit Werten in  $\mathcal{X}$ , so gilt für  $Z := f(X_1, \dots, X_n)$ :

$$\text{Var}(Z) \leq \frac{1}{4} \sum_{i=1}^n c_i^2.$$

*Beweis.* ESU besagt

$$\text{Var}(Z) \leq \inf_{Z_1, \dots, Z_n} \sum_{i=1}^n E((Z - Z_i)^2).$$

In das Infimum dürfen wir die quadratintegrierbare und  $X^{(i)}$ -messbare 'Intervallmitte'

$$\begin{aligned} Z_i &= \frac{1}{2} \left( \sup_{x \in \mathcal{X}} f(X_1, \dots, X_{i-1}, x, X_{i+1}, \dots, X_n) \right. \\ &\quad \left. + \inf_{y \in \mathcal{X}} f(X_1, \dots, X_{i-1}, y, X_{i+1}, \dots, X_n) \right) \\ &\leq \frac{c_i}{2} \end{aligned}$$

einsetzen. Es folgt:

$$(Z - Z_i)^2 \leq \frac{c_i^2}{4} = (\text{ halbe Intervalllänge } )^2$$

□

*Beispiel 3.5* (Längste gemeinsame Teilfolge). Seien  $X_1, \dots, X_n, Y_1, \dots, Y_n \sim \text{Ber}(\frac{1}{2})$  unabhängige ZVen. Betrachte

$$Z := f(X_1, \dots, X_n, Y_1, \dots, Y_n)$$

$$:= \max\{k \mid X_{i_1} = Y_{j_1}, \dots, X_{i_k} = Y_{j_k}, 1 \leq i_1 < \dots < i_k \leq n, 1 \leq j_1 < \dots < j_k \leq n\}.$$

$Z$  ist die Länge der längsten Teilfolge, die in beiden Folgen auftaucht. Beispiel:

$$(i) \ X = (0, 0, 0, 0, 0), Y = (1, 1, 1, 1, 1) \Rightarrow Z = 0$$

$$(ii) \ X = (1, 0, 1, 0, 1), Y = (0, 1, 0, 1, 0) \Rightarrow Z = 4$$

Es ist bekannt, dass  $\lim_{n \rightarrow \infty} E(Z)/n$  f.s. gegen eine Konstante  $\gamma$  konvergiert, deren Wert allerdings unbekannt ist. Vermutung:

$$\gamma := \lim_{n \rightarrow \infty} \frac{E(Z)}{n} = \frac{2}{1 + \sqrt{2}}$$

Wie im Kontext vom üblichen Gesetz der großen Zahlen wollen wir das Konzentrationsphänomen über die  $\text{Var}(Z)$  verstehen. Ändert man nur eine einzige Ziffer in der Folge  $X_1, \dots, X_n, Y_1, \dots, Y_n$ , kann sich  $Z$  um höchstens Eins ändern, also hat  $Z$  beschränkte Differenzen mit  $c_{i,x} = c_{i,y} = 1$ .

$$\Rightarrow \text{Var}(Z) \leq \frac{1}{4} \sum_{i=1}^{2n} 1 = \frac{n}{2}$$

Beachte:  $EZ$  wächst proportional zu  $n$ . Die Čebyšev-Ungleichung liefert

$$P\left(\left|\frac{Z - E(Z)}{n}\right| \geq t\right) = P(|Z - E(Z)| \geq nt) \leq \frac{\frac{n}{2}}{n^2 t^2} = \frac{1}{2t^2 n}$$

bzw.

$$P(|Z - E(Z)| \geq s\sqrt{n}) \leq \frac{\frac{n}{2}}{ns^2} = \frac{1}{2s^2}$$

Mit hoher Wahrscheinlichkeit fällt  $Z$  in ein am Erwartungswert zentriertes Intervall der Länge  $c \times \sqrt{n}$ .

**Definition 3.6.** Eine Folge  $(Z_n)_{n \in \mathbb{N}}$  nichtnegativer ZVen heißt *relativ stabil*, wenn  $\frac{Z_n}{E(Z_n)} \rightarrow 1$  in Wahrscheinlichkeit (Fluktuationen von  $Z_n$  um  $E(Z_n)$  werden klein).

Um relative Stabilität z.z., ist es oft nützlich Schranken an  $\text{Var}(Z_n)$  z.z., denn

$$P\left(\left|\frac{Z_n}{E(Z_n)} - 1\right| \geq \varepsilon\right) = P(|Z_n - E(Z_n)| \geq \varepsilon E(Z_n)) \leq \frac{\text{Var}(Z_n)}{\varepsilon^2 E(Z_n)^2}.$$

*Beispiel 3.7.* (Schätzen von Dichten durch Glättungskerne)

Sei  $X_1, \dots, X_n$  unabhängig, identisch verteilt gemäß unbekannter Dichte  $\phi: \mathbb{R} \rightarrow [0, \infty)$ . Intuitiv würde man anhand von  $n$  Daten  $x_1, \dots, x_n$  die Verteilung von  $X_1$  durch das empirische Maß  $\mu_n := \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$  schätzen. Dieses Maß hat allerdings keine Dichte. Um dem Rechnung zu tragen, ersetzt man das Punktmaß durch eine *Approximation der Eins*, eine stark konzentrierte W'Dichte. Wir schätzen  $\phi$  durch

$$\phi_n(x) := \phi_{n;(x_1, \dots, x_n)}(x) := \frac{1}{nh_n} \sum_{i=1}^n K\left(\frac{x - x_i}{h_n}\right),$$

wobei  $h_n > 0$  *Glättungsparameter* heißt und  $K : \mathbb{R} \rightarrow [0, \infty)$ ,  $\int_{\mathbb{R}} K(x) dx = 1$  die *Approximation der Eins* ist. Wir interessieren uns für die  $\mathcal{L}^1$ -Norm des Fehlers, d.h.

$$Z = f(X_1, \dots, X_n) = \int_{\mathbb{R}} |\phi(x) - \phi_{n;(X_1, \dots, X_n)}(x)| dx.$$

Bei folgender Differenz kürzen sich  $2(n-1)$  Summanden weg:

$$\begin{aligned} |f(x_1, \dots, x_n) - f(x_1, \dots, x'_i, \dots, x_n)| &= \int_{\mathbb{R}} \left| |\phi(x) - \phi_n(x)| - |\phi(x) - \phi_n(x)| \right| dx \\ &\leq \frac{1}{nh_n} \int_{\mathbb{R}} \left| K\left(\frac{x-x_i}{h_n}\right) - K\left(\frac{x-x'_i}{h_n}\right) \right| \\ &\leq \frac{1}{nh_n} \int_{\mathbb{R}} \left( K\left(\frac{x-x_i}{h_n}\right) + K\left(\frac{x-x'_i}{h_n}\right) \right) = \frac{2}{n}. \end{aligned}$$

Mit Korollar 3.4 folgt unmittelbar

$$\text{Var}(Z) \leq \frac{n \cdot 2^2}{4 n^2} = \frac{1}{n}$$

*Beispiel 3.8.* (Supremum eines Rademacher-Prozesses)

Seien  $(\alpha_{i,t})_{i=1, \dots, n, t \in \mathcal{T}}$  reelle Zahlen und  $X_1, \dots, X_n$  unabhängig, identisch verteilte Rademacher-ZVen. Definiere

$$Z := \sup_{t \in \mathcal{T}} \sum_{i=1}^n \alpha_{i,t} X_i.$$

$Z$  heißt *Rademacher-Mittel*.  $E(Z)$  hängt stark von der Wahl der  $\alpha_{i,t}$  ab. Aber ändert man ein  $X_i$ , dann ändert sich  $Z$  um höchstens  $2 \sup_{t \in \mathcal{T}} |\alpha_{i,t}|$ . Korollar 3.4 impliziert dann:

$$\text{Var}(Z) \leq \sum_{i=1}^n \sup_{t \in \mathcal{T}} |\alpha_{i,t}|^2.$$

Dies Schranke kann mit der ESU noch verbessert werden.

Seien  $X'_1, \dots, X'_n$  unabhängige Kopien von  $X_1, \dots, X_n$  und definiere

$$Z'_i := \sup_{t \in \mathcal{T}} \left( \left( \sum_{j:j \neq i} X_j \alpha_{j,t} \right) + X'_i \alpha_{i,t} \right).$$

Sei  $t^*$  der (zufällige) Index in  $\mathcal{T}$ , sodass

$$\sum_{j=1}^n X_j \alpha_{j,t^*} = \sup_{t \in \mathcal{T}} \sum_{j=1}^n X_j \alpha_{j,t}.$$

Nach Definition des Supremus:

$$Z'_i = \sup_{t \in \mathcal{T}} \left( \left( \sum_{j:j \neq i} X_j \alpha_{j,t} \right) + X'_i \alpha_{i,t} \right) \geq \left( \sum_{j:j \neq i} X_j \alpha_{i,t^*} \right) + X'_i \alpha_{i,t^*}$$

Dann gilt für jedes  $i$ :

$$Z - Z'_i = \sum_{j=1}^n X_j \alpha_{j,t^*} - Z'_i \leq (X_i - X'_i) \alpha_{i,t^*}$$

Fallunterscheidung ergibt:

$$(Z - Z'_i)_+^2 \leq (X_i - X'_i)^2 \alpha_{i,t^*}^2.$$

Turmeigenschaft und Unabhängigkeit ergeben

$$\begin{aligned} E((Z - Z'_i)_+^2) &\leq E(E((X_i - X'_i)^2 \alpha_{i,t^*}^2 \mid X_1, \dots, X_n)) \\ &\leq E(\alpha_{i,t^*}^2 E(X_i^2 - 2X_i X'_i + X_i'^2 \mid X_1, \dots, X_n)) \\ &\leq E(\alpha_{i,t^*}^2 E(1 - 2X_i X'_i + 1 \mid X_1, \dots, X_n)) \\ &= 2E(\alpha_{i,t^*}^2) \end{aligned}$$

Da  $T^*$  nicht von dem Laufindex  $i$  abhängt, impliziert ESU:  $\text{Var}(Z) \leq 2E\left(\sum_{i=1}^n \alpha_{i,t^*}^2\right) \leq 2E\left(\sup_{t \in \mathcal{T}} \sum_{i=1}^n \alpha_{i,t}^2\right)$ . Nun steht das Supremum ausserhalb der Summe!

### 3.3. Selbstbeschränkende Funktionen.

**Definition 3.9.** Sei  $\mathcal{X} \neq \emptyset$  eine Menge und  $n \in \mathbb{N}$ . Eine Funktion  $f: \mathcal{X}^n \rightarrow \mathbb{R}_+$  heißt *selbstbeschränkend* oder *self-bounding*, falls  $\forall i \in \{1, \dots, n\} \exists f_i: \mathcal{X}^{n-1} \rightarrow \mathbb{R}$  mit  $\forall x_1, \dots, x_n \in \mathcal{X}$

(1)

$$0 \leq f(x_1, \dots, x_n) - f_i(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) \leq 1.$$

(Approximation der  $n$ -dim. Funktion mit  $(n-1)$ -dim. Funktion mit maximalem Fehler 1)

(2)  $\sum_{i=1}^n (f(x_1, \dots, x_n) - f_i(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)) \leq f(x_1, \dots, x_n)$   
( $\hat{=}$  self-bounding).

Welche Folgerungen ergeben sich?

Für  $i \in \{1, \dots, n\}$ ,  $x_1, \dots, x_n, x'_i \in \mathcal{X}$  gilt:

$$\begin{aligned} & |f(x_1, \dots, x_n) - f(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)| \\ & \leq |f(x_1, \dots, x_n) - f_i(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)| \\ & + |f_i(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) - f(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)| \\ & \leq 2, \text{ also hat } f \text{ beschränkte Differenzen.} \end{aligned}$$

Erwartungshaltung: Solche  $f$  liefern gute Konzentrationsungleichungen.

$$\begin{aligned} & \sum_{i=1}^n (f(x_1, \dots, x_n) - f_i(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n))^2 \\ & \stackrel{(1)}{\leq} \sum_{i=1}^n (f(x_1, \dots, x_n) - f_i(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)) \cdot 1 \\ & \stackrel{(2)}{\leq} f(x_1, \dots, x_n) \end{aligned}$$

Die vorherige Rechnung impliziert folgendes Korollar:

**Korollar 3.10.**

Seien  $X_1, \dots, X_n : \Omega \rightarrow \mathcal{X}$  unabhängige ZVen. Sei  $Z = f(X_1, \dots, X_n)$  für ein selbstbeschränkendes  $f : \mathcal{X}^n \rightarrow \mathbb{R}_+$ . Dann

$$\text{Var}(Z) \leq E(Z).$$

*Beweis.* Hier ist  $Z_i$  quadratintbar und  $X^{(i)}$  messbar:

$$\begin{aligned} \Rightarrow \text{Var}(Z) &\stackrel{\text{Satz 3.1}}{\leq} \sum_{i=1}^n \inf_{Z_1, \dots, Z_n} E((Z - Z_i)^2) \\ &\leq \sum_{i=1}^n E((Z - f_i(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n))^2) \\ &\stackrel{\text{s.o.}}{\leq} E(f(X_1, \dots, X_n)) = E(Z) \end{aligned}$$

□

→ self-bounding-Eigenschaft schließt schwere Ränder (Konzentration des Maßes an den Rändern) aus. Kontrast mit Rademacher ZV.

**Anwendung bei verschiedenen Modellen.** Hat für jedes  $n \in \mathbb{N}$   $Z(n) = f_n(X_1, \dots, X_n)$  die selbstbeschränkende Eigenschaft, so folgt  $\forall \varepsilon > 0$ :

$$\begin{aligned} P\left(\left|\frac{Z(n)}{E(Z(n))} - 1\right| > \varepsilon\right) &= P(|Z(n) - E(Z(n))| > \varepsilon E(Z(n))) \\ &\leq \frac{\text{Var}(Z(n))}{\varepsilon^2 E(Z(n))^2} \leq \frac{1}{\varepsilon^2 E(Z(n))}. \end{aligned}$$

Falls  $E(Z(n)) \rightarrow \infty$  für  $n \rightarrow \infty$ , so liegt relative Stabilität vor.

Es ergeben sich also zwei Vorgehensweisen, um die self-bounding property auszunutzen: Zwei Szenarien beim Abschätzen:

- $E(Z(n))$  von oben abschätzen.
- $\frac{1}{E(Z(n))}$  von oben abschätzen.

Greife eine Klasse von Funktionen heraus, die selbstbeschränkend sind:

**Definition 3.11.** Sei  $\mathcal{X} \neq \emptyset$  Menge,  $n \in \mathbb{N}$ ,  $\Pi_i \subset \mathcal{X}^i$  für  $i = 1, \dots, n$  und  $\Pi$  das  $n$ -Tupel:  $(\Pi_1, \dots, \Pi_n)$ . Für  $m \leq n$  sagen wir, dass der Vektor  $(x_1, \dots, x_m) \in \mathcal{X}^m$  die Eigenschaft  $\Pi$  hat  $:\Leftrightarrow (x_1, \dots, x_m) \in \Pi_m$ .

$\Pi$  heißt *erblich*  $:\Leftrightarrow$  Hat  $(x_1, \dots, x_n)$  die Eigenschaft  $\Pi$ , dann auch jeder Teilvektor  $(x_{i_1}, \dots, x_{i_k}) = (x_j)_{j \in I}$  mit  $I = \{i_1, \dots, i_k\} \subset \{1, \dots, n\}$

Ist  $\Pi$  eine erbliche Eigenschaft, so heißt die Funktion  $f$ , die jedem  $(x_1, \dots, x_m)$ ,

$m \leq n$ , die Mächtigkeit der längsten Teilfolge  $(x_{i_1}, \dots, x_{i_k})$  von  $(x_1, \dots, x_m)$  zuordnet, die die Eigenschaft II hat, die *Konfigurationsfunktion* von II. Eigentlich umfaßt  $f$  gleich  $n$  Funktionen

$$f_1: \mathcal{X} \rightarrow \mathbb{N}_0, f_2: \mathcal{X}^2 \rightarrow \mathbb{N}_0, \dots, f_n: \mathcal{X}^n \rightarrow \mathbb{N}_0$$

Bei der VC-Theorie in Kapitel und Bsp. 3.14 spielt die Konfigurationseigenschaft eine tragende Rolle.

**Korollar 3.12.**

Sei II eine erbliche Eigenschaft auf  $\mathcal{X} \neq \emptyset$ ,  $f$  die Konfigurationsfunktion von II,  $X_1, \dots, X_n$  unabhängige ZVen mit Werten in  $\mathcal{X}$  und  $Z = f(X_1, \dots, X_n)$ .

$$\Rightarrow \text{Var}(Z) \leq E(Z).$$

*Beweis.* Nach Korollar 3.10 genügt es zu zeigen, dass  $f$  selbst-beschränkend ist. Wegen Erbllichkeit gilt

$$Z_i = f_{n-1}(X^{(i)}) \leq f_n(X) \text{ und } Z \leq Z_i + 1 \Rightarrow 0 \leq Z - Z_i \leq 1$$

Andererseits: Gilt für ein  $(x_1, \dots, x_n) \in \mathcal{X}^n: Z = f(x_1, \dots, x_n) = k$ , so existiert nach Definition eine Teilfolge  $(x_{i_1}, \dots, x_{i_k})$  von  $(x_1, \dots, x_n)$  mit  $f(x_{i_1}, \dots, x_{i_k}) = k$ . Sei  $I = \{i_1, \dots, i_k\}$ . Für  $i \notin I$  gilt:  $Z = Z_i$ , denn beide haben  $(x_j)_{j \in I}$  als Teilfolge.

$$\Rightarrow \sum_{i=1}^n (Z - Z_i) = \sum_{i \in I} (Z - Z_i) \leq \sum_{i \in I} 1 = |I| = k = Z.$$

Also ist  $f$  selbst-beschränkend. □

Zusammenfassung:

$$\begin{aligned} f \text{ Konfigurationsfunktion} &\Rightarrow f \text{ ist self-bounding} \\ &\Rightarrow \text{Var}(Z) \leq E(Z) \text{ für } Z = f(X_1, \dots, X_n). \end{aligned}$$

*Beispiel 3.13* (Anzahl verschiedener Werte in einer Stichprobe).

Seien  $X_1, \dots, X_n : \Omega \rightarrow \mathbb{N}$  unabhängig, identisch verteilte ZVen mit Verteilung  $P(X_1 = k) = p_k \in [0, 1]$ , sodass  $\sum_{k \in \mathbb{N}} p_k = 1$ . Ferner betrachte

$$Z_n = \# \text{ unterschiedliche Werte in } (x_1, \dots, x_n).$$

Anders geschrieben heißt das

$$\begin{aligned}
 Z_n &= 1 + \sum_{i=2}^n \mathbb{1}_{\{x_i \neq x_1, \dots, x_i \neq x_{i-1}\}} \\
 &= \sum_{i=1}^n \mathbb{1}_{\{x_i \neq x_1, \dots, x_i \neq x_{i-1}\}} \\
 &= \sum_{i=1}^n \prod_{j=1}^{i-1} \mathbb{1}_{\{x_i \neq x_j\}}. \\
 \Rightarrow E(Z_n) &= \sum_{i=1}^n E \left( \prod_{j=1}^{i-1} \mathbb{1}_{\{X_i \neq X_j\}} \right) \\
 &= \sum_{i=1}^n \sum_{k=1}^{\infty} E(\mathbb{1}_{\{X_i \neq X_1\}} \dots \mathbb{1}_{\{X_i \neq X_{i-1}\}} \mid X_i = k) P(X_i = k) \\
 &= \sum_{i=1}^n \sum_{k=1}^{\infty} (1 - p_k)^{i-1} p_k
 \end{aligned}$$

Es gilt:

$$\lim_{n \rightarrow \infty} \frac{E(Z_n)}{n} = 0 \text{ (Übung),}$$

d.h. Wiederholungen sind nicht vernachlässigbar, bremsen Wachstum von  $Z_n$ .

Welche Konzentrationsungleichungen können wir für  $Z_n$  zeigen?

(A)  $Z_n = f(X_1, \dots, X_n)$  ist Funktion mit beschränkten Differenzen mit Schranken  $c_i = 1 \forall i = 1, \dots, n$ .

$$\stackrel{\text{Kor. 3.4}}{\Rightarrow} \text{Var}(Z_n) \leq \frac{n}{4} \Rightarrow \frac{\text{Var}(Z_n)}{4n} \text{ beschränkt}$$

$$\begin{aligned}
 P \left( \left| \frac{Z_n}{n} \right| > \varepsilon \right) &\leq \frac{E(Z_n)^2}{\varepsilon^2 n^2} \\
 &= \frac{1}{n^2 \varepsilon^2} (E((Z_n - E(Z_n))^2) + E(Z_n)^2) \\
 &= \frac{\text{Var}(Z_n)}{\varepsilon^2 n^2} + \frac{1}{\varepsilon^2} \left( \frac{E(Z_n)}{n} \right)^2 \\
 &\leq \frac{1}{4\varepsilon^2 n} + o(n) = o(n)
 \end{aligned}$$

Somit liegt stochastische Konvergenz gegen 0 vor:

$$\frac{Z_n}{n} \rightarrow 0.$$

Ist das optimal?

(B) Sei  $f(x_1, \dots, x_n)$  eine Konfigurationsfunktion zur Eigenschaft

$$(x_1, \dots, x_m) \in \Pi_m \Leftrightarrow \forall x_1, \dots, x_m \text{ unterschiedlich}$$

$\Pi$  ist erblich.

$$\begin{aligned} &\stackrel{\text{Kor. 3.12}}{\Rightarrow} \text{Var}(Z_n) \leq E(Z_n) \\ &\Rightarrow \frac{\text{Var}(Z_n)}{n} \leq \frac{E(Z_n)}{n}. \end{aligned}$$

Also ist der Quotient  $\left(\frac{\text{Var}(Z_n)}{n}\right)_{n \in \mathbb{N}}$  Nullfolge und nicht nur beschränkt.

(C) Kann man hier Janson-Ungleichung anwenden?

*Beispiel 3.14* (Vapnik-Chervonenkis-Dimension). aus der statistischen Lerntheorie. Sei  $\mathcal{X} \neq \emptyset$  beliebige Menge,  $\mathcal{A} \subset \mathcal{P}(\mathcal{X})$  (z.B.  $\sigma$ -Algebra) und  $H \subset \mathcal{X}$  endlich ( $\hat{=}$  sample). Sei  $G \subset H$ . Wir sagen

$$\begin{aligned} \mathcal{A} \text{ identifiziert } G &:\Leftrightarrow \exists A \in \mathcal{A} : A \cap H = G \\ s(\mathcal{A}, H) &:= \#\{G \subset H \mid \mathcal{A} \text{ identifiziert } G\} \\ &= \#\{A \cap H \mid A \in \mathcal{A}\} =: \#\text{Spur}_{\mathcal{A}}(H) \leq 2^{\#\mathcal{A}} \end{aligned}$$

$H$  ist vollständig identifiziert von  $\mathcal{A} :\Leftrightarrow s(\mathcal{A}, H) = 2^{\#H}$ ,

d.h. jede Teilmenge von  $H$  wird identifiziert. Betrachte *VC-Wachstumsfunktion* von  $\mathcal{A}$

$$n \mapsto s_n(\mathcal{A}) := \sup_{H \in \mathcal{F}_n} s(\mathcal{A}, H),$$

wobei  $\mathcal{F}_n = \{H \in \mathcal{P}(\mathcal{X}) \mid \#H = n\}$ . Die *VC-Dimension* von  $\mathcal{A}$  ist

$$D_{\mathcal{A}} := \sup\{n \in \mathbb{N} \mid s_n(\mathcal{A}) = 2^n\}.$$

Die *VC-Dimension* von  $\mathcal{A}$  bzgl.  $H$  ist

$$\begin{aligned} D_{\mathcal{A}}(H) &:= \sup\{n \in \mathbb{N} \mid \text{es existiert } G \subset H \text{ mit } \#G = n, \\ &\quad \text{das von } \mathcal{A} \text{ vollständig identifiziert wird}\}. \end{aligned}$$

Sei nun  $\mathcal{A}$  fixiert. Falls  $X_1, \dots, X_n : \Omega \rightarrow \mathcal{X}$  ZVen, so bezeichnen wir mit

$$D(X) = D(X_1, \dots, X_n) = D(H) \quad \text{mit} \quad H := \{X_1(\omega), \dots, X_n(\omega)\}.$$

Das ist eine Konfigurationsfunktion zur Eigenschaft II „vollständig identifizierbar“. Diese ist erblich. Unter der Annahme das  $X_1, \dots, X_n$  unabhängig sind, folgt also mit Korollar 3.12

$$\text{Var}(D(X)) \leq E(D(X)).$$

**3.4. Exkurs: Ursprung der VC-Theorie:** Seien  $F: \mathbb{R} \rightarrow [0, 1]$  eine Verteilungsfunktion und  $X_i \sim F$  für  $i = 1, \dots, n$ , unabhängig. Sei  $F_n(t) = \frac{1}{n} \#\{i \mid X_i \leq t\}$  die empirische Verteilungsfunktion von  $X_1, \dots, X_n$ .

$$(3.3) \quad \stackrel{\text{Hoeffding}}{\Rightarrow} \forall t \in \mathbb{R}, \varepsilon > 0, n \in \mathbb{N} : P(|F_n(t) - F(t)| \geq \varepsilon) \leq 2e^{-n\varepsilon^2/2}.$$

Satz von *Glivenko-Cantelli* dagegen besagt:

$$(3.4) \quad \forall n \in \mathbb{N}, \varepsilon > 0 : P(\sup_{t \in \mathbb{R}} |F_n(t) - F(t)| \geq \varepsilon) \rightarrow 0,$$

Man sieht, dass dies eine ‚Uniformisierung‘ von (3.3) ist. Als Einstieg in die weitere Diskussion ist es günstig, letztere Aussage zu

$$\forall n \in \mathbb{N}, \varepsilon > 0 : P(\|P_n - P\|_{\mathcal{A}} \geq \varepsilon) \rightarrow 0,$$

umzuschreiben, wobei

$$P_n(A) = \frac{1}{n} \sum_{j=1}^n \mathbb{1}_{\{X_j \in A\}}.$$

das *empirische Maß*,  $\|P_n - P\|_{\mathcal{A}} := \sup_{A \in \mathcal{A}} \|P_n(A) - P(A)\|$  und  $\mathcal{A}$  das Mengensystem  $\mathcal{A} := \{(-\infty, t] \mid t \in \mathbb{R}\}$  ist.

Es stellt sich die Frage, ob solch eine Konvergenz auch für andere Mengensysteme  $\mathcal{A}$  gilt? Z.B. für  $\mathcal{A} = \mathcal{B}(\mathbb{R})$ ? Durch welche Größe wird dann  $\|F_n - F\|_{\infty}$  ersetzt? Konvergiert die Differenz auch mit dem neuen Konvergenzbegriff gegen 0?

**Satz 3.15** (Theorem von Vapnik-Chervonenkis). *Sei  $\mathcal{A} \subset \mathcal{P}(\mathcal{X})$ . Für alle  $n \in \mathbb{N}$  und  $\varepsilon > 0$  gilt*

$$P(\sup_{A \in \mathcal{A}} |P_n(A) - P(A)| > \varepsilon) \leq 8 \underbrace{s_n(\mathcal{A})}_{\text{siehe oben}} e^{-\frac{n\varepsilon^2}{32}}$$

Allerdings: Was weiß man überhaupt über  $s_n(\mathcal{A})$ ?

**Satz 3.16** (Theorem von Sauer).

$$\forall n \geq D_{\mathcal{A}} \text{ gilt: } s_n(\mathcal{A}) \leq (n+1)^{D_{\mathcal{A}}}$$

Das Theorem von Sauer besagt: Sobald das maximal mögliche exponentielle Wachstum abbricht, ist es sogar nur noch polynomial und der Exponent ist die VC-Dimension.

Theoreme von Vapnik-Chervonenkis und Sauer kombiniert ergeben:

$$\forall n \geq D_{\mathcal{A}} : P(\|P_n - P\|_{\mathcal{A}} > \varepsilon) \leq 8s_n(\mathcal{A})e^{-\frac{n\varepsilon^2}{32}} \leq 8(n+1)^{D_{\mathcal{A}}}e^{-\frac{n\varepsilon^2}{32}}$$

*Bemerkung 3.17* (Besinnen uns nochmal): Was ist eigentlich der formale Rahmen vom Satz von Glivenko-Cantelli?

$$\|F_n - F\|_{\infty} = \sup_{t \in \mathbb{R}} |F_n(t) - F(t)| \xrightarrow[n \rightarrow \infty]{} 0 \text{ f.s.}$$

für i.i.d.  $X_1, \dots, X_n \sim F$  und die zugehörige empirische Verteilungsfunktion

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{(-\infty, X_i]}(t).$$

Frage: Fast sicher bezüglich welchem Maß? ( $P$  oder  $P_X$  oder  $P_n$ ?)

O.E. setzt man  $\Omega = \mathbb{R}$ ,  $\mathcal{A} = \mathcal{B}$ ,  $X = id$  und  $P = P_X$ .

Allgemeiner: Sei  $(S, \mathcal{S})$  Messraum,  $\mathcal{M}_1(S, \mathcal{S})$  Raum der W'maße auf  $(S, \mathcal{S})$  und i.i.d. ZVen  $X_1, \dots, X_n : \Omega \rightarrow S$  mit  $P_X$  als Verteilung.

(Häufig wird zur Vereinfachung  $\Omega = S$  und  $P = P_X$  gesetzt.)

empirisches Maß:  $P_n(A) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_A(X_i)$  für  $A \in \mathcal{S}$ .

Sei

$$\mathcal{C} \subset \mathcal{S}$$

oder

$$\mathcal{F} \subset \{f : S \rightarrow \mathbb{R} \text{ messbar, beschränkt}\}$$

und definiere ZV (falls messbar)

$$\|P_n - P_X\|_{\mathcal{C}} := \sup_{A \in \mathcal{C}} |P_n(A) - P_X(A)|$$

bzw.

$$\|P_n - P_X\|_{\mathcal{F}} := \sup_{A \in \mathcal{C}} |P_n(f) - P_X(f)|$$

$$\text{wobei } \mu(f) = \int f \, d\mu$$

Vorsicht! Nicht notwendigerweise messbar, da  $\mathcal{C}, \mathcal{F}$  überabzählbar sein kann. Man braucht im Allgemeinen Zusatzannahmen an  $\mathcal{C}$  bzw.  $\mathcal{F}$ , wie z.B. Separabilitätseigenschaften!

**Satz 3.18.** Sei  $\mathcal{C} \subset \mathcal{S}$  wie oben und  $P_X \in \mathcal{M}_1(S, \mathcal{S})$ . Dann sind äquivalent:

- (i)  $\|P_n - P_X\|_{\mathcal{C}} \xrightarrow[n \rightarrow \infty]{} 0$  f.s.
- (ii)  $\|P_n - P_X\|_{\mathcal{C}} \xrightarrow[n \rightarrow \infty]{} 0$  stoch.
- (iii)  $E(\|P_n - P_X\|_{\mathcal{C}}) \xrightarrow[n \rightarrow \infty]{} 0$ .

**Definition 3.19.** Wir definieren  $\mathcal{C} \subset \mathcal{S}$  als *Glivenko-Cantelli-Klasse* (GC-Klasse) für  $P_X$ , falls eine (und damit alle) der drei obigen Bedingungen erfüllt ist.

$\mathcal{C} \subset \mathcal{S}$  heißt *universelle GC-Klasse*  $:\Leftrightarrow \mathcal{C}$  ist eine GC-Klasse für jedes  $P_X \in \mathcal{M}_1$ .

$\mathcal{C} \subset \mathcal{S}$  heißt *uniforme GC-Klasse*  $:\Leftrightarrow \sup_{P_X \in \mathcal{M}_1} E(\|P_n - P_X\|_{\mathcal{C}}) \xrightarrow[n \rightarrow \infty]{} 0$  (Konvergenz glm. bzgl. Verteilungen).

$\mathcal{C} \subset \mathcal{S}$  heißt *VC-Klasse*  $:\Leftrightarrow$  VC-Dimension  $D_{\mathcal{C}} < \infty$ .

**Satz 3.20** (Vapnik-Chervonenkis).  $\mathcal{C} \subset \mathcal{S}$  ist genau dann eine *uniforme GC-Klasse*, wenn es eine *VC-Klasse* ist.

Bemerkenswert: Zusammenhang zwischen Wahrscheinlichkeits- und Maßtheorie zu Komplexitäts- und Mengentheorie. (Etwas anschaulich formuliert: Maßtheoretische Strukturen für Suprema von überabzählbaren Mengen instabil.)

*Beispiel 3.21.* Sei  $\mathcal{A} := \{F \sqsubset S\} = \{\text{endliche Teilmengen von } S\}$  und  $P$  ein W-Mass auf  $\mathcal{S}$  ohne Atome. Für  $A_n := \{X_1, \dots, X_n\}$  gilt  $P(A_n) = 0$  und  $P_n(A_n) = P_n^\omega(A_n^\omega) = 1$ .  $\mathcal{A}$  ist keine GC-Klasse für  $P$ .

*Beispiel 3.22.*  $S = \mathbb{R}$ ,  $\mathcal{C} = \{(-\infty, t] \mid t \in \mathbb{R}\}$ .  
GC-Theorem impliziert:  $\mathcal{C}$  ist GC-Klasse.

In der Statistik ist folgende Aussage im Kontext des Kolmogorov-Smirnov-Test wichtig:

**Satz 3.23** (Satz von Kolmogorov:).

$$\sup_{P_X \in \mathcal{M}_1(\mathbb{R}, \mathcal{B})} \|P_n - P_X\|_{\mathcal{C}} \sim \frac{1}{\sqrt{n}}$$

Insbesondere ist  $\mathcal{C}$  eine uniforme GC-Klasse, sogar mit expliziter Fehlerrate!

Verschärfung:

$$\sqrt{n} \|P_n - P_X\|_{\mathcal{C}} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \sup_{t \in [0,1]} |B(F(t))|,$$

wobei  $B$  die Brownsche Brücke ist.

Insbesondere gilt für stetiges  $F$  via Umparametrisierung durch das Simulationslemma/Quantiltransformation: Invarianzprinzip von Donsker für empirische Verteilungen

$$\sqrt{n} \|P_n - P_X\|_{\mathcal{C}} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \sup_{t \in [0,1]} |B(t)|.$$

**Definition 3.24.** Ist  $\mathcal{C} \subset \mathcal{S}$  und  $\mathcal{M} \subset \mathcal{M}_1(S, \mathcal{S})$ , so heißt  $(\mathcal{C}, \mathcal{M})$  ein GC-Paar, falls

$$E(\|P_n - P_X\|_{\mathcal{C}}) \xrightarrow[n \rightarrow \infty]{} 0$$

für jedes  $P_X \in \mathcal{M}$ .

Falls  $\mathcal{C}$  keine universelle GC-Klasse ist, macht es Sinn zu prüfen, ob man sich beider konkreten Anwendung auf eine Klasse  $\mathcal{M}$  von Maßen beschränken kann, so dass  $(\mathcal{C}, \mathcal{M})$  ein GC-Paar sind. Manchmal hat man ja a-priori Informationen über die in Frage kommenden Maße.

Analog für  $\mathcal{F} \subset \{f: \rightarrow \mathbb{R} \text{ messbar, beschränkt}\}$ .

*Beispiel 3.25* (keine universelle GC-Klasse). Seien  $X_j \sim \mathcal{N}(0, 1)$ ,  $j \in \mathbb{N}$ , i. i. d. standardnormalverteilte ZV und  $Y_j := -X_j$ . Betrachten die ZVektoren  $Z_j = (X_j, Y_j)^\top \in \mathbb{R}^2$  und deren empirische Verteilungen

$$P_n := \frac{1}{n} \sum_{j=1}^n \delta_{Z_j} = \frac{1}{n} \sum_{j=1}^n \delta_{(X_j, -X_j)}, n \in \mathbb{N}$$

Die Testfunktion

$$g := g_\omega := \mathbb{1}_{\{(x,y) \in \mathbb{R}^2 | x+y < 0\}} + \mathbb{1}_{\{(X_j(\omega), Y_j(\omega)) | j \in \mathbb{N}\}} : \mathbb{R}^2 \rightarrow [0, 1]$$

ist in jeder Koordinate antiton. Es gilt

$$P(g) - P_n(g) = P(g) - \frac{1}{n} \sum_{j=1}^n g(Z_j) = 0 - 1.$$

Insbesondere gilt für

$$\mathcal{F}_M := \{f: \mathbb{R}^2 \rightarrow \mathbb{R}, \text{ antiton in beiden Koordinaten und } \|f\|_\infty < M\}$$

(3.5)  $\|P - P_n\|_{\mathcal{F}} = \sup_{f \in \mathcal{F}} \|P_n(f) - P_X(f)\| = 1 \not\rightarrow 0.$

**Definition 3.26.** • Eine Teilmenge  $G \subset \mathbb{R}^2$  heißt *strikt monotoner Graph*, falls es eine strikt monotone Funktion  $g: \mathbb{R} \rightarrow \mathbb{R}$  gibt mit  $G = \{(x, g(x)) \mid x \in \mathbb{R}\}$ .

• Sei  $\mathcal{M} \subset \mathcal{M}_1(\mathbb{R}^2, \mathcal{B})$  die Menge der Maße  $P$  mit

$$(3.6) \quad P_c(G) = 0 \quad \text{für jeden strikt monotonen Graphen } G,$$

wobei  $P_c$  den kontinuierlichen Anteil von  $P$  bezeichnet, also

$$P_c := P - \sum_{\substack{x \in \mathbb{R}^2, \\ P(\{x\}) > 0}} \delta_x$$

**Satz 3.27** (De Hardt, Wright).

(a)  $(\mathcal{F}_M, \mathcal{M})$  sind ein GC-Paar.

(b) Gibt es für  $P \in \mathcal{M}_1$  einen strikt monotonen Graphen mit  $P_c(G) > 0$ , dann gilt  $\|P - P_n\|_{\mathcal{F}_M} \not\rightarrow 0$ .

(a) beschreibt eine hinreichende Bedingung für Konvergenz,

(b) eine notwendige.

*Bemerkung 3.28.* Satz zeigt, dass im höher dimensionalen Fall Phänomene auftreten, die im eindimensionalen unmöglich sind. Ab Dimension 2 bleibt alles aber im Prinzip gleich. Es gibt analoge Resultate für Dimension  $> 2$ .

Vergleiche mit Portemanteau-Theorem:

**Satz 3.29** (Portemanteau-Theorem). Seien  $n \in \mathbb{N}$ ,  $S$  topologischer Raum mit Borel-Sigma-Algebra  $\mathcal{S}$ , und  $\mu, \mu_n, n \in \mathbb{N}$ , Wahrscheinlichkeitsmaße auf  $(S, \mathcal{S})$ . Dann sind äquivalent:

(i)  $\mu_n \xrightarrow{w} \mu$ ;

(ii) für alle gleichmäßig stetigen und beschränkten  $f$  gilt:  $\int f d\mu_n \xrightarrow{n \rightarrow \infty} \int f d\mu$ ;

(iii) für alle Lipschitz-stetigen und beschränkten  $f$  gilt:  $\int f d\mu_n \xrightarrow{n \rightarrow \infty} \int f d\mu$ ;

(iv) für alle messbaren und beschränkten  $f$  mit  $\mu(\text{Punkte, an denen } f \text{ unstetig ist}) = 0$  gilt:  $\int f d\mu_n \xrightarrow{n \rightarrow \infty} \int f d\mu$ ;

(v) für alle abgeschlossenen Mengen  $A \subset S$  gilt:  $\limsup_{n \rightarrow \infty} \mu_n(A) \leq \mu(A)$ ;

(vi) für alle offenen Teilmengen  $U \subset S$  gilt:  $\liminf_{n \rightarrow \infty} \mu_n(U) \geq \mu(U)$ ;

(vii) für alle  $B \in \mathcal{B}(S)$  mit  $\mu(\partial B) = 0$  gilt:  $\lim_{n \rightarrow \infty} \mu_n(B) = \mu(B)$ .

Strikt monoter Graph  $G$  spielt die Rolle von  $\partial B$ .

### 3.5. Weitere Anwendungen von self-bounding.

*Beispiel 3.30* (Bedingte Rademacher-Mittelwerte). Seien  $n \in \mathbb{N}$ ,  $X_1, \dots, X_n$ ,  $\varepsilon_1, \dots, \varepsilon_n$  unabhängig mit  $X_i \in [-1, 1]^d$  i.i.d.,  $X_i = (X_{i1}, \dots, X_{id})^\top$  und  $\varepsilon_i \in \{-1, 1\}$  Rademacher-ZVen.

Folgende Größe wird in der Lerntheorie benutzt, um die Komplexität eines Modells zu beschreiben:

$$(3.7) \quad Z = E \left( \max_{j=1, \dots, d} \sum_{i=1}^n \varepsilon_i X_{ij} \mid X_1, \dots, X_n \right).$$

Aus dem *Faktorisierungslemma* folgt, dass es Abbildung  $f: [-1, 1]^d \rightarrow \mathbb{R}$  gibt, so dass

$$(3.8) \quad Z = f(X_1, \dots, X_n) \quad \text{ist.}$$

**Satz 3.31** (Faktorisierungslemma).

Seien  $(S, \mathcal{S})$  Messraum,  $\Omega \neq \emptyset$  und  $f: \Omega \rightarrow S$ ,  $g: \Omega \rightarrow \mathbb{R}$  Abbildungen.

Dann gilt:

$g$  ist  $\sigma(f)$ - $\mathcal{B}$ -messbar  $\Leftrightarrow$  Es existiert ein messbares  $\varphi: (S, \mathcal{S}) \rightarrow \mathbb{R}, \mathcal{B}(\mathbb{R})$  mit  $g = \varphi \circ f$ .

$Z$  hat beschränkte Differenzen:

Ersetzt man  $X_i$  durch unabhängige Kopie  $X'_i \in [-1, 1]^d$ , so kann sich  $\sum_{i=1}^n \varepsilon_i X_{ij}$  um höchstens 2 ändern (für  $j = 1, \dots, d$ ). Setze also  $c_1 = \dots = c_j := 2$ .

Korollar 3.4 impliziert:  $\text{Var}(Z) \leq n$ .

Die ZV  $Z$  ist aber auch selbst-beschränkend:

Um dies z.z., müssen wir neben  $f$  auch  $f_1, \dots, f_n$  identifizieren:

$$Z = f(X_1, \dots, X_n) \text{ wie oben ,}$$

$$Z_i = E \left( \max_{j=1, \dots, d} \sum_{k=1, k \neq i}^n \varepsilon_k X_{kj} \mid X^{(i)} \right) = f_i(X^{(i)}).$$

Als Übungsaufgabe zeigt man:

$$0 \leq Z - Z_i \leq 1 \text{ und } \sum_{i=1}^n (Z - Z_i) \leq Z$$

Mit Korollar 3.10 folgt  $\text{Var}(Z) \leq E(Z)$ .

In vielen Anwendungen gilt  $E(Z) \leq C\sqrt{n}$  und somit  $\text{Var}(Z) \leq C\sqrt{n}$ .

*Beispiel 3.32* (First passage percolation:). Sei  $(V, E)$  ein Graph mit (abzählbarer) Kantenmenge  $E = (e_i)_{i \in I}$  und uniform beschränkten Vertexgrad. Betrachte Kantengewichte

$$X_i: \Omega \rightarrow [0, \infty), \quad E(X_i^2) \leq \sigma^2, \quad i \in I \quad \text{unabhängige ZV}$$

Für einen Pfad  $\gamma$  zwischen Vertices  $x$  und  $y$  in  $V$  definiere Gewicht/Gesamtlänge als

$$\ell(\gamma) = \sum_{i \in I, e_i \subset \gamma} X_i$$

und den Abstand zwischen  $x$  und  $y$  als

$$(3.9) \quad Z := Z(x, y) := \inf\{\ell(\gamma) \mid \gamma \text{ verbindet } x \text{ und } y\}$$

Dies ist offensichtlich auch ein ZV. Ersetzt man im obigen Ausdruck  $X_i$  durch eine unabhängige Kopie  $X'_i$ , wird die entsprechende ZV mit  $Z'_i$  bezeichnet.

Stillschweigend nehmen wir im folgenden an, dass obige  $\inf$  tatsächlich  $\min$  sind. Dies ist z.B. der Fall, falls der Graph endlich viele Kanten enthält.  $Z(x, y)$  kann man als die kürzeste Zeit ansehen, die nötig ist, um von  $x$  nach  $y$  zu kommen oder umgekehrt. Daher der Name *first passage percolation*. Sei  $\gamma^*$  ein Pfad von  $x$  nach  $y$ , das das Infimum in (3.9) realisiert. Ersetzen wir das Gewicht  $X_i$  der  $i$ -ten Kante  $e_i$  durch  $X'_i$ , so ändert sich die Gesamtlänge des Pfades  $\gamma^*$  nur, falls  $e_i \subset \gamma^*$ , also

$$(Z_i - Z'_i)_- = (Z'_i - Z_i)_+ \leq (X'_i - X_i)_+$$

sowie

$$((Z'_i - Z_i)_+)^2 \leq ((X'_i - X_i)_+)^2 \mathbf{1}_{\{e_i \subset \gamma^*\}} \leq (X'_i)^2 \mathbf{1}_{\{e_i \subset \gamma^*\}}$$

Daher folgt unter Nutzung der Unabhängigkeit von  $X_i, X'_i$  aus der ESU:

$$\begin{aligned} \text{Var } Z &\leq \sum_{i \in I} E \left[ ((Z_i - Z'_i)_-)^2 \right] \\ &\leq \sum_{i \in I} E \left[ (X'_i)^2 \mathbf{1}_{\{e_i \subset \gamma^*\}} \right] \\ &\leq E \left[ (X'_i)^2 \right] E \left[ \sum_{i \in I} \mathbf{1}_{\{e_i \subset \gamma^*\}} \right] \\ &\leq \sigma^2 E \left[ \#\text{Kanten im optimalen Pfad } \gamma^* \right] \end{aligned}$$

Übung: Für den Graphen  $\mathbb{Z}^d$  und iid Gewichte  $X_i$  mit Werten in  $[a, b]$ , wobei  $0 < a < b < \infty$ , gilt:

$$\text{Var } Z(x, y) \leq \frac{b}{a} \|y - x\|_1 := \frac{b}{a} \sum_{j=1}^d |y_j - x_j|$$

Dies ist eine deterministische Schranke, die nur von der Verteilung der ZV und dem  $\ell^1$ -Abstand der betrachteten Punkte abhängt.

Tatsächlich wird  $\text{Var } Z(0, n \cdot e_1) \sim n^{2/3}$  vermutet.

*Beispiel 3.33* (Größter Eigenwert einer zufälligen symmetrischen Matrix).

Sei  $A := (X_{i,j})_{i,j}$  symmetrische Matrix mit Koeffizienten

$$\begin{aligned} X_{i,j}, (1 \leq i \leq j \leq n) & \quad \text{unabh. ZV,} \\ X_{i,j} \in [a, b], (1 \leq i \leq j \leq n) & \quad \text{und} \\ X_{i,j} = X_{j,i}, (1 \leq i \leq n, 1 \leq j \leq n). & \end{aligned}$$

Damit sind alle Eigenwerte von  $A$  reell. Sei  $Z := \lambda_{\max}(A)$  der maximale Eigenwert von  $A$  und  $v = (v_1, \dots, v_n) \in \mathbb{R}^n$  ein zugehöriger normierter Eigenvektor. Dann gilt

$$\lambda_{\max} = v^T A v = \sup_{u \in \mathbb{R}^n, \|u\|_2=1} u^T A u = \max_{u \in \mathbb{R}^n, \|u\|_2=1} u^T A u$$

Ersetzt man den Eintrag  $X_{i,j}$  durch unabhängig Kopie  $X'_{i,j}$  (und entsprechend  $X_{j,i}$ ) erhält man  $A'_{i,j}$  und  $Z'_{i,j} = \lambda_{\max}(A'_{i,j})$ . Dann ergeben sich

$$Z'_{i,j} = \max_{u \in \mathbb{R}^n, \|u\|_2=1} u^T A'_{i,j} u \geq v^T A'_{i,j} v$$

sowie

$$\begin{aligned} Z - Z'_{i,j} & \leq v^T A v - v^T A'_{i,j} v \\ & \leq 2 |v_i (X_{i,j} - X'_{i,j}) v_j| \\ & \leq 2(b-a) |v_i v_j| \end{aligned}$$

Um  $(Z - Z'_{i,j})_+$  abzuschätzen interessieren wir uns nur für den Fall  $0 \leq Z - Z'_{i,j} \leq 2(b-a) |v_i v_j|$ . Also folgt

$$\begin{aligned} \sum_{1 \leq i \leq j \leq n} ((Z - Z'_{i,j})_+)^2 & \leq 4(b-a)^2 \sum_{1 \leq i \leq j \leq n} |v_i v_j|^2 \\ & = 4(b-a)^2 \left( \sum_{1 \leq i \leq n} v_i^2 \right)^2 \\ & = 4(b-a)^2 \end{aligned}$$

da  $v$  normiert. Einsetzen in die ESU ergibt

$$\text{Var}(Z) \leq 4(b-a)^2$$

ein Schranke, die nicht von der Matrixgröße abhängt, und von der Verteilung der Einträge nur über den Durchmesser  $(b-a)$ . Die Einträge müssen unabh., nicht aber notwendig identisch verteilt sein.

### 3.6. Eine konvexe Poincaré-Ungleichung.

**Definition 3.34.**  $f: [0, 1]^n \rightarrow \mathbb{R}$  ist *separat konvex*, falls  $\forall i = 1, \dots, n$  :  $\forall x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n$  die Funktion

$$x_i \mapsto f(x_1, \dots, x_i, \dots, x_n) \text{ konvex ist.}$$

**Satz 3.35** (Konvexe Poincaré-Ungleichung).

Seien  $X_1, \dots, X_n \in [0, 1]$  unabhängige ZVen. Sei  $f: [0, 1]^n \rightarrow \mathbb{R}$  separat konvex (und partiell differenzierbar). Dann gilt für  $Z = f(X) = f(X_1, \dots, X_n)$

$$\text{Var}(f(X)) \leq E(\|\nabla f(X)\|^2).$$

Prinzip:  $\mathcal{L}^2$ -Norm von  $f$  durch  $\mathcal{L}^2$ -Norm von  $\nabla f$  abschätzen.

Beachte: Konvexität impliziert bereits Differenzierbarkeit fast überall. Die Aussage lässt sich abschwächen für schwache Gradienten mit den schwachen partiellen Ableitungen.

Ähnliche Ungleichungen werden in Kapitel 3.8 und 5.3 bewiesen. Diese sind unabhängig von Satz 3.35.

*Beweis.* Nutze variationelle Version der Darstellung von  $\nu$  in der Efron-Stein-Ungleichung

$$\text{Var}(Z) \leq \sum_{i=1}^n E[(Z - Z_i)^2]$$

wobei  $Z_i = \inf \tilde{Z}_i$  über alle  $(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$ -messbaren  $\mathcal{L}^2$ -ZVen  $\tilde{Z}_i$  läuft. Insbesondere ist  $\tilde{Z}_i = \inf_{x_i \in [0,1]} f(X_1, \dots, x_i, \dots, X_n)$  eine zulässige Wahl. Sei  $x' = x_i(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$  ein Minimierer des Infimums (Existenz klar wegen Stetigkeit und Kompaktum).

Sei  $\widetilde{X}^{(i)} = (X_1, \dots, X_{i-1}, x', X_{i+1}, \dots, X_n)$ .

$$\begin{aligned} \Rightarrow \sum_{i=1}^n (Z - Z_i)^2 &= \sum_{i=1}^n (f(X) - f(\widetilde{X}^{(i)}))^2 \\ &\stackrel{\text{konvex}}{\leq} \sum_{i=1}^n \left( \frac{\partial f(X)}{\partial x_i} \right)^2 \underbrace{(X - \widetilde{X}^{(i)})^2}_{\leq 1} \leq \sum_{i=1}^n \left( \frac{\partial f(X)}{\partial x_i} \right)^2 = \|\nabla f(X)\|^2 \end{aligned}$$

ABBILDUNG 4. **Konvexität**

□

*Beispiel 3.36* (Größter Singulärwert einer zufälligen Matrix). Sei  $A \in \mathbb{R}^{m \times n}$  Matrix mit unabhängigen zufälligen Koeffizienten  $X_{ij} \in [0, 1]$ . Ähnlich wie bei einer symmetrischen Matrix in Beispiel 3.33 wollen wir die Konzentration des größten Singulärwerts  $Z$  von  $A$  untersuchen. Dieser ist gegeben durch

$$Z = Z(A) = Z((X_{ij})_{ij}) = \sqrt{\lambda_{\max}(A^T A)}$$

Da  $A^T A$  symmetrisch ist, ist  $\lambda_{\max}(A^T A)$  wohldefiniert und reell. Wieder haben wir eine variationelle Darstellung

$$Z = \sqrt{\sup_{u \in \mathbb{R}^n, \|u\|_2=1} u^T A^T A u} = \sqrt{\max_{u \in \mathbb{R}^n, \|u\|_2=1} \|Au\|_2} =: \|A\|$$

wobei  $\|A\|$  *Operator-* oder *Spektralnorm* von  $A$  genannt wird.

Für festes  $u \in \mathbb{R}^n$  ist

$$[0, 1] \ni X_{ij} \mapsto f_{X,u}(X_{ij}) := \|Au\| = \sqrt{\sum_{l=1}^m \left( \sum_{k=1}^n X_{l,k} u_k \right)^2}$$

eine konvexe Funktion von  $X_{ij}$ . Als Maximum konvexer Funktionen ist

$$[0, 1] \ni X_{ij} \mapsto \max_{u \in \mathbb{R}^n, \|u\|_2=1} f_{X,u}(X_{ij})$$

ebenfalls konvex. Nun verwenden wir einen Satz aus der Linearen Algebra der die Störungstheorie von Matrizen betrifft.

**Satz 3.37** (Satz von Liidiski). *Seien  $A$  und  $B \in \mathbb{R}^{m \times n}$  mit Singulärwerten*

$$s_{\max}(M) = s_1(M) \geq s_2(M) \geq \dots \geq s_n(M), \quad M \in \{A, B\}$$

Dann gilt:

$$\begin{aligned}
 (s_1(B) - s_1(A))^2 &\leq \sum_{i=1}^n (s_i(B) - s_i(A))^2 \\
 &\leq \sum_{i=1}^n s_i (B - A)^2 \\
 &= \text{Spur}((B - A)^T (B - A)) \\
 &= \sum_{l=1}^m \sum_{l=1}^n (b_{k,l} - a_{k,l})^2
 \end{aligned}$$

was drei verschiedenen Darstellungen der Hilbert-Schmidt-Norm (oder Frobenius-Norm) zum Quadrat sind.

Gilt nun

$$b_{k,l} = a_{k,l} + \varepsilon \delta_{k,i} \delta_{l,j}$$

so folgt

$$\frac{(s_1(B) - s_1(A))^2}{\varepsilon^2} \leq 1$$

und damit ist  $s_1(A)$  als Funktion  $X_{ij}$  Lipschitz-stetig mit Lipschitz-Konstante gleich Eins. Daraus folgt, dass für Lebesgue-fast alle Werte in  $[0, 1]$

$$a_{ij} \mapsto s_1(A)$$

differenzierbar ist mit  $\left| \frac{\partial}{\partial a_{ij}} s_1(A) \right| \leq 1$ . In der Tat reicht diese abgeschwächte Bedingung, um eine konvexe Poincare-Ungleichung zu beweisen (Übung).

Da wir den Gradienten einer Funktion auf dem  $\mathbb{R}^{m \cdot n}$  untersuchen müssen, erweitern wir unsere 1-dimensionalen Überlegungen. Satz von Liidski impliziert

$$(s_1(B) - s_1(A)) \leq \sqrt{\sum_{l=1}^m \sum_{l=1}^n (b_{k,l} - a_{k,l})^2} = \|B - A\|_{\mathbb{R}^{m \cdot n}}$$

wobei im letzten Ausdruck die Euklidische Norm auf  $\mathbb{R}^{m \cdot n}$  sowie  $A$  und  $B$  als Vektoren in  $\mathbb{R}^{m \cdot n}$  betrachtet werden. Mit anderen Worten ist

$$\mathbb{R}^{m \cdot n} \ni A \mapsto s_1(A)$$

Lipschitz-stetig mit Lipschitz-Konstante gleich Eins, woraus man (als Übung) folgert, das (Lebesgue fast überall)

$$\|\nabla f(X)\| = \left\| \left( \frac{\partial f}{\partial x_{ij}}(X) \right)_{ij} \right\| \leq 1$$

Somit folgt aus der konvexen Poincare-Ungleichung

$$\text{Var}(Z) = \text{Var}(f(X)) \leq E(\|\nabla f(X)\|^2) \leq E(1) = 1$$

Dies ist wieder von den Dimensionen  $n$  und  $m$  unabhängig!

**3.7. Anwendung der Efron-Stein-Ungleichung auf Tail-Events.** Bei der allgemeinen Markov-Ungleichung kann man statt der quadratischen — sofern exponentielle Momente existieren — auch exponentielle Funktionen einsetzen, um schärfere Schranken an das Abfallverhalten zu bekommen. In diesem Sinne wollen wir ESU verbessern! Hier nehmen wir für die Funktion  $f: \mathcal{X}^n \rightarrow \mathbb{R}$  etwas weniger als die beschränkte Differenz-Bedingung an: Es existiere  $\nu > 0$  mit

$$(3.10) \quad \sum_{i=1}^n ((Z - Z'_i)_+)^2 \leq \nu \text{ f.s., wobei wieder} \\ Z'_i = f(X_1, \dots, X_{i-1}, X'_i, X_{i+1}, \dots, X_n).$$

wobei  $X'_i$  eine unabh. Kopie von  $X_i$  ist. Es liegt also eine Art *kummulative* oder *gemittelte* beschränkte Differenz-Bedingung vor.

Beispiel:  $\lambda_{\max}$  ist EW von symm. Matrix  $\Rightarrow$  (3.10) gilt mit  $\nu = 4(b-a)^2$ . Insbesondere ist die Konstante unabhängig von der Konfiguration  $x \in \mathcal{X}^n$  und sogar von der Dimension  $n$ . Dann folgt mit ESU:

$$\text{Var}(Z) \leq E\left(\sum_{i=1}^n ((Z - Z'_i)_+)^2\right) \leq \nu.$$

Für beliebige  $\alpha \in (0, 1)$  sei  $Q_\alpha$  das  $\alpha$ -Quantil von  $Z = f(X) = f(X_1, \dots, X_n)$ , d.h.:

$$Q_\alpha := \inf\{t \in \mathbb{R} \mid P(Z \leq t) \geq \alpha\} = \inf\{t \in \mathbb{R} \mid F_Z(t) \geq \alpha\}.$$

Insbesondere ist  $\mathcal{M}(Z) = Q_{\frac{1}{2}}$  der Median. Zu gegebener Funktion  $f: \mathcal{X}^n \rightarrow \mathbb{R}$  betrachte Modifizierung  $g = g_{a,b}: \mathcal{X}^n \rightarrow \mathbb{R}$  für  $a < b \in \mathbb{R}$ :

$$g(x) = \begin{cases} b & f(x) \geq b \\ f(x) & f(x) \in (a, b) \\ a & f(x) \leq a \end{cases}$$

Betrachte den Fall  $\mathcal{M}(Z) \leq a$ , dann gilt:  $p := P(Z \leq a) \geq \frac{1}{2}$ . Also ist

$$\begin{aligned}
 E(g(X)) &= E(g(X)\mathbf{1}_{\{f(X) \leq a\}}) + E(g(X)\mathbf{1}_{\{f(X) > a\}}) \\
 &\leq aE(\mathbf{1}_{\{f(X) \leq a\}}) + bE(\mathbf{1}_{\{f(X) > a\}}) \\
 &= aP(Z \leq a) + bP(Z > a) \\
 &= ap + b(1 - p) = b + p \underbrace{(a - b)}_{\text{negativ}} \\
 &\leq b + \frac{1}{2}(a - b) = \frac{a + b}{2}.
 \end{aligned}$$

Dann folgt für  $Y = g(X)$ :

$$\begin{aligned}
 \text{Var}(Y) &= E((Y - E(Y))^2) \geq E((Y - E(Y))^2 \mathbf{1}_{\{Y \geq b\}}) \\
 &\geq (b - E(Y))^2 E(\mathbf{1}_{\{Y \geq b\}}) \geq \left(b - \frac{a + b}{2}\right)^2 P(Y \geq b) \\
 &= \frac{(b - a)^2}{4} P(Z \geq b).
 \end{aligned}$$

Komplementäre Abschätzung der  $\text{Var}(Y)$  von oben mit ESU. Zu einem Punkt  $x \in \mathcal{X}^n$  setze

$$\tilde{x}^{(i)} = (x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n).$$

Wir nutzen

$$(3.11) \quad f(x) \leq a \quad \Rightarrow \quad a = g(x) \leq g(\tilde{x}^{(i)}).$$

in der folgenden Ungleichungskette

$$\begin{aligned}
 \text{Var}(g(X)) &\leq \sum_{i=1}^n E\left((g(X) - g(\tilde{X}^{(i)}))^2\right) \stackrel{\text{Symmetrie}}{=} 2 \sum_{i=1}^n E\left((g(X) - g(\tilde{X}^{(i)})_+)^2\right) \\
 &\stackrel{(3.11)}{=} 2E\left(\sum_{i=1}^n \mathbf{1}_{\{f(X) > a\}} (g(X) - g(\tilde{X}^{(i)})_+)^2\right) \\
 &\leq 2\nu P(Z > a).
 \end{aligned}$$

wegen der Annahme (3.10). Kombiniere untere und obere Schranke:

$$\begin{aligned}
 \frac{(b - a)^2}{4} P(Z \geq b) &\leq \text{Var}(Y) \leq 2\nu P(Z > a) \\
 \Rightarrow b - a &\leq \sqrt{8\nu \frac{P(Z > a)}{P(Z \geq b)}}
 \end{aligned}$$

Ziel: Schätze Abstände zwischen Quantilen ab:

Seien  $0 \leq \delta \leq \gamma \leq \frac{1}{2}$ . Setze:

$$a = Q_{1-\gamma}, b = Q_{1-\delta} \quad \text{insbesondere} \quad P(Z > a) \leq \gamma, P(Z > b) \leq \delta.$$

Aus dem gezeigten folgt für die Abstände zwischen Quantilen rechts von  $\mathcal{M}(Z)$ :

$$Q_{1-\delta} - Q_{1-\gamma} \leq \sqrt{\frac{8\nu\gamma}{\delta}}.$$

Um das *tail*-Verhalten zu verstehen, ist es sinnvoll  $\gamma = 2^{-k} > \delta = 2^{-(k+1)}$  für  $k \in \mathbb{N}$  zu setzen. Dann:

$$Q_{1-2^{-k-1}} - Q_{1-2^{-k}} \leq \sqrt{\frac{8\nu}{1/2}} = 4\sqrt{\nu}$$

D.h. die Abstände zwischen aufeinander folgenden Quantilen von exponentiell fallenden Wahrscheinlichkeiten sind beschränkt. Insbesondere:

$$Q_{1-2^{-m-1}} - Q_{1/2} = \sum_{k=1}^m (Q_{1-2^{-k-1}} - Q_{1-2^{-k}}) \leq 4m\sqrt{\nu},$$

also  $Q_{1-2^{-m-1}} \leq \mathcal{M}(Z) + 4m\sqrt{\nu}$  und somit:

$$\forall m \in \mathbb{N} : P(Z > \mathcal{M}(Z) + 4m\sqrt{\nu}) \leq P(Z > Q_{1-2^{-m-1}}) \leq 2^{-m-1}$$

Umformung ergibt:

$$\forall t \geq 0 : P(Z \geq \mathcal{M}(Z) + t) \leq 2^{-\frac{t}{4\sqrt{\nu}}}.$$

Optimal wäre sub-gaußscher Abfall  $P(\dots) \leq 2^{-\frac{t^2}{\nu}}$ . Ist erreichbar mit anderen Methoden.

Nun: zweite Methode um aus ESU Abschätzungen für Wahrscheinlichkeiten von Tail-Events zu erhalten: Führe wie bei Anwendung der Markov-Ungleichung Substitution mit exponentielle Funktion durch.

Für  $\lambda > 0$  betrachte ZVe  $Y := e^{\frac{\lambda Z}{2}}$ .

$$\begin{aligned} \text{Var}(Y) &= E(e^{\lambda Z}) - E(e^{\frac{\lambda Z}{2}})^2 \\ &\stackrel{\text{ESU}}{\leq} E\left(\sum_{i=1}^n \left(e^{\frac{\lambda Z}{2}} - e^{\frac{\lambda Z'_i}{2}}\right)_+^2\right) \leq (*) \end{aligned}$$

wobei  $Z'_i$  unabhängige, identische Kopien von  $Z_i$  sind. Anwendung des Mittelwertsatzes liefert

$$e^{\frac{\lambda z}{2}} - e^{\frac{\lambda w}{2}} \leq \frac{\lambda}{2} e^{\frac{\lambda \xi}{2}} (z - w)$$

für einen Wert  $\xi$  zwischen  $w$  und  $z$ . Da  $x \mapsto \exp(x)$  isoton ist, gilt:

$$\begin{aligned} (e^{\frac{\lambda z}{2}} - e^{\frac{\lambda w}{2}})_+ &= \frac{\lambda}{2} e^{\frac{\lambda \xi}{2}} (z - w)_+ \leq \frac{\lambda}{2} e^{\frac{\lambda z}{2}} (z - w)_+ \\ \Rightarrow (e^{\frac{\lambda z}{2}} - e^{\frac{\lambda w}{2}})_+^2 &\leq \frac{\lambda^2}{4} e^{\lambda z} (z - w)_+^2 \end{aligned}$$

Daraus ergibt sich

$$(*) \leq \frac{\lambda^2}{4} E\left(e^{\lambda Z} \underbrace{\sum_{i=1}^n (Z - Z'_i)_+^2}_{\leq \nu}\right) \leq \frac{\nu \lambda^2}{4} E\left(e^{\lambda Z}\right)$$

wobei wir auch hier die globale Annahme (3.10) verwendet haben. Insgesamt liefert ESU

$$\begin{aligned} \text{Var}(Y) &= E\left(e^{\lambda Z}\right) - E\left(e^{\frac{\lambda Z}{2}}\right)^2 \leq \frac{\nu \lambda^2}{4} E\left(e^{\lambda Z}\right) \\ \Leftrightarrow \left(1 - \frac{\nu \lambda^2}{4}\right) E\left(e^{\lambda Z}\right) &\leq \left(E\left(e^{\frac{\lambda Z}{2}}\right)\right)^2 \\ \Leftrightarrow \left(1 - \frac{\nu \lambda^2}{4}\right) E\left(e^{\lambda Z} e^{-\lambda E(Z)}\right) &\leq \left(E\left(e^{\frac{\lambda Z}{2}}\right)\right)^2 e^{-\lambda E(Z)} \\ \Leftrightarrow \left(1 - \frac{\nu \lambda^2}{4}\right) E\left(e^{\lambda(Z-E(Z))}\right) &\leq \left(E\left(e^{\frac{\lambda}{2}(Z-E(Z))}\right)\right)^2 \\ \Leftrightarrow \left(1 - \frac{\nu \lambda^2}{4}\right) M(\lambda) &\leq M\left(\frac{\lambda}{2}\right)^2, \end{aligned}$$

wobei  $F$  die MEF von  $Z - E(Z)$  ist. Diese Ungleichung für die MEF lässt sich mittels folgenden Lemmas aus der Analysis ausnutzen:

**Lemma 3.38.** (*Lemma aus der Analysis*)

$g : (0, 1) \rightarrow (0, \infty)$  erfülle folgende (sanfte) Regularitätsannahme

$$\lim_{x \rightarrow 0} \frac{g(x) - 1}{x} = 0.$$

Dann impliziert

$$g(x)(1 - x^2) \leq g\left(\frac{x}{2}\right)^2 \text{ für alle } x \in (0, 1),$$

$$\text{dass } g(x) \leq \left(\frac{1}{1 - x^2}\right)^2 \text{ gilt.}$$

*Beweis.*

$$\forall x \in (0, 1) : g(x)(1 - x^2) \leq g\left(\frac{x}{2}\right)^2 \Leftrightarrow \forall x \in (0, 1) : g(x) \leq \frac{1}{1 - x^2} g\left(\frac{x}{2}\right)^2.$$

Wegen  $x \in (0, 1) \Rightarrow \frac{x}{2^n} \in (0, 1)$  kann man dies rekursiv anwenden

$$\forall k \in \mathbb{N} : g(x) \leq (g(x2^{-k}))^{2^k} \prod_{j=0}^{k-1} \left(1 - (x2^{-j})^2\right)^{-2^j}$$

Für  $k \rightarrow \infty$  folgt:

$$g(x) \leq \lim_{k \rightarrow \infty} (g(x2^{-k}))^{2^k} \prod_{j=0}^{k-1} \left(1 - (x2^{-j})^2\right)^{-2^j}$$

Betrachte Limes für ersten Faktor separat.

Regularitätsannahme  $g(x) = 1 + o(x)$  impliziert

$$\ln(g(x2^{-k})^{2^k}) = 2^k \ln(g(x2^{-k})) = 2^k \ln(1 + o(x2^{-k})) \leq 2^k o(x2^{-k}) \xrightarrow[k \rightarrow \infty]{} 0,$$

wobei in der letzten Ungleichung  $\ln(1 + x) \leq x$  eingeht. Es folgt  $g(x2^{-k})^{2^k} \leq \exp(2^k o(x2^{-k})) \xrightarrow[k \rightarrow \infty]{} 1$ . Damit folgt

(3.12)

$$\ln(g(x)) \leq \lim_{k \rightarrow \infty} \sum_{j=0}^{k-1} 2^j (-\ln(1 - (x2^{-j})^2)) = \sum_{j=0}^{\infty} 2^j (-\ln(1 - (x2^{-j})^2)).$$

Da  $x \mapsto \ln(x)$  konkav, ist  $-x \mapsto -\ln(x)$  konvex. Also:

$$u \mapsto -\frac{1}{u} \ln(1 - u) \text{ isoton für } u \in (0, 1)$$

$$\Rightarrow -\frac{1}{x^2 2^{-2j}} \ln(1 - x^2 2^{-2j}) \leq -\frac{1}{x^2} \ln(1 - x^2), \text{ da } x^2 \geq x^2 2^{-2j}.$$

Umformen liefert:

$$-\ln(1 - (x2^{-j})^2) \leq 2^{-2j} (-\ln(1 - x^2))$$

Einsetzen in (3.12) ergibt:

$$\ln(g(x)) \leq \sum_{j=0}^{\infty} 2^j 2^{-2j} (-\ln(1-x^2)) = \left( \sum_{j=0}^{\infty} 2^{-j} \right) (-\ln(1-x^2)) = -2 \ln(1-x^2).$$

Exponentialfunktion anwenden, liefert

$$g(x) = \exp(-2 \ln(1-x^2)) = \exp(\ln((1-x^2)^{-2}))$$

die Behauptung. □

Wende das Lemma auf  $M(\lambda) = E(e^{-\lambda(Z-E(Z))})$  an.

$$M(0) = 1, M'(\lambda) = E\left((Z-E(Z))e^{\lambda(Z-E(Z))}\right) \stackrel{\text{zentriert}}{\Rightarrow} M'(0) = 0,$$

also ist  $M(\lambda) = 1 + o(\lambda)$  für  $\lambda \rightarrow 0$ . Setzen wir  $g(x) = M\left(\frac{2x}{\sqrt{\nu}}\right)$ , dann gilt für alle  $\lambda \in \left(0, \frac{2}{\sqrt{\nu}}\right)$ .

$$(3.13) \quad M(\lambda) = g\left(\lambda \frac{\sqrt{\nu}}{2}\right) \leq \left(\frac{1}{1 - \left(\frac{\lambda\sqrt{\nu}}{2}\right)^2}\right)^2 = \left(1 - \frac{\lambda^2\nu}{4}\right)^{-2}.$$

Insbesondere sehen wir, dass die ESU unter der Annahme (3.10) gewährleistet, dass  $Z$  exponentiell integrierbar ist. Weiterhin erhalten wir eine Abschätzung an W'keiten, dass  $Z$  große Werte annimmt:

$$\begin{aligned} \forall t > 0 : P(Z - E(Z) \geq t) &\stackrel{\text{Markov}}{\leq} E(e^{(Z-E(Z))\nu^{-\frac{1}{2}}})e^{-t\nu^{-\frac{1}{2}}} \\ &\leq e^{-\frac{t}{\sqrt{\nu}}} M\left(\frac{1}{\sqrt{\nu}}\right) \leq e^{-\frac{t}{\sqrt{\nu}}} \underbrace{\left(1 - \frac{1}{4}\right)^{-2}}_{\leq 2} \leq 2e^{-\frac{t}{\sqrt{\nu}}}. \end{aligned}$$

Die Schranke hat die gleiche Struktur wie beim ersten Anwendung der EFU auf *tail*-Abschätzungen. Aber:

- hier haben wir das Ereignis: Überschreitung des Erwartungswerts
- dort haben wir das Ereignis: Überschreitung des Medians!

Nutze (3.13) etwas anders mit Hilfe von  $-\ln(1-u) \leq u(1-u)^{-1}$ :

$$\begin{aligned} \Rightarrow \forall \lambda \in \left(0, 2\nu^{-\frac{1}{2}}\right) : \ln(M(\lambda)) &\leq \ln\left(1 - \frac{\lambda^2\nu}{4}\right) (-2) \leq 2\frac{\lambda^2\nu}{4} \left(1 - \frac{\lambda^2\nu}{4}\right)^{-1} \\ \Rightarrow M(\lambda) &\leq e^{\frac{\lambda^2\nu}{2} \left(1 - \frac{\lambda^2\nu}{4}\right)^{-1}}. \end{aligned}$$

Also:  $Z - E(Z)$  ist eine sub- $\Gamma$  ZV (vgl. Kapitel 2.4) mit Varianzfaktor  $\nu$  und Skalenparameter  $c = \frac{\sqrt{\nu}}{2}$ .

Aus Kapitel 2.4 folgt wieder eine Tail-Abschätzung:

$$(3.14) \quad P(Z - E(Z) \geq \sqrt{2\nu t} + ct) \leq e^{-t} \text{ für } t > 0.$$

Entscheidend: Wie groß ist der Wert  $c$ ? Bei uns  $c = \frac{\sqrt{\nu}}{2}$ . Für nicht zu kleine  $t$  dominiert  $ct = \frac{\sqrt{\nu}}{2}t$  den Term  $\sqrt{2\nu t}$ . Also haben wir dort nur exponentielles Verhalten und (3.14) ist keine sub-gaußsche Abschätzung.

### 3.8. Gaußsche Poincaré-Ungleichung.

Grundprinzipien

- Annahmen an  $f$  abschwächen
- Annahmen an ZVe dafür stärker

**Satz 3.39.** (*Poincaré-Ungleichung*)

Sei  $X = (X_1, \dots, X_d) \sim \mathcal{N}(0, I_d)$ , also sind  $X_1, \dots, X_d$  unabhängig, identisch verteilt. Sei  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  mit  $f \in \mathcal{C}_0^2(\mathbb{R}^d)$ .

$$\text{Dann: } \text{Var}(f(X)) \leq E(\|\nabla f(X)\|^2).$$

Insbesondere impliziert die Annahme, dass  $\text{supp } f$  kompakt ist.

**Satz 3.39** wird in Kapitel 5.3 in Theorem 5.4 allgemeiner bewiesen.

*Beweis.* O.E.  $E(\|\nabla f(X)\|^2) < \infty$ . Zuerst 1-dimensional; höhere Dimensionen liefert ESU. Außerdem betrachte erst  $\varepsilon_1, \dots, \varepsilon_n$  unabhängig, identisch Rademacher-verteilte ZVen, also  $P(\varepsilon_i = 1) = P(\varepsilon_i = -1) = \frac{1}{2}$ . Dann gilt nach ZGS

$$S_n := \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon_i \xrightarrow{\mathcal{D}} X \sim \mathcal{N}(0, \nu).$$

Also  $\forall g \in \mathcal{C}_b(\mathbb{R})$ :

$$\int g(t) P_{S_n}(dt) \xrightarrow{n \rightarrow \infty} \int g(t) dP_X(dt),$$

insbesondere  $\text{Var}(f(S_n)) \rightarrow \text{Var}(f(X)) = \text{Var}(Z)$ , falls man  $g = f^2$  wählt.

Idee: Leite Schranke für linke Seite her und betrachte den Limes. Für jedes  $i$  betrachte Varianz bezüglich  $\varepsilon_i$ :

$$\text{Var}^{(\varepsilon_i)}(f(S_n)) = \frac{1}{4} \left( \underbrace{f\left(S_n + \frac{1 - \varepsilon_i}{\sqrt{n}}\right) - f\left(S_n - \frac{1 + \varepsilon_i}{\sqrt{n}}\right)}_{=: D_{n,i}} \right)^2.$$

Die ESU liefert mit  $Z = \tilde{f}(\varepsilon_1, \dots, \varepsilon_n) = f(S_n)$ :

$$\text{Var}(f(S_n)) \leq \frac{1}{4} \sum_{i=1}^n E(D_{n,i}^2).$$

Setze  $k := \sup_{t \in \mathbb{R}} |f''(t)| < \infty$  nach Annahme  $f \in \mathcal{C}_0^2$ . Taylorentwicklung ergibt

$$\begin{aligned} |D_{n,i}| &\leq \frac{2}{\sqrt{n}} |f'(S_n)| + \frac{1}{2} \left( \frac{2}{\sqrt{n}} \right)^2 k, \text{ also} \\ \frac{1}{4} \sum_{i=1}^n D_{n,i}^2 &\leq \frac{1}{4} \sum_{i=1}^n \left( \frac{4}{n} f'(S_n)^2 + \frac{4k^2}{n^2} + \frac{8k}{n\sqrt{n}} |f'(S_n)| \right) \\ &\leq f'(S_n)^2 + \frac{k^2}{n} + \frac{2k}{\sqrt{n}} \underbrace{|f'(S_n)|}_{\leq \|f'\|_\infty} \end{aligned}$$

Mit dem ZGS folgt

$$\begin{aligned} \limsup_{n \rightarrow \infty} \frac{1}{4} \sum_{i=1}^n E(D_{n,i}^2) &\leq \limsup_{n \rightarrow \infty} E(f'(S_n)^2) + \limsup_{n \rightarrow \infty} E \left( \frac{k^2}{n} + \frac{4k}{\sqrt{n}} \|f'\|_\infty \right) \\ &\stackrel{\text{ZGS}}{=} E(f'(X)^2). \end{aligned}$$

Insgesamt:  $\text{Var}(f(X)) = \lim_{n \rightarrow \infty} \text{Var}(f(S_n)) \leq \lim_{n \rightarrow \infty} \frac{1}{4} \sum_{i=1}^n E(D_{n,i}^2) \leq E(f'(X)^2)$ .

Nun:  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  und  $X = (X_1, \dots, X_d) \sim \mathcal{N}(0, I_d)$ .

ESU liefert für  $Z = f(X_1, \dots, X_n)$ :

$$\begin{aligned} \text{Var}(f(X)) &\leq \sum_{i=1}^d E \left( (Z - E^{(i)}(Z))^2 \right) = \sum_{i=1}^d \text{Var}^{(i)}(f(X)) \\ &\stackrel{\text{Beweisteil 1}}{\leq} \sum_{i=1}^d E \left( \left( \frac{\partial f}{\partial x_i}(X_1, \dots, X_d) \right)^2 \right) = E \left( \sum_{i=1}^d \left( \frac{\partial f}{\partial x_i}(X_1, \dots, X_d) \right)^2 \right) \\ &= E(\|\nabla f(X)\|^2) \end{aligned}$$

□

### 3.9. Beweis für die ES-Ungleichung mittels Dualität.

Wir nutzen nun die Sichtweise der Dualität, um einen alternativen Beweis der ESU zu entwickeln. **Diese Beweismethode wird später im allgemeineren Kontext angewendet.**

Für  $Y, T \in \mathcal{L}(\Omega, \mathcal{A}, P)$  gilt

$$0 \leq \text{Var}(Y - T) = \text{Var}(Y) - 2 \text{Cov}(Y, T) + \text{Var}(T)$$

und umgestellt

$$\text{Var}(Y) \geq 2 \text{Cov}(Y, T) - \text{Var}(T)$$

Dies Ungleichung kann nicht verschärft werden, denn für  $T = Y$  gilt:

$$\text{Var}(Y) = 2 \text{Cov}(Y, T) - \text{Var}(T)$$

Also haben wir bewiesen:

**Proposition 3.40.** Für jedes  $Y \in \mathcal{L}(\Omega, \mathcal{A}, P)$  gilt:

$$(3.15) \quad \text{Var}(Y) = \max_{T \in \mathcal{L}(\Omega, \mathcal{A}, P)} [2 \text{Cov}(Y, T) - \text{Var}(T)]$$

Werden dies i. F. nutzen um Varianz von  $Z = f(X_1, \dots, X_n)$  mit unabhängigen  $X_1, \dots, X_n$  abzuschätzen. Verwenden zunächst wieder Teleskopsumme

$$Z^2 - (EZ)^2 = \sum_{i=1}^n [(E_i Z)^2 - (E_{i-1} Z)^2], \quad \text{wobei}$$

*Bemerkung 3.41* (Erinnerung).

$$\begin{aligned} E_i(Z) &= \int_{\mathcal{X}^{n-i}} f(X_1, \dots, X_i, \xi_{i+1}, \dots, \xi_n) P_{X_{i+1}}(\xi_{i+1}), \dots, P_{X_n}(\xi_n) \\ E^{(i)}(Z) &= \int_{\mathcal{X}} f(X_1, \dots, X_{i-1}, \xi_i, X_{i+1}, \dots, X_n) P_{X_i}(\xi_i) \\ \Delta_i &= E_i(Z) - E_{i-1}(Z) \end{aligned}$$

Damit folgt

$$\text{Var}(Z) = E [Z^2 - (EZ)^2] = \sum_{i=1}^n E [(E_i Z)^2 - (E_{i-1} Z)^2]$$

Erst im folgenden Schritt nutzen wir Orthogonalität. Da  $E_{i-1}(Z) \perp \Delta_i$ , folgt aus Pythagoras

$$\begin{aligned} E [(E_i Z)^2] &= E [(E_{i-1} Z)^2] + E [\Delta^2] \\ \Rightarrow E [\Delta^2] &= E [(E_i Z)^2] - E [(E_{i-1} Z)^2] + \end{aligned}$$

Wie schon vorher bemerkt, gilt für unabh. ZV  $X_1, \dots, X_n$ :

$$\forall i \in \{2, \dots, n\} : E_{i-1}(Z) = E^{(i)}(E_i Z)$$

also auch

$$\begin{aligned} E [(E_i Z)^2 - (E_{i-1} Z)^2] &= E [(E_i Z)^2 - (E^{(i)}(E_i Z))^2] \\ &= E [E^{(i)} [(E_i Z)^2 - (E^{(i)}(E_i Z))^2]] = E [\text{Var}^{(i)}(E_i Z)] \end{aligned}$$

Also haben wir (diesmal ohne Nutzung der Struktureigenschaft (3.1) der Martingaldifferenzen) bewiesen:

$$\text{Var}(Z) = \sum_{i=1}^n E \left[ \text{Var}^{(i)}(E_i Z) \right]$$

Falls wir es nun schaffen, den  $E_i$  Operator an  $\text{Var}^{(i)}$  vorbeizuziehen, verschwindet er wegen der Turmeigenschaft. Mit einer Ungleichheit ist dies wegen Lemma 3.42 in der Tat öglich und wir erhalten

$$\sum_{i=1}^n E \left[ \text{Var}^{(i)}(E_i Z) \right] \leq \sum_{i=1}^n E \left[ \text{Var}^{(i)}(Z) \right]$$

also die ESU.

**Lemma 3.42.** *Ist  $Z = f(X_1, \dots, X_n)$  quadratintegrierbar mit unabhängigen  $X_1, \dots, X_n$ , so gilt*

$$\forall i \in \{1, \dots, n\} : \quad E \left[ \text{Var}^{(i)}(E_i Z) \right] \leq E \left[ \text{Var}^{(i)}(Z) \right]$$

*Beweis.* Setze  $\text{Cov}^{(i)}(Z, T) = E^{(i)} \left[ (Z - E^{(i)}(Z))(T - E^{(i)}(T)) \right]$ . Dann folgt wie bei obiger Proposition

$$(3.16) \quad \text{Var}^{(i)}(Y) \geq 2 \text{Cov}^{(i)}(Y, T) - \text{Var}^{(i)}(T) \quad \text{für bel. } Y, T \in \mathcal{L}(\Omega, \mathcal{A}, P)$$

Andererseits gilt für jede  $\sigma(X_1, \dots, X_i)$ -messbare ZV  $T \in \mathcal{L}^2$ :

$$\begin{aligned} E \left[ \text{Cov}^{(i)}(Y, T) \right] &= E \left[ E^{(i)} \left[ (Z - E^{(i)}(Z))(T - E^{(i)}(T)) \right] \right] \\ \text{Turm-E.} &= E \left[ (Z - E^{(i)}(Z))(T - E^{(i)}(T)) \right] \\ \text{Turm-E.} &= E \left[ E_i \left[ (Z - E^{(i)}(Z))(T - E^{(i)}(T)) \right] \right] \\ T \text{ messbar} &= E \left[ (T - E^{(i)}(T)) E_i \left[ (Z - E^{(i)}(Z)) \right] \right] \\ E \text{ linear} &= E \left[ (T - E^{(i)}(T)) \left[ E_i(Z) - E_i(E^{(i)}(Z)) \right] \right] \\ \text{unabh.} &= E \left[ (T - E^{(i)}(T)) \left[ E_i(Z) - E^{(i)}(E_i(Z)) \right] \right] \\ &= E \left[ \text{Cov}^{(i)}(E_i(Z), T) \right] \end{aligned}$$

insbesondere für  $T = E_i(Z) \in \mathcal{L}^2(\sigma(X_1, \dots, X_i), P)$  mit (3.16)

$$\begin{aligned} E \left[ \text{Cov}^{(i)}(Y, E_i(Z)) \right] &= E \left[ \text{Cov}^{(i)}(E_i(Z), E_i(Z)) \right] = E \left[ \text{Var}^{(i)}(E_i(Z)) \right] \\ &\geq E \left[ 2 \text{Cov}^{(i)}(E_i(Z), E_i(Z)) - \text{Var}^{(i)}(E_i(Z)) \right] \\ &= E \left[ \text{Var}^{(i)}(E_i(Z)) \right] \end{aligned}$$

□

## 4. ENTROPIE

### 4.1. Shannon-Entropie und relative Entropie.

Zunächst einfaches Setting:

- $X: \Omega \rightarrow \mathcal{X}$  ZVe.
- $\mathcal{X}$  ist abzählbar (diskret).
- Gewichtsfunktion  $p(x) = P(X = x)$ ,  $x \in \mathcal{X}$ .

Definition: Die *Shannon-Entropie* der ZVe  $X$  (bzw. Verteilung  $P_X$ ):

$$H(X) := E(-\ln(p(X))) = - \sum_{x \in \mathcal{X}, p(x) > 0} \ln(p(x)) \geq 0$$

Die Abbildung  $[0, \infty) \ni x \mapsto -x \ln(x)$  ist konkav.

Sei  $q: \mathcal{X} \rightarrow [0, 1]$ ,  $\sum_{x \in \mathcal{X}} q(x) = 1$  eine weitere Gewichtsfunktion (z.B. Verteilung einer ZVe  $Y$ ).

Sei  $Z$  eine ZVe mit Gewichtsfunktion  $\frac{1}{2}(p + q)$ . Dann gilt:

$$H(Z) \geq \frac{1}{2}H(X) + \frac{1}{2}H(Y).$$

Das Mischen von Gewichtsfunktionen vergrößert als die Entropie. Allgemeiner gilt für jedes  $t \in (0, 1)$  und für eine ZVe  $Z$  mit Gewichtsfunktion  $tp + (1 - t)q$ :

$$H(Z) \geq tH(X) + (1 - t)H(Y).$$

*Beweis.* Für beliebiges  $t \in (0, 1)$  und  $x \in \mathcal{X}$  gilt:

$$\begin{aligned} & - \sum_{x \in \mathcal{X}} \left( (tp(x) + (1 - t)q(x)) \cdot \ln(tp(x) + (1 - t)q(x)) \right) \\ & \geq - \left( t \sum_{x \in \mathcal{X}} p(x) \ln(p(x)) + (1 - t) \sum_{x \in \mathcal{X}} q(x) \ln(q(x)) \right) \text{ wegen Konkavität.} \end{aligned}$$

□

Definition: Seien  $\mathcal{X}$  höchstens abzählbar,  $P, Q$  Wahrscheinlichkeitsmaße auf  $\mathcal{P}(\mathcal{X})$ ,  $p, q$  dazugehörige Gewichtsfunktionen. Die *Kulback-Leibler-Divergenz* oder *relative Entropie* von  $P$  und  $Q$  ist definiert durch:

$$D(P\|Q) := \begin{cases} \sum_{x \in \mathcal{X}} p(x) \log \left( \frac{p(x)}{q(x)} \right) & \text{falls } P \ll Q \\ \infty & \text{sonst} \end{cases}$$

Die relative Entropie erinnert an einen Abstands begriff. Allerdings ist sie keine echte Metrik, da keine Symmetrie vorliegt. Offensichtlich gilt aber:

$$\begin{aligned} D(P\|Q) &= - \sum_{x \in \mathcal{X}, p(x) > 0} p(x) \ln \left( \frac{p(x)}{q(x)} \right) \\ &\geq - \sum_{x \in \mathcal{X}, p(x) > 0} p(x) \left( \frac{p(x)}{q(x)} - 1 \right) \\ &\geq - \sum_{x \in \mathcal{X}, p(x) > 0} q(x) - p(x) = Q\{x \mid p(x) > 0\} + 1 \geq 0. \end{aligned}$$

Falls  $D(P\|Q) = 0$ , so ist  $\text{supp}(P) = \text{supp}(Q)$ , denn

für  $\subseteq$ : da  $P \ll Q$

für  $\supseteq$ : wegen  $Q\{x \mid p(x) > 0\} = 1$

Es gilt sogar  $P = Q$ , denn es muss  $\forall x \in \mathcal{X}$  gelten:

$$\ln \left( \frac{q(x)}{p(x)} \right) = \frac{q(x)}{p(x)} - 1, \text{ gilt genau dann, wenn } \frac{q(x)}{p(x)} = 1 \text{ ist.}$$

Beispiel: Sei  $\mathcal{X}$  endlich und  $Q$  die Gleichverteilung auf  $\mathcal{X}$ ,  $X: \Omega \rightarrow \mathcal{X}$  eine ZVe mit Verteilung  $P_X$ . Dann gilt:

$$\begin{aligned} 0 \leq D(P\|Q) &= - \sum_{x \in \mathcal{X}, p(x) > 0} p(x) \ln \left( \frac{\frac{1}{|\mathcal{X}|}}{p(x)} \right) \\ &= - \ln \left( \frac{1}{|\mathcal{X}|} \right) \sum_{x \in \mathcal{X}, p(x) > 0} p(x) + \sum_{x \in \mathcal{X}, p(x) > 0} p(x) \ln(p(x)) \\ &= \ln |\mathcal{X}| - H(X) \\ &\Leftrightarrow H(X) \leq \ln |\mathcal{X}|. \end{aligned}$$

Man überlege sich, dass Gleichheit genau dann gilt, wenn  $X$  gleichverteilt auf  $\mathcal{X}$  ist.

Fazit:  $0 \leq H(X) \leq \ln |\mathcal{X}|$  für beliebige Zufallsvariablen auf  $\mathcal{X}$ .

Interpretation: Linke Seite wird angenommen, wenn  $X \sim \delta_x$  z.B. ist. In diesem Fall liegt eine perfekte Information vor. Für Gleichverteilungen liegt die schlechteste Informationstrennung (größtes Chaos) vor.