

Grundlagen zur Wahrscheinlichkeitsrechnung und Pseudo-Zufallszahlenerzeugung

Alexander Asteroth



**Hochschule
Bonn-Rhein-Sieg**
University of Applied Sciences

Fachbereich Informatik
Department of Computer Science

Grundbegriffe/Notationen (Glossar)

Schriftarten

Skalar Ein Skalar $x \in \mathbb{R}$ wird mit einem Kleinbuchstaben geschrieben.

Vektor Ein Vektor $\boldsymbol{x} \in \mathbb{R}^n$ wird mit einem fetten Kleinbuchstaben geschrieben.

Menge Mengen (z.B. $\mathcal{A} \subseteq \mathbb{R}$) werden mit kaligraphischen Buchstaben geschrieben.

Funktion Funktionsnamen (z.B. Zufallsvariable $X : \mathcal{S} \rightarrow \mathbb{R}$) beginnen mit einem Großbuchstaben (oder einem griechischen Kleinbuchstaben). Falls sie vektor-wertig sind werden sie fett geschrieben, ansonsten nicht.

Matrix Für Matrizen $\boldsymbol{A} \in \mathbb{R}^{n \times m}$ verwenden wir, wie für Vektoren, fette Schrift. Allerdings werden für Matrizen Großbuchstaben verwendet, für Vektoren Kleinbuchstaben.

- Der Datentyp ergibt sich so aus der Schriftart.
- **Ausnahme:** Vektorwertige Funktionen und Matrizen werden beide mit fetten Großbuchstaben geschrieben. Wenn möglich werden für Matrizen daher immer die Buchstaben $\boldsymbol{A} - \boldsymbol{D}$ verwendet.

Wahrscheinlichkeiten

Symbol	Begriff	Definition/Signatur/Beispiel
\mathcal{S}	Menge von Stichproben (Stichprobenraum/engl. sample space)	$\mathcal{E} \subseteq \mathcal{P}(\mathcal{S})$ $P : \mathcal{E} \rightarrow [0, 1]$
\mathcal{E}	Menge von Ereignissen	
P	Wahrscheinlichkeitsmaß	
$(\mathcal{S}, \mathcal{E}, P)$	Wahrscheinlichkeitsraum	
	sicheres Ereignis	$P(\mathcal{S}) = 1$
	unmögliches Ereignis	$P(\emptyset) = 0$
	\mathcal{A} und \mathcal{B} unabhängig	$P(\mathcal{A} \cap \mathcal{B}) = P(\mathcal{A}) P(\mathcal{B})$
$P(\mathcal{A} \mathcal{B})$	bedingte Wahrscheinlichkeit	$P(\mathcal{A} \mathcal{B}) = \frac{P(\mathcal{A} \cap \mathcal{B})}{P(\mathcal{B})}$

Zufallsvariablen

X	reelwertige Zufallsvariable	$X : \mathcal{S} \rightarrow \mathbb{R}$
F_X	Wahrscheinlichkeitsverteilung	$F_X(\mathcal{O}) = P(X^{-1}(\mathcal{O}))$ $= P(\{a \in \mathcal{O} \mid X(a) \in \mathcal{O}\})$
f_X	Wahrscheinlichkeitsdichte ... und falls F diff.bar	$\int_{\mathcal{S}} f_X(x) dx = F_X(\mathcal{S})$ $f_X(x) = F'_X(x)$
$X \sim f$	Kurzschreibweise für	$f = f_X$
$X \sim Y$	Kurzschreibweise für	$f_X = f_Y$
f_{XY}	gemeinsame Dichte	$\int_{\mathcal{A}} \int_{\mathcal{B}} f_{XY}(x, y) dx dy = P(X^{-1}(\mathcal{A}) \cap Y^{-1}(\mathcal{B}))$
E, μ	Erwartungswert	$E(X) = \int_{\mathbb{R}} x f(x) dx$ $E(X) = \mu_X = \bar{x}$
Var, σ^2	Varianz	$\text{Var}(X) = \sigma_X^2 = E((X - \bar{x})^2)$
Cov	Kovarianz	$\text{Cov}(X, Y) = E((X - \bar{x})(Y - \bar{y}))$
Corr	Korrelation	$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$

Randdichte (marginal density)	$f(x) = \int_{\mathbb{R}} f(x, y) \, dy$
bedingte Dichte (conditional density)	$f(x \mid y) = \frac{f(x, y)}{f(y)}$
X, Y bedingt unabhängig geg. Z	$f(x, y \mid z) = f(x \mid z) f(y \mid z)$
Satz von Bayes	$f(y \mid x) = \frac{f(x y) f(y)}{f(x)}$

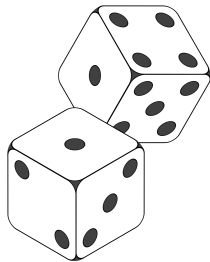
Noch einmal ... langsam — an einem Beispiel

Alea iacta est! (sorry für das langweilige Beispiel)

Würfeln: Wahrscheinlichkeitsmaß

Ausgegangen wird von einem fairen Würfel. Wenn wir als Stichproben eines Zufallsexperiments die Anzahl der Punkte oben auf dem Würfel betrachten erhalten wir als Menge von Stichproben:

$$\mathcal{S} = \{1, 2, 3, 4, 5, 6\}$$



Wenn wir zwei Würfel unabhängig voneinander werfen, besteht der zugehörige Stichprobenraum aus allen Paaren (x, y) wo $x, y \in \mathcal{S}$. Daraus ergibt sich folgender Stichprobenraum:

$$\mathcal{S} \times \mathcal{S}$$

Würfeln: Ereignisraum

Ein Ereignis ist eine Untermenge von \mathcal{S} . Der Ereignisraum \mathcal{E} ist eine Menge von Ereignissen

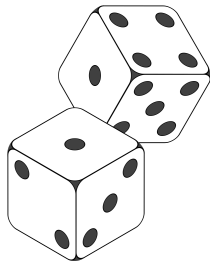
- (i) \mathcal{E} beinhaltet \mathcal{S} ($\mathcal{S} \in \mathcal{E}$)
- (ii) ist unter dem Komplement abgeschlossen
($\mathcal{A} \in \mathcal{E} \implies \overline{\mathcal{A}} \in \mathcal{E}$)
- (iii) ist unter Vereinigung abgeschlossen
($\mathcal{A}, \mathcal{B} \in \mathcal{E} \implies (\mathcal{A} \cup \mathcal{B}) \in \mathcal{E}$)

Wenn die Menge die Elementarereignisse $\{o_i\}$ für alle $o_i \in \mathcal{S}$ beinhaltet, dann folgt durch (i)-(iii) dass

$$\mathcal{E} = \mathcal{P}(\mathcal{S})$$

Im Fall von zwei Würfeln:

$$\mathcal{E} = \mathcal{P}(\mathcal{S}), \quad \text{wobei } \mathcal{S} = \{1, 2, 3, 4, 5, 6\} \times \{1, 2, 3, 4, 5, 6\}$$



Würfeln: Wahrscheinlichkeitsmaß

Ein Wahrscheinlichkeitsmaß ist eine Funktion

$$P : \mathcal{E} \rightarrow [0, 1]$$

sodass für $\mathcal{A}, \mathcal{B} \in \mathcal{E}$

(i) $P(\mathcal{S}) = 1$

(ii) $P(\overline{\mathcal{A}}) = 1 - P(\mathcal{A})$

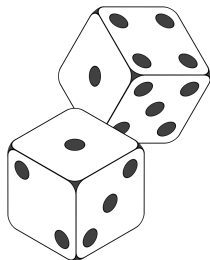
(iii) wenn $\mathcal{A} \cap \mathcal{B} = \emptyset$ dann $P(\mathcal{A} \cup \mathcal{B}) = P(\mathcal{A}) + P(\mathcal{B})$

Da von einem fairen Würfel ausgegangen wird, hat jedes Elementarereignis die selbe Wahrscheinlichkeit:

$$P(\{1\}) = \dots = P(\{6\}) = \frac{1}{6}$$

Wahrscheinlichkeiten von zusammengesetzten Ereignissen wie $\{2, 3\}$ (zwei oder drei Punkte) können, da diese disjunkt sind, von Elementarereignissen abgeleitet werden (iii):

$$P(\{2, 3\}) = P(\{2\}) + P(\{3\}) = \frac{1}{3}$$



Würfeln: Zufallsvariablen und Verteilungen

Eine reellwertige Zufallsvariable X ist eine Funktion

$$X : \mathcal{S} \rightarrow \mathbb{R}$$

sodass das Urbild von Untermengen von \mathbb{R} den Ereignissen entspricht:

$$\mathcal{A} \subseteq \mathbb{R} \implies X^{-1}(\mathcal{A}) \in \mathcal{E}$$

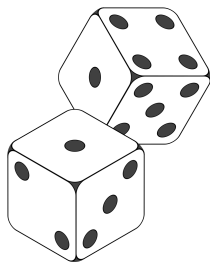
Dadurch kann die Wahrscheinlichkeitsverteilung F_X von X über \mathbb{R} definiert werden durch

$$F_X(\mathcal{A}) = P(X^{-1}(\mathcal{A}))$$

Wenn das Bild $\text{Im}(X)$ endlich/diskret ist nennen wir X eine endliche/diskrete Zufallsvariable (ZV).

Im Fall vom Werfen zweier Würfel können wir eine ZV X definieren indem jedes Paar auf die Summe der Punkte abgebildet wird.

$$X : \mathcal{S} \rightarrow \mathbb{R}, \quad \text{Im}(X) = \{2, \dots, 12\}, \quad \mathcal{S} = \{1, \dots, 6\} \times \{1, \dots, 6\}$$



Würfeln: Zufallsvariablen, Verteilungen, Notation

Betrachten wir wieder den Wurf von zwei Würfeln, so definieren wir eine ZV X durch Abbilden jedes Paares auf die Summe der Punkte.



$$X : \mathcal{S} \rightarrow \mathbb{R}, \quad \text{Im}(X) = \{2, \dots, 12\}, \quad \mathcal{S} = \{1, \dots, 6\} \times \{1, \dots, 6\}$$

Die Verteilung $F_X(\mathcal{A})$ misst die Wahrscheinlichkeit von X durch das Generieren eines Wertes aus \mathcal{A} . Das wird häufig beschrieben durch:

$$P(X \in \mathcal{A})$$

$$\text{z.B. } P(X \in [5, 8])$$

Wenn \mathcal{A} ein Intervall $\mathcal{A} = [a, b]$ ist, dann wird statt $F_X(\mathcal{A})$ oft die Notation

$$P(a \leq X \leq b)$$

$$\text{z.B. } P(5 \leq X \leq 8)$$

verwendet. Wenn $\mathcal{A} = \{a\}$ eine einelementige Menge ist, dann wird $F_X(\mathcal{A})$ beschrieben als:

$$P(X = a)$$

$$\text{z.B. } P(X = 8)$$

Würfeln: Zufallsvariablen und Verteilungsfunktion

Oft ist die Verteilung einer Zufallsvariablen nicht für beliebige Intervalle, sondern für Intervalle der Form definiert:

$$\mathcal{A} = (-\infty, x]$$

In diesem Fall definiert die rechte Grenze von \mathcal{A} die gesamte Menge \mathcal{A} und wir können F_X definieren durch

$$F_X(x) = F_X((-\infty, x]) = P(X^{-1}((-\infty, x]))$$

Wenn F so definiert ist wird die Verteilung als *kumulative Verteilungsfunktion* (CDF) von X bezeichnet.

In unserem Beispiel:

$$F_X(8) \approx 0.72$$

gibt die Wahrscheinlichkeit an, dass die Summe von zwei Würfeln einen Wert kleiner oder gleich 8 annimmt.

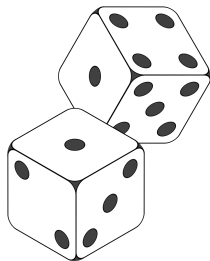
Würfel: Zufallsvariablen und Dichte

Die Dichte einer Zufallsvariablen X ist eine Funktion f_X so dass

$$F_X(\mathcal{A}) = \int_{\mathcal{A}} f_X(x) \, dx \quad (1)$$

oder äquivalent

$$F_X(x) = \int_{-\infty}^x f_X(x') \, dx' \quad (2)$$



Anmerkung

Da das Integral immer zu einer stetigen Funktion führt, existiert f_X **nicht** wenn F_X nicht stetig ist.

Würfeln: Zufallsvariablen und Dichte

Es gilt wieder $X : \mathcal{S} \rightarrow \mathbb{R}$ die
Summe der Augenzahl von zwei
Würfeln

$$\mathcal{S} = \{1, \dots, 6\}^2$$

es gibt 36 mögliche Ausgänge und
für alle Stichproben $(x, y) \in \mathcal{O}$

$$P((x, y)) = \frac{1}{36}$$

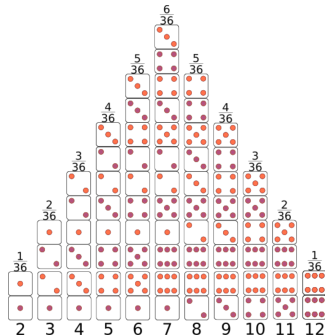
Würfeln: Zufallsvariablen und Dichte

Es gilt wieder $X : \mathcal{S} \rightarrow \mathbb{R}$ die Summe der Augenzahl von zwei Würfeln

$$\mathcal{S} = \{1, \dots, 6\}^2$$

es gibt 36 mögliche Ausgänge und für alle Stichproben $(x, y) \in \mathcal{O}$

$$P((x, y)) = \frac{1}{36}$$



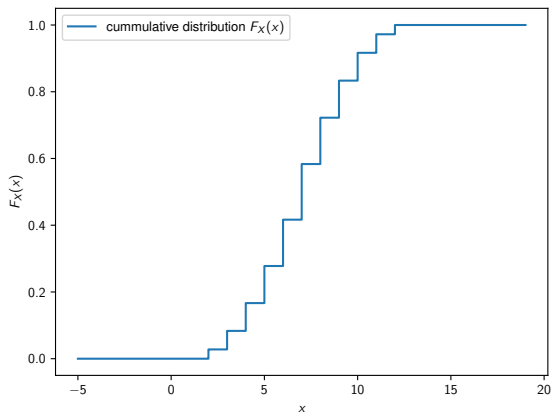
Würfeln: Zufallsvariablen und Dichte

Es gilt wieder $X : \mathcal{S} \rightarrow \mathbb{R}$ die Summe der Augenzahl von zwei Würfeln

$$\mathcal{S} = \{1, \dots, 6\}^2$$

es gibt 36 mögliche Ausgänge und für alle Stichproben $(x, y) \in \mathcal{O}$

$$P((x, y)) = \frac{1}{36}$$



Die kumulative Verteilungsfunktion (CDF) $F_X(x)$ ist rechts abgebildet. **Es handelt sich nicht um eine kontinuierliche Funktion, daher hat X keine Dichte f_X .**

Würfel: Wahrscheinlichkeitsvektor

X hat keine Dichte, die Wahrscheinlichkeiten sind verteilt über eine endliche Menge von diskreten Werten:

$$E = \text{Im}(X) = \{2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$$

Jeder der Werte entspricht einem Elementarergebnis, dessen Wahrscheinlichkeit gemessen werden kann:

$$p_i = F_X(\{e_i\}) = P(X^{-1}(\{e_i\})), \quad e_i \in E$$

(hier haben wir die erste Version der Verteilung F angewendet, nicht die CDF)

Da $E = \text{Im}(X)$ eine endliche Menge ist, weisen wir X den Wahrscheinlichkeitsvektor zu

$$\mathbf{p}_X = (p_1, \dots, p_n)^T$$

In unserem Beispiel

$$\begin{aligned} \mathbf{p}_X &= \frac{1}{36} (1, 2, 3, 4, 5, 6, 5, 4, 3, 2, 1)^T \\ &\approx (0.028, 0.056, 0.083, 0.111, 0.139, 0.167, 0.139, 0.111, 0.083, 0.056, 0.028)^T \end{aligned}$$

Würfeln: Erwartungswert und Varianz

Der Erwartungswert ist definiert als:

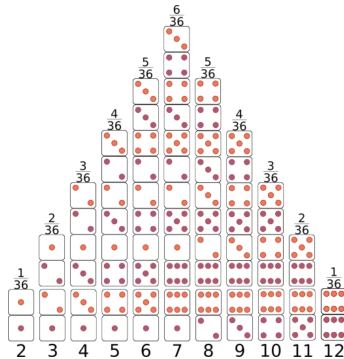
$$E(X) = \int_{\mathbb{R}} x f(x) dx$$

im Fall von diskreten Variablen wird aus dem Integral eine Summe, so ergibt sich für unser Beispiel:

$$\mu = E(X) = \sum_{j=1}^{11} e_j p_j$$

Umsortieren der Terme

$$\begin{aligned} &= (2 + 12) \frac{1}{36} + (3 + 11) \frac{2}{36} + \dots \\ &= 7 \sum p_j = 7 \end{aligned}$$



Würfel: Erwartungswert und Varianz

Varianz ist definiert als:

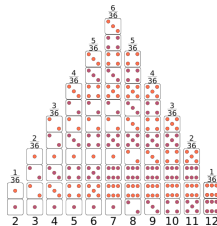
$$\text{Var}(X) = E((X - \mu)^2)$$

in unserem Beispiel:

$$\sigma^2 = \text{Var}(X) = \sum_{j=1}^{11} p_j (e_j - 7)^2$$

Umsortieren der Terme

$$\begin{aligned} &= 2 \cdot 5^2 \frac{1}{36} + 2 \cdot 4^2 \frac{2}{36} + 2 \cdot 3^2 \frac{3}{36} + 2 \cdot 2^2 \frac{4}{36} + 2 \cdot 1^2 \frac{5}{36} + 2 \cdot 0^2 \frac{6}{36} \\ &= \frac{50 + 64 + 54 + 32 + 10}{36} = \frac{210}{36} \approx 5.8 \end{aligned}$$



Würfel: Varianz und Standardabweichung

i

$$\sigma^2 = \frac{210}{36} \approx 5.8$$

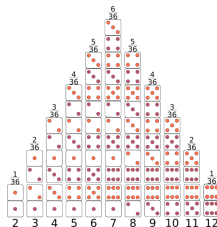
Standardabweichung:

$$\sigma = \sqrt{210/36} \approx 2.4$$

Generell gilt 68% der Stichproben X werden in das Intervall $[\mu - \sigma, \mu + \sigma] = [4.6, 9.4]$ fallen.

Wir können von der Abbildung rechts ableiten, dass $24/36 = 2/3 \approx 0.67$ tatsächlich in das oben angegebene Intervall fallen.

Darüberhinaus fallen 95% der Stichproben von X in ein 2σ Intervall und 99.7% in ein 3σ Intervall.



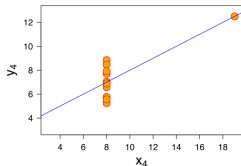
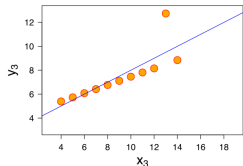
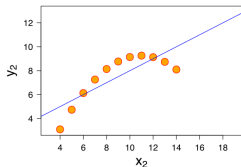
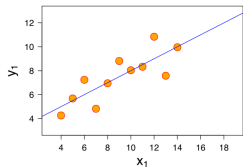
Was machen wir nun mit diesen Konzepten?

1. Gegeben Daten: bestimme passende Verteilung (z.B. maximum likelihood/max. a posteriori)
2. Gegeben einen randomisierten Algorithmus: bestimme Verteilung der Ergebnisse (probabilistische Analyse)
3. Gegeben eine Verteilung: bestimme deren Eigenschaften (z.B. Mittelwert Varianz, ...)
4. Gegeben eine Verteilung: ziehe Stichproben

Warum sollte man Stichproben ziehen?

Um einen randomisierten Algorithmus zufällige Entscheidungen treffen zu lassen.

Warum sollte man Stichproben ziehen?



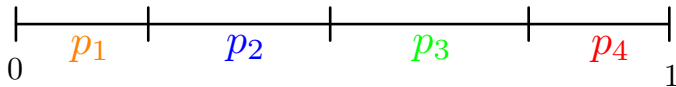
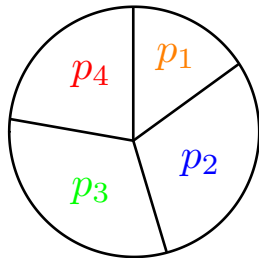
Property	Value	Accuracy
Mean of x	9	exact
Sample variance of x	11	exact
Mean of y	7.50	to 2 decimal places
Sample variance of y	4.125	± 0.003
Correlation between x and y	0.816	to 3 decimal places
Linear regression line	$y = 3 + 0.5x$	to 2 decimal places
R^2	0.67	to 2 decimal places

(info about Anscombe's quartet from Wikipedia)

(Anscombe, Francis J. (1973) Graphs in statistical analysis. American Statistician, 27, 17–21.)

Würfeln: Stichproben (roulette wheel selection)

Angenommen eine RV X und ihr Wahrscheinlichkeitsvektor $\mathbf{p}_X = (p_1, \dots, p_n)^T$ sind gegeben und wir wollen eine Stichprobe x von X ziehen.



Die Gleichverteilung $U(a, b)$

Wenne eine Zufallsvariable X auf ein Intervall $[a, b]$ mit konstanter Dichte auf dem ganzen Intervall abbildet (und Null außerhalb des Intervalls), nennen wir X gleichverteilt auf $[a, b]$ und schreiben dies als:

$$X \sim U(a, b)$$

Example 1.

Oft gilt $a = 0, b = 1$. In diesem Fall gilt

$$f_X(x) = 1 \text{ für alle } x \in [0, 1] \text{ und } f_X(x) = 0 \text{ sonst}$$

Die Verteilung einer solchen Zufallsvariablen auf dem Intervall $[0, a]$ für ein $a \in [0, 1]$ ist linear in $[0, 1]$

$$F_X([0, a]) = a$$

Allgemein gilt, dass die Dichte von $X \sim U(a, b)$ konstant auf $[a, b]$ und die Verteilung linear auf $[a, b]$ ist.

Die Gleichverteilung $U(a, b)$

Formaler:

$$f(x) = U(a, b)(x) = \begin{cases} \frac{1}{b-a} & \text{if } x \in [a, b] \\ 0 & \text{sonst} \end{cases}$$

Würfeln: Ziehen von Stichproben (roulette wheel selection)

Angenommen, eine Zufallsvariable X und der passende Wahrscheinlichkeitsvektor $\mathbf{p}_X = (p_1, \dots, p_n)^T$ seien gegeben, und wir wollen eine Stichprobe x aus X ziehen. Ein Ansatz das zu tun ist:

1. teile $[0, 1]$ in n Unterintervalle I_j der Länge p_j
2. ziehe eine Stichprobe $T \sim U(0, 1)$ aus der Gleichverteilung¹ auf $[0, 1]$,
Bezeichne die Stichprobe als t
3. finde den Index j sodass $t \in I_j$
4. setze $x = e_j$

Die Aufteilung von $[0, 1]$ in Unterintervalle und das anschließende gleichmäßige Ziehen aus $[0, 1]$, sind ähnlich wie das Konzept eines Roulette-Rades. Diese Methode ist daher auch bekannt als *roulette wheel selection*.

¹Pseudo-Zufallszahlengeneratoren für diese Verteilung gibt es in allen Programmiersprachen

Würfel: Ziehen von Stichproben (roulette wheel selection)

In unserem fortlaufenden Beispiel:

$$E = \text{Im}(X) = \{2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$$

$$p_X = \frac{1}{36}(1, 2, 3, 4, 5, 6, 5, 4, 3, 2, 1)^T$$

1. teile $[0, 1]$ in Unterintervalle

$$\begin{aligned} [0, 1] &= I_1 \cup I_2 \cup I_3 \cup \dots \cup I_{11} \\ &= \left[0, \frac{1}{36}\right) \cup \left[\frac{1}{36}, \frac{3}{36}\right) \cup \left[\frac{3}{36}, \frac{6}{36}\right) \cup \dots \cup \left[\frac{35}{36}, 1\right] \end{aligned}$$

2. ziehe eine zufällige Zahl aus $T \sim U(0, 1) \longrightarrow t = 0.3145$

3. finde den Index j sodass $t \in I_j$:

$$0.3145 \in \left[\frac{6}{36}, \frac{10}{36}\right) = I_4 \implies j = 4$$

4. setze

$$x = e_j = e_4 = 5$$

Jetzt wieder zurück zu kontinuierlichen Verteilungen, für die Dichten existieren.

Ziehen von Stichproben bei nicht-diskreten ZV

Angenommen eine Zufallsvariable X mit Dichte $f(x)$ ist gegeben und wir wollen eine Stichprobe daraus ziehen. Beginnen wir mit folgender Beobachtung – führe eine neue ZV T ein

$$T = F(x)$$

wie ist T auf seinem Bild $\text{Im}(T) = [0, 1]$ verteilt²? Sei $a \in [0, 1]$:

$$\begin{aligned} F_T(a) &= P(T \leq a) = P(F(x) \leq a) \\ &= P(X \leq F^{-1}(a)) = F(F^{-1}(a)) = a \end{aligned}$$

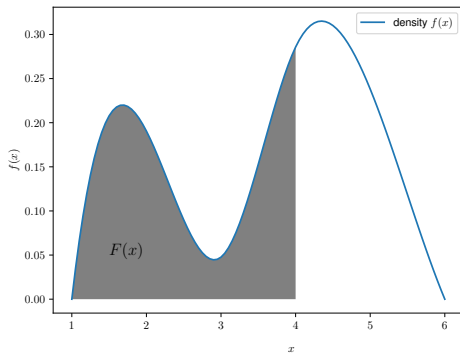
Daher ist T (die Verteilung von X) gleichverteilt in $[0, 1]$.

²unter der Annahme, dass F eine streng ansteigende Funktion in $[0, 1]$ so ist F^{-1} wohl definiert auf $[0, 1]$ (andernfalls wird der Beweis nur viel komplizierter)

Ziehen von Stichproben bei nicht-diskreten ZV

Dies führt zu der folgenden Methode:

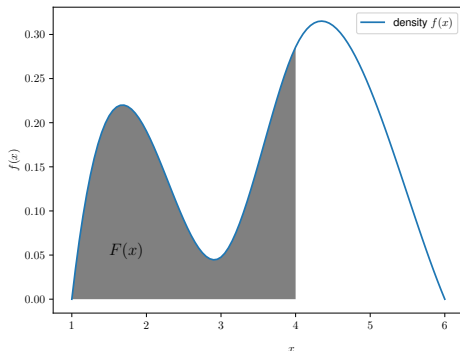
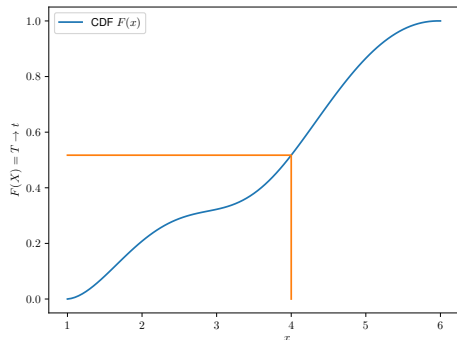
- berechne aus der Dichte $f(x)$ die Verteilung $F(x)$



Ziehen von Stichproben bei nicht-diskreten ZV

Dies führt zu der folgenden Methode:

- berechne aus der Dichte $f(x)$ die Verteilung $F(x)$



- ziehe eine Stichprobe aus $U(0, 1)$, nenne sie t
- finde x sodass $t = F(x)$

How to sample $T \sim \text{U}(0, 1)$? — PRNG's

Let's start with an example of a *pseudo random number generators*³, (*PRNG*):
Ansi C function `rand()` has the following implementation

$$s_0 = 12345$$

$$s_{i+1} = 1103515245s_i + 12345 \mod 2^{31}, \quad i = 0, 1, \dots$$

such a PRNG is calld a *linear congruential generator* (LCG).

³It is also possible to create *real* random numbers (TRNG). (see `/dev/random`, `/dev/urandom`).

Linear congruential generators

LCG's in general have the form

$$s_{i+1} = as_i + b \mod c$$

For efficiency reasons c is often chosen as a power of 2.

Let's consider a minimalistic example:

$$s_0 = 3$$

$$s_{i+1} = as_i \mod 8$$

and investigate the influence of a on the generated sequence:

$$a = 2 \qquad (s_i) = 3, 6, 4, 0, 0, \dots$$

$$a = 3 \qquad (s_i) = 3, 1, 3, 1, 3, \dots$$

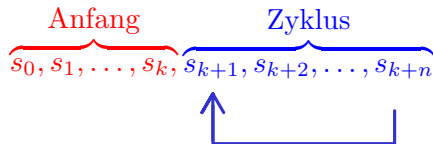
$$a = 4 \qquad (s_i) = 3, 4, 0, 0, 0, \dots$$

$$a = 5 \qquad (s_i) = 3, 7, 3, 7, 3, \dots$$

$$a = 6 \qquad (s_i) = 3, 2, 4, 0, 0, \dots$$

$$a = 7 \qquad (s_i) = 3, 5, 3, 5, 3, \dots$$

In general the generated sequence of random numbers has the following structure:



(Anfang = beginning, Zyklus=cycle)

The quality of a PRNG depends of the length and structure of the cycle. In particular:

- a) length of the cycle
- b) distrubution of numbers in cycle
- c) autocorrelation of the sequence
- d) cryptographically secure RNG's (will not be considered here)

While a) and b) should be more or less obvious, c) needs some explanation. Assume a LCG generates integer valued random numbers from $[0, c - 1] \pmod{c}$. Then the longest possible cycle is c elements long. But if this cycle is

$$0, 1, 2, 3, \dots, c - 1$$

we would not consider it very random. Therefore the autocorrelation of the sequence can help to find out about this kind of structure in the sequence.

If noncyclic sequences are of interest, approximations to irrational numbers are an option.

Pseudo Random Number Generators

... have the structure:

$$s_0 = \text{seed}, \quad s_{i+1} = f(s_i, s_{i-1}, \dots)$$

Linearly Congruential Generators (LCG's)

An LCG generates the next element of the sequence s_{i+1} solely based on the previous s_i :

$$s_{i+1} = as_i + b \mod c$$

often $c = 2^n$

For LCG's it holds that:

Lemma 2.

An LCG has a maximal cycle length of c , if

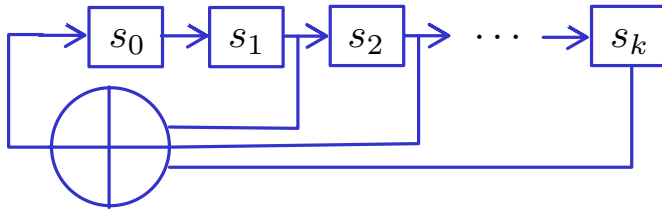
- ▶ *c and b are relatively prime*
- ▶ *every prime factor of c is also a prime factor of $a - 1$*

If $c = 2^n$, $a = 4k + 1$ for $k \in \mathbb{N}$ and b odd, above mentioned requirements are fulfilled
Let

$$s_{i+1} = 13s_i + 5 \mod 2^{10}$$

then (s_i) becomes cyclic after exactly 1024 steps.

Linear Recursive Shift Register



Definition 3 (Linear Recursive Shift Register (LRSR)).

A LRSR (linear recursive shift register) with feedback weights $\mathbf{b} = (b_0, \dots, b_k)$ and seed $s_0 = (s_0^0, s_1^0, \dots, s_k^0)$ is defined by it's register assignments at time $t + 1$ according to:

$$s_{i+1}^{t+1} = s_i^t, \quad i = 0 \dots k-1, \quad s_0^{t+1} = \bigoplus_{i=0}^k b_i s_i^t$$

The register can be considered either as a PRNG of numbers in the range $[0, 2^{k+1})$ or as a generator of a bit sequence (using an arbitrary bit).

Remark 1.

Often global behaviour of shift registers is good while local behaviour is “bad”. Autocorrelation is usually nearly 0 but it is assumed that higher order autocorrelation can be bad.

Kurzschreibweise

Als Variable für die Werte der Zufallsvariablen X, Y, Z verwenden wir die Buchstaben x, y, z . Dies ermöglicht es uns, einige Indizes weg fallen zu lassen:

$$f(x) = f_X(x), \quad F(x) = F_X(x), \quad f(y) = f_Y(y), \quad \dots$$

Unabhängig und identisch verteilte Zufallsvariablen (iid)

Oft beobachten wir das Ergebnis ein und der selben Zufallsvariable mehrere male. x_1, \dots, x_n . Wenn wir z.B. etwas über die Summe dieser x_i herausfinden möchten, dann ist diese selbst wieder eine Zufallsvariable

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

Wobei alle Zufallsvariablen X_i hierbei die gleiche Verteilung haben wie die ursprüngliche Zufallsvariable X . Beachte, dass wir so viele Zufallsvariablen einführen mussten, wie wir Beobachtungen haben und nicht einfach $\sum_{i=1}^n X$ schreiben können, weil es sonst nicht eindeutig wäre, ob X hier einmal oder mehrere Male ausgewertet wurde.

Wir notieren den Sachverhalt, dass alle ZVn die gleiche Verteilung haben wie X als

$$X_i \sim X$$

und nennen X_i identisch verteilt.

Wenn zusätzlich alle ZVn X_i unabhängig sind, nennt man sie unabhängig und identisch verteilt (engl. *independent and identically distributed*) und schreibt dies als **iid**.

Theorem 4 (Gesetz der großen Zahlen).

Seien X_1, \dots, X_n unabhängig und identisch verteilt (iid) mit Mittelwert μ . Dann

$$\lim_{n \rightarrow \infty} \frac{1}{n} (X_1 + \dots + X_n) = \mu$$

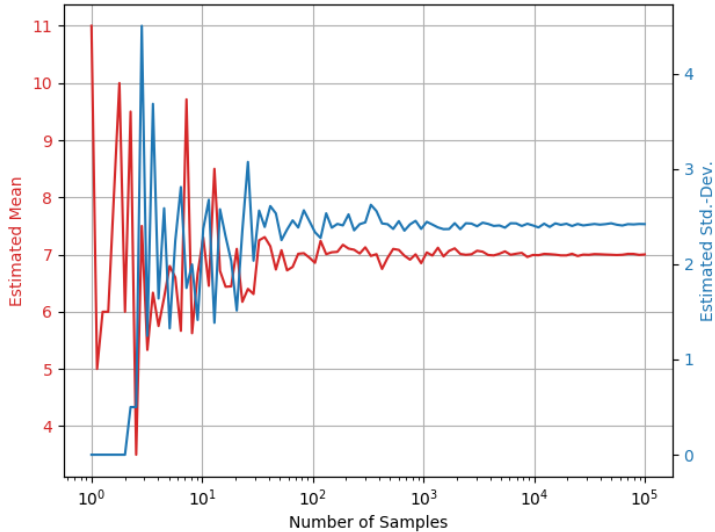
*nahezu sicher*⁴.

Example 5.

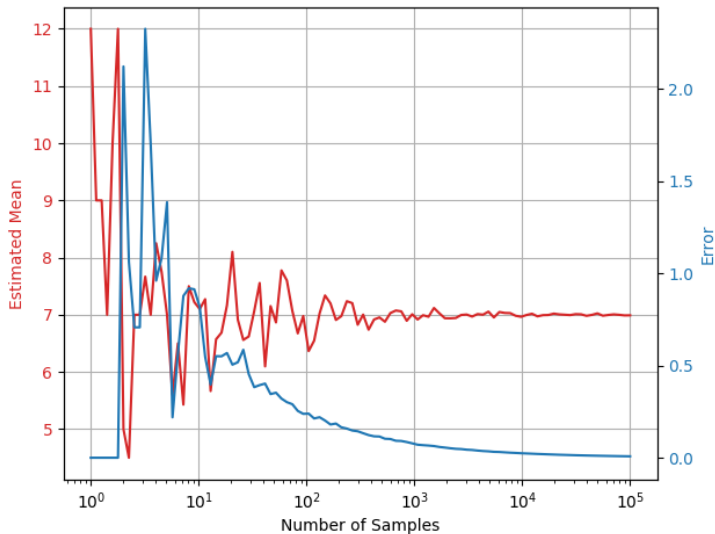
Das Gesetz der großen Zahlen wird angewendet, wenn μ, σ^2 bekannt sind und nicht analytisch abgeleitet werden können. Stichproben werden von f gezogen und μ wird abgeleitet. Wenn μ bekannt ist, kann σ^2 auf die gleiche Weise abgeleitet werden.

⁴ "nahezu sicher" bedeutet mit Wahrscheinlichkeit 1.

Example 6 (Schätzung von μ mit Hilfe des Gesetzes der großen Zahlen).



Example 6 (Schätzung von μ mit Hilfe des Gesetzes der großen Zahlen).



Abstandsmaße für Zufallsvariablen

Definition 7 (Kullback Leibler Divergenz).

Gegeben seien zwei diskrete Zufallsvariablen X und Y , dann ist die Kullback-Leibler Divergenz (d_{KL}) definiert als

$$d_{KL}(X||Y) = \sum_i p_i \log \frac{p_i}{q_i}$$

wobei p, q die Wahrscheinlichkeitsvektoren der ZV X und Y sind und es gilt

$$0 \log \frac{0}{q} = 0, \text{ und wenn } p \neq 0, p \log \frac{p}{0} = \infty$$

im kontinuierlichen Fall wird sie definiert als:

$$d_{KL}(X||Y) = \int_{-\infty}^{\infty} f_X(z) \log \frac{f_X(z)}{f_Y(z)} dz$$

Anmerkungen:

- ▶ Wenn $p_i, q_i \neq 0$ ist die Berechnung von d_{KL} leicht:

$$\text{k1} = (\text{p} * \text{numpy}.\log(\text{p}/\text{q})).\text{sum}()$$

Häufig wird ein kleiner Wert $\varepsilon > 0$ zu allen p_i, q_i hinzugefügt um die Behandlung von Spezialfällen zu vermeiden und eine effiziente Berechnung von d_{KL} zu erlauben.

- ▶ Beachte, dass d_{KL} keine symmetrische Funktion ist, wie man es von einer richtigen Abstandsmetrik erwarten würde. Daher wird oft die folgende Abstandsmetrik verwendet:

$$d_{2KL}(X||Y) := d_{KL}(X||Y) + d_{KL}(Y||X)$$

- ▶ Wenn X, Y kontinuierliche ZV's sind und wenn das Integral nicht gelöst werden kann, wird gewöhnlich f_X, f_Y gesampled und d_{KL} genauso bestimmt wie für diskrete ZV's.

Die Standardnormalverteilung

Definition 8 (Verteilung einer Zufallsvariablen).

Wir notieren eine Zufallsvariable mit der Dichtefunktion f durch

$$X \sim f(x) \quad (\text{d.h. } f_X = f)$$

Definition 9 (Normalverteilung).

Eine *normal verteilte Zufallsvariable* X hat eine Dichte

$$f(x) = N(x_0, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-x_0)^2}$$

also ist X normal verteilt (mit Mittelwert x_0 und Varianz σ^2) genau dann wenn

$$X \sim N(x_0, \sigma^2)$$

Was ist “normal” an der Normalverteilung

Theorem 10 (Zentraler Grenzwertsatz).

Seien X_1, \dots, X_n iid Zufallsvariablen mit der gleichen Dichte $X_i \sim f$, Erwartungswert μ und Varianz σ^2 . Die Zufallsvariable

$$Z_n := \frac{1}{\sigma\sqrt{n}} (X_1 + \dots + X_n - n\mu)$$

konvergiert gegen eine standard normal verteilte Zufallsvariable:

$$Z = \lim_{n \rightarrow \infty} Z_n, \quad Z \sim N(0, 1)$$

alternativ: der Durchschnittswert von n Erkenntnissen der selben Verteilung $X_i \sim f$ konvergiert gegen eine normalverteilte ZV

$$\frac{1}{n} \sum_{i=1}^n X_i \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

Zentraler Grenzwertsatz: Sensor Beispiel

Betrachten wir einen Roboter, der die Eingaben, welche durch eine Zufallsvariable Y realisiert werden, von einem verrauschten Sensor liest. Während der Roboter sich bewegt, liest er zu jedem Zeitpunkt t eine Realisierung von y_t von Y .

Es gibt unterschiedliche Arten von Rauschen (additiv, multiplikativ, Gauß'sch, nicht-Gauß'sch, etc.). Meistens wird Rauschen als additiv mit Mittelwert 0 angenommen, was wir vorerst tun werden:

$$Y = x + E$$

wobei x der echte Wert ist und E additives Rauschen unabhängig von x mit Mittelwert 0 und endlicher Varianz σ^2 .

Nehmen wir an, wir können jede Sensormessung zum Zeitpunkt t n -mal wiederholen

$$y_t^{(j)} = x_t + e_t^{(j)}, \quad 1 \leq j \leq n$$

und betrachten den Mittelwert

$$\bar{y}_t = \frac{1}{n} \sum_{j=1}^n y_t^{(j)} = \frac{1}{n} \sum_{j=1}^n x_t + e_t^{(j)} = x_t + \frac{1}{n} \sum_{j=1}^n e_t^{(j)} = x_t + \bar{e}$$

Zentraler Grenzwertsatz: Robotersensor Beispiel

$$\bar{y}_t = x_t + \bar{e}$$

Das Rauschen \bar{e} ist eine Realisierung der Zufallsvariable

$$E_t = \frac{1}{n} \left(E_t^{(1)} + E_t^{(e)} + \dots + E_t^{(n)} \right)$$

Da $E_t^{(j)} \sim E$ iid mit Mittelwert 0 und Varianz σ^2 sind

$$E_t \sim N \left(0, \frac{\sigma^2}{n} \right)$$

Daher kann der Rauschlevel σ^2 durch Wiederholung auf $\sqrt{\frac{\sigma^2}{n}}$ reduziert werden, also reduzieren n wiederholte Messungen den Rauschlevel um einen Faktor von $\frac{1}{\sqrt{n}}$.

Um additives zero mean Sensorrauschen auf 10% zu reduzieren brauchen wir 100 unabhängige Messungen. Für eine Reduktion auf 50% werden nur 4 Wiederholungen benötigt.

Definition 11 (geometrische Verteilung).

Ein Zufallsexperiment mit den Ausgängen “Erfolg” und “Misserfolg” wird so lange durchgeführt, bis “Erfolg” eintritt. Es resultiert ein Zufallsexperiment mit dem Stichprobenraum

$$\mathcal{S} = \{m^n e \mid n \in \mathbb{N}_0, m = \text{“Misserfolg”}, e = \text{“Erfolg”}\}$$

Die Wahrscheinlichkeiten für “Erfolg” sei konstant p . Definiere eine reellwertige Zufallsvariable $X(m^n e) = n + 1$. Dann gilt für das Bild $\text{Im}(X) = \mathbb{N}$, und es gilt:

$$P(X = t) = (1 - p)^{t-1} p$$

$P(X = t)$ ist also eine geometrische Folge $(a_t) = (a_0 q^t)$ mit $a_0 = \frac{p}{1-p}$ und $q = (1 - p)$. Die Verteilung heißt daher *geometrische Verteilung*.

Satz 1.

Sei X eine geometrisch verteilte Zufallsvariable und $p \in (0, 1)$, dann gilt:

$$E(X) = \frac{1}{p}$$

Beweis.

Sei $q = 1 - p$

$$\begin{aligned} E(X) &= \sum_k k \cdot P(X = k) = \sum_k k q^{k-1} p \\ &= p \sum_k k q^{k-1} \\ &= p \cdot (q^0 + 2q^1 + 3q^2 + 4q^3 + \dots) \\ &= p \cdot (q^0 + 2q^1 + 3q^2 + 4q^3 + \dots) \\ &= p \cdot (q^0 + q^1 + q^2 + q^3 + \dots) \\ &\quad + (q^1 + q^2 + q^3 + \dots) \\ &\quad + (q^2 + q^3 + \dots) \\ &\quad + \dots \\ &= p \cdot \sum_{1 \leq k < \infty} \sum_{i=k}^{\infty} q^{i-1} = p \cdot \sum_{1 \leq k < \infty} q^{k-1} \sum_{i=0}^{\infty} q^i \\ &= p \cdot \frac{1}{1-q} \sum_{k=0}^{\infty} q^k = p \cdot \frac{1}{1-q} \cdot \frac{1}{1-q} = \frac{1}{p} \end{aligned}$$

Satz 2 (Markoff'sche Ungleichung).

Sei $W = (\mathcal{S}, \mathcal{E}, P)$, $\mathcal{E} = \{e_1, e_2, \dots\}$ und X eine reellwertige Zufallsvariable, $X \geq 0$.
Dann gilt für alle $t > 0$

$$P(X \geq t) \leq E(X)/t$$

Beweis.

Wir definieren eine zweite Zufallsvariable:

$$Y(e_i) = \begin{cases} t & \text{if } X(e_i) \geq t \\ 0 & \text{otherwise} \end{cases}$$

expaction of this variable is:

$$\begin{aligned} E[Y] &= 0 \cdot P(X < t) + t \cdot P(X \geq t) = t \cdot P(X \geq t) \\ E(X) &\geq E(Y) = t \cdot P(X \geq t) \end{aligned}$$



Satz 3 (Chebyscheff'sche Ungleichung).

Für alle $t > 0$ ist

$$P(|X - E(X)| \geq t) \leq \frac{\sigma_X^2}{t^2}$$

Beweis.

Sei $Y = |X - E(X)|^2$. Mit der Markow'schen Ungleichung folgt:

$$P(Y \geq t^2) \leq \frac{E(Y)}{t^2}$$

Da $t > 0$ gilt:

$$Y = |X - E(X)|^2 \geq t^2 \quad \Leftrightarrow \quad |X - E(X)| \geq t \quad (3)$$

somit folgt die Behauptung

$$P(|X - E(X)| \geq t) \stackrel{(3)}{=} P(Y \geq t^2) \leq \frac{|X - E(X)|^2}{t^2} = \frac{\sigma_X^2}{t^2}$$



Satz 4 (Chernoff'sche Ungleichung).

Sei $0 < p < 1$ und $X = X_1 + \cdots + X_n$ für unabhängige Indikatorvariablen $X_i \in \mathbb{B}$ mit $P(X_i = 1) = p$.

Dann ist $E(X) = np$ und für alle $\delta \in (0, 1)$ gilt:

$$P(X \leq (1 - \delta) E(X)) \leq e^{-E(X)\delta^2/2}$$

Beweis.

ohne Beweis (siehe Skript).

