



ulm university universität
uulm

Angewandte Stochastik

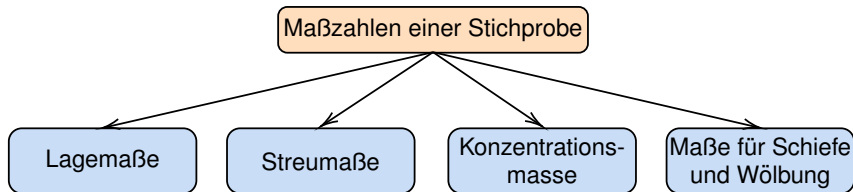
Prof. Dr. Evgeny Spodarev | Vorlesungskurs |

3. Thema

Heutiges Thema

► Beschreibung von Verteilungen

Beschreibung von Verteilungen



- ▶ Es sei eine konkrete Stichprobe (x_1, \dots, x_n) gegeben.
- ▶ Im Folgenden werden Kennzahlen (die sogenannten Maße) dieser Stichprobe betrachtet, welche die wesentlichen Aspekte der der Stichprobe zugrundeliegenden Verteilung wiedergeben:

Beschreibung von Verteilungen

1. Wo liegen die Werte x_i (Mittel, Ordnungsstatistiken, Quantile)? \implies Lagemaße
2. Wie stark streuen die Werte x_i (Varianz)? \implies Streuungsmaße
3. Wie stark sind die Werte x_i in gewissen Bereichen von \mathbb{R} konzentriert? \implies Konzentrationsmaße
4. Wie schief bzw. gewölbt ist die Verteilung von X \implies Maße für Schiefe und Wölbung?

Lagemaße

Man unterscheidet folgende wichtige Lagemaße:

1. Mittelwerte
2. Ordnungsstatistiken und Quantile
3. Modus

1. Mittelwerte

Für eine Stichprobe x_1, \dots, x_n kann man folgende Mittelwerte definieren:

Arithmetisch: $\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i, x_i \in \mathbb{R} \forall i$

Geometrisch: $\bar{x}_n^g = \sqrt[n]{x_1 \dots x_n}, x_i > 0 \forall i$

Harmonisch: $\bar{x}_n^h = \left(\frac{1}{n} \sum_{i=1}^n \frac{1}{x_i} \right)^{-1}, x_i > 0 \forall i$

Gewichtet: $\bar{x}_n^w = \sum_{i=1}^n w_i x_i, w_i \geq 0 \forall i \text{ mit } \sum_{i=1}^n w_i = 1$

Getrimmt: $\bar{x}_n^k = \frac{1}{n-2k} \sum_{i=k+1}^{n-k} x_{(i)}, k \in \mathbb{N}, x_i \in \mathbb{R} \forall i$

1. Mittelwerte

Interessant sind hierbei vor allem:

1. Arithmetisches Mittel:

$\sum_{i=1}^n (x_i - \bar{x}_n) = \sum_{i=1}^n x_i - n\bar{x}_n = 0$, also \bar{x}_n ist der Schwerpunkt des Systems $\{x_i, i = 1, \dots, n\}$ versehen mit Einheitsmaßen.

2. Geometrisches Mittel:

In der Ökonometrie sei B_n ein Faktor der Entwicklung des Marktes (z.B. Zins, Inflationsrate usw.) im Jahr n , B_0 der Ursprungsfaktor und $x_i = \frac{B_i}{B_{i-1}}$ die Veränderungsrate.

Dann gilt

$$B_n = (\sqrt[n]{x_n \dots x_1})^n B_0 = (\bar{x}_n^g)^n B_0.$$

1. Mittelwerte

Außerdem gilt, dass $x_{(1)} \leq \bar{x}_n^h \leq \bar{x}_n^g \leq \bar{x}_n \leq x_{(n)}$ für $x_j > 0$, $j = 1, \dots, n$.

3. Harmonisch:

Sei x_i die Geschwindigkeit der Bewegung des Teils i auf der Produktionslinie der Länge l , für $i = 1, \dots, n$. Dann ist $\frac{l}{x_i}$ die Produktionszeit des Teils i . Die mittlere Produktionszeit ist durch $\frac{1}{n} \sum_{i=1}^n \frac{l}{x_i}$ gegeben und die mittlere Produktionsgeschwindigkeit durch

$$\underbrace{\frac{l}{\frac{1}{n} \sum_{i=1}^n \frac{l}{x_i}}}_{\bar{x}_n^h} \quad \text{definiert.}$$

2. Ordnungsstatistiken und Quantile

Für eine Verteilungsfunktion F sei $F^{-1}(x) = \inf\{y : F(y) \geq x\}$, $x \in (0, 1]$ ihre Quantilfunktion. Das α -Quantil der Verteilung F ist für $\alpha \in [0, 1]$ gegeben durch:

$$\begin{cases} \alpha = 0.25 : F^{-1}(0.25) \text{-unteres Quartil} \\ \alpha = 0.5 : F^{-1}(0.5) \text{-Median} \\ \alpha = 0.75 : F^{-1}(0.75) \text{-oberes Quartil} \end{cases}$$

2. Ordnungsstatistiken und Quantile

Definition

Sei (x_1, \dots, x_n) eine konkrete Stichprobe mit $x_i \in \mathbb{R}$ für $i = 1, \dots, n$. Die *i -te Ordnungsstatistik* der Stichprobe wird definiert durch

$$x_{(i)} = \max \{x_j : \# \{k = 1, \dots, n : x_k \leq x_j\} \geq i\}$$

für $i = 1, \dots, n$.

Hierbei gilt $x_{(1)} = \min_i x_i$, $x_{(n)} = \max_i x_i$.

2. Ordnungsstatistiken und Quantile

Aus obiger Definition lassen sich die empirischen Quantile wie folgt definieren:

1. Empirischer Median:

$$x_{med} = \begin{cases} x_{(\frac{n+1}{2})} & , \text{ falls } n \text{ ungerade} \\ \frac{1}{2} \left(x_{(\frac{n}{2})} + x_{(\frac{n+1}{2})} \right) & , \text{ falls } n \text{ gerade} \end{cases}$$

Insbesondere gilt $x_{med} \approx F^{-1} \left(\frac{1}{2} \right)$ für $n \rightarrow \infty$.

2. Ordnungsstatistiken und Quantile

2. Empirisches α -Quantil:

$$x_{\alpha} = \begin{cases} x_{([n\alpha]+1)} & , \text{ falls } n\alpha \notin \mathbb{N} \\ \frac{1}{2} (x_{([n\alpha])} + x_{([n\alpha]+1)}) & , \text{ falls } n\alpha \in \mathbb{N} \end{cases}$$

für $\alpha \in (0, 1)$. Insbesondere gilt $x_{\alpha} \approx F^{-1}(\alpha)$ für $n \rightarrow \infty$.
 x_{med} ist eine robuste Möglichkeit der Mittelwertbildung, da dieser nicht sensibel bezüglich Ausreißer ist.

Die verschiedenen Quantile lassen sich durch sogenannte **Boxplots** veranschaulichen.

3. Modus

Definition .

1. Sei X eine absolut stetige Zufallsvariable mit Dichte f , welche unimodal ist. Dann ist $x_{mod} = \underset{x}{\operatorname{argmax}} f(x)$ der Modus der Verteilung von X .
2. $\hat{x}_{mod} = \frac{1}{2} (c_{i-1} + c_i)$, wobei das Intervall (c_{i-1}, c_i) die höchste relative Häufigkeit f_i aufweist, heißt der **empirische Modus** der konkreten Stichprobe (x_1, \dots, x_n) . Die relative Häufigkeit f_i ist definiert durch

$$f_i = \frac{\#\{j \in \{1, \dots, n\} : x_j \in (c_{i-1}, c_i)\}}{n}.$$

Lagemasse als Lösungen von Optimierungsaufgaben

Es gilt

$$1. \bar{x}_n = \operatorname{argmin}_a \sqrt{\sum_{i=1}^n (x_i - a)^2}$$

$$2. \bar{x}_{med} = \operatorname{argmin}_a \sum_{i=1}^n |x_i - a|$$

$$3. \hat{x}_{mod} = \frac{1}{2}(c_{i-1} + c_i), \text{ wobei}$$

$$i = \operatorname{argmax}_j \sum_{k=1}^n I(x_k \in (c_{j-1} + c_j]) = \operatorname{argmin}_j \sum_{k=1}^n I(x_k \notin (c_{j-1} + c_j])$$

Streuungsmaße

Bekannte Streuungsmaße einer konkreten Stichprobe (x_1, \dots, x_n) sind die folgenden Größen:

- ▶ **Spannweite** $x_{(n)} - x_{(1)}$,
- ▶ **empirische Varianz** $\bar{s}_n^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2$,
- ▶ **Stichprobenvarianz** $s_n^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2 = \frac{n}{n-1} \bar{s}_n^2$,
- ▶ **empirische Standardabweichungen** $\bar{s}_n = \sqrt{\bar{s}_n^2}$, $s_n = \sqrt{s_n^2}$,
- ▶ **empirischer Variationskoeffizient** $\gamma_n = \frac{s_n}{\bar{x}_n}$, falls $\bar{x}_n > 0$.

Streuungsmaße

- ▶ Die Spannweite zeigt die *maximale Streuung* in den Daten, wobei sich die empirische Varianz mit der *mittleren quadratischen Abweichung* vom Stichprobenmittel auseinandersetzt.
- ▶ Hier sind einige Eigenschaften von \bar{s}_n^2 (bzw. s_n^2 , da sie sich nur durch einen Faktor unterscheiden):

Lemma.

1. Für jedes $b \in \mathbb{R}$ gilt

$$\sum_{i=1}^n (x_i - b)^2 = \sum_{i=1}^n (x_i - \bar{x}_n)^2 + n(\bar{x}_n - b)^2$$

und somit für $b = 0$

$$\bar{s}_n^2 = \frac{1}{n} \sum_{i=1}^n (x_i^2 - \bar{x}_n^2)$$

bzw.
$$s_n^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i^2 - \bar{x}_n^2) .$$

Lemma.

2. Transformationsregel:

Falls die Daten (x_1, \dots, x_n) linear transformiert werden, d.h. jedes y_i lässt sich darstellen als $y_i = ax_i + b$, $a \neq 0$, $b \in \mathbb{R}$, dann gilt

$$\bar{s}_{n,y}^2 = a^2 \bar{s}_{n,x}^2 \quad \text{bzw.} \quad \bar{s}_{n,y} = |a| \bar{s}_{n,x},$$

wobei

$$\bar{s}_{n,y}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y}_n)^2, \quad \bar{s}_{n,x}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2.$$

- ▶ Der Skalierungsunterschied zwischen \bar{s}_n^2 und s_n^2 ist den Eigenschaften der **Erwartungstreue** von s_n^2 zu verdanken, die besagt, dass für eine Zufallsstichprobe (X_1, \dots, X_n) mit X_i unabhängig identisch verteilt, $X_i \sim X$, $\text{Var } X = \sigma^2 \in (0, \infty)$ gilt, dass $\text{Es}_n^2 = \sigma^2$, wobei
$$\text{E}\bar{s}_n^2 = \frac{n}{n-1}\sigma^2 \xrightarrow{n \rightarrow \infty} \sigma^2.$$
- ▶ Das heißt, während bei der Verwendung von s_n^2 zur Schätzung von σ^2 kein Fehler "im Mittel" gemacht wird, ist diese Aussage für \bar{s}_n^2 nur asymptotisch (für große Datenmengen n) richtig.
- ▶ Aufgrund von $\sum_{i=1}^n (x_i - \bar{x}_n) = 0$ ist z.B. $x_n - \bar{x}_n$ durch $x_i - \bar{x}_n$, $i = 1, \dots, n-1$ bestimmt.

- ▶ Somit verringert sich die **Anzahl der Freiheitsgrade** in der Summe $\sum_{i=1}^n (x_i - \bar{x}_n)^2$ um 1 und somit scheint die Normierung $\frac{1}{n-1}$ plausibel zu sein.
- ▶ Die **Standardabweichungen** \bar{s}_n und s_n werden verwendet, damit man die selben Einheiten (und nicht ihre Quadrate, also z.B. Euro und nicht Euro²) erhält. Für normalverteilte Stichproben ($X \sim N(\mu, \sigma^2)$) liefert \bar{s}_n auch die "k-Sigma-Regel", die besagt, dass in den Intervallen

$[\bar{x}_n - \bar{s}_n, \bar{x}_n + \bar{s}_n]$	ca.	68% ,
$[\bar{x}_n - 2\bar{s}_n, \bar{x}_n + 2\bar{s}_n]$	ca.	95% ,
$[\bar{x}_n - 3\bar{s}_n, \bar{x}_n + 3\bar{s}_n]$	ca.	99%

aller Daten liegen.

- ▶ Der Vorteil vom *empirischen Variationskoeffizienten* ist, dass er *maßstabsunabhängig* ist und somit den Vergleich von Streuungseigenschaften unterschiedlicher Stichproben zulässt.

Konzentrationsmaße

- ▶ Insbesondere in den Wirtschaftswissenschaften interessiert man sich oft für die Konzentration von Merkmalsausprägungen in der Stichprobe, z.B. wie sich das Familieneinkommen einer demographischen Einheit auf unterschiedliche Einkommensbereiche (Vielverdiener, Mittelstand, Wenigverdiener) aufteilt, oder wie sich der Markt auf Marktanbieter aufteilt (Marktkonzentration).
- ▶ Dabei ist es wünschenswert, diese Relation mit Hilfe weniger Zahlen oder einer Grafik zum Ausdruck zu bringen.

Konzentrationsmaße

Dies ist mit Hilfe folgender Stichprobenfunktionen möglich:

- *Lorenzkurve L* ,
- *Gini-Koeffizient G* ,
- *Konzentrationsrate CR_g* ,
- *Herfindahl-Index H* .

Lorenzkurve

- ▶ Die Lorenzkurve wurde von M. Lorenz am Anfang des 20. Jahrhunderts für die Charakterisierung der Vermögenskonzentration benutzt.
- ▶ Sei (x_1, \dots, x_n) eine Stichprobe, die in aufsteigender Reihenfolge geordnet werden muss: $(x_{(1)}, \dots, x_{(n)})$.
- ▶ Die **Lorenzkurve** verbindet Punkte

$$(0, 0), (u_1, v_1), \dots, (u_n, v_n), (1, 1)$$

durch Liniensegmente, wobei $u_j = j/n$ der Anteil der j kleinsten Merkmalsträger und $v_j = \sum_{i=1}^j x_{(i)} / \sum_{i=1}^n x_i$ die kumulierte relative Merkmalssumme ist.

Lorenzkurve

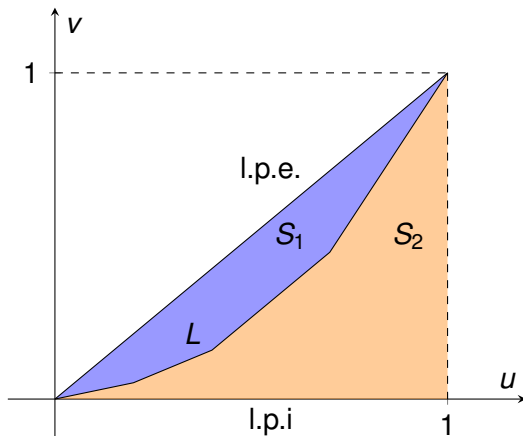


Figure: Abbildung einer typischen Lorenzkurve

Lorenzkurve

- ▶ Der Grundgedanke ist darzustellen, welcher Anteil des Merkmalsträgers auf welchen Anteil der Gesamtmerkmalssumme entfällt.
- ▶ Zum Beispiel lassen sich dadurch Aussagen wie etwa "Auf 20% aller Haushalte im Land entfällt 78% des Gesamteinkommens" machen.
- ▶ Eine Interpretation der Lorenzkurve L ist nur an den Knoten (u_j, v_j) möglich: "Auf $u_j \cdot 100\%$ der kleinsten Merkmalsträger konzentrieren sich $v_j \cdot 100\%$ der Merkmalssumme".
- ▶ Dabei liegt L auf $[0, 1]^2$ immer zwischen der "line of perfect equality" (l.p.e.) $v_i = u_i \quad \forall i$ (Einkommen ist absolut gleichmäßig—also "gerecht"—verteilt) und "line of perfect inequality" (l.p.i.) $v = 0, u \in [0, 1)$ und $(1, 1)$ (das Gesamteinkommen besitzt nur die reichste Familie) und ist immer monoton und konvex.

Lorenzkurve

- ▶ Auf Modellebene gibt es ein Analogon der Lorenzkurve.
- ▶ Dieses ist

$$L = \left\{ (u, v) \in [0, 1]^2 : v = \frac{\int_0^u F_X^{-1}(t) dt}{\int_0^1 F_X^{-1}(t) dt}, \quad u \in [0, 1] \right\},$$

wobei

$$EX = \int_0^1 F_X^{-1}(t) dt.$$

- ▶ Dementsprechend können die Knoten (u_i, v_j) der oben eingeführten empirischen Lorenzkurve als

$$v_j = \frac{\sum_{i=1}^j \frac{x_{(i)}}{n}}{\bar{x}_n}$$

interpretiert werden.

Gini-Koeffizient

Der **Gini-Koeffizient** G ist gegeben durch $G = S_1/S_2$, wobei S_1 die Fläche zwischen der Lorenzkurve L und der Diagonalen $v = u$, S_2 die Fläche zwischen der Diagonalen und der u -Achse ($= 1/2|[0, 1]^2| = 1/2$) ist.

Satz. (Darstellung des Gini-Koeffizienten)

Es gilt

$$G = 2S_1 = \frac{2 \sum_{i=1}^n ix_{(i)}}{n \sum_{i=1}^n x_i} - \frac{n+1}{n}.$$

Gini-Koeffizient

- Es gilt $G \in [0, \frac{(n-1)}{n}]$, wobei

$$G_{\min} = 0 \quad \text{bei } x_1 = x_2 = \dots = x_n \quad \text{"p.e."},$$

$$G_{\max} = \frac{n-1}{n} \quad \text{bei } x_1 = \dots = x_{n-1} = 0, x_n \neq 0 \quad \text{"p.i."}$$

- Somit hängt G_{\max} vom Datenumfang ab.
- Um dies zu vermeiden, betrachtet man oft den normierten Gini-Koeffizienten

$$G^* = \frac{G}{G_{\max}} = \frac{n}{n-1} G \in [0, 1]$$

(Lorenz-Münzner-Koeffizient).

Konzentrationsrate CR_g :

- ▶ Die Lorenzkurve und der Gini-Koeffizient betrachten die *relative Konzentration*, wie etwa bei der Fragestellung "Wieviel % der Familien teilen sich wieviel % des Gesamteinkommens?".
- ▶ Dabei beantwortet die Konzentrationsrate die Frage "Wieviele Familien haben wieviel Prozent des Gesamteinkommens?" für die g reichsten Familien.
- ▶ Somit wird auch die absolute Anzahl aller Familien berücksichtigt.

Konzentrationsrate CR_g :

- ▶ Sei $g \in \{1, \dots, n\}$ und seien $x_{(1)} \leq \dots \leq x_{(n)}$ die Ordnungsstatistiken der Stichprobe (x_1, \dots, x_n) .
- ▶ Für $i \in \{1, \dots, n\}$ sei

$$p_i = \frac{x_{(i)}}{\sum_{j=1}^n x_j} = \frac{x_{(i)}}{n\bar{x}_n}$$

der Merkmalsanteil der i -ten Einheit.

- ▶ Dann gibt die **Konzentrationsrate** $CR_g = \sum_{i=n-g+1}^n p_i$ wider, welcher Anteil des Gesamteinkommens von den g reichsten Familien gehalten wird.

Herfindahl-Index

- ▶ Der **Herfindahl-Index** ist definiert durch $H = \sum_{i=1}^n p_i^2$, wobei der Merkmalsanteil p_i wie oben definiert ist.
- ▶ Bei der gleichen Verteilung des Einkommens ($x_1 = x_2 = \dots = x_n$) gilt $H_{min} = 1/n$, bei völlig ungerechter Verteilung ($x_1 = \dots = x_{n-1} = 0, x_n \neq 0$) gilt $H_{max} = 1$.
- ▶ Sonst gilt $H \in [H_{min}, H_{max}]$, also $1/n \leq H \leq 1$.
- ▶ H ist umso kleiner, je gerechter das Gesamteinkommen verteilt ist.

Maße für Schiefe und Wölbung

- ▶ der Verteilung einer Zufallsvariable X sind:
Schiefe oder Symmetriekoeffizient:

$$\gamma_1 = \frac{\mu'_3}{\sigma^3} = E(\tilde{X}^3),$$

wobei

$$\mu'_k = E(X - EX)^k, \quad \sigma^2 = \mu'_2 = \text{Var } X, \quad \tilde{X} = \frac{X - EX}{\sigma}.$$

Wölbung (Exzess):

$$\gamma_2 = \frac{\mu'_4}{\sigma^4} - 3 = E(\tilde{X}^4) - 3,$$

vorausgesetzt, dass $E(X^4) < \infty$.

Maße für Schiefe und Wölbung

- Falls nun das Merkmal X statistisch in einer Stichprobe (x_1, \dots, x_n) beobachtet wird, wie können γ_1 und γ_2 aus diesen Daten geschätzt und interpretiert werden?
- Als Schätzer für das k -te zentrierte Moment $\mu'_k = E(X - EX)^k$, $k \in \mathbb{N}$ schlagen wir

$$\hat{\mu}'_k = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^k$$

vor, die Varianz σ^2 wird durch

$$\bar{s}_n^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2$$

geschätzt.

Maße für Schiefe und Wölbung

- Somit bekommt man den Momentenkoeffizienten der Schiefe (engl. "skewness")

$$\hat{\gamma}_1 = \frac{\hat{\mu}_3'}{\bar{s}_n^3} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^3}{\left(\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2 \right)^{\frac{3}{2}}}.$$

- Falls die Verteilung von X linksschief ist, überwiegen negative Abweichungen im Zähler und somit gilt $\hat{\gamma}_1 < 0$ für linksschiefe Verteilungen.
- Analog gilt $\hat{\gamma}_1 \approx 0$ für symmetrische und $\hat{\gamma}_1 > 0$ für rechtsschiefe Verteilungen.

Maße für Schiefe und Wölbung

- Das **Wölbungsmaß von Fisher** (engl. "kurtosis") ist gegeben durch

$$\hat{\gamma}_2 = \frac{\hat{\mu}'_4}{\bar{s}_n^4} - 3 = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^4}{\left(\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2\right)^2} - 3.$$

- Falls $\hat{\gamma}_2 > 0$ so ist die Verteilung von X steilgipflig, für $\hat{\gamma}_2 < 0$ ist sie flachgipflig. Falls $X \sim N(\mu, \sigma^2)$, so gilt $\hat{\gamma}_2 \approx 0$.
- Ursache: flachgipflige Verteilungen haben schwerere Tails als die steilgipfligen.
- Als Maß dient dabei die Normalverteilung, für die $\gamma_1 = \gamma_2 = 0$ und somit $\hat{\gamma}_1 \approx 0$, $\hat{\gamma}_2 \approx 0$.
- So definiert sind $\hat{\gamma}_1$ und $\hat{\gamma}_2$ nicht resistent gegenüber Ausreißern.

Maße für Schiefe und Wölbung

- Eine robuste Variante von $\hat{\gamma}_1$ ist beispielsweise durch den sogenannten *Quantilkoeffizienten der Schiefe*

$$\hat{\gamma}_q(\alpha) = \frac{(x_{1-\alpha} - x_{med}) - (x_{med} - x_{\alpha})}{x_{1-\alpha} - x_{\alpha}}, \quad \alpha \in (0, \frac{1}{2})$$

gegeben.

Maße für Schiefe und Wölbung

- ▶ Für $\alpha = 0,25$ erhält man den Quartilkoeffizienten.
- ▶ $\hat{\gamma}_q(\alpha)$ misst den Unterschied zwischen der Entfernung des α - und $(1 - \alpha)$ -Quantils zum Median.
- ▶ Bei linkssteilen (bzw. rechtssteilen) Verteilungen liegt das (untere) x_α -Quantil näher an (bzw. weiter entfernt von) dem Median.
- ▶ Somit gilt
 - $\hat{\gamma}_q(\alpha) > 0$ für linkssteile Verteilungen,
 - $\hat{\gamma}_q(\alpha) < 0$ für rechtssteile Verteilungen,
 - $\hat{\gamma}_q(\alpha) = 0$ für symmetrische Verteilungen.
- ▶ Durch das zusätzliche Normieren (Nenner) gilt $-1 \leq \hat{\gamma}_q(\alpha) \leq 1$.