

Statistik (Business Statistics) / Studium der Verteilung

Deskriptive Statistik: Darstellung + Zusammenfassung beobachteter Daten

Induktive Statistik: (Inferenz) Analyse mittels math. Modellen
Wahrscheinlichkeitstheorie: mathematisches Fundament ↗

Variable: spezielle Eigenschaft eines Individuums

⇒ Variablenwerte: • Kategorisch/qualitativ

-Name

-Werte

-Masseinheit


⇒ vorgegebene Liste von Werte

• numerisch/quantitativ

⇒ eine Zahl mit der man umgeht

Diagramme

Blockdiagramm 

Kuchendiagramm 

Histogramm ⇒ • symmetrisch / asymmetrisch

• Unimodal (nur 1 MaxWert)

• Bimodal (zwei MaxWerte)

• Rechtsschief (links > rechts)

• Linksschief (rechts > links)

• Uniform/Gleichverteilt (alles gleich hoch)

Zeit-Diagramm (x: Zeit, y: Wert der Stichprobe)

Kennzahlen • Lage (Mittelpunkt, der Verteilung)

• Streuung (Abweichung vom Mittelpunkt)

Arithmetisches Mittel: $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

Median: Wert in der Mitte der sortierten Daten

$\bar{x} < m$: Verteilung linksschief

$\bar{x} > m$: Verteilung rechtsschief

$\bar{x} \approx m$: Verteilung \approx symmetrisch

Perzentile: Verallgemeinerung des Medians

$\Rightarrow Q_1$, 25%, erstes Quartil

Q_2 , 50%, Median (zweites Quartil)

Q_3 , 75%, drittes Quartil

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


Standardabweichung: $s = \sqrt{s^2}$

Interquartilabstand: $IQA = Q_3 - Q_1$

Ausreisser: Wert ausserhalb von

$[Q_1 - 1,5 \times IQA, Q_3 + 1,5 \times IQA]$

Fünf-Kennzahlen-Zusammenfassung: $FKZ = (Min, Q_1, m, Q_3, Max)$

\Rightarrow Boxplot: Visualisierung FKZ 

Sammeln von Daten

Stichprobendesign:

- Bevölkerung: Gesamtheit der Individuen, über die eine statistische Aussage gemacht werden soll
- Stichprobe: Teil der Bevölkerung, über den Informationen gesammelt wurden

\Rightarrow Stichprobe muss repräsentativ sein!

- verzerrt (biased): bevorzugt/benachteiligt Segmente der VS. Bevölkerung systematisch
- • unverzerrt (unbiased): keine Segmente bevorzugt/benachteiligt

Schlechte Stichproben:

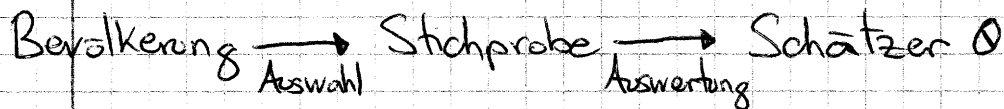
- Freiwillige-Antwort: nur interessierte Teilnehmer mit Zugang zur Stichprobe antworten, Negatives fällt oft aus (da keine Lust zu beantworten usw.)
- Ich-mach's-mir-leicht: nicht zufällig, sondern Situationsabhängig ausgewählt, um die Arbeit zu erleichtern (z.B. vor Ort, spezielles Medium, ...)

Datenherkunft:

- Beobachtungsdaten (durch Befragung gesammelt ohne die Individuen vorher zu manipulieren)
- Experimentelle Daten (durch Erzeugung spezieller Situation gewonnen)

Parameter (Kennzahl der gesamten Bevölkerung)

Statistik (Kennzahl der einen Stichprobe)



Stichprobenverteilung (\Rightarrow Aussagekraft der Statistik)

- Verzerrung des Schätzers: Abstand des Schätzermittels $E(\hat{\theta})$ zum tatsächlichen θ
 \Rightarrow unverzerrt: $E(\hat{\theta}) = \theta$, sonst verzerrt
- Variabilität des Schätzers: Abweichung vom Parameter
 \uparrow Stichproben = \downarrow Variabilität

Wahrscheinlichkeitstheorie

- Zufallsexperiment: Vorgang der nach einer bestimmten Vorschrift ausgeführt wird, beliebig oft wiederholbar ist und zufallsabhängig ein Ergebnis liefert
 - Elementarereignis/Ergebnis: möglicher Ausgang des Experiments
 - Ergebnisraum: Menge aller möglichen Ergebnisse des Experiments
 - Ereignis: Teilmenge des Ergebnisraums
- unmögliches Ereignis: $A = \emptyset$
 - sicheres Ereignis: $A = S$
 - Komplementäre Ereignis: $A^c = S \setminus A$
 - Durchschnitt: $A \cap B$
 - Vereinigung: $A \cup B$
 - Disjunkte Ereignisse: $A, B, A \cap B = \emptyset$
 - Ausschöpfende Ereignisse: $A, B, A \cup B = S$
 - Partition: des Ergebnisraums, eine Folge A_1, \dots, A_n von ausschöpfen und paarweise disjunkten Ereignissen

Wahrscheinlichkeit $P(A)$, P : Wahrscheinlichkeitsfunktion, A : Ereignis

- LAPLACE: alle Ergebnisse gleich wahrscheinlich \Rightarrow Zählregel
$$P(A) = \frac{\text{Anzahl der für } A \text{ günstigen Fälle}}{\text{Anzahl aller möglichen Fälle}} = \frac{\#A}{\#S}$$
- STATISTISCH: Grenzwert der Häufigkeiten wenn die Anzahl Versuche gross wird
$$P(A) = \lim_{\text{Versuche} \rightarrow \infty} \frac{\text{Anzahl der für } A \text{ günstigen Versuche}}{\text{Anzahl Versuche}}$$
- AXIOMATISCH/KOLMOGOROFF: abstrakte Funktion mit:
 - 1) $P(A) \in \mathbb{R}, P(A) \geq 0$
 - 2) $P(S) = 1$
 - 3) disjunkt A, B gilt: $P(A \cup B) = P(A) + P(B)$

• Normierung: $0 \leq P(A) \leq 1$

• unmögliches Ereignis $P=0$, sicheres Ereignis $P=1$

• $P(A^c) = 1 - P(A)$

• Monotonie: $A \subseteq B \Rightarrow P(A) \leq P(B)$

• Zählregel: $P(A) = \sum_{\omega \in A} P(\{\omega\})$

• Zerlegung: $P(A) = P(A \cap B) + P(A \cap B^c)$

• Additionssatz: $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

Bedingte Wahrscheinlichkeit: $P(B|A) = \frac{P(A \cap B)}{P(A)}$

Wahrscheinlichkeitsbaum

\Rightarrow Multiplikationssatz: $P(A \cap B) = P(A) \cdot P(B|A)$

Bayes'sche Formel: $P(A|B) = \frac{P(A) \cdot P(B|A)}{P(B)}$

Unabhängige Ereignisse: A, B unabhängig falls:

$P(A \cap B) = P(A) \cdot P(B)$ (Multiplikationssatz)

Zufallsvariable: Funktion $X: S \rightarrow \mathbb{R}$, die jedem Elementarereignis eine Zahl zuordnet

Realisationen: Werte einer Zufallsvariable

\Rightarrow Wahrscheinlichkeit $P(X=x)$, dass die ZV X den Wert x annimmt

• diskrete ZV: endlich oder abzählbar viele Werte

• stetige ZV: alle Werte in einem Intervall

Wahrscheinlichkeitsfunktion $P(X=x_i)$ von ZV X

$0 \leq P(X=x_i) \leq 1$ zwischen 0 und 1

$\sum_i P(X=x_i) = 1$ Summe aller immer 1

Verteilungsfunktion: $F(x) = P(X \leq x) = \sum_{x_i \leq x} P(X = x_i)$

Mittelwert und Varianz charakterisieren eine ZV.

⇒ genau wie bei statistischen Variablen, aber mit der Wahrscheinlichkeit dass ZV X den Wert x_i annimmt

Erwartungswert: $\mu_x = E(X) = \sum_i x_i \cdot P(X = x_i)$

Varianz: $\sigma_x^2 = \text{Var}(X) = E((X - \mu)^2) = \sum_i (x_i - \mu_x)^2 \cdot P(X = x_i)$

Standardabweichung: $\sigma_x = \sqrt{\sigma_x^2}$ $\hookrightarrow \sigma^2 = E(X^2) - (E(x))^2$

Gesetz der grossen Zahlen: (stat. Var. durch ZV ersetzen)

geht der Stichprobenumfang n gegen ∞ ,
so wird $\bar{x} \approx \mu_x$ und $s \approx \sigma_x$

• Bernoulli-Verteilung: Münzwurf, abhängig von p (W. Erfolg)

$$X \sim \text{Bern}(p)$$

$$P(X=1) = p, P(X=0) = 1-p$$

$$\mu_x = p$$

$$\sigma_x^2 = p(1-p)$$

$$\sigma_x = \sqrt{p(1-p)}$$

• Binomial-Verteilung: Bernoulli, repetiert n Mal

$$X \sim \text{Bin}(n, p)$$

$$P(X=k) = \binom{n}{k} \cdot p^k \cdot (1-p)^{n-k} \quad \left(\binom{n}{k} = \frac{n!}{k! \cdot (n-k)!} \right)$$

$$\mu_x = np$$

$$\sigma_x^2 = np(1-p)$$

$$\sigma_x = \sqrt{np(1-p)}$$

• Poissonverteilung: Erfolg pro Einzelversuch gering, aber Versuche unendlich oft wiederholt werden

$$X \sim P_0(\mu)$$

$$P(X=x) = \frac{\mu^x \cdot e^{-\mu}}{x!}$$

($\mu = \frac{n}{p}$ um Binomial zu ersetzen)

$$\mu_x = \mu$$

$$\sigma_x^2 = \mu$$

$$\sigma_x = \sqrt{\mu}$$

Stetige Wahrscheinlichkeitsverteilung (Intervall)

stetige Verteilung \Rightarrow Dichtefunktion

$$f(x) \geq 0 \text{ und } \int_{-\infty}^{\infty} f(x) dx = 1$$

$$P(a \leq X \leq b) = \int_a^b f(x) dx = F(b) - F(a)$$

Verteilungsfunktion: $F(x) = P(X \leq x) = \int_{-\infty}^x f(t) dt$

$$\Rightarrow f(x) = F'(x) \Leftrightarrow \text{dichte} = d \text{ Verteilung}$$

Erwartungswert: $\mu_x = E(x) = \int_{-\infty}^{\infty} x \cdot f(x) dx$

Varianz: $\sigma_x^2 = \text{Var}(x) = E((X-\mu)^2) = \int_{-\infty}^{\infty} (x-\mu_x)^2 \cdot f(x) dx$

Standardabweichung: $\sigma_x = \sqrt{\sigma_x^2}$

Perzentil: $p\% = \int_{-\infty}^{x_{p\%}} f(x) dx$

$$P(X \leq x_{p\%}) = F(x_{p\%}) = p\%$$

Median: 50% Perzentil mit $F(m) = \frac{1}{2}$

Gleichverteilung:

$$X \sim U[a, b]$$

$$\mu_x = \frac{1}{2} (b+a)$$

$$\sigma_x^2 = \frac{1}{12} (b-a)^2$$

$$\sigma_x = \sqrt{\frac{1}{12}} (b-a)$$

Rechenregeln: X, Y sind ZV

$$E(X+Y) = E(X) + E(Y) \quad \text{addition Erwartungswerte}$$

$$E(a+b \cdot X) = a + b \cdot E(X) \quad \text{skalierung Erwartungswerte}$$

$$\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) \quad \text{addition Varianzen (unabhängig)}$$

$$\text{Var}(a+b \cdot X) = b^2 \cdot \text{Var}(X) \quad \text{skalierung Varianzen (quadratisch)}$$

$$\sigma_{x+y} = \sqrt{\sigma_x^2 + \sigma_y^2}, \quad \sigma(a+b \cdot X) = |b| \cdot \sigma(X)$$

Standardisierung: $Z = \frac{X - \mu_x}{\sigma_x}, \quad E(Z) = 0, \sigma_z = 1$

• Gleichverteilung: $X \sim U[a, b]$

$$f(x) = \frac{1}{b-a} \quad \text{im Intervall } [a, b]$$

$$\mu_x = \frac{1}{2} (b+a)$$

$$\sigma_x^2 = \frac{1}{12} (b-a)^2$$

$$\sigma_x = \sqrt{\frac{1}{12}} (b-a)$$

• Normalverteilung: $X \sim N(\mu, \sigma^2)$

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2} \Rightarrow \text{Gauss'sche Glockenkurve}$$

$$\mu_x = \mu$$

$$\sigma_x^2 = \sigma^2$$

$$\sigma_x = \sigma = \sqrt{\sigma^2}$$

$$\Phi(x) = \frac{1}{\sigma \sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{1}{2} \left(\frac{t-\mu}{\sigma}\right)^2} dt$$

\Rightarrow Tabellen benutzen, zuerst aber standardisieren!

$$P(X \leq a) = P\left(Z \leq \frac{a-\mu}{\sigma}\right) = \Phi\left(\frac{a-\mu}{\sigma}\right)$$

$$P(X \geq a) = 1 - P(X \leq a) = 1 - \Phi\left(\frac{a-\mu}{\sigma}\right)$$

$$P(b \leq X \leq a) = \Phi\left(\frac{a-\mu}{\sigma}\right) - \Phi\left(\frac{b-\mu}{\sigma}\right)$$

Perzentile: α -te Perzentile, z_α

$$P(X \leq x) = 10\% = 0,1$$

X standardisieren $\Rightarrow Z$ (z.B. $N(3,4)$)

$$z_\alpha = z_{0,1} = -z_{0,9} = -1,28$$

$$0,1 = 10\% = P(Z \leq z_{0,1}) = P(Z \leq -1,28)$$

standardisierung rückgängig

$$P\left(\frac{X-3}{2} \leq -1,28\right) = P(X \leq 0,44)$$

Approximation der Binomialverteilung:

gute Approximation bei $n \cdot p \geq 10$, $np(1-p) \geq 10$

$\text{Bin}(n, p)$ mit Stetigkeitskorrektur:

$$X \sim \text{Bin}(np) \Rightarrow Y \sim N(np, np(1-p))$$

$$\Rightarrow P(X=x) \approx P(x-0,5 \leq Y \leq x+0,5) = P\left(\frac{x-0,5-\mu}{\sigma} \leq Z \leq \frac{x+0,5-\mu}{\sigma}\right)$$

Stichprobenverteilung:

Stichprobenerhebung \Rightarrow Erkenntnisse über die Verteilung

\Rightarrow Schätzen von μ, σ

Verteilung des Mittelwerts:

Zentraler Grenzwertsatz: $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ strebt mit grossen

gegen eine Normalverteilung mit: $\mu_{\bar{X}} = \mu$

$$\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}$$

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

Induktive Statistik (Stichprobe \Rightarrow Bevölkerung)

Parameter θ , Schätzer $\hat{\theta}$

\Rightarrow • Verzerrung: $\text{Bias} = E(\hat{\theta}) - \theta$

unverzerrt bei $E(\hat{\theta}) = \theta$

• Variabilität: $\text{Var}(\hat{\theta})$

• Erwartete quadratische Abweichung:

$$(\hat{\theta} - \theta)^2 = \text{Bias}^2 + \text{Var}(\hat{\theta})$$

Verteilungsmittelwert $\mu \Rightarrow$ Stichprobenmittel $\bar{X} \Rightarrow \text{Var}(\bar{X}) = \frac{\sigma^2}{n}$

Verteilungsvarianz $\sigma^2 \Rightarrow$ Stichprobenvarianz $\bar{\sigma}^2 \Rightarrow \text{Var}(\bar{\sigma}^2) = \frac{\sigma^2}{n}$

Eigenschafts-Anteil $p \Rightarrow$ Anteil in der Stichprobe $\hat{p} \Rightarrow \text{Var}(\hat{p}) = \frac{p(1-p)}{n}$

$$\Rightarrow \text{Var}(\hat{p}) = \frac{p(1-p)}{n} \left(\frac{1}{b^2}\right) \text{SA}(\hat{p}) = \sqrt{\frac{p(1-p)}{n}} \left(\frac{1}{b}\right)$$

wenn man SA um b verkleinern will, muss n um b^2 multipliziert werden!

Konfidenzintervalle

\Rightarrow Intervall um Schätzer $\hat{\theta}$, indem Parameter θ mit einer gewissen Konfidenz enthalten ist.

Konfidenzniveau $1-\alpha$

Konfidenzgleichung:

\Rightarrow Konfidenzintervall = Schätzer \pm FM = Schätzer $\pm z_{1-\frac{\alpha}{2}} \times \text{SF}$

FM = Fehlermarge = $z_{1-\frac{\alpha}{2}} \times \text{SF}$

SF = Standardfehler = $\frac{s}{\sqrt{n}}$

s = Standardabweichung

n = Stichprobenzahl

z = Perzentilfunktion der Standardnormalverteilung

$z = z_{1-\frac{\alpha}{2}}$, so dass $P(-z \leq Z \leq z) = 1-\alpha$

α = Signifikanz

Schätzung des Mittelwerts:

Gegeben $1-\alpha, n, s, \bar{x}$

$$\text{SF} = \frac{s}{\sqrt{n}}$$

FM = $z_{1-\frac{\alpha}{2}} \times \text{SF}$, Perzentil z der Standardnormalverteilung

$$[a, b] = \bar{x} \pm z_{1-\frac{\alpha}{2}} \cdot \frac{s}{\sqrt{n}}$$

Schätzung eines Anteils einer Eigenschaft:

Gegeben $1-\alpha, n, \hat{p}$

$$SF = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$$[a, b] = \hat{p} \pm z_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

Schätzung einer Differenz von Mittelwerten:

$$\Delta = \mu_1 - \mu_2$$

$$SF = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$[a, b] = \bar{x}_1 - \bar{x}_2 \pm z_{1-\frac{\alpha}{2}} \cdot \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

Schätzung einer Differenz von Anteilen:

$$\Delta = p_1 - p_2$$

$$SF = \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

$$[a, b] = \hat{p}_1 - \hat{p}_2 \pm z_{1-\frac{\alpha}{2}} \cdot SF$$

Frage nach der Stichprobenzahl:

Gegeben $FM, s, 1-\alpha$

$$n \geq \left(\frac{z_{1-\frac{\alpha}{2}} \cdot s}{FM} \right)^2$$

Hypothesentests \Rightarrow Aussage über Wert des Parameters bestätigen

Hypothese: H_0 , Nullhypothese

H_A , Alternativhypothese

	H_0 korrekt	H_0 falsch
H_0 verworfen	Fehler Typ 1	OK
H_0 beibehalten	OK	Fehler Typ 2

Signifikanz α kontrolliert nur Fehler Typ 1.

Güte eines Tests: $1 - P(\text{Fehler Typ 2})$

- grosses $\alpha \Rightarrow$ grosse Güte
- Je weiter μ von μ_0 entfernt, grössere Güte
- Je kleiner Standardabweichung s , grössere Güte
- Je grösser Probenanzahl n , grössere Güte
- einseitige Tests grössere Güte als zweiseitige

Hypothesentests für Mittelwert μ

$$H_0: \mu = \mu_0$$

Testvariable: $Z = \frac{\text{Schätzer-Vermutung}}{\text{Standardfehler}} = \frac{\bar{X} - \mu_0}{s/\sqrt{n}}$

Teststatistik: $z = \frac{\text{Konkretes Ergebnis-Vermutung}}{\text{Standardfehler}} = \frac{\bar{z} - \mu_0}{s/\sqrt{n}}$

p-Wert:
$$PW = \begin{cases} P(Z \geq z) & \text{falls } H_A: \mu > \mu_0 \\ P(Z \leq z) & \text{falls } H_A: \mu < \mu_0 \\ 2P(Z \geq |z|) & \text{falls } H_A: \mu \neq \mu_0 \end{cases}$$

Testentscheidung: $PW \leq \alpha \Rightarrow H_0$ verwerfen zugunsten H_A
 $PW > \alpha \Rightarrow H_0$ behalten

Hypothesentests für Anteil p $H_0: p = p_0$

Testvariable:
$$Z = \frac{\bar{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$$

Teststatistik:
$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$$

p -Wert: analog zu "Hypothesentests für Mittelwert μ "

Testentscheidung: analog zu "Hypothesentests für Mittelwert μ "

Hypothesentests für Differenzen von Mittelwerten $H_0: \Delta = \Delta_0$

Testvariable:
$$Z = \frac{\bar{X}_1 - \bar{X}_2 - \Delta_0}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}$$

Teststatistik:
$$z = \frac{\bar{x}_1 - \bar{x}_2 - \Delta_0}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}$$

p -Wert: analog zu "Hypothesentests für Mittelwert μ "

Testentscheidung: analog zu "Hypothesentests für Mittelwert μ "

Hypothesentests für Differenz von Anteilen $H_0: \Delta = \Delta_0$

$$\text{Testvariable: } z = \frac{\bar{p}_1 - \bar{p}_2 - \Delta_0}{\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}}$$

$$\text{Teststatistik: } z = \frac{\hat{p}_1 - \hat{p}_2 - \Delta_0}{\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}}$$

p-Wert: analog zu "Hypothesentests für Mittelwert μ "

Testentscheidung: analog zu "Hypothesentests für Mittelwert μ "

Anwendungen der χ^2 -Verteilung

H_0 : X hat eine vorgegebene Verteilung

⇒ endlich viele Klassen (je mehr, je genauer, aber komplizierter)

h_j : absolute Häufigkeit der Stichprobenelemente in der Klasse s

e_j : erwartete Häufigkeit laut H_0

$$\chi^2 = \sum_{s=1}^k \frac{(h_s - e_s)^2}{e_s}$$

wenn χ^2 zu gross, lehnt man H_0 ab!

⇒ Freiheitsgrade, Tabelle