

Estimation Theory

- ▶ Random variables and parameterized theoretical models build the theoretical foundation of modern statistics
- ▶ If we accept a family of distributions as a realistic model of a data generating process which shows variable (uncertain) outcomes, we need to know in specific situations which value of the parameter should be used to solve real world questions
- ▶ Estimation theory deals with estimating the values of the actual parameters based on empirical data. Substitution of the parameter estimates into the model distribution formulas gives a specific shape of the distribution reflecting the specific physical scenario and enables us to calculate any probability or gives us answers to any probability related question.

Estimator (Schätzfunktion; Schätzer)

- ▶ An estimator (Schätzer) is a rule for calculating an estimate of a given quantity based on observed data
- ▶ An estimate (Schätzung) is the result of the application of an estimator for a given data set
- ▶ This chapter deals firstly with point estimators which yield single-valued results
- ▶ Latter we will talk about interval estimator, where the result will be a range of plausible values
- ▶ Statistical theory is concerned with the properties of estimators for choosing “good” estimators

Properties of Estimators

- ▶ Consider a fixed parameter θ characterizing a statistical model.
- ▶ We start our estimation task from a given sample $\mathbf{x} = \{x_1, x_2, \dots, x_n\}$
 $\mathbf{x} \in X$, where X is the sample space.
- ▶ An estimator $T(X)$ is a function that maps the sample space to a set of sample estimates. The resulting estimate for any given sample can be written

$$\hat{\theta} = T(x)$$

- ▶ For a given sample \mathbf{x} the error of the estimate is defined

$$e(x) = T(x) - \theta = \hat{\theta} - \theta$$

Mean Squared Error

- ▶ The “Mean Squared Error” is defined as the expected value of the squared errors

$$MSE(T) = E\left[(T(X) - \theta)^2\right]$$

- ▶ It is used to indicate how far, on average, the collection of resulting estimates are away from the single true parameter which is estimated.
- ▶ Note: Expectation is taken over the whole sample space

Bias and Variance

- ▶ The *bias* of an estimator is the distance between the average of the collection of estimates, and the single parameter being estimated. It is just the expected value of the error.

$$BIAS(T) = E(T(X) - \theta)$$

- ▶ The *variance* of an estimator is the expected value of the squared sampling deviations. It is used to indicate how far, on average, the collection of estimates are from the *expected value* of the estimates.

$$VAR(T) = E\left[T(X) - E(T(X))\right]^2$$

MSE and Variance

- ▶ Note the difference between Mean Squared Error and variance. The MSE takes also a possible Bias into account
- ▶ The following relationship holds

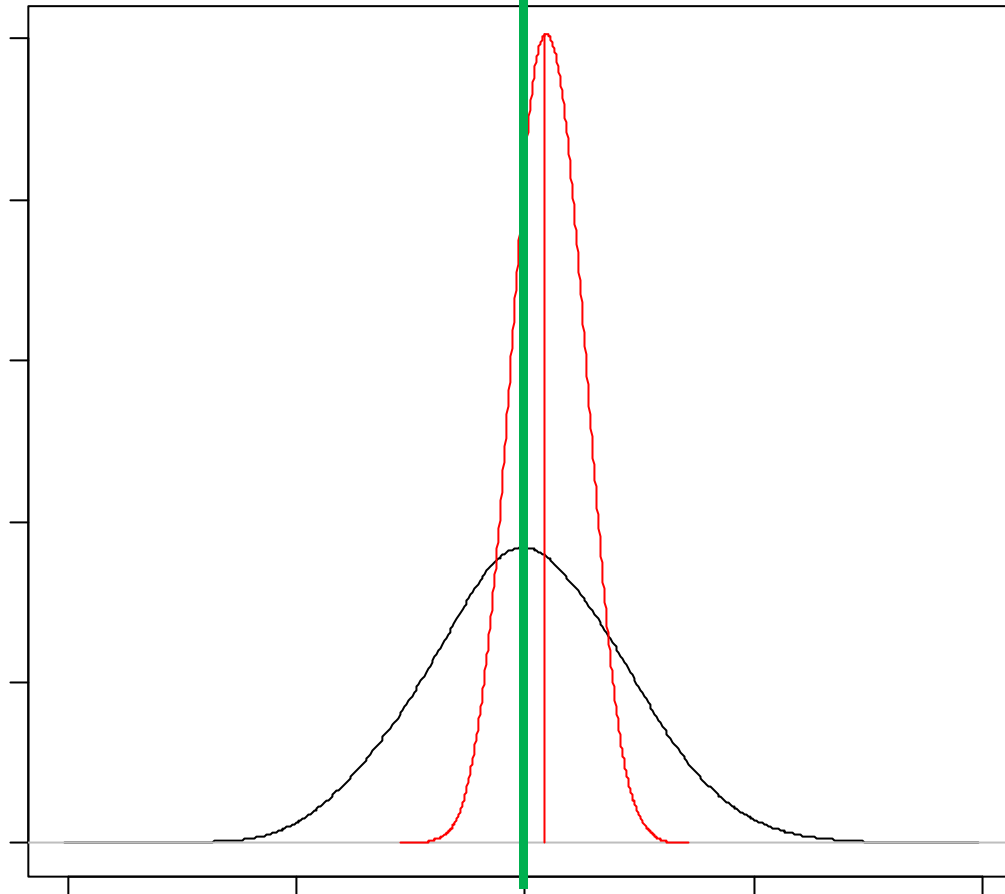
$$MSE(T) = VAR(T) + (BIAS(T))^2$$

“Good” Estimators

- ▶ 2 desirable properties of estimators:
 - ▶ Unbiasedness
 - ▶ Minimal Mean Squared Error
- ▶ In general both criteria cannot be satisfied simultaneously: a biased estimator may have lower MSE than any unbiased estimator, since despite having bias, the estimator variance may be sufficiently smaller than that of any unbiased estimator, and it may therefore be preferable to use, despite the bias.
- ▶ Among unbiased estimators, there often exists one with the lowest variance, which is called the minimum variance unbiased estimator

Bias versus Variance

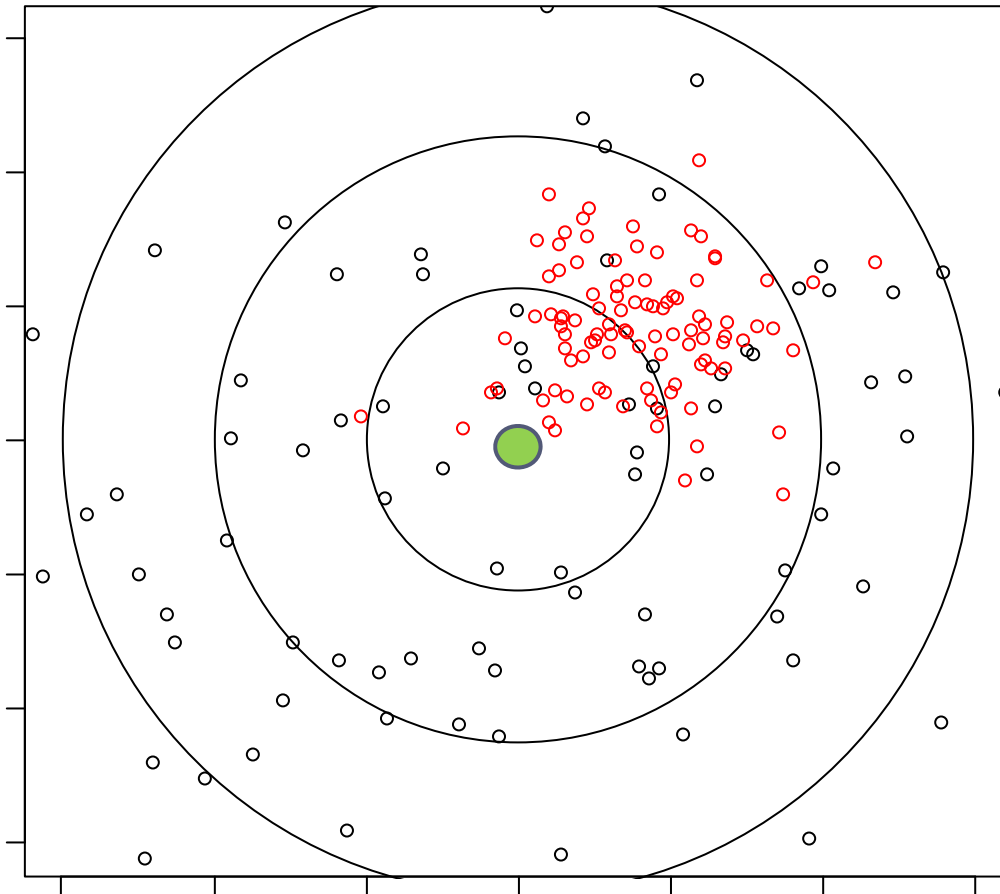
θ unknown, true parameter value



The black estimator is obviously unbiased but produces estimates with a larger variance than the red one, which is slightly biased

MSE of the red estimator is smaller than that of the black estimator – also it is not unbiased

Bias versus Variance



Again the black estimator is unbiased but produces estimates with a larger variance than the red one, which is systematically biased, but all points are quite close to the target

θ unknown, true parameter value

Übungsbeispiel (a)

- ▶ Es liegen fünf unabhängige Beobachtungen X_1, \dots, X_5 einer ZV mit Erwartungswert μ und Varianz σ^2 vor.
- ▶ Wir betrachten die folgenden 4 Schätzfunktionen für μ :
- ▶ $T1 = 1/5 (X_1 + X_2 + X_3 + X_4 + X_5)$
- ▶ $T2 = 1/3 (X_1 + X_2 + X_3)$
- ▶ $T3 = 1/6 (X_1 + X_2 + X_3 + X_4 + X_5)$
- ▶ $T4 = 1/6 (X_1 + X_2 + X_3 + X_4 + 2X_5)$
- ▶ (a) Welche dieser Schätzer sind erwartungstreu für μ ?
- ▶ $E(T1) = E(T2) = E(T4) = \mu$... sind erwartungstreu
- ▶ $E(T3) = 1/6 E(X_1 + X_2 + X_3 + X_4 + X_5) = 5/6 \mu$... nicht erwartungstreu

Übungsbeispiel (b)

- ▶ Es liegen fünf unabhängige Beobachtungen X_1, \dots, X_5 einer ZV mit Erwartungswert μ und Varianz σ^2 vor.
- ▶ Wir betrachten die folgenden 4 Schätzfunktionen für μ :
- ▶ $T1 = 1/5 (X_1 + X_2 + X_3 + X_4 + X_5)$
- ▶ $T2 = 1/3 (X_1 + X_2 + X_3)$
- ▶ $T3 = 1/6 (X_1 + X_2 + X_3 + X_4 + X_5)$
- ▶ $T4 = 1/6 (X_1 + X_2 + X_3 + X_4 + 2X_5)$

▶ (b) Welche dieser Schätzer hat die kleinste Varianz?

▶ $\text{Var}(T1) = 1/25 \text{Var}(X_1 + X_2 + X_3 + X_4 + X_5) = 5/25 \sigma^2 = \sigma^2/5$

▶ $\text{Var}(T2) = 1/9 \text{Var}(X_1 + X_2 + X_3) = 3/9 \sigma^2 = \sigma^2/3$

▶ $\text{Var}(T3) = 1/36 \text{Var}(X_1 + X_2 + X_3 + X_4 + X_5) = 5/36 \sigma^2$

▶ $\text{Var}(T4) = 1/36 \text{Var}(X_1 + X_2 + X_3 + X_4 + 2X_5) = 8/36 \sigma^2 = 2\sigma^2/9$

Innerhalb der Teilmenge der erwartungstreuen Schätzer hat der Schätzer T1 (~ arithmetische Mittel) die geringste Varianz

T3 hat zwar eine kleinere Varianz ist aber nicht erwartungstreu

Estimation of the Variance

- ▶ If we like to estimate the variance of a population from a random sample, the formula dividing the sum of squared differences of the mean by $n-1$ provides an unbiased estimate of the variance of the population

$$S_1^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2.$$

$$\begin{aligned} E(S_1^2) &= \frac{1}{n} \sum_{i=1}^n E((X_i - \bar{X})^2) = \frac{1}{n} E\left(\sum_{i=1}^n (X_i - \mu + \mu - \bar{X})^2\right) \\ &= \frac{1}{n} E\left(\sum_{i=1}^n ((X_i - \mu)^2 - 2(X_i - \mu)(\bar{X} - \mu) + (\bar{X} - \mu)^2)\right) \\ &= \frac{1}{n} E\left(\sum_{i=1}^n (X_i - \mu)^2 - 2\sum_{i=1}^n (X_i - \mu)(\bar{X} - \mu) + n(\bar{X} - \mu)^2\right) \\ &= \frac{1}{n} E\left(\sum_{i=1}^n (X_i - \mu)^2 - 2n(\bar{X} - \mu)(\bar{X} - \mu) + n(\bar{X} - \mu)^2\right) \\ &= \frac{1}{n} E\left(\sum_{i=1}^n (X_i - \mu)^2 - n(\bar{X} - \mu)^2\right) \\ &= \frac{1}{n} \left(\sum_{i=1}^n E((X_i - \mu)^2) - n E((\bar{X} - \mu)^2)\right) \\ &= \frac{1}{n} (n \operatorname{Var}(X) - n \operatorname{Var}(\bar{X})) \\ &= \operatorname{Var}(X) - \operatorname{Var}(\bar{X}) = \sigma^2 - \frac{\sigma^2}{n} = \frac{n-1}{n} \sigma^2, \end{aligned}$$

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

$$E(\hat{\sigma}^2) = \sigma^2$$

Maximum Likelihood Estimation

- ▶ **Maximum likelihood estimation (MLE)** is a very popular statistical method to fit a theoretical model to empirical data.
- ▶ Modeling real world data by estimating maximum likelihood offers a way of tuning the free parameters of the model to provide an optimum fit.
- ▶ The basic idea is to search for that parameter values which maximize the likelihood of the observed sample, i.e. MLEs are chosen so that the in reality observed sample arises very likely
- ▶ The method was developed by geneticist and statistician Sir R.A. Fisher almost a century ago.



Maximum Likelihood Principle (1)

- ▶ Consider a family D_θ of probability distributions parameterized by an unknown parameter θ (which could be vector-valued), associated with either a known probability density function (continuous distribution) or a known probability mass function (discrete distribution), denoted as f_θ .
- ▶ We have a sample of n values x_1, x_2, \dots, x_n which we assume are independent realizations from this family of distributions (iid...independently identically distributed)
- ▶ Thus each x_i is a realization of a RV $X \sim f_\theta$
- ▶ The joint probability density associated with our observed data is then

$$f_\theta(x_1, \dots, x_n | \theta)$$

Maximum Likelihood Principle (2)

- ▶ Due to independence of the observations the joint probability density of the sample may be written

$$\prod_{i=1}^n f_{\theta}(x_i | \theta)$$

- ▶ If we consider the observed sample as fixed and interpret this expression as a function of θ we derive the so called likelihood-function

$$\mathcal{L}(\theta) = \prod_{i=1}^n f_{\theta}(x_i | \theta)$$

- ▶ The method of maximum likelihood estimates θ by finding that value of θ that maximizes $L(\theta)$. The result is called the **maximum likelihood estimator (MLE)** of θ .

Maximum Likelihood Principle (3)

- ▶ Since maxima are unaffected by monotone transformations, one can take the logarithm of this expression to turn it into a sum, which makes derivation of the maximum easier

$$\mathcal{L}^*(\theta) = \sum_{i=1}^n \log f_{\theta}(x_i | \theta).$$

- ▶ The maximum of this expression can then be found by taking the derivative and find either an analytical solution with a closed formula or find values numerically using various optimization algorithms

ML-Estimate for an Exponential RV

$$L(\lambda) = \prod_{i=1}^n \lambda \exp(-\lambda x_i) = \lambda^n \exp\left(-\lambda \sum_{i=1}^n x_i\right) = \lambda^n \exp(-\lambda n\bar{x}),$$

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

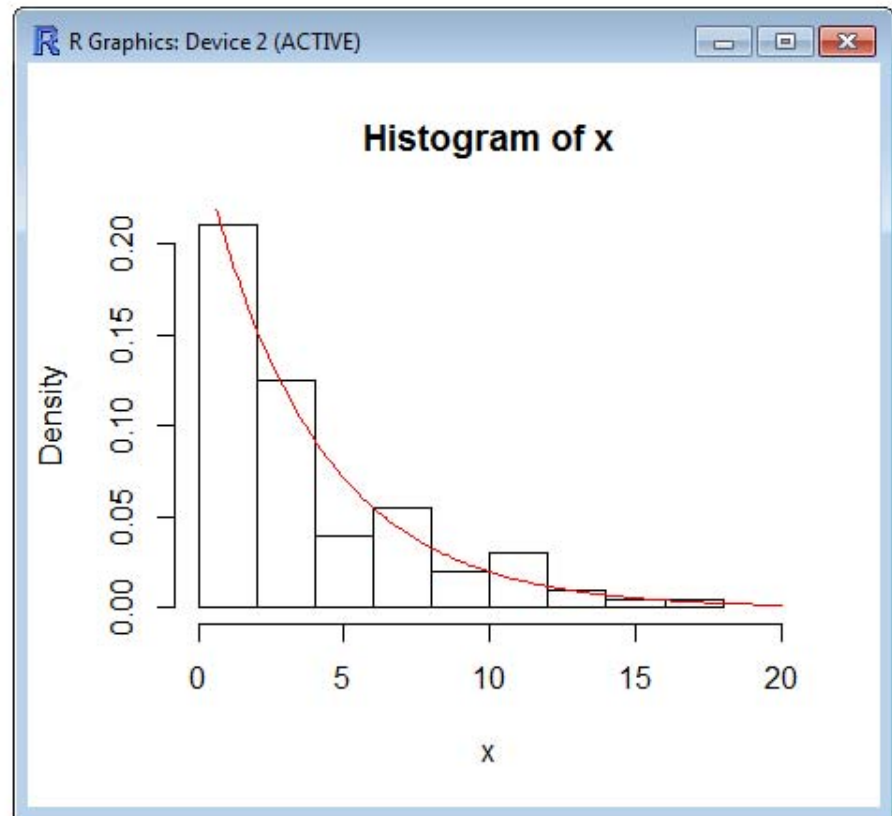
$$\frac{d}{d\lambda} \ln L(\lambda) = \frac{d}{d\lambda} (n \ln(\lambda) - \lambda n\bar{x}) = \frac{n}{\lambda} - n\bar{x} \begin{cases} > 0 & \text{if } 0 < \lambda < 1/\bar{x}, \\ = 0 & \text{if } \lambda = 1/\bar{x}, \\ < 0 & \text{if } \lambda > 1/\bar{x}. \end{cases}$$

$$\hat{\lambda} = \frac{1}{\bar{x}}.$$

Praktische Anwendung

R Console

```
> # Example
>
> # We generate 100 realisations from an exponential distribution
> # with parameter lambda = 0.25
>
> set.seed(1234)
> x <- rexp(100, 0.25)
> hist(x, prob=T, xlim=c(0,20), )
>
> # Maximum Likelihood Estimate of Lambda
> print(paste("ML-estimate of lambda", round(1/mean(x),2)))
[1] "ML-estimate of lambda 0.26"
>
> jj <- seq(from=0, to=20, length=200)
> lines(jj, dexp(jj, 1/mean(x)), col=2)
> |
```



Method of Moments

- ▶ The method of moments in general provides estimators which are consistent but not as efficient as the Maximum likelihood ones. (i.e. they have larger variances)
- ▶ They are often used, because they lead to very simple computations, unlike ML method which can become very cumbersome. (Sometimes the results derived by the method of moments are used as starting value for numerical optimization)
- ▶ It is based on the calculation of empirical moments (mean, variance) and equating them with the theoretical formula and solving for the unknown parameters.

Example with Pareto Distribution

- ▶ We use the following parameterization:

$$f(x; \alpha, \lambda) = \frac{\alpha \lambda^\alpha}{(\lambda + x)^{\alpha+1}}$$

$$E(X) = \frac{\lambda}{\alpha - 1} \quad V(X) = \frac{\alpha \lambda^2}{(\alpha - 1)^2(\alpha - 2)}$$

- ▶ The method of moments simply equates the empirical moments into the formula

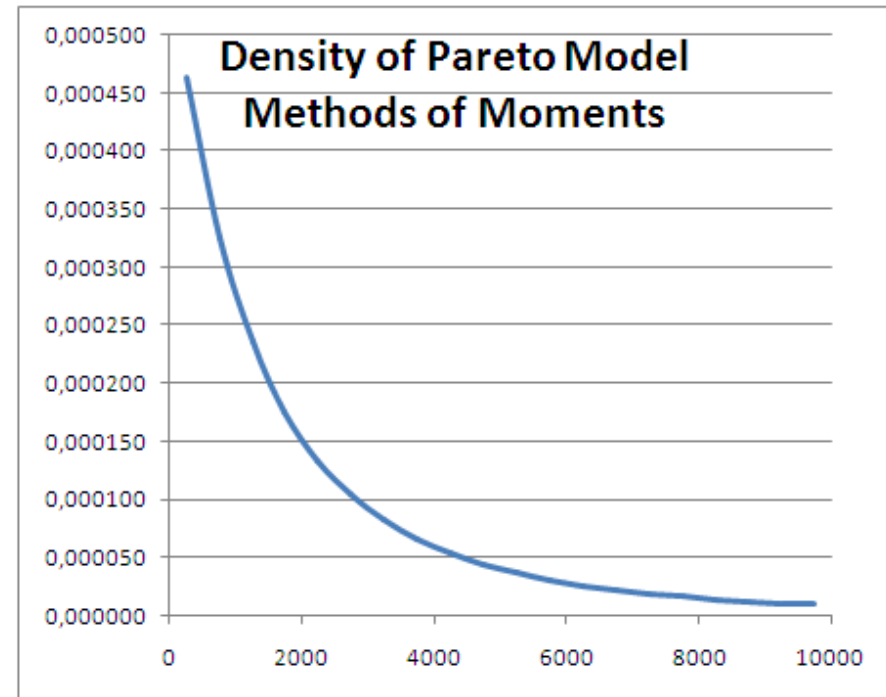
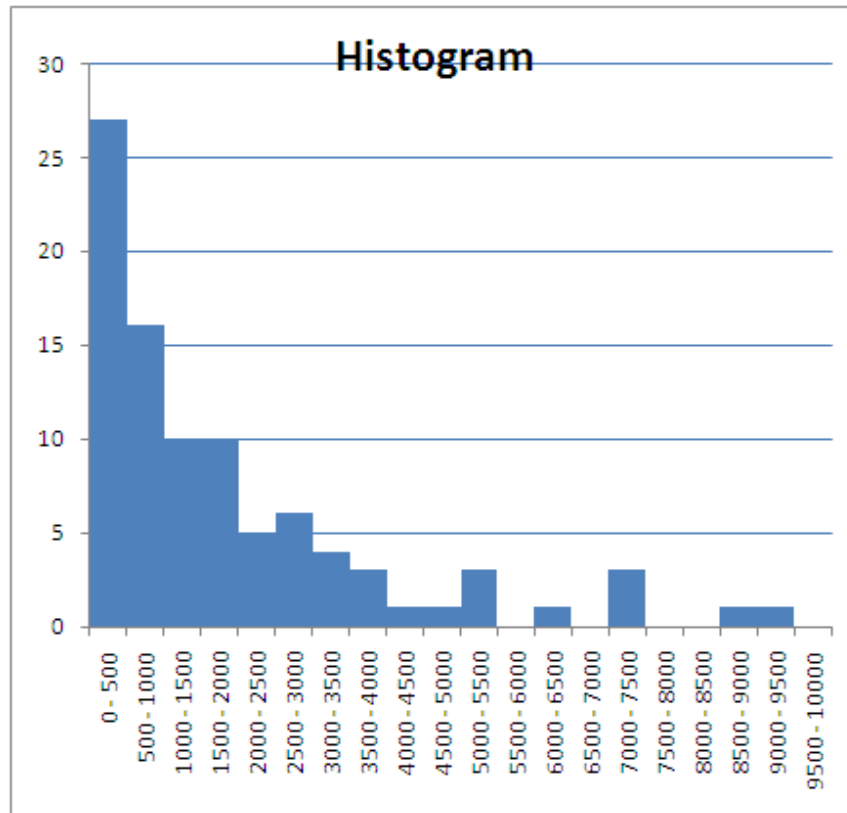
$$\bar{x} = E(X) = \frac{\lambda}{\alpha - 1} \quad s^2 = V(X) = \frac{\alpha \lambda^2}{(\alpha - 1)^2(\alpha - 2)}$$

$$\hat{\alpha} = \frac{2s^2}{s^2 - \bar{x}^2} \quad \hat{\lambda} = (\hat{\alpha} - 1) \cdot \bar{x}$$

Worked-out Example



Pareto.xls

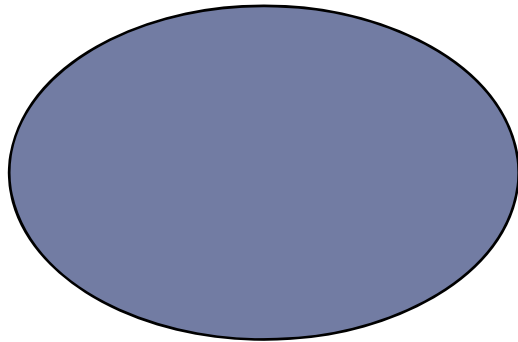


Empirical Moments	
mean	2,989.83
var	47,006,239
Moment Estimate	
lambda	4394
alpha	2.470



Grundproblem der Inferenzstatistik

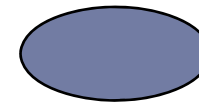
Grundgesamtheit



Stichprobenziehung



Zufalls-
Stichprobe



π ... "wahre", unbekannte Anteil
nicht zufällig

p ... beobachtete Anteil
zufällig



?

Inferenzschluss

Stichprobenziehung

- ▶ Teilerhebung (sample survey)
versus
Vollerhebung (census)
- ▶ Gründe für Stichprobenerhebung
 - ▶ Kostenersparnis
 - ▶ Zeitgewinn
 - ▶ Praktische Unmöglichkeit einer Vollerhebung

Stichprobentechniken

Arten der Stichprobenziehung

▶ Zufallsauswahlverfahren

Jedes Element der Grundgesamtheit besitzt eine bestimmte, von null verschiedene Wahrscheinlichkeit in die Stichprobe zu gelangen

▶ Verfahren der bewußten Auswahl

Vorgabe von Quotenmerkmalen, durch die die Stichprobenstruktur in wichtigen Variablen der Struktur der Grundgesamtheit entspricht

Problem: Verbleibender subjektiver Spielraum läßt keine wahrscheinlichkeits theoretisch abgesicherten Aussagen über die Zuverlässigkeit der Ergebnisse zu

Zufallsauswahlverfahren

- ▶ Einfachste Variante:
Jedes Element besitzt die gleiche Wahrscheinlichkeit gezogen zu werden
Uneingeschränkte Zufallsauswahl (simple random sampling)
- ▶ In der Praxis:
Geschichtete Zufallsstichprobe (stratified random sampling)
ermöglicht genauere Aussagen in heterogenen Populationen
Klumpenstichprobe (cluster sampling)
reduziert Erhebungskosten; oft aus praktischen Gründen erforderlich

Voraussetzung für echte Zufallsauswahl

- ▶ "sampling frame"
Zuverlässiges Register (Datenbank) aller Elemente der Grundgesamtheit (sampling units)
- ▶ Sampling frame
 - ▶ ermöglicht die Operationalisierung von Zufallsauswahlen
 - ▶ eröffnet die Möglichkeit zur Durchführung der Erhebung (Adressen, Tel.Nr., etc)
 - ▶ Enthält zusätzliche Informationen zur Erhöhung der Präzision bei der Stichprobenziehung bzw. bei der Schätzung/Hochrechnung der Ergebnisse

Typische Probleme

- ▶ "Selection Bias"

1936 US-President Election 2,4 Mio Fragebogen Adressen aus Telefonbuch; KfZ-Registration; Mitglieder eines Buchklubs

Prognose: Landon 57% Roosevelt 43%

Ergebnis: Roosevelt > 60%

- ▶ "Household Bias"

population units: Person

sampling units: Haushalt

Pro Haushalt wird ein Mitglied in die Stichprobe aufgenommen ==> Mitglieder von Großfamilien sind systematisch unterrepräsentiert

Andere Probleme

- ▶ **Non-Response Bias**
 - ▶ Nicht-Antwörter können sich von den Antwortern systematisch unterscheiden
- ▶ **Response Bias**
 - ▶ Befragte wollen sich nicht deklarieren
- ▶ **Gestaltung der Frage kann Antwort beeinflussen**
 - ▶ z.B. hat in einer experimentellen Studie das Vertauschen der Reihenfolge von Antwortalternativen zu einer 5%-igen Veränderung des Ergebnisses geführt
 - ▶ Formulierung !
 - ▶ Skalierung (Anzahl der Antwortalternativen gerade versus ungerade)

Fehlerstruktur bei Stichprobenerhebungen

Stichprobenschätzung

=

Wahre Parameter

+

Erhebungs-Bias (Verzerrung aufgrund der
Befragungstechnik)

+

Stichprobenfehler (Unsicherheit aufgrund der Teilerhebung)

Anwendungskontext

- ▶ Wir betrachten zunächst eine Grundgesamtheit mit einem binärem Merkmal (homogruader Fall)
- ▶ Wie kann man von der Stichprobe auf die Grundgesamtheit schließen ?
- ▶ Bei Kenntnis der Parameter der Grundgesamtheit
<Anzahl interessierender Ereignisse (M) und Umfang der Grundgesamtheit (N) bzw. des Anteils $\pi=M/N$ >
wissen wir bereits, wie Aussagen über zentrale Schwankungsintervalle für die Anzahl (X) bzw. den Anteil (p) in der Stichprobe gemacht werden können.

Theoretisches Vorwissen

- ▶ Ziehen ohne Zurücklegen
Grundgesamtheit mit N Elementen
davon M mit der interessierenden Eigenschaft
Stichprobe vom Umfang n
 1. Exakte Bestimmung der Wahrscheinlichkeiten aller möglichen Stichprobenergebnisse mittels der Hypergeometrischen Verteilung
 2. Näherung der Hypergeometrischen Verteilung durch die Binomialverteilung
 3. Approximation durch die Normalverteilung mit Varianzformel der Hypergeometrischen oder der Binomialverteilung

Theoretisches Vorwissen (1 Hypergeom. Vert.)

- ▶ Ziehen ohne Zurücklegen aus einer Grundgesamtheit mit N Elementen davon M interessierende Elemente

Stichprobe vom Umfang n

X ... Anzahl der interessierenden Elemente in der Stichprobe

p ... Anteil der interessierenden Elemente in der Stichprobe

$$E(X) = n \cdot \frac{M}{N} = n \cdot \pi$$

$$E(p) = \frac{M}{N} = \pi$$

$$V(X) = n\pi(1-\pi) \cdot \frac{N-n}{N-1}$$

$$V(p) = \frac{\pi(1-\pi)}{n} \cdot \frac{N-n}{N-1}$$

$$\frac{X - n\pi}{\sqrt{n \cdot \pi(1-\pi) \cdot \frac{N-n}{N-1}}} \approx N(0;1)$$

$$\frac{p - \pi}{\sqrt{\frac{\pi(1-\pi)}{n} \cdot \frac{N-n}{N-1}}} \approx N(0;1)$$

Theor. Vorwissen (2 Binom. Verteilung)

- ▶ Ziehen einer Stichprobe vom Umfang n
- ▶ Vereinfachte Berechnung der Varianz, falls das Ziehen mit Zurücklegen erfolgt oder ein kleiner Auswahlsatz (n/N) gegeben ist.

$$E(X) = n \cdot \frac{M}{N} = n \cdot \pi$$

$$V(X) = n\pi(1 - \pi)$$

$$\frac{X - n\pi}{\sqrt{n \cdot \pi(1 - \pi)}} \approx N(0;1)$$

$$E(p) = \pi$$

$$V(p) = \frac{\pi(1 - \pi)}{n}$$

$$\frac{p - \pi}{\sqrt{\frac{\pi(1 - \pi)}{n}}} \approx N(0;1)$$

Beispiel

Parameter der Grundgesamtheit

$$N = 10.000$$

$$M = 6.000$$

$$\pi = 0,60$$

Parameter der Stichprobe

$$n = 100$$

$$E(X) = 60$$

Mit Stichprobenkorrektur

$$V(X) = 23,7624$$

$$s(X) = 4,8747$$

Ohne Stichprobenkorrektur

$$V(X) = 24,0000$$

$$s(X) = 4,8990$$

$$E(p) = 0,60$$

Mit Stichprobenkorrektur

$$V(p) = 0,0024$$

$$s(p) = 0,0487$$

Ohne Stichprobenkorrektur

$$V(p) = 0,0024$$

$$s(p) = 0,0490$$

X ... Anzahl der Erfolge

p ... Anteil der Erfolge

Ergebnisvergleich

	Hy	Bi	Norm mit	Norm ohne
Prob($p < 0,5$)=	0,0163	0,0168	0,0156	0,0160
Prob($p < 0,6$)=	0,4565	0,4567	0,4592	0,4594
Prob($p = 0,6$)=	0,0816	0,0812	0,0817	0,0813
Prob($0,55 < p < 0,65$)	0,6440	0,6416	0,6441	0,6417
Prob($p > 0,6$)=	0,4618	0,4621	0,4592	0,4594

Exakte Ergebnisse



Approximation in der Praxis



Inklusions- bzw. Repräsentationsschluß

- ▶ Bei bekanntem Parameter der Grundgesamtheit haben wir bislang Aussagen über die Verteilung der Stichprobengrößen getroffen (Inklusion; direkter Schluß)
- ▶ durch Umkehrung gelangen wir zur neuen Aufgabe: ausgehend von einem Stichprobenergebnis soll auf die Parameter der Grundgesamtheit geschlossen werden (Repräsentations- oder Inferenz-Schluss; indirekter Schluss)

Punktschätzung

- ▶ Da die Stichprobe zufällig ist, ist auch das Ergebnis aus der Stichprobe zufällig
- ▶ Da wir wissen, dass der Erwartungswert für den Anteilswert der Stichprobe (p) gleich dem wahren Wert für den Parameter π in der Grundgesamtheit ist, erscheint es bei Vorliegen einer konkreten Stichprobe sinnvoll, den Stichprobenanteil p als Schätzung für π zu verwenden.
- ▶ Weiters können wir den Standardfehler zur Ermittlung der Präzision dieser Schätzung nutzen.

Konzept der Konfidenzintervalle - Beispiel

- ▶ Grundgesamtheit mit binärem Merkmal
z.B.: „Kandidat-A“ ... Erfolg Kandidat-B“ ... Misserfolg
- ▶ Stichprobe mit $n=500$
- ▶ Angenommen der wahre Wert in der Grundgesamtheit sei $\pi=0.5$ und wir wählen eine Sicherheitswahrscheinlichkeit von $(1-\alpha)=0,99$:
- ▶ X sei die Anzahl der Erfolge in der Stichprobe
- ▶ p sei der Anteil der Erfolge in der Stichprobe

Zentrales Schwankungsintervall - Beispiel

- ▶ Für das zentrale Schwankungsintervall von X ergibt sich dann
- ▶ $P(n \cdot \pi - 2,58 \cdot \sigma < X < n \cdot \pi + 2,58 \cdot \sigma) = 0,99$
- ▶ Mit $\sigma^2 = n \cdot \pi \cdot (1 - \pi)$
- ▶ $n \cdot \pi = 500 \cdot 0,5 = 250$ $\sigma^2 = n \cdot \pi \cdot (1 - \pi) = 500 \cdot 0,5 \cdot 0,5 = 125$
- ▶ $\sigma = 11,18$
- ▶ $P(250 - 2,58 \cdot 11,18 < X < 250 + 2,58 \cdot 11,18) = 0,99$
- ▶ $P(221,15 < X < 278,85) = 0,99$
- ▶ Falls der wahre Wert in der Grundgesamtheit $\pi = 0,5$ beträgt wird die Anzahl der Befragten für Kandidat A mit einer Wahrscheinlichkeit von 99% im Intervall von 221 bis 279 liegen.

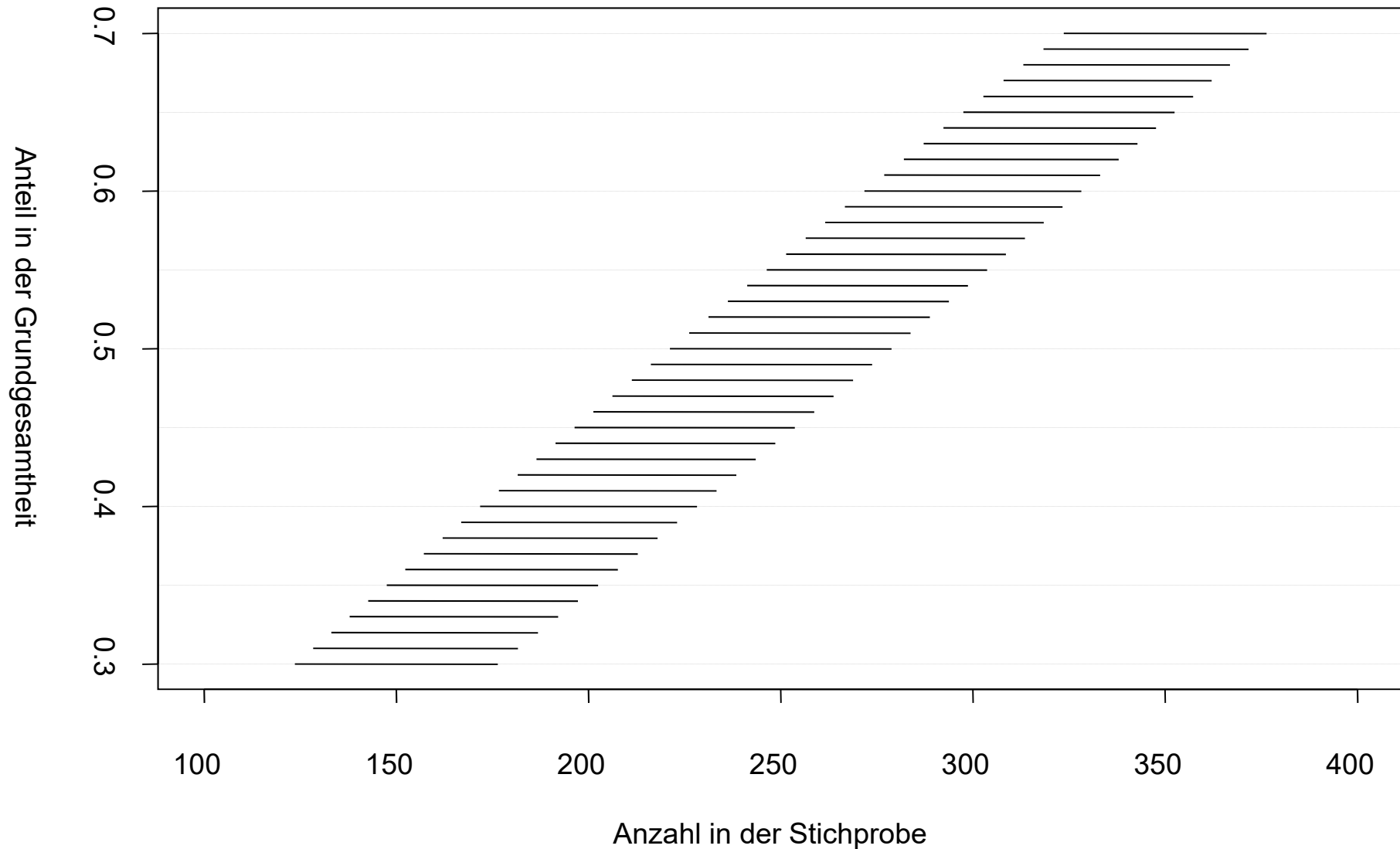
Konzept der Konfidenzintervalle

- Für unterschiedliche Werte von π ergeben sich natürlich verschiedene zentrale Schwankungsintervalle

n= 500
 $\alpha= 0,01$ 2,5758 Quantilswert

Zentrales Schwankungsintervall für die Anzahl					Zentrales Schwankungsintervall für den Anteil			
π	E(X)	Var(X)	UG(X)	OG(X)	E(p)	Var(p)	UG(p)	OG(p)
0,10	50	45,0	33	67	0,10	0,0002	6,5%	13,5%
0,15	75	63,8	54	96	0,15	0,0003	10,9%	19,1%
0,20	100	80,0	77	123	0,20	0,0003	15,4%	24,6%
0,25	125	93,8	100	150	0,25	0,0004	20,0%	30,0%
0,30	150	105,0	124	176	0,30	0,0004	24,7%	35,3%
0,35	175	113,8	148	202	0,35	0,0005	29,5%	40,5%
0,40	200	120,0	172	228	0,40	0,0005	34,4%	45,6%
0,45	225	123,8	196	254	0,45	0,0005	39,3%	50,7%
0,50	250	125,0	221	279	0,50	0,0005	44,2%	55,8%
0,55	275	123,8	246	304	0,55	0,0005	49,3%	60,7%
0,60	300	120,0	272	328	0,60	0,0005	54,4%	65,6%
0,65	325	113,8	298	352	0,65	0,0005	59,5%	70,5%
0,70	350	105,0	324	376	0,70	0,0004	64,7%	75,3%
0,75	375	93,8	350	400	0,75	0,0004	70,0%	80,0%
0,80	400	80,0	377	423	0,80	0,0003	75,4%	84,6%
0,85	425	63,8	404	446	0,85	0,0003	80,9%	89,1%
0,90	450	45,0	433	467	0,90	0,0002	86,5%	93,5%

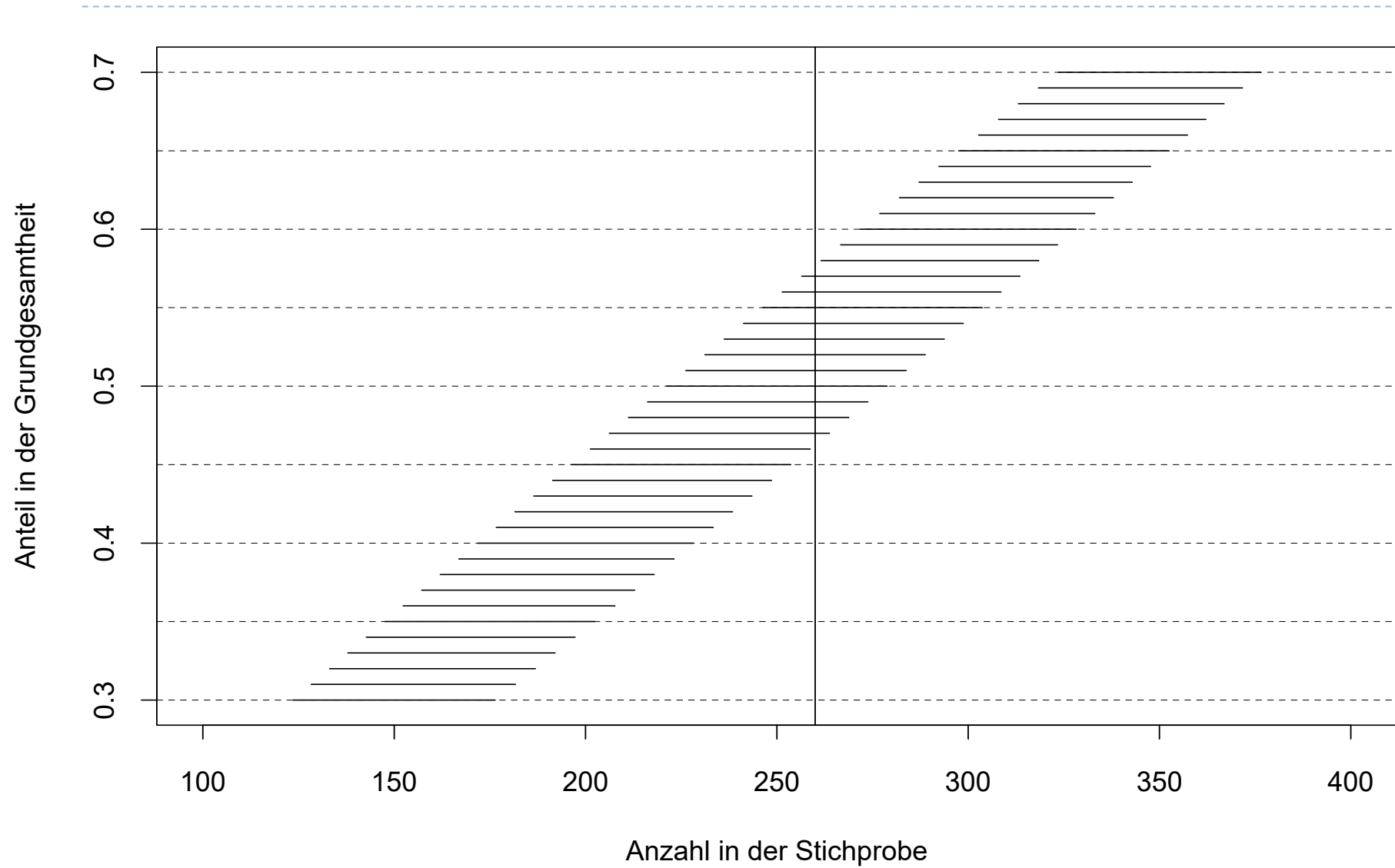
Die Graphik stellt für unterschiedliche Werte von π die zentralen Schwankungsintervalle durch horizontale Linien dar.



Konfidenzintervalle

- ▶ Trägt man in diese Graphik den konkreten real beobachteten Stichprobenwert (z.B. 260) mittels einer vertikalen Linie ein, so kann man ablesen, dass ein solches Stichprobenergebnis bei einer Sicherheitswahrscheinlichkeit von 99% mit einem π -Wert in der Grundgesamtheit von 0,46 bis 0,58 konform geht.
- ▶ Hingegen erscheinen Werte von beispielsweise $\pi=0,7$ bzw. $\pi=0,4$ für den Anteil in der Grundgesamtheit mit dem Stichprobenergebnis nicht verträglich.

Konfidenzintervall



Formel für das Konfidenzintervall

Ausgangspunkt:

$$P\left(-z \leq \frac{p - \pi}{\sqrt{\frac{\pi(1-\pi)}{n}}} \leq z\right) \approx \Phi(z) - \Phi(-z) = 1 - \alpha$$

Auflösung nach π

$$P\left(p - z\sqrt{\frac{\pi(1-\pi)}{n}} \leq \pi \leq p + z\sqrt{\frac{\pi(1-\pi)}{n}}\right) = 1 - \alpha$$

$$\sigma_p = \sqrt{\frac{\pi(1-\pi)}{n}} \quad \text{ersetzen durch} \quad s_p = \hat{\sigma}_p = \sqrt{\frac{p(1-p)}{n-1}}$$

Approximatives Konfidenzintervall für Anteile

- ▶ Liefert eine Stichprobe den empirischen Anteilswert p , so überdeckt das folgende Konfidenzintervall den wahren Parameter π mit einer Wahrscheinlichkeit von $1-\alpha$. (z sei in der Formel das $(1-\alpha/2)$ -Quantil der Standard-Normalverteilung)

$$P\left(p - z\sqrt{\frac{p(1-p)}{n-1}} \leq \pi \leq p + z\sqrt{\frac{p(1-p)}{n-1}}\right) = 1 - \alpha$$

bzw. bei großen n

$$P\left(p - z\sqrt{\frac{p(1-p)}{n}} \leq \pi \leq p + z\sqrt{\frac{p(1-p)}{n}}\right) = 1 - \alpha$$

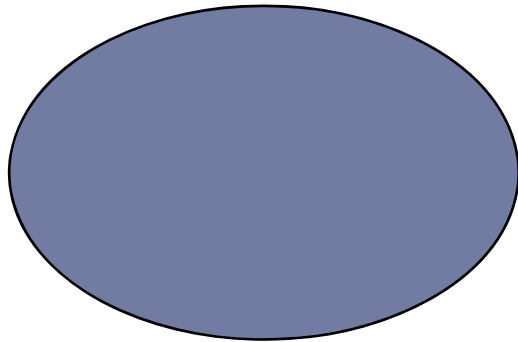
Korrekturfaktor

- ▶ Bei einem relativ großen Auswahlsatz und Ziehen ohne Zurücklegen ist der Korrekturfaktor zu berücksichtigen:

$$P\left(p - z\sqrt{\frac{p(1-p)}{n-1} \frac{N-n}{N-1}} \leq \pi \leq p + z\sqrt{\frac{p(1-p)}{n-1} \frac{N-n}{N-1}}\right) = 1 - \alpha$$

Grundproblem der Inferenzstatistik

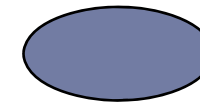
Grundgesamtheit



Stichprobenziehung



Zufalls-
Stichprobe



π ... "wahre", unbekannte Anteil
nicht zufällig

p ... beobachtete Anteil
zufällig

$$P \left(p - z_{1-\alpha/2} \sqrt{\frac{p(1-p)}{n}} \leq \pi \leq p + z_{1-\alpha/2} \sqrt{\frac{p(1-p)}{n}} \right) = 1 - \alpha$$

Inferenzschluss

Beispiel

- ▶ N=10.000 Haushalte
- ▶ Stichprobe:
n=100 davon 30 Haushalte mit mehr als 1 Auto
p=0,3 ... Punktschätzung für π
- ▶ Gesucht 95% ($\alpha=0,05$) Konfidenzintervall für den unbekanntem Anteil π der Haushalte mit mehr als einem Auto in der Grundgesamtheit
- ▶ Varianz der Punktschätzung: $s_p^2=0,3*0,7/99=0,0021$
- ▶ Standardfehler: $s_p=0,046$
- ▶ Für 95% ($\alpha=0,05$) Konfidenzintervall: $z=1,96$
- ▶ $0,3 - 1,96*0,046 < \pi < 0,3 + 1,96*0,046$
- ▶ $P(0,21 < \pi < 0,39) = 0,95$
- ▶ **Die Wahrscheinlichkeit, dass der unbekanntem Anteilswert der Grundgesamtheit durch ein Intervall von 21% bis 39% überdeckt wird, beträgt approximativ 95%.**

Lösung mit R

```
# Example Confidence Interval for Proportion
ci.prop <- function(k, n, alpha=0.05)
# ad-hoc function
{
  P <- k/n
  SE <- sqrt(P*(1-P)/n)      # standard error
  LB <- P - qnorm(1-alpha/2)*SE
  UB <- P + qnorm(1-alpha/2)*SE
  list(sample.size=n, proportion.observed=P,
        confidence.level=1-alpha, lower.bound=LB,
        upper.bound=UB)
}
```

```
> ci.prop(30, 100)
$sample.size
[1] 100

$proportion.observed
[1] 0.3

$confidence.level
[1] 0.95

$lower.bound
[1] 0.2101832

$upper.bound
[1] 0.3898168
```

KonfidenzAnteil.xls

n=	100	
X=	30	
p=	30,00%	= Anteil in der Stichprobe
(1-p)=	0,7000	
var(p)=	0,0021	
sigma(p)=	0,0458	
α =	0,05	
Tab=	1,9600	
Tab*sigma=	8,98%	=emax
UG	21,02%	
OG	38,98%	
I=	17,96%	=Länge des Konfidenzintervalls

Beispiel

- ▶ Selbes Beispiel aber höhere Sicherheit der Aussage wird gewünscht:
- ▶ Gesucht 99% ($\alpha=0,01$) Konfidenzintervall für π
- ▶ $p=0,3$ $s_p^2=0,3*0,7/99=0,0021$ $s_p=0,046$
- ▶ $z=2,58$
- ▶ $0,3 - 2,58*0,046 < \pi < 0,3 + 2,58*0,046$
- ▶ $P(0,18 < \pi < 0,42) = 0,99$
- ▶ **Die Wahrscheinlichkeit, dass der unbekannte Anteilswert der Grundgesamtheit durch ein Intervall von 18% bis 42% überdeckt wird, beträgt 99%.**
- ▶ **→** Höhere Sicherheit bedingt eine weniger Präzise Aussage bzw. Hohe Präzision impliziert häufige Fehlaussagen

```
> ci.prop(30, 100, alpha=0.01)
$sample.size
[1] 100

$proportion.observed
[1] 0.3

$confidence.level
[1] 0.99

$lower.bound
[1] 0.1819607

$upper.bound
[1] 0.4180393
```

Beispiel

- ▶ Selbes Beispiel aber vierfach so große Stichprobe
- ▶ $N=10.000$ Haushalte
- ▶ $n=400$ mit 120 Haushalten mit mehr als 1 Auto
- ▶ Gesucht 95% ($\alpha=0,05$) Konfidenzintervall für π
- ▶ $p=0,3$ $s_p^2=0,3*0,7/399=0,0005$ $s_p=0,023$
- ▶ $z=1,96$
- ▶ $0,3 - 1,96*0,023 < \pi < 0,3 + 1,96*0,023$
- ▶ $P(0,26 < \pi < 0,34) = 0,95$
- ▶ Gesucht 99% ($\alpha=0,01$) Konfidenzintervall für π
- ▶ $0,3 - 2,58*0,023 < \pi < 0,3 + 2,58*0,023$
- ▶ $P(0,24 < \pi < 0,36) = 0,99$
- ▶ Vierfache Stichprobe halbiert die Länge des Konfidenzintervalls

Länge des Konfidenzintervalls

$$P\left(p - z\sqrt{\frac{p(1-p)}{n}} \leq \pi \leq p + z\sqrt{\frac{p(1-p)}{n}}\right) = 1 - \alpha$$

$$L = 2 \cdot z_{1-\alpha/2} \sqrt{\frac{p(1-p)}{n}}$$

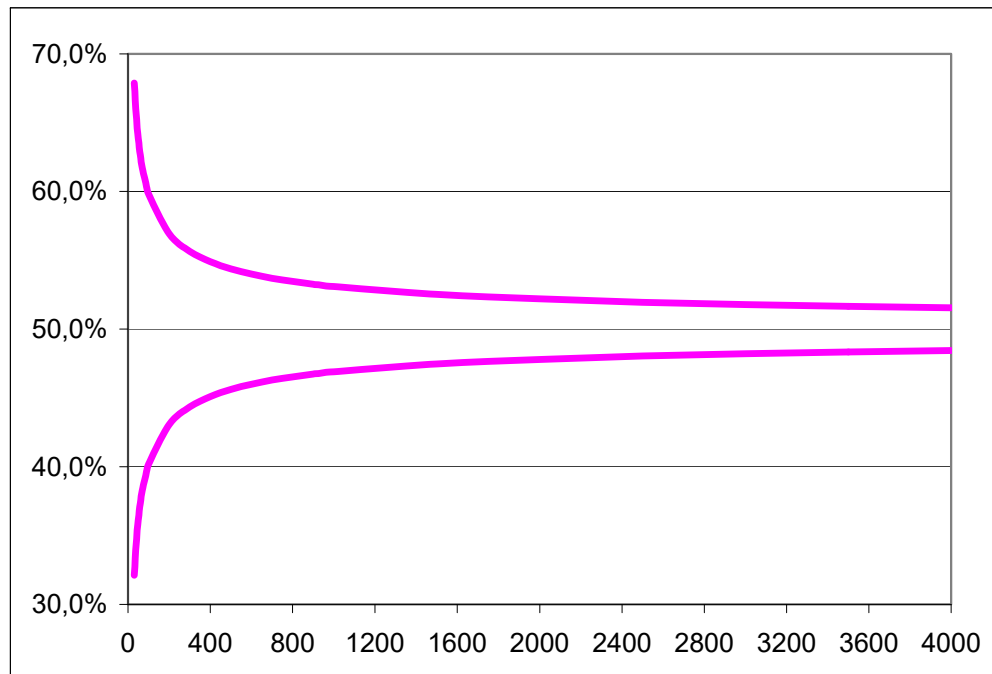
Hängt ab:

- ▶ von der Wahl von α (mit größerem α [~Irrtumswahrscheinlichkeit] wird Länge kleiner)
- ▶ vom Stichprobenumfang n (mit größerem n wird Länge kleiner; Wurzelgesetz!)
- ▶ von der Größe von p bzw. π (bei $p=1/2$ maximal)

Konfidenzintervall mit wachsendem Stichprobenumfang

Beobachtete Responserate $p = 50\%$
 Irrtumswahrscheinlichkeit $\alpha = 0,05$
 Konfidenzniveau $1-\alpha = 0,95$

Fallzahl	Standardfehler	KONFIDENZINTERVALL		Maximaler Fehler
		Untergrenze	Obergrenze	
30	0,091	32,1%	67,9%	17,9%
40	0,079	34,5%	65,5%	15,5%
50	0,071	36,1%	63,9%	13,9%
60	0,065	37,3%	62,7%	12,7%
70	0,060	38,3%	61,7%	11,7%
80	0,056	39,0%	61,0%	11,0%
90	0,053	39,7%	60,3%	10,3%
100	0,050	40,2%	59,8%	9,8%
200	0,035	43,1%	56,9%	6,9%
300	0,029	44,3%	55,7%	5,7%
400	0,025	45,1%	54,9%	4,9%
500	0,022	45,6%	54,4%	4,4%
600	0,020	46,0%	54,0%	4,0%
700	0,019	46,3%	53,7%	3,7%
800	0,018	46,5%	53,5%	3,5%
900	0,017	46,7%	53,3%	3,3%
1000	0,016	46,9%	53,1%	3,1%
1500	0,013	47,5%	52,5%	2,5%
2000	0,011	47,8%	52,2%	2,2%
2500	0,010	48,0%	52,0%	2,0%
3000	0,009	48,2%	51,8%	1,8%
3500	0,008	48,3%	51,7%	1,7%
4000	0,008	48,5%	51,5%	1,5%



Bestimmung des Stichprobenumfanges

$$n \geq \frac{4 \cdot z_{1-\alpha/2}^2 \cdot p(1-p)}{L^2}$$

Falls keine a-priori Kenntnis bezüglich p besteht, geht man vom „worst case“ $p=1/2$ aus [$p(1-p)$ wird dann maximal], wodurch sich die Formel wie folgt vereinfacht:

$$n \geq \frac{z_{1-\alpha/2}^2}{L^2}$$

Bestimmung des Stichprobenumfanges

$$L = 2 e_{\max}$$

e_{\max} ... Maximaler Fehler, des Konfidenzintervalls; bezeichnet bei vorgegebenem Signifikanzniveau, die maximale plus/minus Abweichung vom wahren Parameter

Beispiel:

Bestimme n , so dass der maximaler Fehler 5 Prozentpunkte beträgt
==> Länge des Konfidenzintervalls also maximal 10%

$$L=0,10$$

Ohne Vorkenntnis von p :

$$\text{Bei } \alpha=0,05: n > 1,96^2/0,01 = 384,1 \implies n=385$$

$$\text{Bei } \alpha=0,01: n > 2,58^2/0,01 = 663,5 \implies n=664$$

Beispiel

- ▶ Umfrage bei $n=2.000$ Wahlberechtigten
- ▶ Wie genau kann ein Anteil bei einem Konfidenzniveau von 95% vorhergesagt werden ?

p	Wurzel $[p(1-p)/n]$	max. Fehler	Länge des KI
0,1	0,0067	$\pm 1,31\%$	2,62%
0,2	0,0089	$\pm 1,75\%$	3,50%
0,3	0,0102	$\pm 2,01\%$	4,02%
0,4	0,0110	$\pm 2,15\%$	4,30%
0,5	0,0112	$\pm 2,19\%$	4,38%

Beispiel

- ▶ Gesucht ist eine Stichprobe vom Umfang n , mit der der Anteil der Ja-Wähler bei einer Volksabstimmung auf 1% genau geschätzt werden kann ($L=0,02$) Sicherheitsniveau 0,95

- ▶ a) bei Vorkenntnis, dass $\pi \sim 0,25$ sei:

$$n = 4 \cdot 3,84 \cdot 0,1875 / 0,0004 = 7.203$$

$$n \geq \frac{4 \cdot z_{1-\alpha/2}^2 \cdot p(1-p)}{L^2}$$

- ▶ b) ohne Vorkenntnis über den Anteil

$$n = 3,84 / 0,0004 = 9.604$$

$$n \geq \frac{z_{1-\alpha/2}^2}{L^2}$$

Beispiel aus den Medien



Interviews with 1,022 adult Americans conducted by telephone by Opinion Research Corporation on September 5-7, 2008. The margin of sampling error for results based on the total sample is plus or minus 3 percentage points.

n=1022

2/2a. If Barack Obama and Joe Biden were the Democratic Party's candidates and John McCain and Sarah Palin were the Republican Party's candidates, who would you be more likely to vote for -- Barack Obama and Joe Biden, the Democrats, or John McCain and Sarah Palin, the Republicans? (IF UNSURE:) As of today, who do you lean more toward? (RANDOM ORDER)

	<u>Obama</u>	<u>McCain</u>	<u>Neither</u> <u>(vol.)</u>	<u>Other</u> <u>(vol.)</u>	<u>No</u> <u>Opinion</u>
September 5-7, 2008	48%	48%	3%	1%	*

n= 1022
 X= 491
 p= **48,0%** = Anteil in der Stichprobe
 (1-p)= 0,5196
 var(p)= 0,0002
 sigma(p)= 0,0156
 α= 0,05
 Tab= 1,9600
 Tab*sigma= 3,06% =emax

UG 44,98%
OG 51,11%

l= 6,13% =Länge des Konfidenzintervalls

Beispiele: Konfidenzintervall für Anteile

- ▶ Sie lesen in einer Publikation, dass in einer Stichprobe von $n=600$ Befragten ein empirischer Anteil von $p=30\%$ erhoben wurde.
- ▶ Welche Aussage können Sie für den Anteil π in der Grundgesamtheit treffen, wenn Sie bereit sind mit einer 5%-Irrtumswahrscheinlichkeit zu argumentieren?

$$P\left(p - z\sqrt{\frac{p(1-p)}{n}} \leq \pi \leq p + z\sqrt{\frac{p(1-p)}{n}}\right) = 1 - \alpha$$

Beispiele: Konfidenzintervall für Anteile

- ▶ Stichprobe mit $p=0,30$ und $n=600$
- ▶ $p*(1-p)=0,3*0,7=0,21$
- ▶ $z=1,96$ $\sigma_p=0,0187$
- ▶ $P(0,30-1,96*0,0187 < \pi < 0,30+1,96*0,0187)=0,95$
- ▶ $P(26,33\% < \pi < 33,67\%)=0,95$

▶ **Die Wahrscheinlichkeit, dass der unbekannte Anteilswert in der Grundgesamtheit durch ein Intervall von 26,33% bis 33,67% überdeckt wird, beträgt 95%.**

- ▶ Falls Sie bereit sind Ihre Irrtumswahrscheinlichkeit auf 10% zu erhöhen, kommen Sie zu folgendem Ergebnis:
- ▶ $z=1,64$ $P(26,92\% < \pi < 33,08\%)=0,90$

Konfidenzintervalle für den Erwartungswert

Wir betrachten eine quantitative Zufallsvariable X mit „wahrem“ Erwartungswert μ

▶ X sei entweder normalverteilt

oder

▶ Der Stichprobenumfang n sei genügend groß (Faustregel: $n > 30$)

Dann gilt für das Stichprobenmittel:

$$\bar{x} \sim N(\mu; \sigma_{\bar{x}}^2)$$

$$\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n} \quad \text{Ziehen mit Zurücklegen oder kleiner Auswahlsatz}$$

$$\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n} \cdot \frac{N-n}{N-1} \quad \text{Ziehen ohne Zurücklegen}$$

Intervalle bei bekannter Varianz

$$i) P\left(\mu - z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \bar{x} \leq \mu + z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

durch Umformung :

$$ii) P\left(\bar{x} - z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

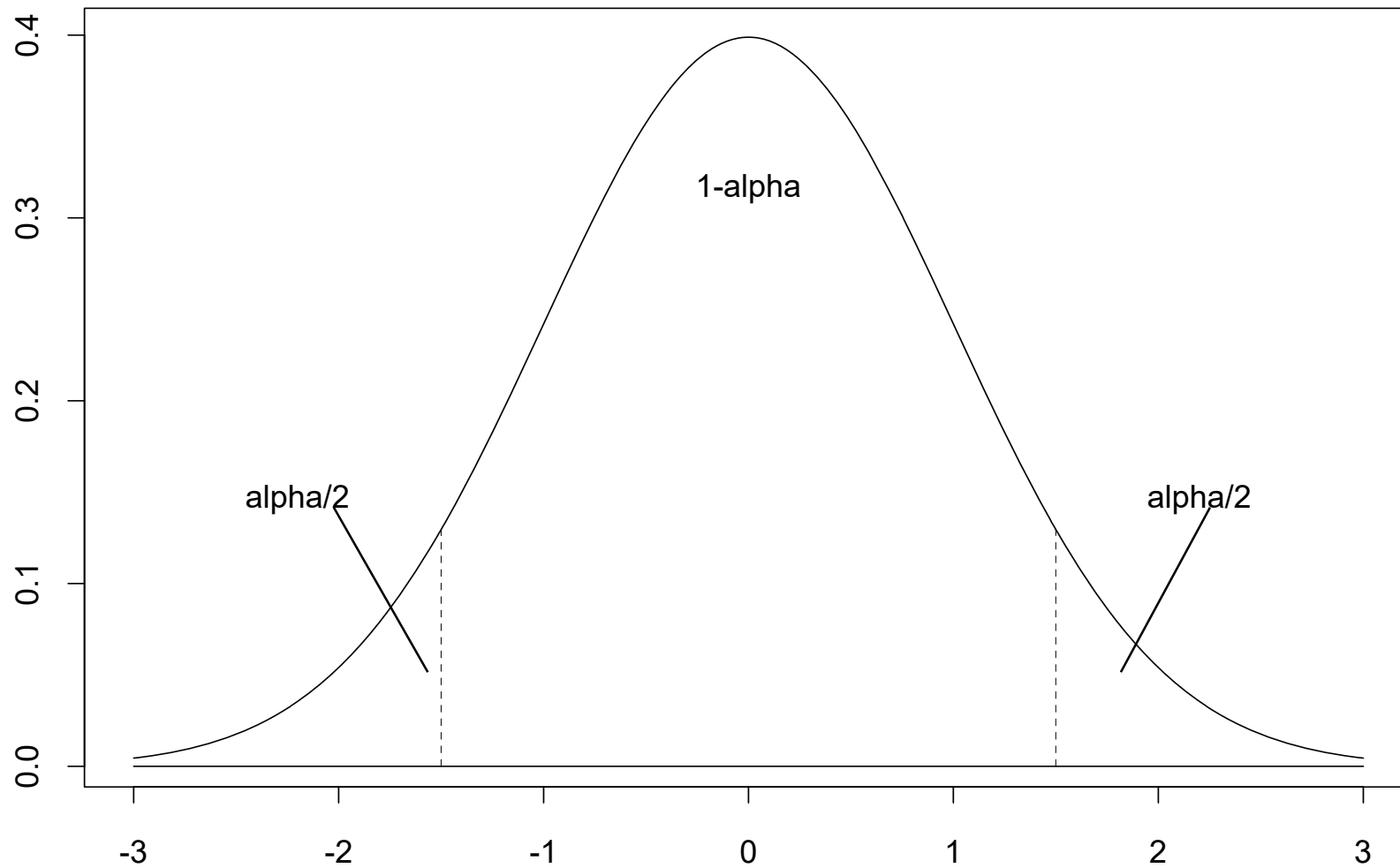
i) zentrales Schwankungsintervall für \bar{x}

ii) Konfidenzintervall für μ

Interpretation

- ▶ Das Zentrale Schwankungsintervall gibt bei Kenntnis des Erwartungswertes μ an, wie viel Prozent der Werte der Stichprobenfunktion „Arithmetisches Mittel“ bei wiederholter Stichprobenziehung in diesem Intervall zu liegen kommen werden.
- ▶ Das Konfidenzintervall gibt aus der Stichprobe abgeleitete Grenzen an, durch die der unbekannte Parameter μ mit einer vorgegebenen Irrtumswahrscheinlichkeit überdeckt oder eingeschlossen wird.

Konzept zentraler Schwankungsintervalle



Zentrale Schwankungsintervalle

- ▶ Sei $X \sim N(\mu, \sigma^2)$ so ergibt sich das zentrale Schwankungsintervall für das arithmetische Mittel einer Stichprobe, welches eine Wahrscheinlichkeit von $1-\alpha$ abdeckt durch:

$$P(\mu - z_{1-\alpha/2} \sigma_{\bar{x}} \leq \bar{x} \leq \mu + z_{1-\alpha/2} \sigma_{\bar{x}}) = 1 - \alpha$$

$$P(\mu - z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \bar{x} \leq \mu + z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}) = 1 - \alpha$$

- ▶ Für $\alpha=0,1$ ($\alpha=0,05$; $\alpha=0,01$) ergibt sich aus der Tabelle für $z_{1-\alpha/2} = 1,6449$ (1,96; 2,5758)
- ▶ d.h.

$P(\mu - 1,6449 < Z < \mu + 1,6449)$	$= 0,9$
$P(\mu - 1,96 < Z < \mu + 1,96)$	$= 0,95$
$P(\mu - 2,5758 < Z < \mu + 2,5758)$	$= 0,99$



Beispiel

- ▶ Sei ein IQ-Test so normiert, dass gilt $IQ \sim N(100, 15^2)$
- ▶ Gesucht ist ein zentrales Schwankungsintervall, für das arithmetische Mittel des IQ einer Stichprobe von $n=30$ Personen, welches eine Wahrscheinlichkeit von 0,95 aufweist.

$$P\left(\mu - z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \bar{x} \leq \mu + z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

- ▶ $\alpha = 0,05$ $1-\alpha/2 = 0,975$ $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{15}{\sqrt{30}} = 2,74$
 - ▶ $P(100 - 1,96 * 2,74 < \bar{x} < 100 + 1,96 * 2,74) = 0,95$
 - ▶ $P(94,63 < \bar{x} < 105,37) = 0,95$
-



Excel-Sheet: Konfidenz Mittelwert

Grundgesamtheit

$$E(X) = 100$$

$$V(X) = 225$$

$$\sigma(X) = 15$$

Legende:

xq...arithmetisches Mittel (x-quer)

Stichprobenumfang

$$n = 30$$

$$E(xq) = 100$$

$$V(xq) = V(X)/30 = 7,5$$

$$\sigma(xq) = 2,74$$

Irrtumswahrscheinlichkeit

$$\alpha = 0,05$$

$$z\text{-Wert} = 1,96$$

Grenzen des zentralen Schwankungsintervalls für das arithmetische Mittel

$$UG = 94,63$$

$$OG = 105,37$$

Konfidenzintervall bei bekannter Varianz

Durch Inversion der Formel für das Schwankungsintervall ergibt sich die Formel für das Konfidenzintervall bei bekannter Varianz.

$$P\left(\bar{x} - z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{x} + z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

Allerdings ist das Anwendungs-Szenario „Ziehen einer Stichprobe aus einer Grundgesamtheit mit bekannter Varianz“ bei praktischen Fragestellungen extrem selten, weshalb wir uns direkt mit dem Fall unbekannter Varianzen beschäftigen.

Unbekannte Varianz

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

- ▶ Ersetzen des unbekanntes Wertes von σ^2 durch die Stichprobenschätzung

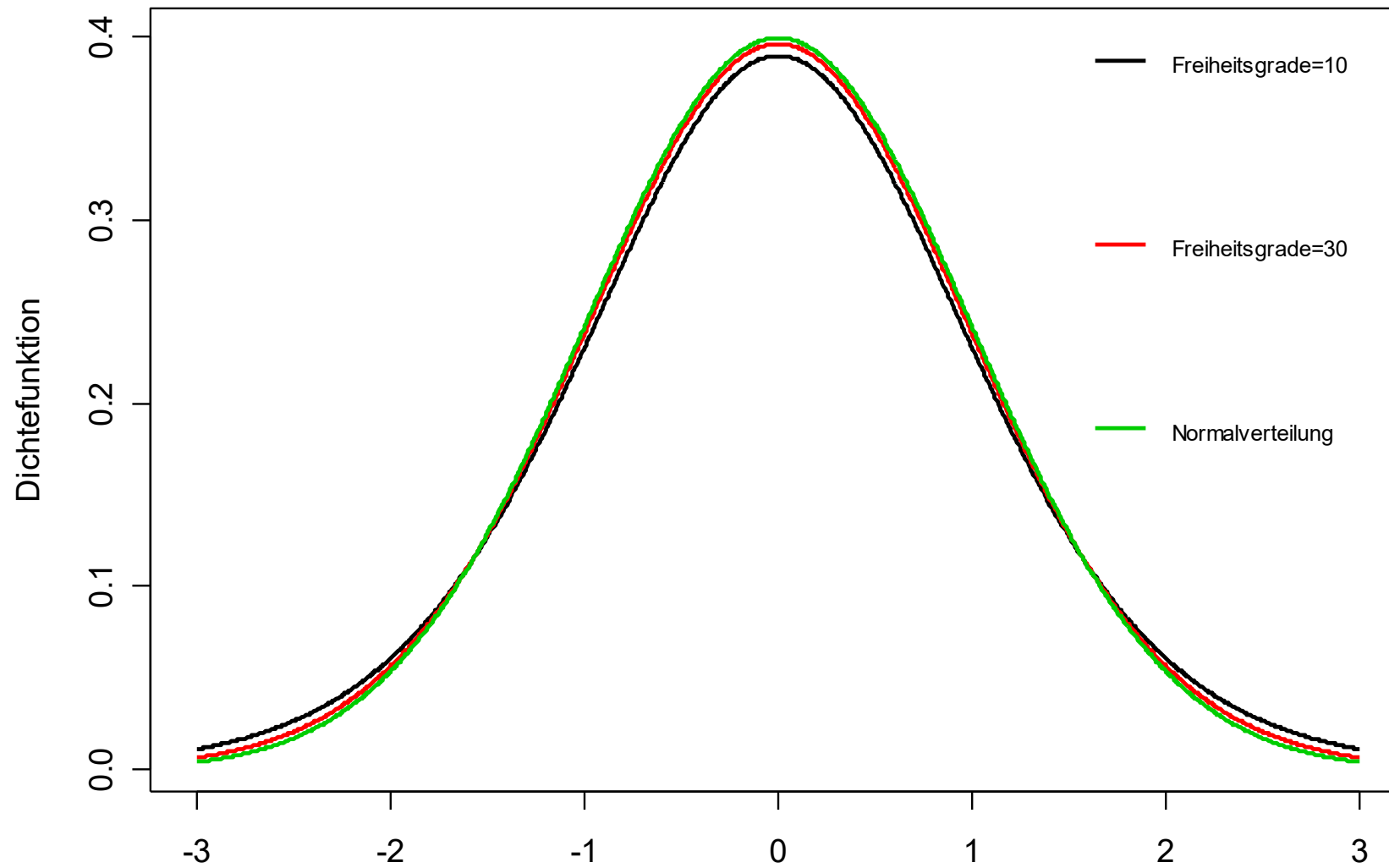
$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Schätzt man aus der Stichprobe, die Varianz der Grundgesamtheit ist die Division durch (n-1) angebracht

- ▶ Die aus der Varianzschätzung resultierende zusätzliche Unsicherheit muss zumindest bei kleinen Stichproben kompensiert werden
-



Exkurs Student-Verteilung (t-Verteilung)



Konfidenzintervall bei unbekannter Varianz

$$P\left(\bar{x} - t_{n-1;1-\alpha/2} \frac{\hat{\sigma}}{\sqrt{n}} \leq \mu \leq \bar{x} + t_{n-1;1-\alpha/2} \frac{\hat{\sigma}}{\sqrt{n}}\right) = 1 - \alpha$$

Konfidenzintervall für μ bei unbekannter Varianz

Falls n groß ist (Faustregel: $n > 30$) können wir auch mit der Normalverteilung arbeiten :

$$P\left(\bar{x} - z_{1-\alpha/2} \frac{\hat{\sigma}}{\sqrt{n}} \leq \mu \leq \bar{x} + z_{1-\alpha/2} \frac{\hat{\sigma}}{\sqrt{n}}\right) = 1 - \alpha$$

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

In Worten

- ▶ Wir erhalten ein Konfidenzintervall für den „wahren“ Mittelwert μ in der Grundgesamtheit durch folgende Formel:
- ▶ Mittelwert der Stichprobe plus/minus dem geschätzten Standardfehler multipliziert mit dem zugehörigen Quantil der Normalverteilung (bei kleinen Stichproben Student-Verteilung)
- ▶ Allgemeines Prinzip von Konfidenzintervallen

Stichprobenmittelwert plus/minus

Tabellenwert mal geschätzter Standardfehler



Konfidenzintervall für Mittelwert

- ▶ Eine Stichprobenuntersuchung unter $n=300$ Angestellten einer bestimmten Branche ergab folgendes Ergebnis:
- ▶ Durchschnittseinkommen = 1.200 €
- ▶ Standardabweichung = 140 €
- ▶ Gesucht ist ein Konfidenzintervall, für das Durchschnittseinkommen in der Grundgesamtheit
- ▶ Unbekannte Varianz; großes n ; zentraler Grenzwertsatz

$$P\left(\bar{x} - z_{1-\alpha/2} \frac{\hat{\sigma}}{\sqrt{n}} \leq \mu \leq \bar{x} + z_{1-\alpha/2} \frac{\hat{\sigma}}{\sqrt{n}}\right) = 1 - \alpha$$

Konfidenzintervall für Mittelwert

▶ $Z=1,96$ $\frac{\hat{\sigma}}{\sqrt{n}} = \frac{140}{\sqrt{300}} = 8,08$

▶ $P(1200-1,96*8,08 < \mu < 1200+1,96*8,08) = 0,95$

▶ $P(1184,16 < \mu < 1215,84) = 0,95$

▶ Bei Verwendung der t-Verteilung:

▶ $P(1184,09 < \mu < 1215,91) = 0,95$

▶ Ein Intervall von $[1184,09; 1215,91]$ überdeckt den unbekanntem Mittelwert der Grundgesamtheit mit 95%-iger Wahrscheinlichkeit.

Excel-Sheet: Konfidenz Mittelwert

Legende:

xq...arithmetisches Mittel (x-quer)

s...Standardabweichung der Stichprobenwerte

nur bei großem n anzuwenden

Stichprobenergebnisse

n= 300

xq= 1200

s= 140

s(xq)= 8,082904

Irrtumswahrscheinlichkeit

alpha= 0,05

t-Wert= 1,9679

z-Wert 1,960

Konfidenzintervall für den Erwartungswert (Mittelwert der Grundgesamtheit)

UG= 1184,09

UG= 1184,16

OG= 1215,91

OG= 1215,84

Implementation mit R

```
# Example Confidence Interval for Mean
ci.mean <- function(n, xq, sd, alpha=0.05)
{
  #ad.hoc function
  SE <- sd/sqrt(n)           # standard error
  LB <- xq - qt(1-alpha/2, n-1)*SE
  UB <- xq + qt(1-alpha/2, n-1)*SE
  list(sample.size=n, mean.observed=xq,
        confidence.level=1-alpha, lower.bound=LB,
        upper.bound=UB)
}
```

```
ci.mean(300, 1200, 140)
ci.mean(300, 1200, 140, alpha=0.01)
```

```
> ci.mean(300, 1200, 140)
$sample.size
[1] 300

$mean.observed
[1] 1200

$confidence.level
[1] 0.95

$lower.bound
[1] 1184.093

$upper.bound
[1] 1215.907
```

```
> ci.mean(300, 1200, 140, alpha=0.01)
$sample.size
[1] 300

$mean.observed
[1] 1200

$confidence.level
[1] 0.99

$lower.bound
[1] 1179.046

$upper.bound
[1] 1220.954
```

Konfidenzintervall für Mittelwert

- ▶ Stichprobe unter 50 Haushalten einer Stadt mit Kategorie-A Wohnungen zwischen 80-100m² ergab:

- ▶ Kaltmiete pro m²:

- ▶ Arithm. Mittelwert: 8,30€

- ▶ Standardabweichung: 2,10€

$$P\left(\bar{x} - z_{1-\alpha/2} \frac{\hat{\sigma}}{\sqrt{n}} \leq \mu \leq \bar{x} + z_{1-\alpha/2} \frac{\hat{\sigma}}{\sqrt{n}}\right) = 1 - \alpha$$

- ▶ 90%-Konfidenzintervall $z=1,645$

- ▶ Standardfehler: $2,10/\sqrt{50} = 0,297$

- ▶ $P(8,30 - 1,645 \cdot 0,297 < \mu < 8,30 + 1,645 \cdot 0,297) = 0,90$

- ▶ $P(7,81 < \mu < 8,79) = 0,90$

- ▶ Ein Intervall von [7,81 bis 8,79] überdeckt den unbekanntem Mittelwert der Grundgesamtheit mit 90%-iger Wahrscheinlichkeit.



Bestimmung der Fallzahl

- ▶ Bestimmung des Stichprobenumfanges, um eine vorgegebene Genauigkeit erzielen zu können
- ▶ Bei der Messung von Reaktionszeiten schätzt ein Psychologe aufgrund der Erfahrung aus früheren Studien die Standardabweichung auf 0,05sec.
- ▶ Wie groß muss die Stichprobe sein, damit er zu A) 95% bzw. B) 99% davon ausgehen kann, dass der maximale Schätzfehler nicht größer als 0,01sec sein wird?

$$e_{\max} = z \cdot \hat{\sigma} / \sqrt{n}$$

$$n > z^2 \cdot \hat{\sigma}^2 / e_{\max}^2$$

- ▶ A) $n > 1,96^2 \cdot 0,05^2 / 0,01^2 = 96,04 \implies n = 97$
- ▶ B) $n > 2,58^2 \cdot 0,05^2 / 0,01^2 = 166,4 \implies n = 167$