

# **Bayesianische Statistik:**

**Welche Probleme löst sie  
und welche nicht?**

Michael Höfler

# Frequentistische Statistik

Interessierender Parameter  $\theta$  (Wahrscheinlichkeit, Erwartungswert, Assoziation, Effekt)

## Bevor man Daten erhebt:

Verteilungsmodell für Daten, gegeben fester, aber unbekannter Wert von  $\theta$ :

$P(\text{Daten} \mid \theta)$

z.B. Wahrscheinlichkeit (Prävalenz) von Depression; Anzahl der Depressionsfälle in der Stichprobe (Daten),  $\mathbf{X}$ , ist binomialverteilt, Wahrscheinlichkeit für  $\mathbf{x}$  Fälle:

$$P(X = x) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$$

## Nach Datenerhebung:

Man hat  $\mathbf{X} = \mathbf{x}$  Fälle beobachtet, nun ist  $\mathbf{x}$  fest. Es handelt sich jetzt um eine Funktion von  $\theta$ , genannt Likelihood-Funktion:

$$L(\theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$$

# Maximum-Likelihood-Prinzip:

beste Schätzung für  $\theta$  ist der Wert, der dann  $L(\theta)$  maximiert;  
hier:  $\theta = x/n$  (Anteil Depressionsfälle in Stichprobe). So ein „Schätzer“  
ist auch Grundlage statistischer Tests über  $\theta$ .

Man macht keine Vorannahmen über  $\theta$ , alle Information über  $\theta$   
kommt aus den Daten.

## Wichtig:

Es gibt **kein** Verteilungsmodell für

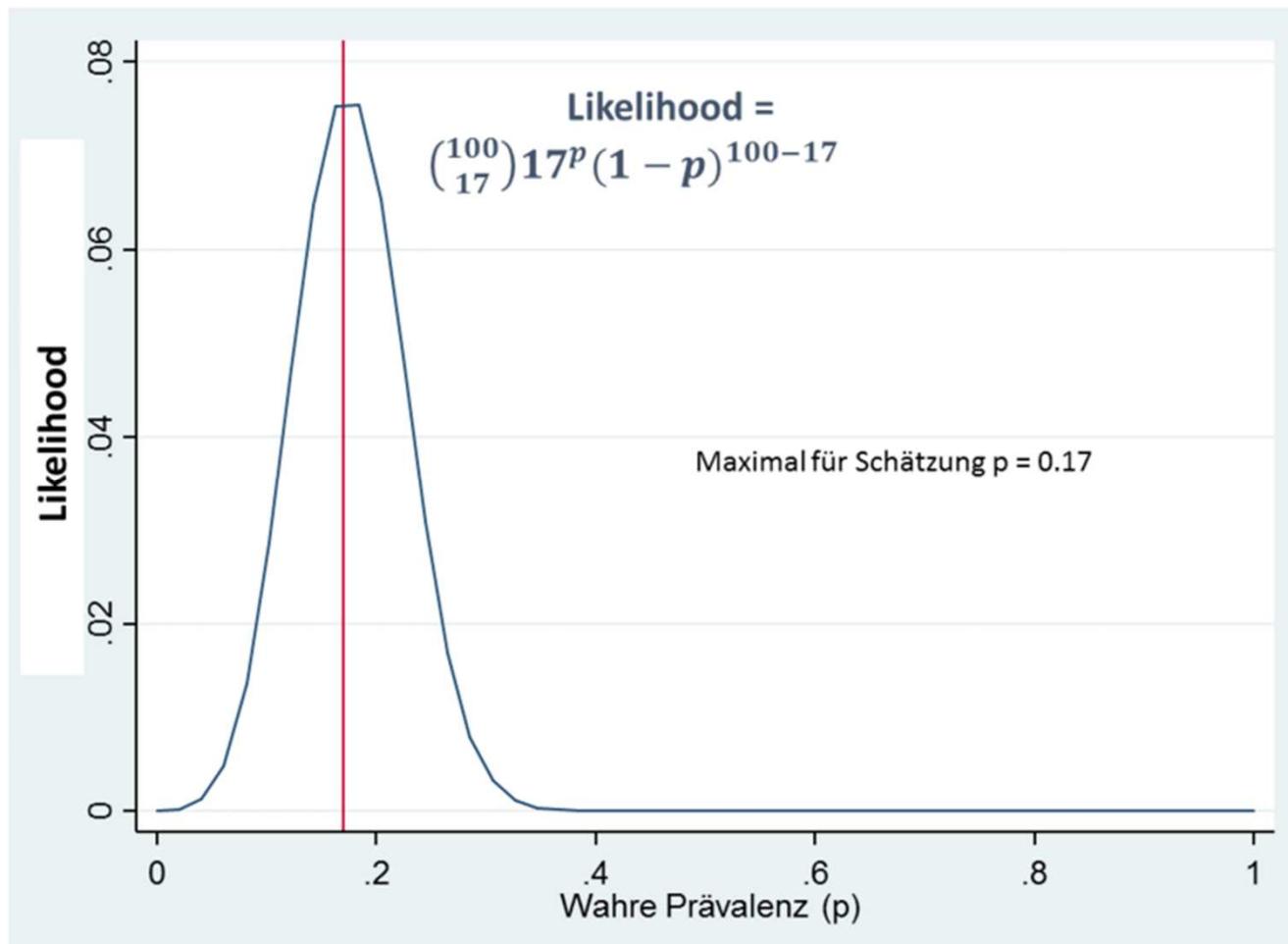
**$P(\theta | \text{Daten})$**

$\theta$  ist ein unbekannter, aber fester Wert.

**Keine Wahrscheinlichkeitsaussagen über  $\theta$  möglich!**

# (Beispiel numerisch)

In einer Stichprobe von 100 Personen erfüllen 17 Personen die Kriterien einer Major Depression. Die Maximum-Likelihood-Schätzung für die Depressions-Wahrscheinlichkeit,  $\theta = p$ , beträgt 0.17 (17%).



**Wahrscheinlichkeit für Zutreffen von  $H_0/H_1$  ist 0 oder 1!**

**Wahrscheinlichkeit, dass ein Konfidenzintervall das wahre  $\theta$  enthält, ist 0 oder 1!**

**Für erwünschte Wahrscheinlichkeitsaussagen ( $0 < P(H_1) < 1$ ) Interpretationskrücke nötig!**



Bildquelle: Pixaby.com

## Frequentistisches Gedankenexperiment

**Häufigkeitsinterpretation:** Vorstellung, man führt ganz oft (z.B. 10.000 mal) mit identischen Methoden eine Studie durch und analysiert (schätzen, testen) jedes mal auf gleiche Art. Zufall durch zufällige Stichprobenziehung/Gruppeneinteilung führt jedes mal zu anderem Ergebnis. In diesem Sinne kann man dann von Wahrscheinlichkeiten sprechen:

Falls  $H_0$  zutrifft, verwirft ein **Test** in  $100 \cdot \alpha$  % der Fälle  $H_0$  zu unrecht (5% falsch Positive, falls  $\alpha = .05$ ). Weiter keine absolute W.keitsaussage über  $H_1$  (nur bedingt auf  $H_0$ )!

Ein  $(1 - \alpha) \cdot 100\%$ -iges **Konfidenzintervall\***, z.B. 95%, für einen Parameter (z.B. Regressionskoeffizient  $\beta$ ), enthält in 95% der Fälle den wahren Parameterwert.

# Problem 1: Gedankenexperiment weit hergeholt

## 1. Funktioniert nur, wenn Studie **Zufallskomponente** hat

- Zufällige Stichprobenziehung
- Oder zufällige Gruppeneinteilung
- (sonst noch mehr Gedankenexperimente nötig, z.B. Vorstellung, dass es hinter einer untersuchten (sehr speziellen) klinischen Stichprobe eine größere Population gebe, die sich nicht systematisch von der Stichprobe unterscheidet).

## (... und Analyse sollte vollständig determiniert sein)

- z.B. immer t-Test durchführen, egal wie Verteilungen in Daten
- Wünschenswert aber: auch Alternativverfahren benutzen (z.B. U-Test), das viel weniger voraussetzt (sonst “analytical bias“)
- → unterschiedliche Analyse in unterschiedlichen Datensätzen, “Garden of forking paths“, frequentistische Interpretation stimmt nicht mehr, v.a. wenn dann **unterschiedliche Schätzer** untersucht werden (U-Test basiert nicht auf Mittelwertsunterschied, sondern Unterschied im Durchschnittsrang).

# Problem 2: falsche Interpretation des p-Werts

**p-Wert** =  $P(X > x \mid H_0)$ : Wahrscheinlichkeitsaussage, aber nur über **X** und bedingt auf  $H_0$

- **1 - p** ist **nicht** die Wahrscheinlichkeit für das Zutreffen von  $H_1$  (gegeben die Daten), **p** ist nicht Wahrscheinlichkeit für  $H_0$
- **1 - p** ist **kein** Maß für die **Replizierbarkeit** von Studienergebnissen
- **p** misst **nicht** die Größe/Bedeutsamkeit eines Zusammenhangs
- u.u.u.:



# Problem 3: „dichotomia“

- Reduziere Ergebnisse ohne Notwendigkeit auf **binäre Entscheidung**,  $p < .05$ ?
- **p** sollte eigentlich gar nicht kategorisiert werden (willkürlich; .05 ist bloß eine soziale Konvention, die mechanisch angewendet wird; s.u.)

PsycARTICLES: Journal Article

The earth is round ( $p < .05$ ).

© Request Permissions

Cohen, Jacob

American Psychologist, Vol 49(12), Dec 1994, 997-1003

After 4 decades of severe criticism, the ritual of null hypothesis significance testing (mechanical dichotomous decisions around a sacred .05 criterion) still persists. This article reviews the problems with this practice, including near universal misinterpretation of  $p$  as the probability that  $H_0$  is false, the misinterpretation that its complement is the probability of successful replication, and the mistaken assumption that if one rejects  $H_0$  one thereby affirms the theory that led to the test. Exploratory data analysis and the use of graphic methods, a steady improvement in and a movement toward standardization in measurement, an emphasis on estimating effect sizes using confidence intervals, and the informed use of available statistical methods are suggested. For generalization, psychologists must finally rely, as has been done in all the older sciences, on replication. (PsycINFO Database Record (c) 2016 APA, all rights reserved)



American Journal of Epidemiology  
© The Author(s) 2017. Published by Oxford University Press on behalf of the Johns Hopkins Bloomberg School of Public Health. All rights reserved. For permissions, please e-mail: journals.permissions@oup.com.

Vol. 186, No. 6  
DOI: 10.1093/aje/kwx259  
Advance Access publication:  
August 21, 2017

Invited Commentary

Invited Commentary: The Need for Cognitive Science in Methodology

Sander Greenland\*

\* Correspondence to Dr. Sander Greenland, Department of Epidemiology, School of Public Health, University of California, Los Angeles, Los Angeles, CA 90095 (e-mail contact only, e-mail: lesdomes@ucla.edu).

Initially submitted June 6, 2017; accepted for publication June 7, 2017.

“[...] degrading continuous measures of evidence into decisive conclusions, **feeding the strong cognitive bias of dichotomania: the compulsion to replace quantities with dichotomies (“black-and-white thinking”)**, even when such dichotomization is unnecessary and misleading for inference.”

# Problem 4: Falschinterpretation von Testergebnissen je nach Stichprobengröße

- **Kleine Stichproben,  $p \geq .05$** : Behauptung, Nullhypothese treffe zu, obwohl geringe Power (Sekundäranalysen!)
- **Große Stichproben,  $p < .05$** : Behauptung, es gebe einen bedeutsamen Unterschied, dabei kann bereits kleiner Unterschied zur Ablehnung von  $H_0$  führen; falsche Praxis von \*, \*\*, \*\*\*-Signifikanz (hat keine methodische Basis; Vermischung der stat. Schulen von K. Pearson und R.A. Fisher; wahrnehmungspsychologisch sehr problematisch; ).

# Problem 5: Statistische Tests sind asymmetrisch

Die Verteilung von Teststatistiken ist oft nur berechenbar, wenn  $H_0$  zutrifft (z.B.  $\chi^2$ -Unabhängigkeitstest) → man kann besser auf  $H_1$  schließen.

$p < \text{vorgegebenes } \alpha$ , z.B.  $\alpha = .05 \rightarrow H_0 \text{ verwerfen}$

$p \geq \text{vorgegebenes } \alpha$ , z.B.  $\alpha = .05 \rightarrow H_0 \text{ beibehalten}$

Korrekt: Es gibt keine Evidenz/keinen Hinweis dafür, dass  $H_0$  verletzt ist.

Falsch:  $H_0$  trifft zu.

# Problem 6: Statistisches Modell kann falsch sein

Falsche Ergebnisse, weil z.B. die Annahmen eines Modells (ausgedrückt in Likelihood-Funktion) verletzt sind („analytical bias“; eher Regel als Ausnahme).



The image shows a screenshot of a journal article page. At the top, there is a header for 'Behaviour Research and Therapy' with the Elsevier logo on the left and a ScienceDirect link on the right. The article title is 'Robust statistical methods: A primer for clinical psychology and experimental psychopathology researchers' by Andy P. Field and Rand R. Wilcox. Below the title, there are footnotes for the authors' affiliations. The page is partially obscured by a vertical scrollbar on the left and a horizontal scrollbar at the bottom.

Contents lists available at [ScienceDirect](#)

**Behaviour Research and Therapy**  
journal homepage: [www.elsevier.com/locate/brat](http://www.elsevier.com/locate/brat)

**Robust statistical methods: A primer for clinical psychology and experimental psychopathology researchers**

Andy P. Field <sup>a,\*</sup>, Rand R. Wilcox <sup>b</sup>

<sup>a</sup> School of Psychology, University of Sussex, Falmer, Brighton, BN1 9QH, UK  
<sup>b</sup> Department of Psychology, University of Southern California, 618 Seeley Mudd Building, University Park Campus, Los Angeles, CA 90089-1061, USA

# Problem 7: Annahmen sind intransparent

- Annahmen im Gedankenexperiment
- Bevor Daten erhoben, sind alle Parameterwerte (z.B. Effekt = -2, 1, 0, +1, +2) „gleichwahrscheinlich“, „Gleichverteilung“ (korrekter: alle gleichbreiten Intervalle haben Wahrscheinlichkeit, z.B. [-2, -1], [-1, 0], [0, 1] ...], stetige Verteilung, Wahrscheinlichkeit für jeden Punktwert = 0)
- Mit 100% Wahrscheinlichkeit keinerlei Bias in einer Schätzung, z.B. durch gemeinsame Ursachen von **X** und **Y**, Messfehler in **X** und **Y**, Selektion, s.u.

# Problem 8: instabile Schätzungen

- Kleine Datensätze, zu viele Einflussvariablen
- Sehr große Zusammenhänge/Effekte (Regressionskoeffizienten  $\beta$ ) aber oft unwahrscheinlich
- In frequentistischer Statistik variieren Schätzungen der  $\beta$  jedoch völlig frei (zwischen  $-\infty$  und  $+\infty$ ).

## Research Methods & Reporting

### Sparse data bias: a problem hiding in plain sight

BMJ 2016 ; 352 doi: <https://doi.org/10.1136/bmj.i1981> (Published 27 April 2016)

Cite this as: BMJ 2016;352:i1981

Sander Greenland, professor of epidemiology and statistics<sup>1</sup>, Mohammad Ali Mansournia, assistant professor of epidemiology<sup>2</sup>, Douglas G Altman, professor of statistics in medicine<sup>3</sup>

---

<sup>1</sup>Department of Epidemiology and Department of Statistics, University of California, Los Angeles, CA, USA

<sup>2</sup>Department of Epidemiology and Biostatistics, School of Public Health, Tehran University of Medical Sciences, PO box 14155-6446, Tehran, Iran

<sup>3</sup>Centre for Statistics in Medicine, Nuffield Department of Orthopaedics, Rheumatology and Musculoskeletal Sciences, University of Oxford, Oxford, UK

---

Correspondence to: M A Mansournia [mansournia\\_ma@yahoo.com](mailto:mansournia_ma@yahoo.com)

Accepted 7 March 2016

# Problem 9: grundsätzlich falscher Gebrauch von Statistik

- Statistik wird **gedankenlos, mechanisch** und als bloßes **Ritual** verwendet.
- **p-hacking**: Missbrauch der Variation in den Ergebnissen je nach Auswertungsmethode, um best. Ergebnis zu erhalten

## Statistical Rituals: The Replication Delusion and How We Got There

Gerd Gigerenzer

First Published June 14, 2018 | Research Article |  Check for updates

<https://doi.org/10.1177/2515245918771329>

Article information 

Altmetric | 347 

### Abstract

The "replication crisis" has been attributed to misguided external incentives gamed by researchers (the *strategic-game hypothesis*). Here, I want to draw attention to a complementary internal factor, namely, researchers' widespread faith in a statistical ritual and associated delusions (the *statistical-ritual hypothesis*). The "null ritual," unknown in statistics proper, eliminates judgment precisely at points where statistical theories demand it. The crucial delusion is that the  $p$  value specifies the probability of a successful replication (i.e.,  $1 - p$ ), which makes replication studies appear to be superfluous. A review of studies with 839 academic psychologists and 991 students shows that the replication delusion existed among 20% of the faculty teaching statistics in psychology, 39% of the professors and lecturers, and 66% of the students. Two further beliefs, the illusion of certainty (e.g., that statistical significance proves that an effect exists) and Bayesian wishful thinking (e.g., that the probability of the alternative hypothesis being true is  $1 - p$ ), also make successful replication appear to be certain or almost certain, respectively. In every study reviewed, the majority of researchers (56%–97%) exhibited one or more of these delusions. Psychology departments need to begin teaching statistical thinking, not rituals, and journal editors should no longer accept manuscripts that report results as "significant" or "not significant."

### Keywords

replication, p-hacking, illusion of certainty, p value, null ritual

- Testergebnisse werden als **Wahrheitsaussage** statt als **Verhaltensregel** verstanden. Ein signifikantes Ergebnis bedeutet nur, dass es eine Auffälligkeit gibt, die erklärt werden muss.

---

 American Journal of Epidemiology  
© The Author(s) 2017. Published by Oxford University Press on behalf of the Johns Hopkins Bloomberg School of Public Health. All rights reserved. For permissions, please e-mail: journals.permissions@oup.com.

Vol. 186, No. 6  
DOI: 10.1093/aje/kwx259  
Advance Access publication:  
August 21, 2017

---

**Invited Commentary**

---

**Invited Commentary: The Need for Cognitive Science in Methodology**

---

**Sander Greenland\***

\* Correspondence to Dr. Sander Greenland, Department of Epidemiology, School of Public Health, University of California, Los Angeles, Los Angeles, CA 90095 (e-mail contact only; e-mail: lesdomes@ucla.edu).

*Initially submitted June 6, 2017; accepted for publication June 7, 2017.*

---

[...] both confidence intervals and  $\alpha$ -level tests were conceived as **decision rules for behavior** but were rapidly misinterpreted as **rules for belief**, and thus fed the false notion that a single study can by itself tell us whether an effect is present or absent.”

# Bayesianische Statistik

Frequentistische Statistik:  $P(\text{Daten} | \theta)$

Wie erhält man daraus das interpretierbare  $P(\theta | \text{Daten})$ ?  
Wahrscheinlichkeit z.B. für  $\theta = \text{Effekt} > 0$ ,  
gegeben die Daten (und das Modell dazu/Likelihood)?

(Man kann auch analog wie im Folgenden vorgehen, nur bedingt darauf, dass eine  $H_1$  über  $\theta$  zutrifft; „prior“ ist dann die Wahrscheinlichkeit für das Zutreffen von  $H_1$ , bevor man die Daten sieht:

$P(\text{Daten} | H_1) \rightarrow P(H_1 | \text{Daten})$ ?

**Zweiseitige Hypothesen** über  $\theta$ , z.B.  $H_0: \theta = 0$  lassen sich damit untersuchen; aber nicht, wenn man (wie auf den folgenden Folien) eine stetige Verteilung („prior“) über  $\theta$  zugrundelegt, dann würde gelten:  $P(H_0) = 0$ .)

# Theorem von Bayes:

## Prior:

Wahrscheinlichkeitsverteilung der Größe des Effekts, bevor man Daten untersucht

## Posterior:

$$P(\text{Effekt}|\text{Daten}) = \frac{P(\text{Daten}|\text{Effekt}) \times P(\text{Effekt})}{\text{Konstante} *}$$

Falls prior = **Gleichverteilung** („flat prior“) = jeder Wert gleichwahrscheinlich

$$= \frac{P(\text{Daten}|\text{Effekt})}{\text{Konstante} *}$$

Prior → „Basisratenproblem“ in der Psychologie des Entscheidens:

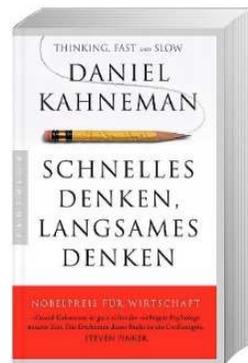


Bild: weltbild.com

= Ergebnis frequentistischer Statistik, aber bessere Interpretation mit **transparenter Annahme** (man weiß nichts über den Effekt, bevor man die Daten sieht)

\* Konstante dient nur dazu, dass das Ergebnis als Verteilung interpretiert werden kann; sorgt dafür, dass „Wahrscheinlichkeit über alle möglichen  $\beta$ -Werte“ = 1

(V.a. für kausale Inferenz erweitert um **Bias-Parameter\*** und Verteilungen, die die Unsicherheit über diese modellieren)

*J. R. Statist. Soc. A (2005)  
168, Part 2, pp. 267–306*

**Multiple-bias modelling for analysis of observational data**

Sander Greenland  
University of California, Los Angeles, USA

# Theorem von Bayes, erweitert:

Verteilungen von Bias-Parametern\*, die Unsicherheiten über diese widerspiegeln

prior

$$P(\text{Effekt} | \text{Daten}, \text{Annahmen}) = \frac{P(\text{Daten} | \text{Effekt}) P(\text{Annahmen}) P(\text{Effekt})}{\text{Konstante}}$$

Wahrscheinlichkeit für überhaupt kein Bias = 1!

prior = Gleichverteilung

\* Z.B. Effekte unberücksichtigten Confounders auf **X** und **Y**, Selektionswahrscheinlichkeiten, Messfehler.

$$= \frac{P(\text{Daten} | \text{Effekt})}{\text{Konstante}}$$

= Ergebnis frequentistischer Statistik

**Posterior** = Wahrscheinlichkeitsverteilung von  $\theta$

(z.B. Ausmaß eines Effekts), **gegeben Daten und prior**  
(und Modell/Likelihood)

**Vorteil:** daraus lässt sich alles bestimmen, was man wissen möchte, z.B. **Wahrscheinlichkeit** für

$H_1 = \text{Effekt} > 0$

Effekt  $> \mathbf{s}$

Effekt liegt in beliebigem Intervall  $[\mathbf{s}_1, \mathbf{s}_2]$ .

$\mathbf{s}$  = minimale Größe für inhaltlich relevanten Effekts,  
z.B. via „klinische Signifikanz“

# Umsetzung: Wie bestimmt man die Posterior?

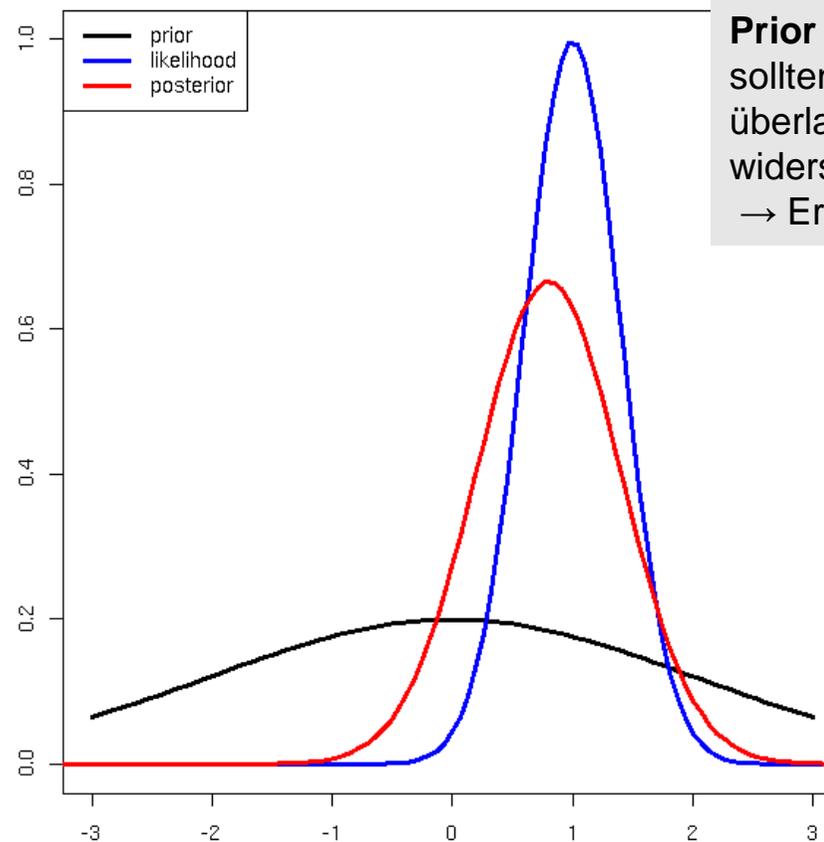
Simuliert man oft mit MCMC = „Monte-Carlo-Markov-Chain“

<https://www.youtube.com/watch?v=OTO1DygELpY&feature=youtu.be>

Beispiel für **prior**,  
**likelihood** und  
**posterior**

**Prior:** Normalverteilung um Erwartungswert 0 = kein Effekt; weder positiven, noch negativen Werten Präferenz gegeben; große Varianz, aber große positive wie negative Effekte unwahrscheinlicher

**Likelihood:** Normalverteilung um aus Daten geschätzten Wert von 1; hier viel weniger Varianz (Schätzung von Regressionskoeff.  $\beta$  ist normalverteilt, zum. In großen Stichproben)



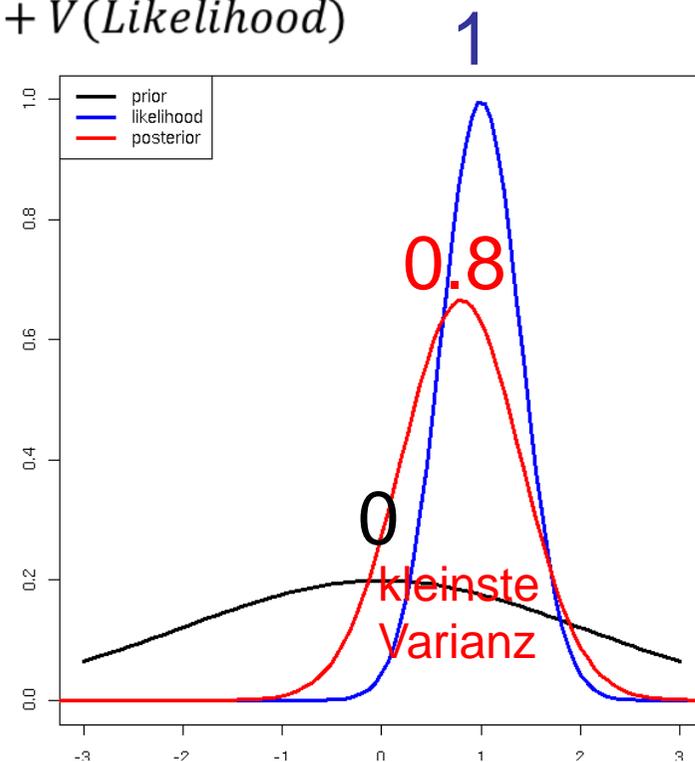
**Prior und likelihood** sollten sich gut überlappen, sonst widersprüchlich, → Erklärungsbedarf)

# Beispiel Normalverteilung

Erwartungswert (E) der Posterior = mit Varianz (V) gewichtetes Mittel von Prior und Likelihood (“inverse variance-weighted“):

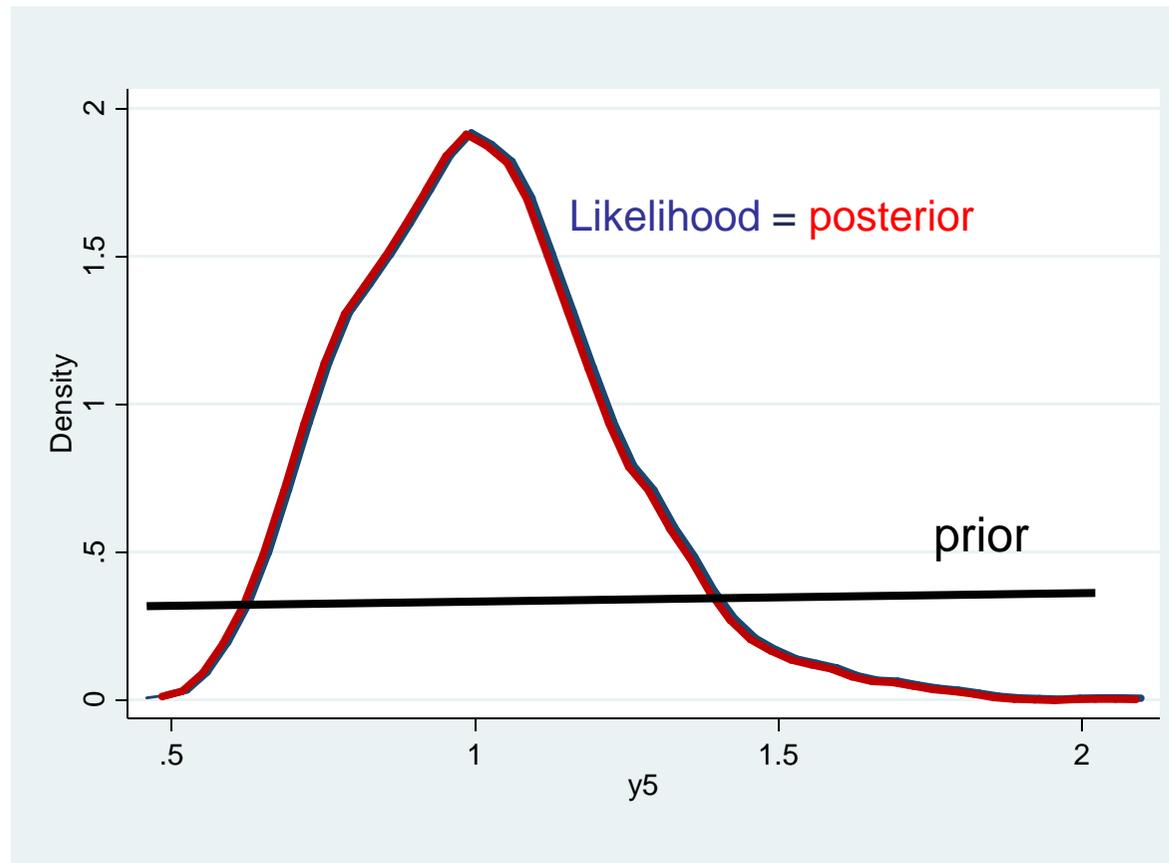
$$E(\text{Posterior}) = \frac{\frac{E(\text{Prior})}{V(\text{Prior})} + \frac{E(\text{Likelihood})}{V(\text{Likelihood})}}{V(\text{Prior}) + V(\text{Likelihood})}$$

**Posterior** hat kleinste Varianz; je kleiner die Varianz der **Prior**, desto größer ihr Einfluss auf die Posterior.

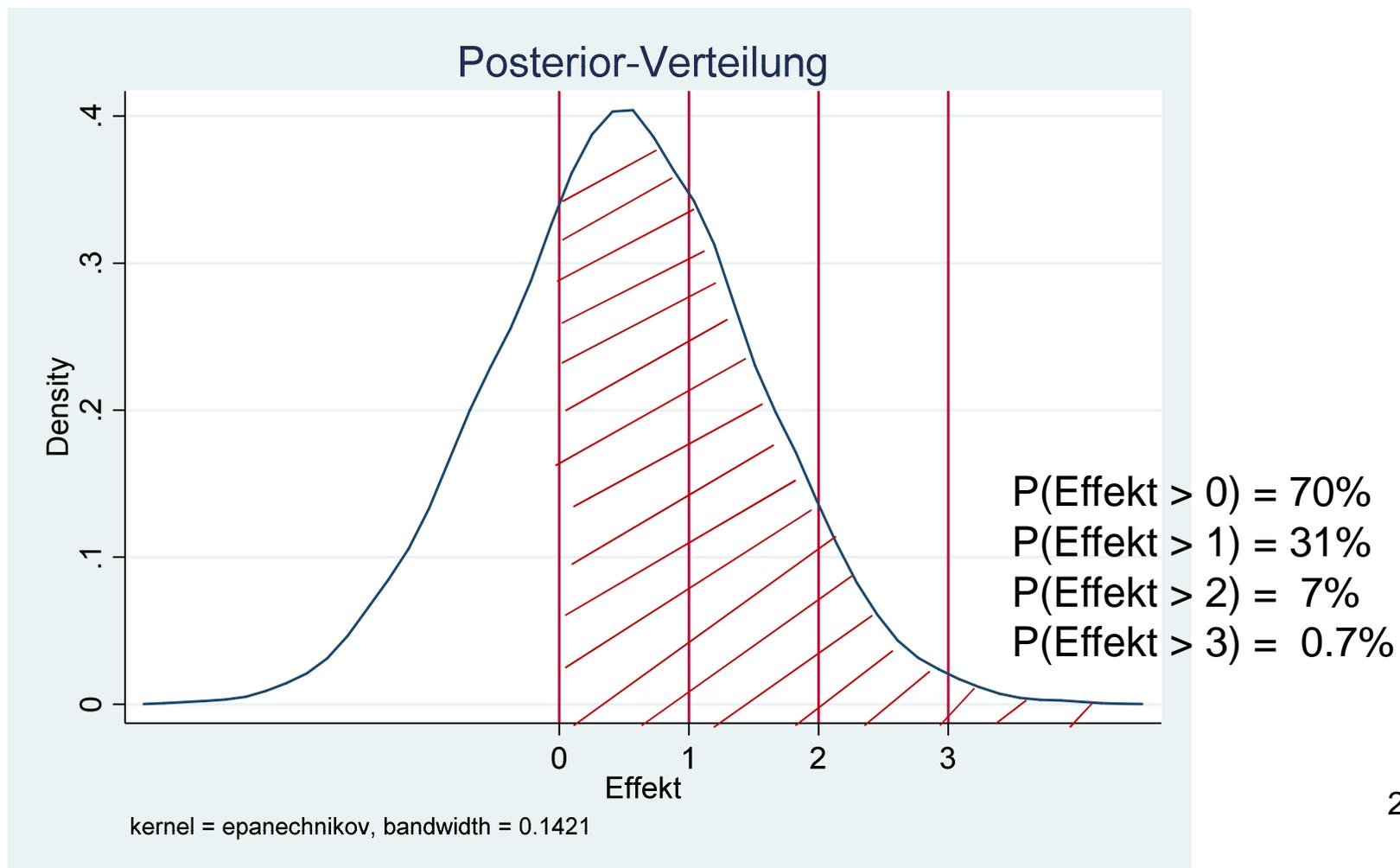


# Frequentische Statistik bayesianisch betrachtet: Gleichverteilung

(„flat prior“, „uninformative prior“)  $\Rightarrow$  **likelihood** = **posterior**



# Allgemein: Aus der **Posterior** ist jede Intervall-Wahrscheinlichkeit berechenbar



# Fiktives Beispiel

In einem Experiment habe man zwei Gruppen à 50 Personen randomisiert in  $X = 0$  und  $X = 1$  eingeteilt.

$X = 1$ : Übungen zum Training der Intelligenz durchgeführt.

$X = 0$ : Magazin zum Lesen gegeben, z.B. „Beef!“.

$Y$  = Wert in IQ-Test einen Tag später.

# Frequentistische Auswertung

```
. regress y x
```

Source	SS	df	MS	Number of obs = 100
Model	916.334659	1	916.334659	F(1, 98)
Residual	25274.2761	98	257.900776	Prob > F
Total	26190.6107	99	264.551624	R-square
				Adj R-sq
				Root MSE

Unterschied zweiseitig nicht signifikant ( $p=0.062$ ), aber unter  $X = 1$  könnte IQ um bis zu 12.4 höher sein!

y	Coef.	Std. Err.	t	P> t	[95% Conf. Intervall]
x	6.054204	3.211858	1.88	0.062	-.3196234 12.42803
_cons	99.18632	2.271126	43.67	0.000	94.67934 103.6933

**Einseitiger p-Wert für  $H_1$ :**

**Effekt > 0 =  $0.062/2 =$**

**0.031** (erlaubt, falls Schätzung von  $\beta$  das erwartete Vorzeichen hat (hier positiv), sonst  $p > 0.5$ )).

# Bayesianische Auswertung

## a. „flat prior“: gleiches Ergebnis

(bis auf zufällige Abweichungen — durch Simulationsmethode in der Berechnung):

```

. bayes , rseed(64674376) prior({ y:x }, normal(0, 1000000)) mcmcsize(100000) burnin(5000) thinning(2) savin
> g(post) : regress y x
note: discarding every other sample observation; using observations 1,3,5,...

Burn-in ...
Simulation ...

file post.dta saved

Model summary
-----
Likelihood:
  y ~ regress(xb_y, {sigma2})

Priors:
  {y:x} ~ normal(0,1000000)          (1)
  {y:_cons} ~ normal(0,10000)       (1)
  {sigma2} ~ igamma(.01,.01)

(1) Parameters are elements of the linear form xb_y.

Bayesian linear regression          MCMC iterations = 204,999
Random-walk Metropolis-Hastings sampling  Burn-in = 5,000
                                          MCMC sample size = 100,000
                                          Number of obs = 100
                                          Acceptance rate = .4402
                                          Efficiency: min = .1509
                                          avg = .2263
                                          max = .3762
Log marginal likelihood = -434.50671
    
```

	Mean	Std. Dev.	MCSE	Median	Equal-tailed [95% Cred. Intervall]	
y	<b>Punktschätzung</b>		<b>SE</b>			
	x	6.109508	3.255388	.026502	6.12945	- .3341289 12.51243
	_cons	99.12405	2.30417	.018705	99.13456	97.73206 100.51604
	sigma2	263.2974	38.48576	.198435	259.4171	198.5666 349.0749

Note: Default priors are used for some model parameters.

technische Details

Bayesianisches Schätzintervall = „credibility interval“

Anmerkung: Es gibt hier keine „p-Werte“, aber man kann zum Hypothesenprüfen vorgehen wie auf der folgenden Folie.

Daraus berechnen, **wie wahrscheinlich Differenz > 0** ist:

```
. bayestest interval {y:x}, lower(0) upper(100)

Interval tests      MCMC sample size =      10,000

prob1 : 0 < {y:x} < 100
```

	Mean	Std. Dev.	MCSE
prob1	.9721	0.16469	.0038601

Gegeben die Daten und die Prior beträgt diese Wahrscheinlichkeit 97.2%!

Theoretisch = 1 - einseitiger p-Wert (.031) vorhin

.. Oder wie wahrscheinlich **praktisch relevante Differenz von mindestens 10** (= 2/3\*SD)

```
. bayestest interval {y:x}, lower(10) upper(100)

Interval tests      MCMC sample size =      10,000

prob1 : 10 < {y:x} < 100
```

	Mean	Std. Dev.	MCSE
prob1	.116	0.32024	.008951

Gegeben die Daten und die Prior beträgt diese Wahrscheinlichkeit 11.6%!

**Andere Prior („informative“):**  $N(0, 56.25)$ , keine Präferenz für eine der beiden Bedingungen (Erwartungswert = 0), aber großer Effekt in beide Richtungen unwahrscheinlich ( $SD = \sqrt{56} = 7.5$ : jeweils 2.5% W.keit, dass eine der beiden Bedingungen um mehr als 15 (SD der IQ-Verteilung) besser ist; bei Normalverteilung liegt 95% W.keitsmasse zwischen Erwartungswert und  $\pm 2 \cdot SD$ ).

```
. bayes , rseed(64674270) prior({y:x} , normal(0,56.25)) mcmcsize(10000) burnin(5000) thinning(2) : regress y x
note: discarding every other sample observation; using observations 1,3,5,...
```

```
Burn-in ...
Simulation ...
```

Model summary

Likelihood:

```
y ~ regress(xb_y, {sigma2})
```

Priors:

```
{y:x} ~ normal(0,56.25) (1)
{y:_cons} ~ normal(0,10000) (1)
{sigma2} ~ igamma(.01,.01)
```

(1) Parameters are elements of the linear form xb\_y.

```
Bayesian linear regression           MCMC iterations =    24,999
Random-walk Metropolis-Hastings sampling  Burn-in           =     5,000
                                           MCMC sample size =   10,000
                                           Number of obs     =     100
                                           Acceptance rate   =    .4444
                                           Efficiency: min   =    .1677
                                           avg               =    .2434
                                           max               =    .3865
Log marginal likelihood = -429.99019
```

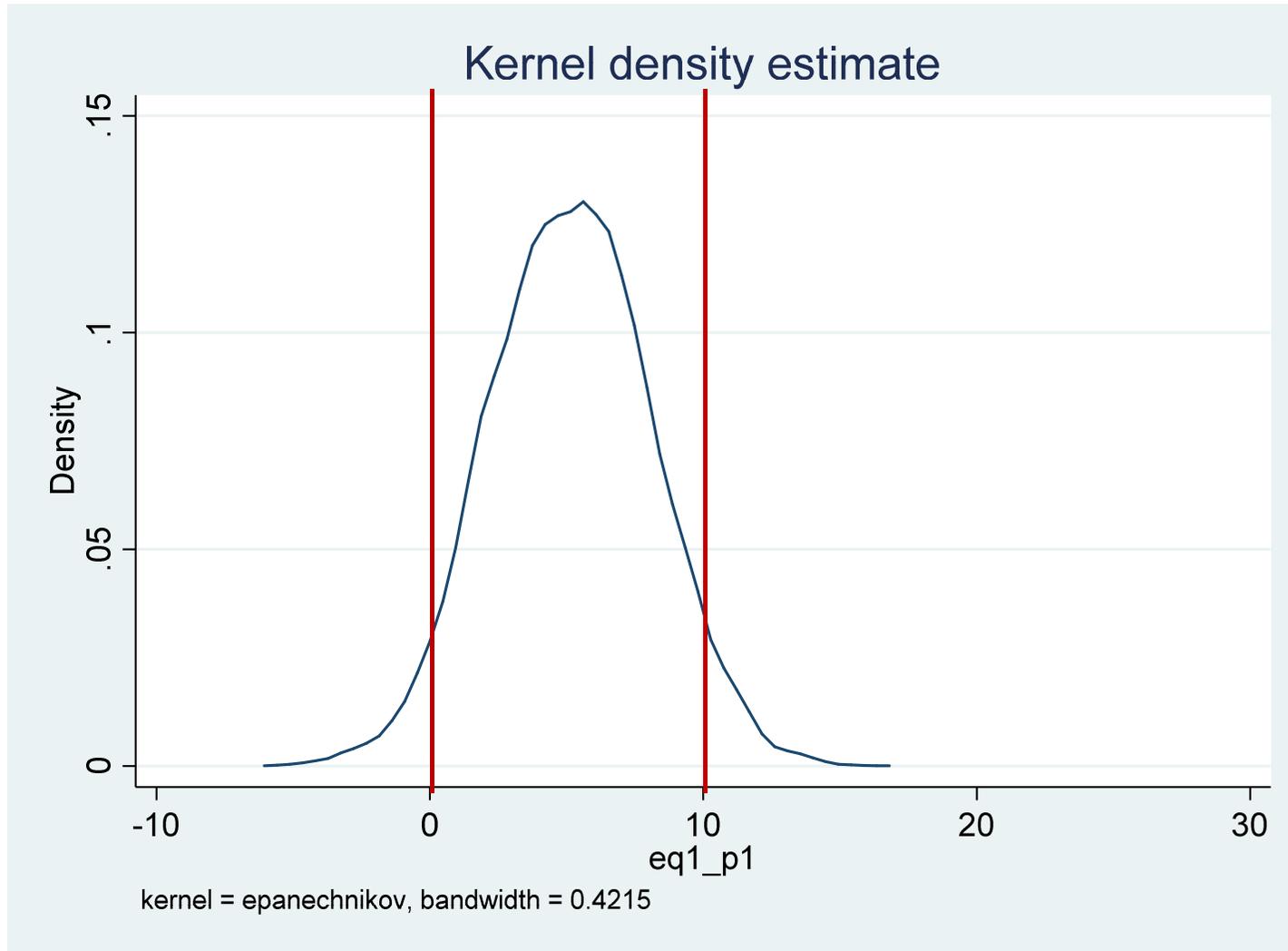
**Punktschätzung zu 0 hin verschoben (a-priori Erwartungswert = 0)**

	Mean	Std. Dev.	MCSE	Median	[95% Cred. Interval]
y					
x	5.153889	2.933098	.069891	5.170395	[-.4623962 10.96649]
_cons	99.81193	2.221009	.054232	99.56019	[95.38456 104.0596]
sigma2	263.2707	38.6592	.621865	259.0131	[199.058 350.758]

Schmaleres Schätzintervall als bei non-informative prior

Note: Default priors are used for some model parameters.

# Posterior



**Effekt von Intelligenztraining auf IQ-Wert**

Aus posterior ausrechnen, **wie wahrscheinlich** Differenz > 0 und Differenz > 10 sind:

```
. bayestest interval {y:x} , lower(0) upper(100)
```

```
Interval tests      MCMC sample size =      10,000
```

```
prob1 : 0 < {y:x} < 100
```

	Mean	Std. Dev.	MCSE
prob1	.9625	0.18999	.0032323

Gegeben die Daten und die prior beträgt diese Wahrscheinlichkeit hier 96% (kaum geringer).

Oder ein praktisch relevanter Unterschied von mindestens 10 IQ-Punkten:

```
. bayestest interval {y:x} , lower(10) upper(100)
```

```
Interval tests      MCMC sample size =      10,000
```

```
prob1 : 10 < {y:x} < 100
```

	Mean	Std. Dev.	MCSE
prob1	.0488	0.21546	.0036812

Gegeben die Daten und die prior beträgt diese Wahrscheinlichkeit hier nur 5%!

# Gretchenfrage: Wie bestimmt man die Prior? „educated guess“

(falls man sie nicht aus anderen Daten bestimmen kann)

Es gibt viele mögliche Verteilungstypen einer Prior, auch für schiefe Verteilungen, z.B. Gamma-Verteilung.

## Als ob man wetten würde

- Regressionskoeffizient  $\beta$  (Assoziation, Effekt) normalverteilt = viele, unabhängige Einflüsse auf Vorwissen
- Wert  $e$ , für den gilt: genauso viel darauf wetten, dass  $\beta > e$  wie dass  $\beta < e \rightarrow$  Erwartungswert =  $e$
- Bestimme Wert  $sd$ , sodass du 2:1 darauf wetten würdest, dass wahres  $\beta$  zwischen  $e - sd$  und  $e + sd$  liegt  
(bei Normalverteilung liegen ca. 2/3 Wahrscheinlichkeitsmasse zwischen  $e -$  Standardabweichung und  $e +$  Standardabweichung)
- Varianz =  $sd^2$
- Ergibt Prior =  $N(e, sd^2)$

# Mittels „data equivalents“

(ausführliches Beispiel über Assoziation im Anhang)

- Beispiel Schätzung unbekannter Prävalenz für Depression
- Nehme an, Vorwissen entspreche einer **fiktiven Studie**, die bei  $n = 100$   $p = 10/100 = 10\%$  Depression geschätzt hat
- Entspricht „Betaverteilung“\*  $\text{beta}(a,b) = \text{beta}(p^* (n-2) , (p-2)^* (1-p))$

\* Verteilungsmodell für Wahrscheinlichkeiten

[Int J Epidemiol. 2006 Jun;35\(3\):765-75. Epub 2006 Jan 30.](#)

**Bayesian perspectives for epidemiological research: I. Foundations and basic methods.**

[Greenland S<sup>1</sup>.](#)

# Spezielle Bayes-Methoden:

- „**Empirical (objective) Bayes**“: Prior wird aus Daten geschätzt. Formel von Bayes dient einfach dazu, Wissen aktuell zu halten (z.B. neue Probanden machen ein Paradigma):

$$P(\text{Effekt}|\text{neue Daten}) = \frac{P(\text{neue Daten}|\text{Effekt}) \times P(\text{Effekt}|\text{bisherige Daten})}{\text{Konstante}}$$

- „**Bayes-Faktor**“: Faktor, um den  $H_1$  wahrscheinlicher wird, nachdem man die Daten analysiert hat

Kritik: die absolute Wahrscheinlichkeit für  $H_1$  ist interessanter. Der Bayes-Faktor wird heuristisch eingeteilt in Stufen der Evidenz → **dichotomia**

# Spezielle Bayes-Methoden:

- **„Semi-Bayes“**: nur Priors für manche Parameter spezifizieren, z.B. für ein  $\beta_1 =$  Zusammenhang von  $\mathbf{X}$  und  $\mathbf{Y}$ , aber nicht für  $\beta_0 =$  Grundniveau von  $\mathbf{Y}$  (hierfür implizit flat prior verwendet; die Beispiele hier sind semi-bayesian)
- **„Reverse Bayes“**: **Wie unwahrscheinlich darf  $H_1$  gerade noch sein**, damit man, gegeben die Daten, den Schluss ziehen kann, dass  $H_1$  zutrifft?

$$P(H_1|\text{neue Daten}) = \frac{P(\text{neue Daten}|H_1) \times P(H_1|\text{bisherige Daten})}{\text{Konstante}}$$

# Welche Probleme löst die bayesianische Statistik nun?

## 1. Gedankenexperiment ✓

Entfällt, da man die gewünschte Wahrscheinlichkeit aus der Posterior enthält

## 2. Falsche Interpretation des p-Werts ✓

Bayesianische Wahrscheinlichkeiten haben die gewünschte Interpretation.

## 3. Dichotomia ✓ —

Solange man Wahrscheinlichkeit aus Posterior nicht wieder dichotomisiert (kategorisiert).

## 4. Fehlinterpretation je nach Stichprobengröße —

$P(H_1|\text{Daten})$  hängt (über Likelihood) von Stichprobengröße ab (jedoch weniger, umso informativer (geringere Varianz) die Prior ist).

## 5. Asymmetrie statistischer Tests ✓

$H_0$  und  $H_1$  werden symmetrisch behandelt

$$P(H_1|\text{Daten}) = 1 - P(H_0|\text{Daten})$$

## 6. Statistisches Modell kann falsch sein —

Analytical bias (in der Likelihood) geht in die Posterior ein, erhält allerdings durch die Prior geringeres Gewicht.

## 7. Annahmen intransparent

Annahmen stecken **explizit** in Prior, dadurch sollte diese Gegenstand der Diskussion über die Ergebnisse werden.

## 8. Instabile Schätzungen

Die Posterior hat höchstens so große Varianz wie die Likelihood. Bei entsprechendem Vorwissen kann man Schätzungen stark stabilisieren (siehe Anhang).

## 9. Falscher Gebrauch von Statistik

Durch die Wahl der Prior noch **eine Ebene mehr**, auf der man Variation in den Ergebnissen erzeugen kann.

→ Prior unbedingt **vorher registrieren!**

# Pro und kontra Bayes

- **Studienergebnisse nicht vergleichbar**, wenn sie unterschiedliche Priors benutzen  
**Ausweg**: immer auch Likelihood-Verteilung angeben.  
Eine Leserin kann diese dann mit einer Prior ihrer Wahl kombinieren und damit ihre **eigene Posterior** erhalten.

- Bayes ist durch die Prior **subjektiv**, frequentistische Statistik objektiv

Die **scheinbare Objektivität frequ. Statistik** basiert auf absurden Annahmen:

1. kein Vorwissen — ohne dem würde man bestimmten Effekt gar nicht untersuchen; Forschung beginnt nie bei null; wenn man trotzdem bei null anfängt, verschwendet man Ressourcen (da unnötig große Stichproben).
2. keinerlei Bias mit 100%iger Wahrscheinlichkeit

**Subjektivität lässt sich verringern**, z.B. durch data equivalents.

# Literatur

Baldwin SA, Larson MJ. An introduction to using Bayesian linear regression with clinical data. *Behaviour Research and Therapy* 2017; 98:58-75

Etz A, Vandekerckhove J. Introduction to Bayesian Inference for Psychology *Psychon Bull Rev* 2018; 25:5–34

Greenland S. Bayesian perspectives for epidemiological research: I. Foundations and basic methods. *Int J Epidemiol* 2006; 35:765-75

Greenland S. Putting background information about relative risks into conjugate priors. *Biometrics* 2001;57:663-70

Quintana DS, Williams DR: Bayesian alternatives for common null hypothesis significance tests in psychiatry: a non-technical guide using JASP. *BMC Psychiatry* 2018; 18:178

# Anhang: Beispiel für data equivalents und Stabilisierung von Schätzungen

- kleine Stichprobe mit großem Zufallsfehler in der Schätzung eines  $\theta$
- Beispiel für binäres  $X$ , binäres  $Y$
- Numerisch sehr instabiles, schlecht zu interpretierendes, aber häufig benutztes Zusammenhangsmaß: odds ratio (OR)

		Faktor X		
		0	1	
Outcome Y	0	49 (98.9%)	14 (90%)	
	1	1 (1.1%)	6 (10%)	
		50 (100%)	20 (100%)	70

- OR geschätzt als 21 (2.3 – 189.3), unrealistisch groß!

```
. logistic Y X [fweight=freq]
```

```
Logistic regression          Number of obs   =          70
                             LR chi2(1)           =          11.27
                             Prob > chi2          =          0.0008
Log likelihood = -17.119242   Pseudo R2      =          0.2477
```

Y	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]
X	21	23.55844	2.71	0.007	2.329783 189.288
_cons	.0204082	.0206154	-3.85	0.000	.0028181 .1477908

Note: \_cons estimates baseline odds.

#### research methods & reporting

### Sparse data bias: a problem hiding in plain sight

*BMJ* 2016 ;352 doi: <https://doi.org/10.1136/bmj.i1981> (Published 27 April 2016)

Cite this as: *BMJ* 2016;352:i1981

Sander Greenland, professor of epidemiology and statistics<sup>1</sup>, Mohammad Ali Mansournia, assistant professor of epidemiology<sup>2</sup>, Douglas G Altman, professor of statistics in medicine<sup>3</sup>

<sup>1</sup>Department of Epidemiology and Department of Statistics, University of California, Los Angeles, CA, USA

<sup>2</sup>Department of Epidemiology and Biostatistics, School of Public Health, Tehran University of Medical Sciences, PO box 14155-6446, Tehran, Iran

<sup>3</sup>Centre for Statistics in Medicine, Nuffield Department of Orthopaedics, Rheumatology and Musculoskeletal Sciences, University of Oxford, Oxford, UK

Correspondence to: M A Mansournia [mansournia\\_ma@yahoo.com](mailto:mansournia_ma@yahoo.com)

Accepted 7 March 2016



Beobachtet

Prior entspricht diesem Auffüllen der Daten

Nachgerechnet, man muss nach Indikator ind adjustieren:

		Häufigkeit in Zelle	Indikatorvariable für aufgefüllte Daten	
X	Y	freq	ind	
0	0	49	0	
0	1	1	0	
1	0	14	0	
1	1	6	0	
0	0	100000	1	
0	1	4.5	1	
1	0	100000	1	
1	1	4.5	1	

```
. logistic Y factor ind [iweight=freq]

Logistic regression               Number of obs   =           8
                                LR chi2(2)       =        94.86
                                Prob > chi2      =         0.0000
Log likelihood = -118.56639       Pseudo R2      =         0.2857
```

Y	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
factor	3.174619	1.78973	2.05	0.040	1.051518	9.584441
ind	.0003	.000164	-14.83	0.000	.0001027	.0008761
_cons	.0710765	.0348868	-5.39	0.000	.0271599	.1860051

Note: \_cons estimates baseline odds.

[Int J Epidemiol](#). 2006 Jun;35(3):765-75. Epub 2006 Jan 30.

**Bayesian perspectives for epidemiological research: I. Foundations and basic methods.**

[Greenland S](#)<sup>1</sup>.