

## 2.4.6 Konditionale Bayes-Inferenz: Begrifflicher Hintergrund und „Erinnerungen“

In Kapitel 1.6.1 wurde davon gesprochen, dass die Lösung/ Darstellung des datengestützten Entscheidungsproblems über das Auswertungsproblem und die damit verbunden Suche nach optimalen Entscheidungsfunktionen durchaus auch auf Kritik stösst. Diese stützt sich v.a. einerseits

- auf die immense computationale Komplexität

und andererseits ganz prinzipiell

- auf die Problematik kontrafaktischer Ereignisse bei der Bewertung von Entscheidungsfunktionen mittels der Risikofunktion.

Dies legt eine konditionale Sicht als mögliche Alternative nahe: Es werden optimale Lösungen für die konkret beobachtete Datenkonstellation  $\{x\}$  gesucht („auf  $\{x\}$  konditionierte Betrachtung“). Dies führt auf die „übliche Bayes-Inferenz“, die mit Hilfe der sogenannten Posteriori-Verteilung gegeben die Daten arbeitet, und in diesem Sinne hier zur Abgrenzung als „*konditionale Bayes-Inferenz*“ bezeichnet werde. Zur Vorbereitung werde an das Theorem von Bayes in seiner allgemeinen Form erinnert:

### Proposition 2.86 (Allgemeines Theorem von Bayes)

Seien  $X$  und  $U$  zwei Zufallsvariablen mit gemeinsamer Wahrscheinlichkeitsfunktion  $f_{X,U}(\cdot)$  bzw. Dichte  $f_{X,U}(\cdot)$  (bezüglich eines dominierenden  $\sigma$ -finiten Maßes  $\nu \otimes \lambda$ ) und den bedingten Wahrscheinlichkeitsfunktionen bzw. bedingten Dichten  $f_{X|U}(\cdot|u)$  und  $f_{U|X}(\cdot|x)$  (bezüglich  $\nu$  bzw.  $\lambda$ ).

Dann gilt:

$$f_{U|X}(u|x) = \frac{f_{X|U}(x|u) \cdot f_U(u)}{f_X(x)} \quad (2.41)$$

mit

$$f_X(x) = \int f_{X|U}(x|u) \cdot f_U(u) d\nu(u). \quad (2.42)$$

**Bem. 2.87**

Bei stetigem  $U$  mit Dichte  $f_U(u)$  erhält man also Proposition 2.86 mit

$$f_X(x) = \int f_{X|U}(x|u) \cdot f_U(u) du. \quad (2.43)$$

Im Fall von diskreten Zufallsvariablen  $X$  und  $U$  – mit  $\mathcal{U}$  als Träger von  $U$  – ergibt sich

$$p(\{U = u\}|\{X = x\}) = \frac{p(\{X = x\}|\{U = u\}) \cdot p(\{U = u\})}{p(\{X = x\})} \quad (2.44)$$

mit

$$p(\{X = x\}) = \sum_{u \in \mathcal{U}} p(\{X = x\}|\{U = u\}) \cdot p(\{U = u\}) \quad (2.45)$$

Betrachtet werde im Folgenden immer einer dieser beiden Spezialfälle .  
Die allgemeinere Formulierung über beliebige Dichten bezüglich geeigneter dominierender Maße ist unproblematisch.

**Bem. 2.88 (Normierungskonstante)**

$f_X(x)$  aus (2.42) spielt die Rolle einer reinen Normierungskonstante, die nicht von  $u$  abhängt. Häufig reicht es daher,  $f_{X|U}(x|u) \cdot f_U(u)$  zu berechnen. Da man weiß, dass sich insgesamt eine Wahrscheinlichkeitsdichte ergeben muss, kennt man implizit auch die Normierungskonstante. Man schreibt dann mit  $\propto$  als Symbol für „proportional zu“

$$f_{X|U}(x|u) \propto f_{X|U}(u|x) \cdot f_U(u).$$

## Bem. 2.89 (Konditionale Bayes-Inferenz: Konzeptionelle Hintergr

Gegeben sei ein datengestütztes Entscheidungsproblem  $((\mathbb{A}, \Theta, l(\cdot)); (\mathcal{X}, \sigma(\mathcal{X}), (p_{\vartheta}(\cdot))_{\vartheta \in \Theta})$ , wobei  $p_{\vartheta}(\cdot)$  die Wahrscheinlichkeitsfunktion bzw. Dichte  $f_{\vartheta}(\cdot)$  besitze, und eine Priori-Verteilung auf einem geeigneten Maßraum  $(\Theta, \sigma(\Theta))$  mit Dichte bzw. Wahrscheinlichkeitsfunktion  $\pi(\cdot)$ . Sei, als gedankliche Hilfskonstruktion,  $U$  („Umwelt“, „Natur“) eine „Zufallsgrösse“ (Zufallsvariable/-element), die das Eintreten des Umweltzustands  $\vartheta$  beschreibt und  $X$  eine Zufallsgrösse, die den Ausgang des Informationbeschaffungsexperiments beschreibt. Dann gelte für die Dichte bzw. Wahrscheinlichkeitsfunktion  $f_{X,U}(x, \vartheta)$  der gemeinsamen Verteilung von  $X$  und  $U$  für alle  $x \in \mathcal{X}$  und  $\vartheta \in \Theta$ :

$$f_{X,U}(x, \vartheta) = \pi(\vartheta) \cdot f_{\vartheta}(x)$$

Das heißt für alle  $x \in \mathcal{X}$  und  $\vartheta \in \Theta$  gilt

$$f_{\vartheta}(x) = f_{X|U}(x|\vartheta) ;$$

die Verteilung der Zufallsgrösse aus der Informationsstruktur wird als bedingte Verteilung von  $X$  gegeben  $U$  interpretiert (!!).

Dann ergibt sich (!!!) für jedes  $x$  aus dem Satz von Bayes gemäß Proposition 2.86 für (eine Version der bzw.) die Dichte bzw. Wahrscheinlichkeitsfunktion  $\pi(\vartheta|x)$  der bedingten Verteilung von  $U$  gegeben  $X = x$ :

$$\pi(\vartheta|x) = c(x) \cdot f_{\vartheta}(x) \cdot \pi(\vartheta) \quad (2.46)$$

mit

$$\frac{1}{c(x)} = f_X(x) = \int f_{X|U}(x|\vartheta)\pi(\vartheta)d\vartheta \quad (2.47)$$



im Falle von stetigem  $X$  und  $U$ , und bei diskretem  $X$  und  $U$

$$\begin{aligned} \frac{1}{c(x)} = p(\{X = x\}) &= \sum_{j=1}^m p(\{X = x\}|\{U = \vartheta_j\}) \cdot \pi(\{U = \vartheta_j\}) \quad (2.48) \\ &= \sum_{j=1}^m p(\{X = x\}|\{U = \vartheta_j\}) \cdot \pi(\vartheta_j). \end{aligned}$$

Für jede Beobachtung  $\in \mathcal{X}$  wird die bedingte Verteilung von  $U$  gegeben  $X = x$  als *Posteriori-Verteilung des Parameters  $\vartheta$  nach der Beobachtung  $x$*  bezeichnet. Die zugehörige Dichte bzw. Wahrscheinlichkeitsfunktion  $\pi(\vartheta|x)$  heißt *Posteriori-Dichte nach der Beobachtung*, und  $f_{\vartheta}(x)$  heißt *Likelihood*.

Die marginale Verteilung von  $X$  mit Dichte  $f_X(x)$  aus (2.47) bzw. Wahrscheinlichkeitsfunktion  $(p(\{X = x\}))_{x \in \mathcal{X}}$  aus (2.48) heißt *Priori-Prädiktive-Verteilung*.

Die Größen  $f_X(x)$  und  $p(\{X = x\})$  sind nicht zu verwechseln mit den als bedingte Verteilungen interpretierten  $f_{\vartheta}(x)$  und  $p_{\vartheta}(\{X = x\})$ .

Analog gibt es auch eine *posteriori-prädiktive Verteilung*, wenn man in analoger Weise über die Posteriori-Verteilung herausintegriert bzw. -summiert. Dies ist dann die Wahrscheinlichkeitsverteilung der nächsten Beobachtung, basierend auf dem aktuellen Wissensstand.

**Bem. 2.90 (Bayes Postulat (nicht entscheidungstheoretisch))**

Nach der Beobachtung der Stichprobe enthält die (klassische) Posteriori-Verteilung die volle Information, d.h. sie beschreibt das Wissen über den unbekannt Parameter vollständig.

Alle statistischen Analysen haben sich ausschließlich auf die Posteriori zu stützen; darauf baut insbesondere auch die Konstruktion von

- Bayesschen-Punktschätzungen: *MPD-Schätzer (Maximum Posteriori Density-Schätzer)*
- Bayesschen-Intervallschätzungen: *HPD-Intervalle (Highest posterior density-Intervalle)*

- Bayes-Tests
  - Suffiziente Statistik: Enthält volle Information der Stichprobe über den Parameter (allgemeine statistische Theorie). Jetzt zwei verschiedene Arten von „enthält volle Information“; Wie passt Suffizienz in diesen Zusammenhang?
  - Posteriori-Verteilung enthält volle Information (vgl. Bem 2.90).
- ⇒ Posteriori hängt tatsächlich nur von suffizienter Statistik ab.

## Proposition 2.91 (Suffizienz und Posteriori-Verteilung)

Ist in der Situation von Bemerkung 2.89  $T$  eine für  $\vartheta$  suffiziente Statistik mit Wahrscheinlichkeitsfunktion bzw. Dichte  $g_{\vartheta}(\cdot)$ , so hängt die Posteriori  $\pi(\vartheta|x)$  nur mehr über  $t = T(x)$  von  $x$  ab. Es gilt <sup>6</sup>

$$\pi(\vartheta|x) \propto g_{\vartheta}(t) \cdot \pi(\vartheta)$$

Beweis:

Gemäß (2.46) ist

$$\pi(\vartheta|x) \propto f_{\vartheta}(x) \cdot \pi(\vartheta)$$

wobei wegen der Suffizienz von  $T$  sich  $f_{\vartheta}(x)$  schreiben lässt als  $f_{\vartheta}(x) = h_{X|T}(x) \cdot g_{\vartheta}(t)$ . Einsetzen liefert die Behauptung.

---

<sup>6</sup> $\propto$  „proportional zu“, vgl. Bemerkung 2.88

## Def. 2.92 (Vorbereitende Erinnerung: Exponentialfamilien)

Sei  $(\mathcal{X}, \sigma(\mathcal{X}), (p_{\vartheta})_{\vartheta \in \Theta})$  ein statistisches Modell mit  $\Theta \subseteq \mathbb{R}^q$ .

- $(p_{\vartheta})_{\vartheta \in \Theta}$  bildet eine (oder  $p_{\vartheta}$  ist für jedes  $\vartheta \in \Theta$  ein Mitglied der)  $q$ -parametrische(n) *Exponentialfamilie* in  $(T_1, \dots, T_q)$  mit *natürlichem Parameter*  $(c_1(\vartheta), \dots, c_q(\vartheta))$ , wenn sich die Dichte  $f_{\vartheta}(x)$  bezüglich eines dominierenden  $\sigma$ -finiten Maßes (also insbesondere Dichte/Wahrscheinlichkeitsfunktion) in die folgende Form bringen läßt: Mit  $t_1 := T_1(\vec{x}), \dots, t_q := T_q(\vec{x})$  ist

$$f_{\vartheta}(x) = h(x) \cdot g(\vartheta) \cdot \exp\left(\sum_{\ell=1}^q c_{\ell}(\vartheta)t_{\ell}\right).$$

- Enthält  $\Theta$  echt innere Punkte und sind  $1, c_1(\vartheta), c_2(\vartheta), \dots, c_q(\vartheta)$  und  $1, T_1(x), T_2(x), \dots, T_q(x)$  (f.-s.) jeweils linear unabhängig, so spricht man von einer *strikt*  $q$ -parametrischen Exponentialfamilie. (Der „natürliche Parameterraum“ hat wirklich die Dimension  $q$ .)



## 2.4.7 Konjugierte Verteilungen, Bayes-Lernen

## a) Ein Motivationsbeispiel

### Bsp. 2.93 (Beta-Binomialmodell)

#### *Absolutes Standardbeispiel*

- Stichprobenmodell: Bernoulliverteilung (allgemein: Binomialverteilung)  
zu Parameter  $\vartheta$

$$p_{\vartheta}(\{X_i = x_i\}) = \vartheta^{x_i}(1 - \vartheta)^{1-x_i}$$

jetzt im Bayes Kontext als bedingte Verteilung schreiben (wieder mit „Hilfsvariable“  $U$ ):

$$p(\{X_i = x_i\} \mid \{U = \vartheta\}) = \vartheta^{x_i}(1 - \vartheta)^{1-x_i}$$

- gebräuchliche Priori-Verteilung:

Betaverteilung, gilt als sehr flexibel, zwei Parameter  $a > 0$ ,  $b > 0$   
hier als Priori verwendet, Bezeichnung  $\pi(\cdot)$

$$\pi(\vartheta) = \frac{\vartheta^{a-1}(1-\vartheta)^{b-1}}{B(a,b)} \cdot I_{[0;1]}(\vartheta) \quad (2.49)$$

$B(a,b)$  ist eine reine Normierungskonstante.

Es gilt:

$$\text{Erwartungswert: } \frac{a}{a+b} \quad \text{Modus: } \frac{a-1}{a+b-2}, \quad a > 1, b > 1$$

## Abbildung 1: Ruger, (1999) Test- und Schatztheorie I, Seite 193

2.4. BAYES-INFERENZ

193

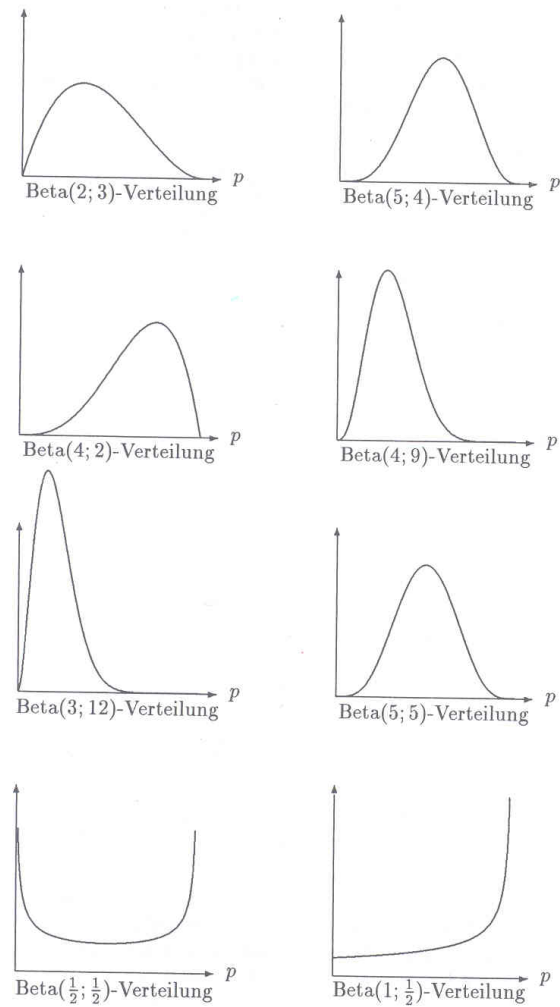


Abbildung 2.17: Einige Beta( $a; b$ )-Verteilungen.  
 Die Beta(1; 1)-Verteilung ist die Gleichverteilung. Die an der Senkrechten im Punkt 0.5 gespiegelte Dichte einer Beta( $a; b$ )-Verteilung ist die Beta( $b; a$ )-Verteilung.

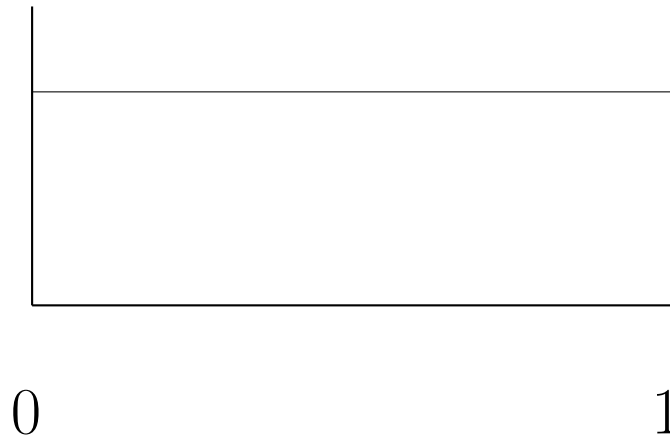
Jetzt Satz von Bayes anwenden: Posteriori nach einer Beobachtung berechnen.

$$\begin{aligned}
 \pi(\vartheta|x_i) &= \frac{\vartheta^{x_i}(1-\vartheta)^{1-x_i} \cdot \vartheta^{a-1}(1-\vartheta)^{b-1}}{\underbrace{\text{Norm.} \cdot B(a,b)}_{\text{Normierung}}} \cdot I_{[0;1]}(\vartheta) \\
 &\propto \vartheta^{x_i+a-1} \cdot (1-\vartheta)^{b-x_i} \cdot I_{[0;1]}(\vartheta) \\
 &= \vartheta^{a'-1} \cdot (1-\vartheta)^{b'-1} \cdot I_{[0;1]}(\vartheta)
 \end{aligned}$$

- Posteriori ist also wieder eine Betaverteilung, nun mit den Parametern

$$a' = a + x_i \quad \text{und} \quad b' = b - x_i + 1 = b + (1 - x_i).$$

Start z.B. mit  $a^{(0)} = 1, b^{(0)} = 1$ :



Gleichverteilung (als Nichtwissen verkaufbar?)

$x_1 = 1$  beobachtet  $\Rightarrow a^{(1)} = a^{(0)} + 1 = 2, b^{(1)} = b^{(0)} + 0 = 1$

*Beta*(2, 1)-Verteilung

$$\pi(\vartheta \mid x_1) \propto \vartheta I_{[0;1]}(\vartheta)$$





- Jetzt weiteres Experiment:

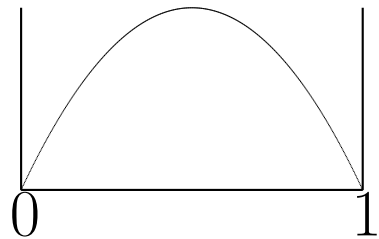
(neue) Priori = (alte Posteriori): Beta(2,1) Verteilung

neue Stichprobe  $x_2$

neue Posteriori:  $Beta(a^{(2)}, b^{(2)}) = Beta(a^{(1)} + x_i, b^{(1)} + (1 - x_i))$

z.B.  $x_2 = 0 \rightarrow a^{(2)} = 2 + 0 = 2, b^{(2)} = 1 + 1 - 0 = 2$

$$\pi_2(\vartheta | (x_1, x_2)') = \frac{\vartheta^{2-1} \cdot (1 - \vartheta)^1}{\text{Norm.}} = \frac{\vartheta^1 \cdot (1 - \vartheta)^1}{\text{Norm.}} = \frac{\vartheta - \vartheta^2}{\text{Norm.}} \quad \text{für } \vartheta \in [0; 1]$$



$$\pi_2(0|x) = \pi_2(1|x) = 0$$

- Weitere Beobachtung  $x_3 = 1$ :

neue Posteriori:  $Beta(a^{(3)}, b^{(3)}) = Beta(a^{(2)} + x_i, b^{(2)} + (1 - x_i))$

$$a^{(3)} = 2 + 1, \quad b^{(3)} = 2$$

$$\pi_3(\vartheta | (x_1, x_2, x_3)') \propto \vartheta^2 (1 - \vartheta)^1 I_{[0;1]}(\vartheta) = \vartheta^2 - \vartheta^3 I_{[0;1]}(\vartheta)$$

- Weiter zusätzliche Beobachtung  $x_4 = 1$ :  
neue Posteriori:  $Beta(a^{(4)}, b^{(4)}) = Beta(a^{(3)} + x_i, b^{(3)} + (1 - x_i))$   
 $a^{(4)} = 3 + 1, b^{(4)} = 2$

$$\pi_4(\vartheta | (x_1, x_2, x_3, x_4)') \propto \vartheta^3 (1 - \vartheta)^1$$

- Allgemein gilt bei  $n$  unabhängigen Wiederholungen:

Die Posteriori  $\pi_n(\vartheta | (x_1, \dots, x_n)')$  ist eine

$$B \left( a^{(0)} + \sum_{i=1}^n x_i; b^{(0)} + n - \sum_{i=1}^n x_i \right) \text{ Verteilung.} \quad (2.50)$$

Man kann zeigen: Dasselbe Ergebnis erhält man, wenn man  $x_1, \dots, x_n$  auf einmal verarbeitet.

- In diesem Beispiel gilt für die posteriori-prädiktive Verteilung

$$p(X_{n+1} = 1 | X_1 = x_1, \dots, X_n = x_n)$$

$$= \mathbb{E}(\pi_n(\vartheta | x_1, \dots, x_n))$$

$$= \frac{a + \sum_{i=1}^n x_i}{a + b + n}.$$

Für die Gleichverteilung (vgl. oben) als Ausgangspriori ergibt sich wegen  $a^{(0)} = b^{(0)} = 1$

$$\begin{aligned} & p(X_{n+1} = 1 | X_n = x_n, \dots, X_1 = x_1) \\ &= \frac{\left( \sum_{i=1}^n x_i \right) + 1}{n + 2}; \end{aligned}$$

## b) Konjugiertheit: Definition und klassische Ergebnisse (vgl. „Schätzen und Testen I“)

### Def. 2.94 (Konjugiertheit)

Eine Verteilungsfamilie  $\Pi$  von Priori-Verteilungen heißt zu einer Menge  $\mathcal{P}$  von Stichprobenverteilungen *konjugiert*, wenn für jede Priori-Verteilung  $\pi(\cdot) \in \Pi$  und jedes  $p(\cdot) \in \mathcal{P}$  die zugehörige Posteriori-Verteilung wieder ein Element von  $\Pi$  ist. Man sagt dann auch, dass jedes Element  $\pi(\cdot) \in \Pi$  zu  $\mathcal{P}$  konjugiert ist.

**Proposition 2.95 (Beispiele für Konjugiertheit: Beta-Binomial/Dirichlet-Multinomial-Modell/Gamma-Poisson-Modell, Selbstkonjugiertheit der Normalverteilung)**



a) Die Menge der Beta-Verteilungen als Priori ist zur Menge der Bernoulli-Verteilungen konjugiert (vgl. Bsp. 2.93).

Allgemeiner gilt:

Ist  $\vec{X} = (X_1, \dots, X_k)$  eine Stichprobe eines zum Parameter  $\vec{\vartheta} = (\vartheta_1, \dots, \vartheta_k)$  multinomial-verteilten Untersuchungsmerkmals, besitzt  $\vec{X}$  also die Wahrscheinlichkeitsfunktion

$$f(x|\vec{\vartheta}) \propto \prod_{j=1}^k \vartheta_j^{x_j}$$

und wählt man die sog. *Dirichlet-Verteilung* zum Parameter  $\vec{\alpha} = (\alpha_1, \dots, \alpha_k)^T$

$$\pi(\vec{\vartheta}) = \prod_{j=1}^k \vartheta_j^{(\alpha_j-1)},$$

so ist die Posteriori-Verteilung eine Dirichlet-Verteilung mit dem Parameter  $\alpha' = (\alpha'_1, \dots, \alpha'_k)^T$ , wobei

$$\alpha'_j = \alpha_j + x_j - 1, \quad j = 1, \dots, k.$$

b) Ist  $\vec{X} = (X_1, \dots, X_n)$  eine i.i.d. Stichprobe eines zum Parameter  $\lambda$  Poisson verteilten Untersuchungsmerkmals, besitzt  $\vec{X}$  also die Wahrscheinlichkeitsfunktion

$$f(x|\lambda) = \frac{\lambda^{\sum_{i=1}^n x_i}}{x_1! x_2! \dots x_n!} e^{-n\lambda},$$

und wählt man als Priori-Verteilung eine Gamma-Verteilung mit Parametern  $a$  und  $b$ , d.h. eine Verteilung mit der Dichte

$$\pi(\lambda) = \frac{b^a}{\underbrace{\Gamma(a)}_{\text{Norm.konst.}}} \lambda^{a-1} e^{-b\lambda}, \quad (2.51)$$

so ist die Posteriori-Verteilung eine Gamma-Verteilung mit den

Parametern

$$a + \sum_{i=1}^n x_i \quad \text{und} \quad b + n.$$

### Bsp. 2.96 (Normalverteilung)

Ist  $\vec{X} = (X_1, \dots, X_n)$  eine i.i.d. Stichprobe eines mit den Parametern  $\mu$  und  $\sigma^2$  normalverteilten Untersuchungsmerkmals, so gilt:

- (i) Ist  $\sigma^2$  bekannt und wählt man als Priori-Verteilung für  $\mu$  eine Normalverteilung mit den Parametern  $\nu$  und  $\rho^2$ , so ist die a posteriori Verteilung  $\pi(\mu|\vec{x})$  eine Normalverteilung mit den Parametern  $\nu'$  und  $\rho'^2$  mit

$$\nu' = \frac{\bar{x}\rho^2 + \nu\frac{\sigma^2}{n}}{\rho^2 + \frac{\sigma^2}{n}} \quad (2.52)$$

und

$$\rho^{2'} = \frac{\rho^2 \cdot \frac{\sigma^2}{n}}{\rho^2 + \frac{\sigma^2}{n}}. \quad (2.53)$$

- (ii) Ist  $\mu$  bekannt, aber  $\sigma^2$  unbekannt, so erhält man die konjugierte Verteilung, indem man  $\frac{1}{\sigma^2}$  als gammaverteilt annimmt. Man sagt dann,  $\sigma^2$  sei *invers gammaverteilt*.

Wie findet man solche konjugierten Paare?

### Satz 2.97 (Zur Konjugiertheit in Exponentialfamilien)

Hat in der Situation von Def. 2.89 jedes Element der Menge  $\mathcal{P}$  der Stichprobenverteilungen eine Dichte bzw. Wahrscheinlichkeitsfunktion  $f(x|\vartheta)$  der Form

$$f(x|\vartheta) \propto h(\vartheta) \exp(T(x) \cdot b(\vartheta)) \quad (2.54)$$

und jedes Element der Menge  $\Pi$ , aus der die Priori-Verteilung stammt, eine Dichte bzw. Wahrscheinlichkeitsfunktion der Form

$$\pi(\vartheta) \propto [h(\vartheta)]^\alpha \exp(b(\vartheta) \cdot \beta), \quad (2.55)$$

so sind  $\Pi$  und  $\mathcal{P}$  konjugiert. Es gilt dann

$$\pi(\vartheta|x) \propto [h(\vartheta)]^{\alpha+1} \cdot \exp((T(x) + \beta) \cdot b(\vartheta)). \quad (2.56)$$



Beweis:

(2.56) ergibt sich unmittelbar durch Anwenden der Formel für die Posteriori-Verteilung auf (2.54) und (2.55). Dann ist (2.56) mit  $\alpha' := \alpha + 1$  und  $\beta' := \beta + T(x)$  von der Form (2.55), also sind tatsächlich  $\Pi$  und  $\mathcal{P}$  konjugiert.

**Bem. 2.98 (zu Satz 2.97)**

- Der Satz kann also direkt zur Konstruktion geeigneter, konjugierter Priori-Verteilungen verwendet werden, indem man die Stichprobenverteilung in die Form (2.54) bringt und dann eine Priori gemäß (2.55) wählt.
- $b(\vartheta)$  spielte in (2.54) und in (2.55) eine ganz unterschiedliche Rolle:  
In (2.54) ist  $b(\vartheta)$  der natürliche Parameter der Exponentialfamilie, aus der die Likelihood / Stichprobenverteilung stammt.  
In (2.55) hingegen ist  $b(\vartheta)$  die suffiziente Statistik für den natürlichen Parameter  $\beta$  der Exponentialfamilie, aus der die Priori stammt. (Bei der Priori ist ja der Wert von  $\vartheta$  „zufällig“.)
- Ähnliches gilt für  $h(\vartheta)$ .

**Bsp. 2.99 (Beispiele zu Satz 2.97)**

Man bestimme in folgenden Situationen unter Verwendung von Satz 2.97 jeweils eine konjugierte Priori-Verteilung:

- a)  $X_1, \dots, X_n$  ist i.i.d. normalverteilt mit unbekanntem  $\mu$  und bekannter Varianz  $\sigma^2$
  
- b)  $X$  ist binomialverteilt zum unbekanntem Parameter  $p$
  
- c)  $X_1, \dots, X_n$  ist i.i.d. Poisson-verteilt mit unbekanntem Parameter  $\lambda$

## 2.4.8 (Reine) Bayes-Punktschätzung

### Def. 2.100 (MPD-Schätzung)

Gegeben eine Beobachtung  $\vec{x}$  und die Posteriori-Verteilung mit Dichte bzw. Wahrscheinlichkeitsfunktion  $\pi(\vartheta|\vec{x})$  heißt  $\hat{\vartheta}$  mit

$$\pi(\hat{\vartheta}|\vec{x}) = \max_{\vartheta \in \Theta} \pi(\vartheta|\vec{x})$$

*(reiner) Bayes-Schätzwert* oder *Maximum (bzw Highest) Posteriori Density Schätzwert* (MPD- (bzw. HPD-) Schätzwert) oder *Posteriori-Modus-Schätzwert*. Die zugehörige Schätzfunktion  $\hat{\vartheta}(\vec{X})$  heißt *reine Bayes-Schätzung* oder *MPD- (bzw. HPD-) Schätzung* bzw. *Posteriori-Modus-Schätzung*.

**Bem. 2.101 (Zur MPD-Schätzung)**

- a) Ist die Posteriori-Verteilung unimodal, so ist  $\hat{\vartheta}$  der Modus der Posteriori.
- b) Ist der Zustandsraum  $\Theta$  beschränkt und liegt dem Schätzproblem als Priori-Verteilung eine Gleichverteilung zugrunde, so gilt

$$\pi(\vartheta|\vec{x}) \propto f(\vec{x}|\vartheta) \cdot \pi(\vartheta) = f(\vec{x}|\vartheta) \cdot \text{Konstante}$$

D.h. der MPD-Schätzer ist dasjenige  $\vartheta$ , das  $f(x|\vartheta)$  maximiert, also der Maximum-Likelihood-Schätzwert.

c) Im Falle  $\Theta = \mathbb{R}^+$  oder  $\Theta = \mathbb{R}$  gibt es keine Gleichverteilung auf  $\Theta$ , denn mit

$$f(x) = c \quad \text{ist} \quad \int_0^{\infty} f(x) dx = \int_0^{\infty} c dx = [x]_0^{\infty} = \infty$$

unabhängig von  $c > 0$ .

Man kann aber zeigen, dass viele der zentralen Ergebnisse der Bayes-Theorie erhalten bleiben, wenn man auch nicht normierbare  $\sigma$ -finite Maße als Prioris zulässt (z.B. Lebesgue Maß  $\lambda(\cdot)$ ;  $\lambda([a, b]) := b - a$ : „*improper priors*“ )

### Bsp. 2.102 (Beta-Binomialmodell)

$\pi(\vartheta|x_1, \dots, x_n)$  ist  $B(a + \sum_{i=1}^n x_i; b + n - \sum_{i=1}^n x_i) =: B(a', b')$ -verteilt.

Hat man bei der Priori  $a=1=b$  gewählt, so ergibt sich mit  $\frac{a' - 1}{a' + b' - 2}$  als

Modus der  $Beta(a', b')$ -Verteilung der MPD-Schätzwert

$$\hat{\vartheta} = \frac{1 + \sum_{i=1}^n x_i - 1}{1 + \sum_{i=1}^n x_i + 1 + n - \sum_{i=1}^n x_i - 2} = \frac{1}{n} \sum_{i=1}^n x_i,$$

also in der Tat der ML-Schätzwert.





## 2.4.9 Der Hauptsatz der Bayes-Entscheidungstheorie

### Def. 2.103 (Posteriori-Verlust-optimale Aktionen, konditionale Bayes-Aktionen)

Gegeben sei ein datenbasiertes Entscheidungsproblem  $((\mathbb{A}, \Theta, l(\cdot)); (\mathcal{X}, A, (p_\vartheta)_{\vartheta \in \Theta}))$  und eine Priori-Verteilung  $\pi(\cdot)$  über  $(\Theta, \sigma(\Theta))$ .

Eine Aktion  $a_x^* \in \mathbb{A}$  heißt *Posteriori-Verlust optimal* zur *Beobachtung*  $x \in \mathcal{X}$  oder *konditionale Bayes-Aktion zu  $x$  und der Priori-Verteilung  $\pi(\cdot)$* , wenn gilt

$$\mathbb{E}_{\pi(\cdot|x)} l(a_x^*, \vartheta) \leq \mathbb{E}_{\pi(\cdot|x)} l(a, \vartheta) \quad \forall a \in \mathbb{A}.$$

$a_x^*$  ist also sozusagen Bayes-Aktion zur Posteriori-Verteilung  $\pi(\cdot|x)$  als „aufdatierter Priori-Verteilung“  $\pi(\cdot|x)$ .

Analog definiert man eine *Posteriori-Nutzen-Optimalität*.

2 Arten, Bayes-Entscheidungstheorie zu betreiben

datengestütztes Entscheidungsproblem  
+  
Informationsbeschaffungsexperiment  
+  
Priori-Verteilung;

Auswertungsproblem + Priori-Verteilung  
komplexer Aktionsraum: alle  
Entscheidungsfunktionen

Bayes-optimale  
Entscheidungsfunktion  $d^* : \mathcal{X} \rightarrow \mathbb{A}$   
( $\rightarrow$  Testfunktion, Schätzfunktion)

konkrete Beobachtung  $x$

Bayes optimale Aktion  
 $a^* = d^*(x)$

Priori-Verteilung

Stichprobenverteilung;  
Informationsbeschaffungsexperiment

konkrete Beobachtung

Posteriori-Verteilung

Bayes Postulat

Posteriori-Verlust optimale Aktion  $a_x^*$ ,  
z.B. reine (optimale) Bayes Schätzung  
 $\hat{\vartheta}_x$   
reiner/optimaler Bayes Test  $\varphi_x$

?

„Priori-Risiko“ optimale Aktion

## Satz 2.104 (Hauptsatz der Bayes-Entscheidungstheorie)

Gegeben sei ein datengestütztes Entscheidungsproblem  $((\mathbb{A}, \Theta, \ell(\cdot)); (\mathcal{X}, \mathcal{A}, (p_\vartheta)_{\vartheta \in \Theta}))$ , bestehend aus einem *datenfreien Entscheidungsproblem*  $(\mathbb{A}, \Theta, \ell(\cdot))$  und einer Informationsstruktur  $(\mathcal{X}, \mathcal{A}, (p_\vartheta)_{\vartheta \in \Theta})$  sowie eine Priori-Verteilung  $\pi(\cdot)$  über  $(\Theta, \sigma(\Theta))$ .

Eine Entscheidungsfunktion

$$\begin{aligned} d^* : \mathcal{X} &\longrightarrow \mathbb{A} \\ x &\longmapsto d^*(x) \end{aligned}$$

ist genau dann Bayes-optimal im zugehörigen Auswertungsproblem, wenn für jedes  $x \in \mathcal{X}$  die zugehörige Aktion  $d^*(x)$  Posteriori-Verlust optimal zur Beobachtung  $x$  ist.

Beweis: Für den diskreten Fall <sup>7</sup>

- Vorneweg eine Hilfsüberlegung: Suche die Lage des Minimums  $\vec{z}_{min}$  einer Funktion  $f(\vec{z})$  mit  $\vec{z} = (z_1, \dots, z_n)$ , wobei  $f(\vec{z}) = \sum_{i=1}^n c_i f_i(z_i)$ , also die  $i$ -te Komponente von  $z$  nur im  $i$ -ten Summanden auftritt.

$$f(\vec{z}) = \sum_{i=1}^n c_i f_i(z_i) \rightarrow \min_{\vec{z}}$$

$\iff f_i(z_i) \longrightarrow \min_{z_i}$  für jedes  $i$  unabhängig von den anderen Summanden.

- Der Deutlichkeit halber wird wieder eine Hilfsvariable  $U$  (vgl. Bem. 2.89 ) eingeführt und  $p_{\vartheta}(\{X = x\})$  wird als  $p(\{X = x\}|\{U = \vartheta\})$  geschrieben.

---

<sup>7</sup>für den allgemeinen Fall: siehe z.B. Rürger (1999, S. 283f.)

Angewendet auf Entscheidungsprobleme mit der Posteriori  $\pi(\vartheta|x)$  ergibt sich mit dieser Notation

$$\pi(\vartheta|x) = \frac{p(\{X = x\}|\{U = \vartheta\}) \cdot \pi(\vartheta)}{p(\{X = x\})}. \quad (2.57)$$

Nun betrachte man Entscheidungsfunktionen  $d(\cdot)$  im Auswertungsproblem:  
Für die Risikofunktion

$$R(d, \vartheta) = \mathbb{E}_{p_\vartheta}(\ell(d(x), \vartheta))$$

gilt hier

$$R(d, \vartheta) = \sum_{x \in \mathcal{X}} \ell(d(x), \vartheta) \cdot p_\vartheta(\{X = x\}).$$

Die optimale Entscheidungsfunktion zur Priori  $\pi(\cdot)$  minimiert unter allen  $d$

$$\mathbb{E}_\pi(R(d, \vartheta)),$$

löst also

$$\sum_{\vartheta \in \Theta} \left( \sum_{x \in \mathcal{X}} \ell(d(x), \vartheta) \cdot p_\vartheta(\{X = x\}) \right) \cdot \pi(\vartheta) \rightarrow \min_d$$

$$\begin{aligned}
&\iff \sum_{\vartheta \in \Theta} \sum_{x \in \mathcal{X}} \left( \ell(d(x), \vartheta) \cdot \underbrace{p(\{X = x\} | U = \vartheta) \cdot \pi(\vartheta)}_{= \pi(\vartheta|x) \cdot p(\{X=x\})} \right) \rightarrow \min_d \\
&\stackrel{(2.57)}{\iff} \underbrace{\sum_{x \in \mathcal{X}}}_{\hat{=} \sum_{i=1}^n} \left( \underbrace{\sum_{\vartheta \in \Theta} \ell(d(x), \vartheta) \cdot \pi(\vartheta|x)}_{\hat{=} f_i(z_i)} \right) \cdot \underbrace{p(\{X = x\})}_{\hat{=} c_i; \text{ priori-prädiktiv, marginal}} \rightarrow \min_d
\end{aligned}$$

- Wegen der Hilfsüberlegung ist dies äquivalent dazu, für jedes feste  $x$

$$\sum_{\vartheta \in \Theta} \ell(d(x), \vartheta) \cdot \pi(\vartheta|x)$$

separat zu minimieren nach  $a_x := d(x)$  für festes  $x$ .

Dies liefert jeweils genau die Posteriori-Verlust optimale Aktion, also die Bayes-Aktion zur Posteriori als neuer Priori.



**Satz 2.105 (Bestimmung von Bayes-optimalen  
Entscheidungsfunktionen, z.B. Rüger (1999, Satz 2.20))**

Gegeben sei das Schätzproblem als datengestütztes Entscheidungsproblem gemäß Kapitel 1.5 sowie eine Priori-Verteilung  $\pi(\cdot)$ .

Dann gilt:

- i) Wählt man die quadratische bzw. absolute Verlustfunktion, so gilt für die Bayes-optimale Entscheidungsfunktion  $d_{quad}^*(\cdot)$  bzw.  $d_{abs}^*(\cdot)$ :  
Für jedes  $x$  ist  $d_{quad}^*(\cdot)$  genau der Erwartungswert und  $d_{abs}^*(\cdot)$  der Median der Posteriori-Verteilung  $\pi(\mathcal{V}|x)$ .

ii) Die HPD-Schätzung ergibt sich näherungsweise für kleine  $\epsilon$ , wenn man die sogenannte Toleranzverlustfunktion zum Grade  $\epsilon$  verwendet:

$$l_{\epsilon}(\hat{\vartheta}, \vartheta) = \begin{cases} 1 & |\hat{\vartheta} - \vartheta| > \epsilon \\ 0 & |\hat{\vartheta} - \vartheta| \leq \epsilon \end{cases}$$

## 2.4.10 „Asymptotische Objektivität“ der konditionalen Bayes-Inferenz

Motivationsbeispiel: Betrachte (2.52) und (2.53) für  $n \rightarrow \infty$

$$\lim_{n \rightarrow \infty} \nu' = \frac{\bar{x}\rho^2 + 0}{\rho^2 + 0} = \bar{x} \quad (2.58)$$

$$\lim_{n \rightarrow \infty} \rho^{2'} = \frac{\rho^2 \cdot 0}{\rho^2 + 0} = 0 \quad (2.59)$$

1. Mal wieder kommt  $\bar{X}$  raus.
2. Viel wichtiger: Die Grenzwerte (2.58) und (2.59) hängen nicht von  $\rho$  und  $\nu$ , also von den Parametern der Priori-Verteilung ab. Hat man

eine sehr große Stichprobe, so „verschwindet der Einfluss der Priori-Verteilung“. Dies gilt ganz allgemein und wird oft als eine pragmatische Rechtfertigung dafür gesehen, bei großem Stichprobenumfängen „angenehme“, z.B. konjugierte, Prioris zu verwenden.

## Satz 2.106 („Asymptotische Objektivität von Bayes-Verfahren“, „Konsistenzsatz“)

Sei  $\Theta = \{\vartheta_1, \dots, \vartheta_m\}$  ein endlicher Parameterraum und  $\vec{X} = (X_1, \dots, X_n)$  eine i.i.d. Stichprobe eines beliebig verteilten (reellwertigen) Untersuchungsmerkmals mit Dichten  $f(x_i | \vartheta_{wahr})$ ,  $\vartheta_{wahr} \in \Theta$ .

Sei  $\pi(\vartheta)$  die Wahrscheinlichkeitsfunktion der Priori-Verteilung auf  $\Theta$  mit  $\pi(\vartheta) > 0$  für alle  $\vartheta$ . Dann gilt für die Wahrscheinlichkeitsfunktion der nach  $n$  Beobachtungen gebildeten Posteriori-Verteilung  $\pi_n(\vartheta | x)$

$$\lim_{n \rightarrow \infty} \pi_n(\vartheta | x) = \begin{cases} 1 & \text{falls } \vartheta = \vartheta_{wahr} \\ 0 & \text{falls } \vartheta \neq \vartheta_{wahr} \end{cases}$$

## Bem. 2.107 (Erneute kritische Diskussion des Bayes-Ansatzes)

## 2.5 Einige alternative Regeln (im Kontext der klassischen Entscheidungstheorie)

## 2.5.1 Die Laplace Regel



Foto:

<http://www.mathematik.de/ger/information/landkarte/gebiete/wahrscheinlichkeitstheorie/wahrscheinlichkeitstheorie.html>

[Stand: 25.06.13]



**Def. 2.108 (Laplace-Regel)**

Gegeben sei das datenfreie Entscheidungsproblem  $(\mathbb{A}, \Theta, u(\cdot))$  mit endlichem  $\Theta = \{\vartheta_1, \dots, \vartheta_m\}$ .

Die Kriteriumsfunktionen

$$\begin{aligned} \tilde{\Phi} : \mathbb{A} &\longrightarrow \mathbb{R} \\ a &\longmapsto \sum_{j=1}^m u(a, \vartheta_j) \end{aligned} \tag{2.60}$$

und

$$\begin{aligned}\Phi : \mathbb{IA} &\longrightarrow \mathbb{R} \\ a &\longmapsto \frac{1}{m} \sum_{j=1}^m u(a, \vartheta_j)\end{aligned}\tag{2.61}$$

heißen *Laplace-Regel*.

**Bem. 2.109**

Die beiden Kriteriumsfunctonen (2.60) und (2.61) liefern dieselbe Ordnung auf der Aktionenmenge.

- a) Die Kriteriumsfuncton (2.61) entspricht einer Bayes-Regel mit Priori-Verteilung  $\pi(\cdot) = (\frac{1}{m}, \frac{1}{m}, \dots, \frac{1}{m})$  (mit  $m = |\Theta|$ ), also einer Gleichverteilung auf  $\Theta$ .

Damit geometrisch: Höhenlinien senkrecht auf Equilibrator-Linie.

- b) Viele Eigenschaften der optimalen Aktion können deshalb aus den in Kapitel 2.4 formulierten Sätzen über Bayes-Regeln abgeleitet werden. (Zulässigkeit, Entbehrlichkeit randomisierter Aktionen,...)

c) Rechtfertigung durch „Prinzip vom unzureichenden Grund“ (Laplace):  
Wenn nichts dafür spricht, dass eines der Elementarereignisse wahrscheinlicher ist als die anderen, dann sind sie gleichwahrscheinlich, also

$$\pi(\{\vartheta_1\}) = \pi(\{\vartheta_2\}) = \dots = \pi(\{\vartheta_m\}).$$

Da

$$\pi(\{\vartheta_1\}) + \pi(\{\vartheta_2\}) + \dots + \pi(\{\vartheta_m\}) = 1$$

ist zwangsläufig

$$\pi(\{\vartheta_j\}) = \frac{1}{m}, \quad j = 1, \dots, m.$$

- d) Verallgemeinerung auf unendliches  $\Theta$ :  
Theorie der nichtinformativen Prior-Verteilungen, siehe Bemerkung 2.111.

### Bem. 2.110 (Beispiel und Kritik)

Abwandlung von Beispiel aus Kapitel 1.3.2, Lotterie  
 Urnen bestehend aus einer unbekanntem Anzahl von grünen, blauen und  
 restlichen (rote, schwarze, violette) Kugeln.

Man kann entweder

$a_1$  nicht spielen,

$a_2$  zum Preis von  $c_g = 60\text{€}$  auf grün setzen oder

$a_3$  zum Preis von  $c_b = 90\text{€}$  auf blau setzen.

Es wird eine Kugel zufällig gezogen. Man erhält  $240\text{€}$ , wenn die Kugel, auf  
 die man gesetzt hat, gezogen wird.

	$\{g\}$	$\{b\}$	$\{\text{rest}\}$	
$a_1$	0	0	0	0
$a_2$	180	-60	-60	20 ← optimale Aktion
$a_3$	-90	150	-90	-10

**Bem. 2.111 („Nichtinformative“ Priori-Verteilung und ihr Inform**

In der konditionalen Inferenz gibt es verschiedene Versuche, ähnlich der Laplace-Regel, „nichtinformative“ Priori-Verteilungen zu definieren und diese dann als Standardbewertungen heranzuziehen.

- z.B. die Gleichverteilung, diese ist aber nicht invariant gegenüber Transformationen des Parameters. Man hat dann also „keine Information“ über  $\vartheta$ , aber eine informative Priori z.B. über eine bijektive, nichtlineare Transformation von  $\vartheta$ .
- z.B. Verteilungen, die invariant bezüglich bijektiver Transformationen des Parameters sind (Jeffrey-Regel).

- z.B. Verteilungen, die die Entropie maximieren (Jaynes-Regel)
- Ganz neue Möglichkeiten ergeben sich beim Übergang zu Credalmengen (siehe Kapitel 3).



## 2.5.2 Die Minimax-Regret-Regel von L.J. Savage, auch Niehans-Savage-Regel genannt

fnewpage

### Def. 2.112 (Minimax-Regret-Aktion)

Gegeben sei ein datenfreies Entscheidungsproblem  $(\mathbb{A}, \Theta, u(\cdot))$  in Nutzenform bzw.  $(\mathbb{A}, \Theta, l(\cdot))$  in Verlustform.

Seien  $\mathbb{A}$  und  $\Theta$  endlich,  $\Theta = \{\vartheta_1, \dots, \vartheta_m\}$ ,  $\mathbb{A} = \{a_1, \dots, a_n\}$ .

Das datenfreie Entscheidungsproblem  $(\mathbb{A}, \Theta, r(\cdot))$  in Verlustform mit

$$\begin{aligned} r : \mathbb{A} \times \Theta &\longrightarrow \mathbb{R} \\ (a_i, \vartheta_j) &\longmapsto r(a_i, \vartheta_j) \end{aligned}$$

und

$$r(a_i, \vartheta_j) = \max_{\ell=1, \dots, n} (u(a_\ell, \vartheta_j)) - u(a_i, \vartheta_j) \quad (2.62)$$

$$\text{bzw.} \quad r(a_i, \vartheta_j) = l(a_i, \vartheta_j) - \min_{\ell=1, \dots, n} (l(a_\ell, \vartheta_j)) \quad (2.63)$$

heißt **induziertes Regret Problem** bzw. **induzierte Regret-Tafel**.

Jedes  $a^* \in \mathbb{A}$  mit

$$\max_{j=1, \dots, m} r(a^*, \vartheta_j) \leq \max_{j=1, \dots, m} r(a, \vartheta_j) \quad \text{für alle } a \in \mathbb{A}$$

heißt **Minimax-Regret-Aktion**.

## Bsp. 2.113 (Beispiel und Kritik)

**Proposition 2.114 (Bayes-Aktionen in Regrettafeln)**

Gegeben sei ein datenfreies Entscheidungsproblem  $(\mathbb{A}, \Theta, u(\cdot))$  bzw.  $(\mathbb{A}, \Theta, l(\cdot))$  mit  $\mathbb{A} < \infty$  und  $|\Theta| < \infty$  und die Priori-Bewertung  $\pi(\cdot)$ . Eine Aktion  $a^*$  ist genau dann Bayes-Aktion zu  $\pi(\cdot)$ , wenn  $a^*$  Bayes-Aktion zu  $\pi(\cdot)$  im induzierten Regret-Problem ist.

### 2.5.3 Das Hurwicz-Kriterium

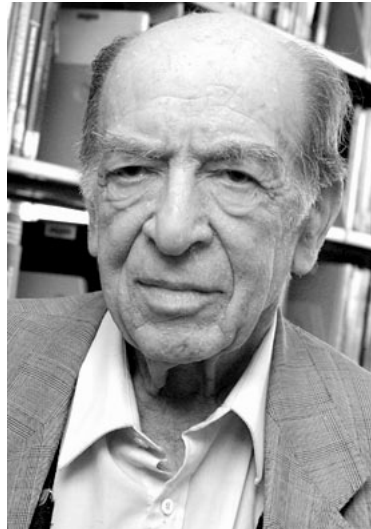


Foto: [http://www.nobelprize.org/nobel\\_prizes/economic-sciences/laureates/2007/hurwicz-facts.html](http://www.nobelprize.org/nobel_prizes/economic-sciences/laureates/2007/hurwicz-facts.html)  
[Stand: 25.06.13]

### Def. 2.115 (Hurwicz-Kriterium)

Gegeben sei ein datenfreies Entscheidungsproblem  $(\mathbb{A}, \Theta, u(\cdot))$  mit  $|\Theta| < \infty$ , sowie  $\alpha \in [0, 1]$ .

$$\begin{aligned} \Phi : \mathbb{A} &\longrightarrow \mathbb{R} \\ a &\longmapsto \Phi(a) \end{aligned}$$

mit

$$\Phi(a) = \alpha \max_j u(a, \vartheta_j) + (1 - \alpha) \min_j u(a, \vartheta_j) \quad (2.64)$$

heißt *Hurwicz-Kriterium* zum *Optimismusparameter*  $\alpha$ .

### 2.5.4 Das Erfahrungskriterium von J.L. Hodges und E.L. Lehmann (1952)

**Def. 2.116 (Erfahrungskriterium von Hodges & Lehmann)**

Gegeben sei ein datenfreies Entscheidungsproblem  $(\mathbb{A}, \Theta, u(\cdot))$ ,  $|\Theta| < \infty$ ,  $\mu \in [0, 1]$  und eine Priori-Bewertung  $\pi(\cdot)$ .

$$\begin{aligned} \Phi : \mathbb{A} &\longrightarrow \mathbb{R} \\ a &\longmapsto \Phi(a) \end{aligned}$$

mit

$$\begin{aligned} \Phi(a) = & \mu \cdot \left( \sum_{j=1}^m u(a, \vartheta_j) \pi(\{\vartheta_j\}) \right) \\ & + (1 - \mu) \cdot \left( \min_j (u(a, \vartheta_j)) \right) \end{aligned} \quad (2.65)$$

heißt *Erfahrungskriterium von Hodges und Lehmann* zum Vertrauensparameter  $\mu$ .



## 2.6 Gleichmäßig beste Verfahren in der statistischen Entscheidungstheorie

### Def. 2.117 (Gleichmäßig beste Verfahren)

Gegeben sei das auf eine Menge  $\mathcal{D}_0$  eingeschränkte Auswertungsproblem  $(\mathcal{D}_0, \Theta, R(\cdot))$  eines datengestütztes Entscheidungsproblem  $((\mathbb{A}, \Theta, l(\cdot)); (\mathcal{X}, \sigma(\mathcal{X}), (p_\vartheta)_{\vartheta \in \Theta}))$

Eine Entscheidungsfunktion  $d^* \in \mathcal{D}_0$  heißt *gleichmäßig bestes Verfahren* aus  $\mathcal{D}_0$ , wenn für die Risikofunktion gilt:

$$R(d^*, \vartheta) \leq R(d, \vartheta) \quad \text{für alle } \vartheta \in \Theta \text{ und alle } d \in \mathcal{D}_0. \quad (2.66)$$

**Bem. 2.118 (zu Def. 2.117)**

- i) Man beachte, dass die Risikofunktion von  $d^*$  für *alle*  $\vartheta \in \Theta$  nicht größer sein soll;  $d^*$  soll also im zugehörigen Auswertungsproblem  $(\mathcal{D}_0, \Theta, R(\cdot))$  alle Elemente von  $\mathcal{D}_0$  dominieren.
- ii) Dies ist für „großes  $\mathcal{D}_0$ “ eine extrem starke Forderung – insbesondere im Lichte der elementaren Beispiele aus Kapitel 1.3 – aber v.a. bei Exponentialfamilie und geeigneter Einschränkung der Menge der betrachteten Entscheidungsfunktionen möglich (UMVU-Schätzer, gleichmäßig bester Test, siehe später).

**Satz 2.119 (Lehmann und Scheffé (1950))**

Gegeben sei ein Schätzproblem im Sinne Beispiel 1.42. Ferner sei

- $(p_{\vartheta})_{\vartheta \in \Theta}$  eine strikt  $q$ -parametrische Exponentialfamilie in  $T(\vec{x}) = (T_1(\vec{x}), \dots, T_q(\vec{x}))$ .
- $\ell(\hat{\vartheta}, \vartheta)$  konvex in  $\hat{\vartheta}$  für alle  $\vartheta$ .

Betrachtet man die Schätzung einer Transformation  $\gamma(\vartheta)$  von  $\vartheta$  und die Klasse  $\mathcal{D}_\gamma$  aller für  $\gamma(\vartheta)$  erwartungstreuen Schätzer, so gilt:

Ist  $\mathcal{D}_\gamma$  nicht leer, so gibt es eine nichtrandomisierte Schätzfunktion der Form  $\eta(T(\vec{X}))$ , die gleichmäßig beste in der Klasse  $\mathcal{D}_\gamma$  ist.

Ist umgekehrt  $\eta(T(\vec{X}))$  eine erwartungstreue Schätzfunktion für  $\gamma(\vartheta)$ , so ist  $\eta(T(\vec{X}))$  gleichmäßig bestes Verfahren in  $\mathcal{D}_\gamma$ .

## Bem. 2.120 (Zur Interpretation des Satzes)

- „Informationsdeutung“:
  
- Konstruktive Anwendung:

**Korollar 2.121**

Gegeben sei eine strikt  $q$ -parametrische Exponentialfamilie in  $T(\vec{x}) = (T_1(\vec{x}), \dots, T_q(\vec{x}))$ .

Dann ist für jede (messbare) Funktion  $\eta(\cdot)$  der Schätzer  $\eta(T_1, \dots, T_q)$  UMVU-Schätzer für  $\mathbb{E}_\vartheta(\eta(T_1, \dots, T_q))$ .

**Bem. 2.122 (Zum Satz von Lehmann-Scheffé)**

Die Beschränkung auf erwartungstreue Schätzer ist wesentlich.



**Bem. 2.123**

Korollar 2.121 wird typischerweise „andersherum“ angewendet. Will man eine Funktion  $\gamma(\vartheta)$  schätzen, so sucht man eine Funktion  $g(T)$ , so dass  $\mathbb{E}g(T) = \gamma(\vartheta)$ . Gemäß Korollar 2.121 weiß man dann, dass  $g(T)$  UMVU für  $\gamma(\vartheta)$  ist.

Üblicherweise wird man versuchen einen Ansatz zu wählen, bei dem sich  $g(\cdot)$  aus „einfachen Grundfunktionen“ zusammensetzt.

**Bsp. 2.124**

Gegeben sei eine i.i.d. Stichprobe eines normalverteilten Merkmals mit unbekanntem Mittelwert  $\mu$ , aber bekannter Varianz  $\sigma^2$ . Man bestimme einen UMVU Schätzer

- a) für  $\mu$  und
- b) für  $\exp(\mu)$ .

## Satz 2.126 (Optimale Tests)

Betrachtet werde das Testproblem als Entscheidungsproblem gemäß Beispiel 1.44 mit

$a_0$  für  $H_0$  entscheiden

$a_1$  für  $H_1$  entscheiden

	$\vartheta$	
$a_0$	0	1
$a_1$	1	0

$$l(a_0, \vartheta) = \begin{cases} 0 & \vartheta \in \Theta_0 \\ 1 & \vartheta \in \Theta_1 \end{cases}$$

$$l(a_1, \vartheta) = \begin{cases} 1 & \vartheta \in \Theta_0 \\ 0 & \vartheta \in \Theta_1 \end{cases}$$

Ferner sei  $\Theta_0 = \{\vartheta | \vartheta \leq \vartheta_0\}$ ,  $\Theta_1 = \{\vartheta | \vartheta \geq \vartheta_1\}$ ,  $\vartheta_1 > \vartheta_0$ , und ein Signifikanzniveau  $\alpha \in (0, 1)$  vorgegeben.

Bildet  $(P_{\vartheta}^{\oplus n})$  eine strikt einparametrische Exponentialfamilie in  $T$  (mit dem natürlichen Parameter  $c(\vartheta)$ , der in eindeutiger Beziehung zu  $\vartheta$  stehe), so gibt es ein  $\kappa \in \mathbb{R}$ , so dass der Test

$$\varphi^*(\vec{x}) = \begin{cases} 1 & T(\vec{x}) > \kappa \\ \gamma & T(\vec{x}) = \kappa \\ 0 & T(\vec{x}) < \kappa \end{cases} \quad (2.67)$$

mit

$$\mathbb{E}_{\vartheta_0} \varphi^* = \alpha \quad (2.68)$$

gleichmäßig bester Test (UMP) ist, d.h. es gilt

$$\mathbb{E}_{\vartheta} \varphi^* \geq \mathbb{E}_{\vartheta} \varphi, \quad \forall \vartheta \in \Theta_1,$$

für alle  $\varphi$  mit  $\sup_{\vartheta \in \Theta_0} \mathbb{E}_{\vartheta} \varphi \leq \alpha$ .

**Bem. 2.127 (Zur Interpretation von Satz 2.126)**

- Informationsdeutung:
  
- konstruktive Anwendung:

**Bem. 2.128**

- Die Aussage gilt allgemeiner für Verteilungen mit *monotonen Dichtequotienten*; bei denen also für alle  $\vartheta_1 \in \Theta_1$ ,  $\vartheta_0 \in \Theta_0$  und eine geeignete Funktion  $T(\vec{X})$  der Quotient  $\frac{f_{T||\vartheta_1}(t)}{f_{T||\vartheta_0}(t)}$  monoton in  $t$  ist.
- Auch bei der Fragestellung  $H_0 : \vartheta = \vartheta_0$  gegen  $H_1 : \vartheta \neq \vartheta_0$ , gibt es bei Exponentialfamilien einen gleichmäßig besten Test, wenn man sich auf *unverfälschte* Tests (d.h. Gütefunktion  $\geq \alpha$  für alle Elemente der Alternative) beschränkt.  
Ähnliches gilt in der Situation  $H_0 : \vartheta \in [\underline{\vartheta}_0, \bar{\vartheta}_0]$  gegen  $H_1 : \vartheta \in [\underline{\vartheta}_1, \bar{\vartheta}_1]$ , wenn man nur unverfälschte und ähnliche Tests (Gütefunktion am Rand der Nullhypothese  $= \alpha$ ) betrachtet.

## Bsp. 2.130 Anwendung von Satz 2.126 auf den Mittelwertstest bei der Normalverteilung und auf das Testen bei der Binomialverteilung

a)  $X_1, \dots, X_n$  *i.i.d.*  $\sim N(\mu, \sigma^2)$

$H_0 : \mu \leq \mu_0$  gegen  $H_1 : \mu \geq \mu_1 > \mu_0$   $\sigma^2$  bekannt

$$\begin{aligned}
 f_{X_1, \dots, X_n | \mu}(x_1, \dots, x_n) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right) \\
 &= \underbrace{\left(\frac{1}{\sqrt{2\pi}}\right)^n}_{\text{konstant (bei bek. } \sigma^2)} \cdot \frac{1}{\sigma^n} \cdot \underbrace{\exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n x_i^2\right)}_{\text{nur von } X \text{ abhängig (bei bekanntem } \sigma^2)} \cdot \underbrace{\exp\left(-\frac{1}{2\sigma^2} n\mu^2\right)}_{\text{nur von } \mu \text{ abh. nicht von } x} \cdot \exp\left(\frac{\mu}{\sigma^2} \cdot \sum_{i=1}^n x_i\right)
 \end{aligned}$$

natürlicher Parameter  $\frac{\mu}{\sigma^2}$

Suffiziente Statistik:  $T = \sum_{i=1}^n X_i$

Die kritische Region  $K(\kappa)$ , d.h. der Bereich aller  $\vec{x}$  mit  $\varphi(\vec{x}) = 1$ , ist wegen (3.5) von der Form

$$K(\kappa) := \{\vec{x} | T(\vec{x}) > \kappa\}$$

mit  $\kappa$  so, dass

$$P_{\mu_0}(K(\kappa)) \stackrel{!}{=} \alpha.$$

(Man weiß ja wegen (2.68), dass die maximale Fehlerwahrscheinlichkeit an der Grenze der Nullhypothese angenommen wird.)



Das heißt, es soll gelten

$$P_{\mu_0}(K(\kappa)) = P_{\mu_0}(\{\vec{x} | T(\vec{x}) > \kappa\}) = P_{\mu_0}(\{\vec{x} | \sum_{i=1}^n X_i > \kappa\}) = \alpha.$$

Zur Berechnung dieser Wahrscheinlichkeit benutzt man

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

d.h.

$$\frac{\sum_{i=1}^n X_i - n\mu_0}{\sigma\sqrt{n}} \sim N(0, 1).$$

Also führt der Ansatz  $P_{\mu_0}(K(\kappa)) \stackrel{!}{=} \alpha$ , d.h.

$$P \left( \frac{\sum_{i=1}^n X_i - n\mu_0}{\sigma \cdot \sqrt{n}} > \underbrace{\frac{\kappa - n\mu}{\sigma \cdot \sqrt{n}}}_{\kappa' = \tau_\alpha} \right) \stackrel{!}{=} \alpha,$$

Fraktile der Normalverteilung

dazu, dass die kritische Region so zu wählen ist, dass

$$\begin{aligned} & \frac{\sum_{i=1}^n X_i - n \cdot \mu}{\sigma \cdot \sqrt{n}} > \tau_\alpha \\ \iff & \frac{\bar{X} - \mu}{\sigma} \cdot \sqrt{n} > \tau_\alpha \\ \iff & \bar{X} > \mu_0 + \frac{\tau_\alpha \cdot \sigma}{\sqrt{n}}. \end{aligned}$$

Da ferner  $\mu_1 > \mu_0$  bei dieser Konstruktion beliebig gewählt werden konnte, gilt die Aussage für *alle*  $\mu > \mu_0$ .

Beachte: es ergibt sich die übliche kritische Region des Gauss Tests, dieser ist also UMP.

b) Bei der Bernoulliverteilung sei zu einem konkreten Beispiel übergegangen.

Man testet  $H_0 : p \leq 0.5$  gegen  $H_1 : p \geq 0.6$ , wobei bei der i.i.d. Stichprobe  $X_1, \dots, X_n$  der Stichprobenumfang  $n = 5$  sei und das Signifikanzniveau auf  $\alpha = 0.1$  gesetzt sei.

Zur Anwendung von Satz 2.126 bringt man zunächst die Wahrscheinlichkeitsfunktion  $f_p(\cdot)$  der Bernoulliverteilung auf

„Exponentialfamilien-Gestalt“.

$$\begin{aligned}
 f_p(x_1, \dots, x_n) &= \prod_{i=1}^n p^{x_i} \cdot (1-p)^{1-x_i} && = \\
 &= (1-p)^n \cdot \left(\frac{p}{1-p}\right)^{\sum_{i=1}^n x_i} && = \\
 &= (1-p)^n \cdot \exp\left(\ln\left(\left(\frac{p}{1-p}\right)^{\sum_{i=1}^n x_i}\right)\right) && = \\
 &= (1-p)^n \cdot \exp\left(\sum_{i=1}^n x_i \cdot \ln\left(\frac{p}{1-p}\right)\right)
 \end{aligned}$$

Man erhält  $T = \sum_{i=1}^n X_i$ , und  $\ln\left(\frac{p}{1-p}\right)$  ist der natürliche Parameter.

Der Ansatz

$$K(\kappa) := \{\vec{x} | T(\vec{x}) > \kappa\}$$

für die kritische Region führt auf ein Problem. Da  $T = \sum_{i=1}^n X_i$  binomialverteilt ist, ergibt sich <sup>8</sup>

$$P_{p_0}(K(\kappa)) = P_{p_0}(\{T > \kappa\}) = \sum_{\substack{j > \kappa \\ j \in \mathbb{N}}}^n \binom{5}{j} \underbrace{0.5^j \cdot 0.5^{5-j}}_{0.5^5}$$

---

<sup>8</sup>Wegen (2.68) setzt man hier wieder den oberen Randwert der Nullhypothese ein, als denjenigen Wert der Nullhypothese, der am schwersten von  $H_1$  zu unterscheiden ist.

Allerdings gibt es kein  $\kappa$ , so dass diese Gleichung erfüllt ist:

$$\kappa \in (4, 5] \quad P(K(\kappa)) = 0$$

$$\kappa = 4 \Rightarrow P(K(\kappa)) = \binom{5}{5} \cdot \left(\frac{1}{2}\right)^5 = \frac{1}{32} < 0.1$$

$$\kappa \in (3, 4] \Rightarrow P(K(\kappa)) = P(K(4)) = \frac{1}{32}$$

$$\kappa = 3 \Rightarrow P(K(\kappa)) = \frac{1}{32} + \binom{5}{4} \cdot 0.5^5 = \frac{1}{32} + \frac{5}{32} = \frac{6}{32} > 0.1$$

Was tun? Klar ist

$\varphi(\vec{x}) = 1$  , d.h.  $H_0$  ablehnen, für alle  $\vec{x}$  mit  $T(\vec{x}) = 5$  und

$\varphi(\vec{x}) = 0$  , d.h.  $H_0$  nicht ablehnen, für alle  $\vec{x}$  mit  $T(\vec{x}) \leq 3$ .

Ist  $T(\vec{x}) = 4$ , dann hat man so zu randomisieren, dass  $\mathbb{E}_{0.5}\varphi = \alpha$ . Man setzt hierzu  $\varphi(\vec{x}) = \gamma$ , falls  $T(\vec{x}) = 4$ , und erhält:

$$\begin{aligned}\mathbb{E}_{0.5}\varphi &= 1 \cdot P(\varphi(x) = 1) + \gamma \cdot P(\varphi(x) = \gamma) + \\ &\quad + 0 \cdot P(\varphi(x) = 0) \stackrel{!}{=} \alpha\end{aligned}$$

also

$$\begin{aligned}\gamma &= \frac{\alpha - P(\varphi(x) = 1)}{P(\varphi(x) = \gamma)} = \\ &= \frac{\alpha - P(T = 5)}{P(T = 4)} = \frac{0.1 - \frac{1}{32}}{\frac{5}{32}} = 0.44 .\end{aligned}$$