

Methoden der Statistik

Mathias Trabs

1. September 2015

Inhaltsverzeichnis

1	Grundbegriffe der Statistik	2
1.1	Drei grundlegende Fragestellungen	3
1.1.1	Schätzprobleme	3
1.1.2	Hypothesentests	6
1.1.3	Konfidenzmengen (Bereichsschätzung)	9
1.2	Minimax- und Bayesansatz	10
1.3	Ergänzungen: Quantile	14
2	Lineares Modell	15
2.1	Regression und kleinste Quadrate	15
2.2	Inferenz unter Normalverteilungsannahme	21
2.3	Varianzanalyse	26
3	Exponentialfamilien and verallgemeinerte lineare Modelle	30
3.1	Die Informationsungleichung	30
3.2	Verallgemeinerte Lineare Modelle	34
3.3	Ergänzung: Numerische Bestimmung des Maximum-Likelihood-Schätzers	37
4	Klassifikation	38
4.1	Logistische Regression	38
4.2	Bayesklassifikation	40
4.3	Lineare Diskriminanzanalyse	42
5	Ausblick	43

Literatur

- Georgii, H.-O.: *Stochastik. Einführung in die Wahrscheinlichkeitstheorie und Statistik*, de Gruyter, 2007
- James, G., Witten, D., Hastie, T., Tibshirani, R.: *An Introduction to Statistical Learning (with Applications in R)*, Springer, 2013
- Lehmann, E.L., Romano, J.P.: *Testing Statistical Hypotheses*, Springer, 2005
- Lehmann, E.L. and G. Casella: *Theory of Point Estimation*, Springer, 2003
- Shao, J.: *Mathematical Statistics*, Springer, 2003
- Venables, W.N., Ripley, B.D.: *Modern Applied Statistics with S-Plus*, Springer, 1997
- Wasserman, L.: *All of Statistics*, Springer, 2003
- Witting, H.: *Mathematische Statistik I*, Teubner, 1985

1 Grundbegriffe der Statistik

Während die Wahrscheinlichkeitstheorie anhand eines gegebenen Modells die Eigenschaften der (zufälligen) Ereignisse untersucht, ist das Ziel der Statistik genau andersherum: Wie kann man aus den gegebenen Beobachtungen Rückschlüsse auf das Modell ziehen?

Beispiel 1.1 (Werbung). Wir verwenden den “Advertising”-Datensatz aus James et al. (2013). Für 200 Märkte haben wir die Anzahl der verkauften Produkte Y sowie das jeweilige Budget für Fernsehwerbung X^F , für Radiowerbung X^R und für Zeitungsannoncen X^Z gegeben.

Betrachten wir das *Modell*

$$Y_i = aX_i^F + b + \varepsilon_i, \quad i = 1, \dots, 200,$$

wobei die zufälligen Störgrößen ε_i Marktunsicherheiten, externe Einflüsse etc. modellieren. Plausible Annahmen an das Modell sind

- (i) (ε_i) sind unabhängig (näherungsweise),
- (ii) (ε_i) sind identisch verteilt,
- (iii) $\mathbb{E}[\varepsilon_i] = 0$ (kein systematischer Fehler)
- (iv) ε_i normalverteilt (wegen ZGWS).

Naheliegende *Ziele/Fragestellungen*:

- (i) Es sollen a, b anhand der Daten ermittelt werden. Ein mögliches Schätzverfahren ist der *Kleinste-Quadrate-Schätzer*

$$(\hat{a}, \hat{b}) := \arg \min_{a, b} \sum_{i=1}^n (Y_i - aX_i - b)^2$$

(wir minimieren die Summe der quadrierten Residuen). Mit \hat{a}, \hat{b} erhalten wir die *Regressionsgrade*

$$y = \hat{a}x^F + \hat{b}.$$

- (ii) Sind die Modellannahmen erfüllt? Histogramm, Boxplot und QQ-Plot (Quantil-Quantil-Plot) der Residuen.
- (iii) Wenn wir die Verteilung von \hat{a} kennen (Verteilungsannahme an ε nötig!), können wir Intervalle der Form $I = [\hat{a} - c, \hat{a} + c]$ für $c > 0$ konstruieren, so dass der tatsächlich Parameter a mit vorgegebener Wahrscheinlichkeit in I liegt.
- (iv) Wir wollen *testen*, ob es einen Effekt gibt, d.h. gilt die Hypothese $H_0 : a = 0$ oder kann sie verworfen werden? Beispielsweise kann man die Hypothese verwerfen, falls $|\hat{a}| > c$ für einen kritischen Wert $c > 0$. Um einen sinnvollen Wert zu bestimmen, benötigen wir wieder Verteilungsannahmen an die Fehler (ε_i) .

Wir können das Modell auf polynomielle Regression $Y_i = a_0 + a_1X_i^F + \dots + a_n(X_i^F)^n + \varepsilon_i$ oder multiple Regression $Y_i = a_0 + a_1X_i^F + a_2X_i^R + a_3X_i^Z + \varepsilon_i$ erweitern. Dies führt auf das Problem der Modellwahl.

Definition 1.2. Ein messbarer Raum $(\mathcal{X}, \mathcal{F})$ versehen mit einer Familie $(\mathbb{P}_\vartheta)_{\vartheta \in \Theta}$ von Wahrscheinlichkeitsmaßen mit einer beliebigen Parametermenge $\Theta \neq \emptyset$ heißt statistisches Experiment oder statistisches Modell. \mathcal{X} heißt Stichprobenraum. Jede $(\mathcal{F}, \mathcal{S})$ -messbare Funktion $Y : \mathcal{X} \rightarrow \mathcal{S}$ heißt Beobachtung oder Statistik mit Werten in $(\mathcal{S}, \mathcal{S})$ und induziert das statistische Modell $(\mathcal{S}, \mathcal{S}, (\mathbb{P}_\vartheta^Y)_{\vartheta \in \Theta})$. Sind die Beobachtungen Y_1, \dots, Y_n für jedes \mathbb{P}_ϑ unabhängig und identisch verteilt (iid.), so nennt man Y_1, \dots, Y_n eine mathematische Stichprobe.

Beispiel 1.3 (mathematische Stichprobe). Für $n \in \mathbb{N}$ seien X_1, \dots, X_n iid. verteilte Zufallsvariablen mit Werten in \mathcal{X} und Randverteilung $X_1 \sim \mathbb{P}_\vartheta$ mit Parameter $\vartheta \in \Theta$. Dann ist der Stichprobenvektor (X_1, \dots, X_n) gemäß dem Produktmaß $\mathbb{P}_\vartheta^n(dx) = \prod_{i=1}^n \mathbb{P}_\vartheta(dx_i)$ auf $(\mathcal{X}^n, \mathcal{F}^{\otimes n})$ verteilt.

Wir werden uns in dieser Vorlesung weitgehend mit (verallgemeinerten) linearen Modellen befassen, d.h. die Abhängigkeit der Zufallsvariablen X_i bzw. deren Verteilung vom unbekanntem Parameter kann durch eine lineare Abbildung dargestellt werden.

1.1 Drei grundlegende Fragestellungen

Die meisten statistischen Fragestellungen kann man einer der drei Grundprobleme *Schätzen*, *Testen* und *Konfidenzintervalle* zuordnen. Diese werden im folgenden kurz umrissen und im Laufe der Vorlesung weiter vertieft.

1.1.1 Schätzprobleme

Ziel ist es, aufgrund der vorhandenen Beobachtungen den unbekanntem Parameter im statistischen Modell $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ zu bestimmen, also einen einzelnen (bestmöglichen) Wert anzugeben (*Punktschätzung*). Damit ist ein Schätzer eine Abbildung, die nur von den Beobachtungen abhängt.

Definition 1.4. Sei $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ ein statistisches Modell, $\rho : \Theta \rightarrow \mathbb{R}^d$ ein (abgeleiteter) d -dimensionaler Parameter, $d \in \mathbb{N}$. Ein Schätzer ist eine messbare Abbildung $\hat{\rho} : \mathcal{X} \rightarrow \mathbb{R}^d$. Gilt $\mathbb{E}_\vartheta[\hat{\rho}] = \rho(\vartheta)$ so heißt $\hat{\rho}$ unverzerrt oder erwartungstreu (engl.: unbiased).

Beispiel 1.5. Seien X_1, \dots, X_n eine Bernoulli-verteilte mathematische Stichprobe mit Parameter $p \in (0, 1)$. Betrachte den Schätzer $\hat{p}_n := n^{-1} \sum_{i=1}^n X_i$. Dann gilt $\mathbb{E}_p[\hat{p}_n] = n^{-1} \sum_{i=1}^n \mathbb{E}[X_i] = p$. Also ist \hat{p}_n erwartungstreu. Um die Streuung des Schätzers um den wahren Parameter p zu messen, berechnen wir

$$\text{Var}_p(\hat{p}_n) = n^{-2} \sum_{i=1}^n \text{Var}_p(X_i) = \frac{p(1-p)}{n}.$$

Für größer werdenden Stichprobenumfang konzentriert sich also \hat{p}_n um p .

Wie gut ein Schätzer ist, wird mithilfe einer Verlustfunktion bestimmt. Diese misst den Abstand zwischen geschätztem und wahren Parameter.

Definition 1.6. Eine Funktion $L : \Theta \times \mathbb{R}^d \rightarrow \mathbb{R}_+$ heißt Verlustfunktion, falls $L(\vartheta, \cdot)$ für jedes $\vartheta \in \Theta$ messbar ist. Der erwartete Verlust $R(\vartheta, \hat{\rho}) := \mathbb{E}_\vartheta[L(\vartheta, \hat{\rho})]$ eines Schätzers $\hat{\rho}$ heißt Risiko. Typische Verlustfunktionen sind

- (i) der 0-1-Verlust $L(\vartheta, r) = \mathbb{1}_{\{r \neq \rho(\vartheta)\}}$,
- (ii) der absolute Verlust $L(\vartheta, r) = |r - \rho(\vartheta)|$ (euklidischer Abstand im \mathbb{R}^p) sowie
- (iii) der quadratische Verlust $L(\vartheta, r) = |r - \rho(\vartheta)|^2$.

Lemma 1.7 (Bias-Varianz-Zerlegung). Sei $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ ein statistisches Modell und $\hat{\rho} : \mathcal{X} \rightarrow \mathbb{R}^d$ ein Schätzer des Parameters $\rho(\vartheta)$ mit $\mathbb{E}_\vartheta[|\hat{\rho}|^2] < \infty$ für alle $\vartheta \in \Theta$. Dann gilt für den quadratischen Verlust

$$\mathbb{E}_\vartheta[|\hat{\rho} - \rho(\vartheta)|^2] = \text{Var}_\vartheta(\hat{\rho}) + \underbrace{|\mathbb{E}_\vartheta[\hat{\rho}] - \rho(\vartheta)|^2}_{\text{Bias}} \quad \text{für alle } \vartheta \in \Theta.$$

Beweis. Es gilt

$$\begin{aligned} \mathbb{E}_\vartheta[|\hat{\rho} - \rho(\vartheta)|^2] &= \mathbb{E}_\vartheta[|\hat{\rho} - \mathbb{E}_\vartheta[\hat{\rho}] + \mathbb{E}_\vartheta[\hat{\rho}] - \rho(\vartheta)|^2] \\ &= \mathbb{E}_\vartheta[|\hat{\rho} - \mathbb{E}_\vartheta[\hat{\rho}]|^2] + 2\mathbb{E}_\vartheta[(\hat{\rho} - \mathbb{E}_\vartheta[\hat{\rho}])^\top (\mathbb{E}_\vartheta[\hat{\rho}] - \rho(\vartheta))] + |\mathbb{E}_\vartheta[\hat{\rho}] - \rho(\vartheta)|^2 \\ &= \text{Var}_\vartheta(\hat{\rho}) + |\mathbb{E}_\vartheta[\hat{\rho}] - \rho(\vartheta)|^2. \end{aligned} \quad \square$$

Beispiel. In der Situation von Beispiel 1.5, betrachten wir den Schätzer $\tilde{p}_n := (\sum_{i=1}^n X_i + 1)/(n + 2)$. Dieser hat den Bias

$$\mathbb{E}[\tilde{p}_n] - p = \frac{1 - 2p}{n + 2}$$

und die Varianz

$$\text{Var}(\tilde{p}_n) = \frac{np(1-p)}{(n+2)^2}.$$

Damit hat \tilde{p}_n einen kleineren quadratischen Fehler als \hat{p}_n , wenn $|p - 1/2| \leq 1/\sqrt{8}$.

Bemerkung 1.8. Ein Schätzproblem, bei dem der interessierende Parameter nur endliche viele Werte annehmen kann, heißt auch *Klassifikationsproblem* und der entsprechende Schätzer heißt *Klassifizierer* (mehr dazu in Kapitel 4).

Obwohl wir in dieser Vorlesung keine Asymptotik, d.h. das Verhalten der Schätzer bei Stichprobenumfängen $n \rightarrow \infty$, behandeln, seien noch zwei weitere wichtige Grundbegriffe erwähnt.

Definition 1.9. Sei $X_1, \dots, X_n \stackrel{iid.}{\sim} \mathbb{P}_\vartheta$ eine mathematische Stichprobe. Dann heißt ein Schätzer $\hat{\rho}_n$ vom abgeleiteten Parameter $\rho(\vartheta)$ konsistent, falls

$$\hat{\rho}_n \xrightarrow{\mathbb{P}_\vartheta} \rho(\vartheta) \quad \text{für } n \rightarrow \infty.$$

Der Schätzer $\hat{\rho}_n$ heißt asymptotisch normalverteilt, falls $\mathbb{E}[|\hat{\rho}_n|^2] < \infty$ und

$$\frac{\hat{\rho}_n - \mathbb{E}_\vartheta[\hat{\rho}_n]}{\sqrt{\text{Var}_\vartheta(\hat{\rho}_n)}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1) \quad \text{unter } \mathbb{P}_\vartheta.$$

Aufgrund des zentralen Grenzwertsatzes sind viele Schätzer asymptotisch normalverteilt, so auch in Beispiel 1.5. Daher kommt der Untersuchung von statistischen Modellen unter Normalverteilungsannahme eine besondere Bedeutung zu.

Zwei wichtige Konstruktionsprinzipien von Schätzern sind die Momentenmethode und Maximum-Likelihood-Schätzer:

Methode 1: Momentenmethode. Sei X_1, \dots, X_n eine mathematische Stichprobe reeller Zufallsvariablen mit $\mathbb{E}[|X_1|^d] < \infty$. Offensichtlich hängen i.A. die Momente einer Verteilung $m_k = m_k(\vartheta) := \mathbb{E}_\vartheta[X_1^k], k \in \mathbb{N}$, von ihrem Parameter $\vartheta \in \mathbb{R}^d$ ab. Aufgrund des Gesetzes der großen Zahlen ist der kanonische Schätzer von m_k gegeben durch das Stichprobenmoment $\hat{m}_k := \frac{1}{n} \sum_{j=1}^n X_j^k$. Der Momentenschätzer $\hat{\vartheta}$ von ϑ ist definiert als die Lösung der d -Gleichungen

$$\begin{aligned} m_1(\hat{\vartheta}) &= \hat{m}_1, \\ m_2(\hat{\vartheta}) &= \hat{m}_2, \\ &\vdots \\ m_d(\hat{\vartheta}) &= \hat{m}_d. \end{aligned}$$

Beispiel 1.10. Sei $X_1, \dots, X_n \stackrel{iid.}{\sim} \mathcal{N}(\mu, \sigma^2)$. Dann ist $m_1 = \mathbb{E}_{\mu, \sigma^2}[X_1] = \mu$ und $m_2 = \mathbb{E}_{\mu, \sigma^2}[X_1^2] = \text{Var}_{\mu, \sigma^2}(X_1) + \mathbb{E}_{\mu, \sigma^2}[X_1]^2 = \sigma^2 + \mu^2$. Folglich müssen wir die Gleichungen

$$\hat{\mu} = \frac{1}{n} \sum_{j=1}^n X_j \quad \text{und} \quad \hat{\sigma}^2 + \hat{\mu}^2 = \frac{1}{n} \sum_{j=1}^n X_j^2$$

lösen. Bezeichnen wir das Stichprobenmittel mit $\bar{X}_n := \frac{1}{n} \sum_{j=1}^n X_j$, erhalten wir die Lösung

$$\hat{\mu} = \bar{X}_n, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{j=1}^n (X_j - \bar{X}_n)^2.$$

Die Momentenmethode kann auf die Erwartungswerte allgemeinerer Funktionale verallgemeinert werden (siehe Übung 4). Für die zweite Methode benötigen wir etwas mehr Struktur, die wir auch im weiteren Verlauf der Vorlesung immer wieder aufgreifen.

Definition 1.11. Ein statistisches Modell $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ heißt dominiert, falls es ein σ -endliches Maß μ gibt, so dass \mathbb{P}_ϑ absolut stetig bzgl. μ ist ($\mathbb{P}_\vartheta \ll \mu$) für alle $\vartheta \in \Theta$. Die durch ϑ parametrisierte Radon-Nikodym-Dichte

$$L(\vartheta, x) := \frac{d\mathbb{P}_\vartheta}{d\mu}(x), \quad \vartheta \in \Theta, x \in \mathcal{X}$$

heißt Likelihoodfunktion, wobei diese meist als durch x parametrisierte Funktion in ϑ aufgefasst wird.

Beispiel 1.12.

- (i) $\mathcal{X} = \mathbb{R}, \mathcal{F} = \mathcal{B}(\mathbb{R}), \mathbb{P}_\vartheta$ ist gegeben durch die Lebesguedichte f_ϑ , beispielsweise $\mathbb{P}_{\mu, \sigma^2} = \mathcal{N}(\mu, \sigma^2)$ oder $\mathbb{P}_\vartheta = \mathcal{U}([0, \vartheta])$. Dann ist $L(\vartheta, x) = f_\vartheta(x)$.
- (ii) Jedes statistische Modell auf dem Stichprobenraum $(\mathbb{N}, \mathcal{P}(\mathbb{N}))$ oder allgemeiner auf einem abzählbaren Raum $(\mathcal{X}, \mathcal{P}(\mathcal{X}))$ ist vom Zählmaß dominiert. Die Likelihoodfunktion ist durch die Zähldichte gegeben.
- (iii) Ist $\Theta = \{\vartheta_1, \vartheta_2, \dots\}$ abzählbar, so ist $\mu = \sum_i c_i \mathbb{P}_{\vartheta_i}$ mit $c_i > 0$ und $\sum_i c_i = 1$ ein dominierendes Maß.

Methode 2: Maximum-Likelihood-Prinzip. Für ein dominiertes statistisches Modell mit Likelihoodfunktion $L(\vartheta, x)$ heißt eine Statistik $\hat{\vartheta}: \mathcal{X} \rightarrow \Theta$ (Θ trage eine σ -Algebra) Maximum-Likelihood-Schätzer (MLE: maximum likelihood estimator), falls

$$L(\hat{\vartheta}, x) = \sup_{\vartheta \in \Theta} L(\vartheta, x) \quad \text{für } \mathbb{P}_\vartheta\text{-f.a. } x \in \mathcal{X} \text{ und alle } \vartheta \in \Theta.$$

Beispiel 1.13. Betrachten wir wieder eine mathematische Stichprobe X_1, \dots, X_n normalverteilter Zufallsvariablen. Dann ist $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), \mathbb{P}_{\mu, \sigma^2}^n)$ mit $\mathbb{P}_{\mu, \sigma^2} = \mathcal{N}(\mu, \sigma^2)$ ein vom Lebesguemaß auf \mathbb{R}^n dominiertes Modell mit Likelihoodfunktion, $x = (x_1, \dots, x_n) \in \mathbb{R}^n$,

$$L(\mu, \sigma^2; x) = (2\pi\sigma^2)^{-n/2} \prod_{j=1}^n \exp\left(-\frac{(x_j - \mu)^2}{2\sigma^2}\right).$$

Um den Maximum-Likelihood-Schätzer zu berechnen, nutzen wir die Monotonie des Logarithmus und betrachten

$$\log L(\mu, \sigma^2; x) = -\frac{n}{2}(\log(2\pi) + \log \sigma^2) - \sum_{j=1}^n \frac{(x_j - \mu)^2}{2\sigma^2} \rightarrow \max_{\mu, \sigma^2}.$$

Ableiten nach μ und σ^2 führt auf die Gleichungen

$$0 = \sigma^{-2} \sum_{j=1}^n (x_j - \mu), \quad \frac{n}{2\sigma^2} = \frac{1}{2\sigma^4} \sum_{j=1}^n (x_j - \mu)^2.$$

Umstellen der ersten Gleichung nach μ liefert $\hat{\mu} = \bar{X}_n$ und Einsetzen in die zweite Gleichung ergibt $\hat{\sigma}^2 = n^{-1} \sum_j (X_j - \bar{X}_n)^2$. Es ist leicht nachzuprüfen, dass $\hat{\mu}$ und $\hat{\sigma}^2$ tatsächlich das Maximierungsproblem lösen (und messbar sind). In diesem Fall stimmt der Maximum-Likelihood-Schätzer also mit dem Momentenschätzer überein.

Beispiel 1.14. Sei X_1, \dots, X_n eine Poisson-verteilte mathematische Stichprobe mit Parameter $\lambda > 0$, d.h. $\mathcal{X} = \mathbb{Z}_+^n, \mathcal{F} = \mathcal{P}(\mathcal{X})$ und $\mathbb{P}_\lambda(X_1 = k) = \frac{\lambda^k e^{-\lambda}}{k!}$. Dann ist die gemeinsame Verteilung gegeben durch

$$\mathbb{P}_\lambda(X_1 = k_1, \dots, X_n = k_n) = \frac{\lambda^{\sum_i k_i} e^{-n\lambda}}{(k!)^n}, \quad k_1, \dots, k_n \in \mathbb{Z}_+.$$

Ableiten nach λ und null setzen führt auf den Maximum-Likelihood-Schätzer $\hat{\lambda} = \bar{X}_n$ (hinreichende Bedingung prüfen!).

1.1.2 Hypothesentests

Häufig interessiert man sich weniger für die gesamte zugrunde liegende Verteilung, als die Frage, ob eine bestimmte Eigenschaft erfüllt ist, oder nicht. Beispielsweise möchte man wissen, ob eine neue Behandlungsmethode I besser ist als die alte bisher genutzte Methode II. Aufgrund einer Beobachtung soll entschieden werden, ob die Hypothese ‘I ist besser als II’ akzeptiert werden kann oder verworfen werden sollte.

Um derartige Fragestellungen in einem statistischen Modell $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ zu formalisieren, wird die Parametermenge in zwei disjunkte Teilmengen Θ_0 und Θ_1 zerlegt, d.h. $\Theta = \Theta_0 \cup \Theta_1$ und $\emptyset = \Theta_0 \cap \Theta_1$. Das *Testproblem* liest sich dann als

$$H_0 : \vartheta \in \Theta_0 \quad \text{versus} \quad H_1 : \vartheta \in \Theta_1.$$

Dabei werden H_0, H_1 als *Hypothesen* bezeichnet, genauer heißt H_0 *Nullhypothese* und H_1 *Alternativhypothese* oder *Alternative*. Ein statistischer Test entscheidet nun zwischen H_0 und H_1 aufgrund einer Beobachtung $x \in \mathcal{X}$.

Definition 1.15. Ein (nicht-randomisierter) statistischer Test ist eine messbare Abbildung $\varphi: (\mathcal{X}, \mathcal{F}) \rightarrow (\{0, 1\}, \mathcal{P}(\{0, 1\}))$, wobei $\varphi(x) = 1$ heißt, dass die Nullhypothese verworfen/ die Alternative angenommen wird und $\varphi(x) = 0$ bedeutet, dass die Nullhypothese nicht verworfen wird/ akzeptiert wird. Die Menge $\{\varphi = 1\} = \{x \in \mathcal{X} : \varphi(x) = 1\}$ heißt Ablehnbereich von φ .

Allgemeiner ist ein randomisierter statistischer Test eine messbare Abbildung $\varphi: (\mathcal{X}, \mathcal{F}) \rightarrow ([0, 1], \mathcal{B}([0, 1]))$. Im Fall $\varphi(x) \in (0, 1)$ entscheidet ein unabhängiges Bernoulli-Zufallsexperiment mit Erfolgswahrscheinlichkeit $p = \varphi(x)$, ob die Hypothese verworfen wird.

Testen beinhaltet mögliche Fehlerentscheidungen:

- (i) Fehler 1. Art (α -Fehler, type I error): Entscheidung für H_1 , obwohl H_0 wahr ist,
- (ii) Fehler 2. Art (β -Fehler, type II error): Entscheidung für H_0 , obwohl H_1 wahr ist.

Definition 1.16. Sei φ ein Test der Hypothese $H_0 : \vartheta \in \Theta_0$ gegen die Alternative $H_1 : \vartheta \in \Theta_1$ im statistischen Modell $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$. Die Gütefunktion von φ ist definiert als

$$\beta_\varphi: \Theta \rightarrow \mathbb{R}_+, \vartheta \mapsto \mathbb{E}_\vartheta[\varphi]$$

Ein Test φ erfüllt das Signifikanzniveau $\alpha \in [0, 1]$ (oder φ ist Test zum Niveau α), falls $\beta_\varphi(\vartheta) \leq \alpha$ für alle $\vartheta \in \Theta_0$. Ein Test φ zum Niveau α heißt unverfälscht, falls $\beta_\varphi(\vartheta) \geq \alpha$ für alle $\vartheta \in \Theta_1$.

Somit hat ein nicht-randomisierten Test das Niveau $\alpha \in (0, 1)$, falls

$$\mathbb{P}_\vartheta(\varphi = 1) \leq \alpha, \quad \text{für alle } \vartheta \in \Theta_0,$$

beschränkt also die Wahrscheinlichkeit des Fehlers 1. Art mit der vorgegeben oberen Schranke α . In der Regel ist es nicht möglich, die Wahrscheinlichkeiten für die Fehler 1. und 2. Art gleichzeitig zu minimieren. Daher werden diese typischerweise asymmetrisch betrachtet:

- (i) Begrenzung der Fehlerwahrscheinlichkeit 1. Art durch ein vorgegebenes Signifikanzniveau α .

(ii) Unter der Maßgabe (i) wird die Wahrscheinlichkeit für Fehler 2. Art minimiert.

Eine zum Niveau α statistisch abgesicherte Entscheidung kann also immer nur zu Gunsten von H_1 getroffen werden. Daraus folgt die Merkregel “Was nachzuweisen ist, stets als Alternative H_1 formulieren”.

Beispiel 1.17 (Einseitiger Binomialtest). Von den 13 Todesfällen unter 55- bis 65-jährigen Arbeitern eines Kernkraftwerkes im Jahr 1995 waren 5 auf einen Tumor zurückzuführen. Die Todesursachenstatistik 1995 weist aus, dass Tumore bei etwa $1/5$ aller Todesfälle die Ursache in der betreffenden Altersklasse (in der Gesamtbevölkerung) darstellen. Ist die beobachtete Häufung von tumorbedingten Todesfällen signifikant zum Niveau 5%?

Bezeichne X die Anzahl der Tumortoten unter $n = 13$ Todesfällen. Dann ist das statistische Modell gegeben durch $\mathcal{X} = \{0, \dots, n\}$, $\mathcal{F} = \mathcal{P}(\mathcal{X})$ und $\mathbb{P}_p = \text{Bin}(13, p)$ mit Parameter $p \in [0, 1]$ und das Testproblem ist gegeben durch

$$H_0 : p \leq 1/5 \quad \text{versus} \quad H_1 : p > 1/5.$$

Ziel ist ein nicht-randomisierter Test zum Niveau $\alpha = 0,05$. Naheliegenderweise konstruieren wir $\varphi(x) = \mathbb{1}_{\{x > c\}}$ wobei der kritische Wert $c > 0$ so gewählt wird, dass $\sup_{p \leq 1/5} \mathbb{P}_p(X > c) \leq \alpha$. Um eine möglichst große Güte zu erreichen, sollte c unter dieser Nebenbedingung möglichst klein gewählt werden. Für $k \in \mathcal{X}$ gilt

$$\mathbb{P}_p(X \leq k) = \sum_{l=0}^k \binom{13}{l} p^l (1-p)^{13-l}.$$

Da $p \mapsto \mathbb{P}_p(X \leq k)$ für alle $k \in \mathcal{X}$ monoton fallend auf $[0, 1]$ ist (ableiten), folgt $\sup_{p \leq 1/5} \mathbb{P}_p(X > c) = \mathbb{P}_{1/5}(X > c)$. Wegen

$$\mathbb{P}_{1/5}(X \leq 4) \approx 0,901 \quad \text{und} \quad \mathbb{P}_{1/5}(X \leq 5) \approx 0,970,$$

wählen wir $c = 5$. Somit kann die Hypothese zum Niveau 0,05 nicht verworfen werden. Die Gütefunktion von φ

$$\beta_\varphi(p) = \mathbb{P}_p(X > 5) = \sum_{l=6}^{13} \binom{13}{l} p^l (1-p)^{13-l}, \quad p \in [0, 1],$$

ist monoton wachsend und somit ist φ auch unverfälscht.

Dieses Beispiel führt uns auf ein allgemeines Konstruktionsprinzip von Tests einer Hypothese $H_0 : \vartheta \in \Theta_0$ vs. $H_1 : \vartheta \in \Theta_1$ mit $\Theta_0 \neq \emptyset$ und $\Theta_1 = \Theta \setminus \Theta_0$.

Methode 3: Teststatistiken. Für Ablehnbereiche $(\Gamma_\alpha)_{\alpha \in (0,1)} \subseteq \mathcal{B}(\mathbb{R})$ und eine Teststatistik $T : (\mathcal{X}, \mathcal{F}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$ sei ein Test gegeben durch

$$\varphi(x) = \mathbb{1}_{\{T(x) \in \Gamma_\alpha\}}, \quad x \in \mathcal{X}. \tag{1.1}$$

Oft werden die Ablehnbereiche als Intervalle $\Gamma_\alpha = (c_\alpha, \infty)$ konstruiert für kritische Werte

$$c_\alpha = \inf \left\{ c \in \mathbb{R} : \sup_{\vartheta \in \Theta_0} \mathbb{P}_\vartheta(T(X) > c) \leq \alpha \right\}, \quad \alpha \in (0, 1). \tag{1.2}$$

Ist $\Theta_0 = \{\vartheta_0\}$ einelementig, dann sind die kritischen Werte genau das $(1-\alpha)$ -Quantil der Verteilung von T unter \mathbb{P}_{ϑ_0} . Ein wichtiges Konzept in der Testtheorie, insbesondere in Anwendungen, sind die p-Werte.

Definition 1.18. Sei $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ ein statistisches Modell und der Test φ der Hypothese $H_0 : \vartheta \in \Theta_0 \neq \emptyset$ gegeben durch (1.1). Dann ist der p-Wert einer Realisierung $x \in \mathcal{X}$ bezüglich φ definiert als

$$p_\varphi(x) = \inf_{\alpha : T(x) \in \Gamma_\alpha} \sup_{\vartheta \in \Theta_0} \mathbb{P}_\vartheta(T(X) \in \Gamma_\alpha).$$

Statt nur zu prüfen, ob ein Test eine Hypothese akzeptiert oder ablehnt, gibt der p-Wert (die Signifikanzwahrscheinlichkeit) das kleinste Signifikanzniveau an, zu dem eine Hypothese abgelehnt würde. Damit gibt der p-Wert Aufschluss darüber “wie stark” die Daten der Hypothese widersprechen.

Satz 1.19. Sei $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ ein statistisches Modell und sei φ ein Test der Hypothese $H_0 : \vartheta \in \Theta_0 \neq \emptyset$ gegeben durch $\varphi = \mathbb{1}_{\{T > c_\alpha\}}$ für eine Teststatistik $T : \mathcal{X} \rightarrow \mathbb{R}$ und kritische Werten $(c_\alpha)_{\alpha \in (0,1)}$ aus (1.2). Dann ist der p-Wert einer Realisierung $x \in \mathcal{X}$ bezüglich φ gegeben durch

$$p_\varphi(x) = \sup_{\vartheta \in \Theta_0} \mathbb{P}_\vartheta(T(X) \geq t^*) \quad \text{mit} \quad t^* := T(x).$$

Sei $\alpha \in (0, 1)$ ein **fest vorgegebenes** Niveau. Ist die Verteilung \mathbb{P}_ϑ^T stetig für alle $\vartheta \in \Theta_0$, gilt

$$\varphi(x) = 1 \iff p_\varphi(x) < \alpha \quad \mathbb{P}_\vartheta - \text{f.s. für alle } \vartheta \in \Theta_0.$$

Ist \mathbb{P}_ϑ^T (topologisch) diskret verteilt für alle $\vartheta \in \Theta_0$, gilt

$$\varphi(x) = 1 \iff p_\varphi(x) \leq \alpha \quad \mathbb{P}_\vartheta - \text{f.s. für alle } \vartheta \in \Theta_0.$$

Beweis. Definiere $\mathbb{P}_0 := \sup_{\vartheta \in \Theta_0} \mathbb{P}_\vartheta$. Da $c \mapsto \mathbb{P}_0(T > c)$ monoton fallend ist, gilt

$$p_\varphi(x) = \inf_{\alpha: t^* > c_\alpha} \mathbb{P}_0(T > c_\alpha) \geq \mathbb{P}_0(T \geq t^*).$$

Da $c_\alpha < t^*$ äquivalent zur Existenz eines $c < t^*$ mit $\mathbb{P}_0(T > c) \leq \alpha$ ist, folgt aus $\mathbb{P}_0(T > c_\alpha) \leq \alpha$ (Rechtsstetigkeit der Verteilungsfunktion), dass

$$p_\varphi(x) \leq \inf\{\alpha : c_\alpha < t^*\} \leq \inf\{\alpha : \mathbb{P}_0(\cap_{c < t^*} \{T > c\}) \leq \alpha\} = \mathbb{P}_0(T \geq t^*).$$

Zusammen erhalten wir $p_\varphi(x) = \mathbb{P}_0(T \geq t^*)$.

Sei nun α fest und T zunächst stetig verteilt. Aus $p_\varphi(x) = \mathbb{P}_0(T \geq t^*) < \alpha$ und $\lim_{c \uparrow t^*} \mathbb{P}_0(T \in (c, t^*)) = 0$ folgt $\mathbb{P}_0(T > c) \leq \alpha$ für ein $c < t^*$. Dann muss aber $\varphi(x) = 1$ gelten. Andersherum gilt

$$\varphi(x) = 1 \implies \exists c < t^* : \mathbb{P}_0(T \geq t^*) \leq \alpha - \mathbb{P}_0(T \in (c, t^*)).$$

Dabei gilt $\mathbb{P}_\vartheta(T \in (c, t^*)) = \mathbb{P}_\vartheta(T \in (c, T(x))) > 0$ für \mathbb{P}_ϑ -f.a. $x \in \mathcal{X}$ und für alle $\vartheta \in \Theta_0$. Ist T diskret verteilt, bleibt zu bemerken, dass $p_\varphi(x) = \mathbb{P}_0(T \geq t^*) = \mathbb{P}_0(T > c)$ für ein $c < t^*$. \square

Bemerkung 1.20.

- (i) Der Vorteil von p-Werten ist, dass sie unabhängig von einem a priori festgesetzten Signifikanzniveau α berechnet werden können. Deshalb werden in allen gängigen Statistik-Softwaresystemen statistische Hypothesentests über die Berechnung von p-Werten implementiert.
- (ii) **Warnung:** Alle Rahmenbedingungen des Experiments, insbesondere also das Signifikanzniveau, müssen vor dessen Durchführung festgelegt werden! Ein Signifikanzniveau darf nicht a posteriori aufgrund der erzielten p-Werte festgelegt werden. Dies widerspricht richtiger statistischer Praxis! Insbesondere wäre α eine Zufallsvariable (als Funktion in den Beobachtungen) und obiger Satz kann nicht angewendet werden.
- (iii) Der p-Wert gibt eine Antwort auf die Frage: “Wie wahrscheinlich sind die gemessenen Daten, gegeben, dass die Nullhypothese stimmt?” (und **nicht** auf die Frage “Wie wahrscheinlich ist es, dass die Nullhypothese wahr ist, gegeben den gemessenen Daten?”)

Beispiel 1.21. Geburten in Berlin:

- (i) *Hypothese*: Es werden genauso viele Jungen wie Mädchen geboren. Sind von $n \in \mathbb{N}$ Geburten $w \leq n$ Mädchen zur Welt gekommen, ist das statistische Modell gegeben durch den Stichprobenraum $\mathcal{X} = \{0, \dots, n\}$ und somit $(\mathcal{X}, \mathcal{P}(\mathcal{X}), (\mathbb{P}_\vartheta)_{\vartheta \in [0,1]})$ mit Binomialverteilungen $\mathbb{P}_\vartheta = \text{Bin}(n, \vartheta)$. Die Hypothese führt auf das zweiseitige Testproblem

$$H_0 : \vartheta = 1/2 \quad \text{versus} \quad H_1 : \vartheta \neq 1/2,$$

wobei $w \in \mathcal{X}$ beobachtet wird. Wir setzen das Niveau $\alpha = 0,05$. Die Teststatistik $T(w) = \left| \frac{w}{n} - \vartheta \right|$ führt auf einen *zweiseitigen Binomialtest*.

- (ii) *Hypothese*: Höchstens die Hälfte der geborenen Kinder hat nicht verheiratete Eltern. Von $n \in \mathbb{N}$ geborenen Kindern haben $v \leq n$ verheiratete Eltern. Mit $(\mathcal{X}, \mathcal{P}(\mathcal{X}), (\mathbb{P}_\vartheta)_{\vartheta \in [0,1]})$ wie oben betrachten wir hier das einseitige Testproblem

$$H_0 : \vartheta \leq 1/2 \quad \text{versus} \quad H_1 : \vartheta > 1/2,$$

wobei $v \in \mathcal{X}$ beobachtet wird. Das Niveau $\alpha = 0,05$ zusammen mit der Teststatistik $T(w) = \frac{w}{n} - \vartheta$ führt auf einen *einseitigen Binomialtest*.

Bemerkung 1.22. Bei großen Stichprobenumfängen ist es sinnvoll, einen *Gauß-Test* für geeignet normalisierter Teststatistik zu verwenden, um Binomialtest zu approximieren: Für $\vartheta \in (0,1)$ normalisieren wir die Beobachtung $X \sim \text{Bin}(n, \vartheta)$ durch $Y := \frac{X - n\vartheta}{\sqrt{n\vartheta(1-\vartheta)}}$. Aus dem Zentralen Grenzwertsatz folgt dann für eine standardnormalverteilte Zufallsvariable $Z \sim \mathcal{N}(0,1)$, dass

$$\begin{aligned} \mathbb{P}_\vartheta(T(X) > c_\alpha) &= \mathbb{P}_\vartheta\left(\frac{|X - n\vartheta|}{\sqrt{n\vartheta(1-\vartheta)}} > \sqrt{\frac{n}{\vartheta(1-\vartheta)}} c_\alpha\right) \\ &\xrightarrow{n \rightarrow \infty} \mathbb{P}\left(|Z| > \sqrt{\frac{n}{\vartheta(1-\vartheta)}} c_\alpha\right) \\ &= \frac{1}{2} \left(1 - \Phi\left(\sqrt{\frac{n}{\vartheta(1-\vartheta)}} c_\alpha\right)\right) \stackrel{!}{=} \alpha, \end{aligned}$$

Mit der Verteilungsfunktion $\Phi(x) = \mathbb{P}(Z \leq x)$. Folglich wählen wir $c_\alpha = \sqrt{\frac{\vartheta_0(1-\vartheta_0)}{n}} q_{1-2\alpha} = \sqrt{\frac{\vartheta_0(1-\vartheta_0)}{n}} \Phi^{-1}(1-2\alpha)$ mit $\vartheta = \vartheta_0$ unter H_0 .

1.1.3 Konfidenzmengen (Bereichsschätzung)

Während ein (Punkt-)Schätzer einen einzelnen Wert angibt, möglichst in der Nähe des wahren Parameters, um Rückschlüsse auf das zugrunde liegende Modell zu ziehen, geben Konfidenzbereiche ein Intervall an, in dem der Parameter mit gegebener Wahrscheinlichkeit liegt.

Definition 1.23. Sei $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ ein statistisches Modell mit abgeleitetem Parameter $\rho: \Theta \rightarrow \mathbb{R}^d$. Eine mengenwertige Abbildung $C: \mathcal{X} \rightarrow \mathcal{P}(\mathbb{R}^d)$ heißt Konfidenzmenge zum Konfidenzniveau $1 - \alpha$ (oder zum Irrtumsniveau α) für $\alpha \in (0,1)$, falls die Messbarkeitsbedingung $\{x \in \mathcal{X} : \rho(\vartheta) \in C(x)\} \in \mathcal{F}$ für alle $\vartheta \in \Theta$ erfüllt ist und es gilt

$$\mathbb{P}_\vartheta(\rho(\vartheta) \in C) = \mathbb{P}_\vartheta(\{x \in \mathcal{X} : \rho(\vartheta) \in C(x)\}) \geq 1 - \alpha \quad \text{für alle } \vartheta \in \Theta.$$

Im Fall $d = 1$ und falls $C(x)$ für jedes $x \in \mathcal{X}$ ein Intervall ist, heißt C Konfidenzintervall.

Beachte, dass $\rho(\vartheta)$ fix ist, während C zufällig ist. Man muss Konfidenzmengen also wie folgt *interpretieren*: Werden in m unabhängigen Experimenten für (verschiedene) Parameter Konfidenzmengen zum Niveau 0,95 konstruiert, dann liegt der unbekannte Parameter in 95% der Fälle in der jeweiligen Konfidenzmenge (für m groß genug; starkes Gesetz der großen Zahlen).

Ein verbreitetes Konstruktionsprinzip für die Konfidenzintervalle ist die Verwendung eines Schätzers und dessen Verteilung, wie im nächsten Beispiel illustriert.

Beispiel 1.24. Im Bernoulli-Experiment von Beispiel 1.5 gilt für $C_n := [\hat{p}_n - \varepsilon_n, \hat{p}_n + \varepsilon_n]$

$$\mathbb{P}_p(p \in C_n) = \mathbb{P}_p(|\hat{p}_n - p| < \varepsilon_n) = \mathbb{P}_p\left(\left|\sum_{i=1}^n (X_i - p)\right| < n\varepsilon_n\right) \stackrel{!}{\geq} 1 - \alpha.$$

Da $\sum_{i=1}^n X_i \sim \text{Bin}(n, p)$ können wir ε_n mithilfe der Quantile der Binomialverteilung bestimmen. Für große n könnte man wieder eine Normalapproximation verwenden. Das resultierende Konfidenzintervall besitzt dann aber nur asymptotisch das Niveau $1 - \alpha$.

Eine alternative Konstruktion von Konfidenzmengen bietet folgender Korrespondenzsatz:

Satz 1.25. Sei $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ ein statistisches Modell und $\alpha \in (0, 1)$. Dann gilt:

- (i) Liegt für jedes $\vartheta_0 \in \Theta$ ein Test φ_{ϑ_0} der Hypothese $H_0 : \vartheta = \vartheta_0$ zum Signifikanzniveau α vor, so definiert $C(x) = \{\vartheta \in \Theta : \varphi_\vartheta(x) = 0\}$ eine Konfidenzmenge zum Konfidenzniveau $1 - \alpha$.
- (ii) Ist C eine Konfidenzmenge zum Niveau $1 - \alpha$, dann ist $\varphi_{\vartheta_0}(x) = 1 - \mathbb{1}_{C(x)}(\vartheta_0)$ ein Niveau- α -Test der Hypothese $H_0 : \vartheta = \vartheta_0$.

Beweis. Nach Konstruktion erhält man in beiden Fällen,

$$\forall \vartheta \in \Theta : \forall x \in \mathcal{X} : \varphi_\vartheta(x) = 0 \iff \vartheta \in C(x).$$

Damit ist φ_ϑ ein Test zum Niveau α für alle ϑ genau dann, wenn

$$1 - \alpha \leq \mathbb{P}_\vartheta(\varphi = 0) = \mathbb{P}_\vartheta(\{x : \vartheta \in C(x)\})$$

und somit ist C eine Konfidenzmenge zum Niveau α . □

Beispiel 1.26. Mit Hilfe des Korrespondenzsatzes können wir ein Konfidenzintervall zum Niveau 0,95 für die Geburtswahrscheinlichkeit von Mädchen in Berlin berechnen. Im Modell aus Beispiel 1.21(i) ist das Konfidenzintervall gegeben durch

$$C(w) = \{\vartheta \in [0, 1] : \left|\frac{w}{n} - \vartheta\right| \leq c_{0,05}\} = \{\vartheta \in [0, 1], p_\varphi(w) > 0,05\},$$

wobei $p_\varphi(w)$ den zu φ gehörigen p-Wert der Realisierung w bezeichnet. Ist C sogar ein Konfidenzintervall? (Übung □).

1.2 Minimax- und Bayesansatz

Wir haben bereits verschiedene Schätzmethoden, wie den Maximum-Likelihood-Schätzer oder die Momentenmethode kennen gelernt. Natürlich gibt es noch viel mehr Konstruktionen. Wie sollte eine Methode anhand des gegebenen Schätzproblems ausgewählt werden? Sei also $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ ein statistisches Modell mit abgeleitetem Parameter $\rho : \Theta \rightarrow \mathbb{R}^d$ und Verlustfunktion L . Als mögliches Vergleichskriterium käme die Risikofunktion $R(\vartheta, \hat{\rho}) = \mathbb{E}_\vartheta[L(\vartheta, \hat{\rho})]$ eines Schätzers $\hat{\rho}$ in Frage. Beachte jedoch folgendes Beispiel:

Beispiel 1.27. Sei $X \sim \mathcal{N}(\mu, 1)$, $\mu \in \mathbb{R}$, und $L(\mu, \hat{\mu}) = (\hat{\mu} - \mu)^2$. Betrachte die zwei Schätzer $\hat{\mu}_1 = X$ und $\hat{\mu}_2 = 5$. Die Risiken sind dann gegeben durch

$$R(\mu, \hat{\mu}_1) = \mathbb{E}_\vartheta[(X - \mu)^2] = 1 \quad \text{und} \quad R(\mu, \hat{\mu}_2) = (5 - \mu)^2.$$

Damit hat $\hat{\mu}_1$ kleineres Risiko als $\hat{\mu}_2$ genau dann, wenn $\mu \notin [4, 6]$.

Definition 1.28. Im statistischen Modell $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ mit abgeleitetem Parameter $\rho : \Theta \rightarrow \mathbb{R}^d$ und Verlustfunktion L , heißt ein Schätzer $\hat{\rho}$ minimax, falls

$$\sup_{\vartheta \in \Theta} R(\vartheta, \hat{\rho}) = \inf_{\tilde{\rho}} \sup_{\vartheta \in \Theta} R(\vartheta, \tilde{\rho}),$$

wobei sich das Infimum über alle Schätzer (d.h. messbaren Funktionen) $\tilde{\rho} : \mathcal{X} \rightarrow \mathbb{R}^d$ erstreckt.

Definition 1.29. Der Parameterraum Θ trage eine σ -Algebra \mathcal{F}_Θ , die Verlustfunktion L sei produktmessbar und $\vartheta \mapsto \mathbb{P}_\vartheta(B)$ sei messbar für alle $B \in \mathcal{F}$. Die a-priori-Verteilung π des Parameters ϑ ist gegeben durch ein Wahrscheinlichkeitsmaß auf $(\Theta, \mathcal{F}_\Theta)$. Das zu π assoziierte Bayesrisiko eines Schätzers $\hat{\rho}$ ist

$$R_\pi(\hat{\rho}) := \mathbb{E}_\pi[R(\vartheta, \hat{\rho})] = \int_\Theta \int_{\mathcal{X}} L(\vartheta, \hat{\rho}(x)) \mathbb{P}_\vartheta(dx) \pi(d\vartheta).$$

Der Schätzer $\hat{\rho}$ heißt Bayesschätzer oder Bayes-optimal (bezüglich π), falls

$$R_\pi(\rho) = \inf_{\tilde{\rho}} R_\pi(\tilde{\rho}),$$

wobei sich das Infimum über alle Schätzer (d.h. messbaren Funktionen) $\tilde{\rho}: \mathcal{X} \rightarrow \mathbb{R}^d$ erstreckt.

Während ein Minimaxschätzer den maximal zu erwartenden Verlust minimiert, kann das Bayesrisiko als ein (mittels π) gewichtetes Mittel der zu erwartenden Verluste angesehen werden. Alternativ wird π als die subjektive Einschätzung der Verteilung des zugrundeliegenden Parameters interpretiert.

Beispiel 1.27 (fortgesetzt). Offensichtlich kann $\hat{\mu}_2$ kein Minimaxschätzer sein. Zunächst ist es aber nicht klar, ob es einen besseren Schätzer als $\hat{\mu}_2$ gibt. Tatsächlich werden wir später beweisen, dass $\hat{\mu}_1$ minimax ist. Unter der a-priori-Verteilung $\mu \sim \pi = \mathcal{U}([4, 6])$ hat jedoch $\hat{\mu}_2$ das kleinere Bayesrisiko $R_\pi(\hat{\mu}_2) = \frac{1}{3} < 1 = R_\pi(\hat{\mu}_1)$.

Das Bayesrisiko kann auch als insgesamt zu erwartender Verlust in folgendem Sinne verstanden werden: Definiere $\Omega := \mathcal{X} \times \Theta$ und die gemeinsame Verteilung von Beobachtung und Parameter $\tilde{\mathbb{P}}$ auf $(\mathcal{X} \times \Theta, \mathcal{F} \otimes \mathcal{F}_\Theta)$ gemäß $\tilde{\mathbb{P}}(dx, d\vartheta) = \mathbb{P}_\vartheta(dx) \pi(d\vartheta)$. Bezeichnen X und T die Koordinatenprojektionen von Ω auf \mathcal{X} bzw. Θ , dann gilt $R_\pi(\hat{\rho}) = \mathbb{E}_{\tilde{\mathbb{P}}}[L(T, \hat{\rho}(X))]$.

Wiederholung: Auf einem Wahrscheinlichkeitsraum $(\Omega, \mathcal{F}, \mathbb{P})$ ist die bedingte Wahrscheinlichkeit eines Ereignisses $A \in \mathcal{F}$ gegeben $B \in \mathcal{F}$ mit $\mathbb{P}(B) > 0$ definiert als $\mathbb{P}(A|B) = \mathbb{P}(A \cap B) / \mathbb{P}(B)$. Sei $\Omega = \bigcup_{i \in I} B_i$ eine abzählbare Zerlegung in paarweise disjunkte Ereignisse $B_i \in \mathcal{F}$, dann besagt die **Bayesformel** für jedes $A \in \mathcal{F}$ mit $\mathbb{P}(A) > 0$ und alle $k \in I$

$$\mathbb{P}(B_k|A) = \frac{\mathbb{P}(B_k) \mathbb{P}(A|B_k)}{\sum_{i \in I} \mathbb{P}(B_i) \mathbb{P}(A|B_i)}.$$

Mittels bedingten Erwartungswerten (Stochastik II) kann diese Formel auf Dichten ausgedehnt werden.

Definition 1.30. Sei $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ ein von μ dominiertes statistisches Modell mit Dichten $f_{X|T=\vartheta} := \frac{d\mathbb{P}_\vartheta}{d\mu}$. Sei π eine a-priori-Verteilung auf $(\Theta, \mathcal{F}_\Theta)$ mit Dichte f_T bzgl. einem Maß ν . Ist $f_{X|T= \cdot}: \mathcal{X} \times \Theta \rightarrow \mathbb{R}_+$ ($\mathcal{F} \otimes \mathcal{F}_\Theta$)-messbar, dann ist die a-posteriori-Verteilung des Parameters gegeben der Beobachtung $X = x$ definiert durch die ν -Dichte

$$f_{T|X=x}(\vartheta) = \frac{f_{X|T=\vartheta}(x) f_T(\vartheta)}{\int_\Theta f_{X|T=t}(x) f_T(t) \nu(dt)}, \quad \vartheta \in \Theta \quad (\tilde{\mathbb{P}}^X\text{-f.ü.}) \quad (1.3)$$

Das a-posteriori-Risiko eines Schätzers $\hat{\rho}$ gegeben $X = x$ ist definiert durch

$$R_\pi(\hat{\rho}|x) = \int_\Theta L(\vartheta, \hat{\rho}(x)) f_{T|X=x}(\vartheta) \nu(d\vartheta).$$

Beachte, dass im Nenner in (1.3) die Randdichte $f_X = \int_\Theta f_{X|T=t}(\cdot) f_T(t) \nu(dt)$ bzgl. μ von X in $(\mathcal{X} \times \Theta, \mathcal{F} \otimes \mathcal{F}_\Theta, \tilde{\mathbb{P}})$ steht, so dass der Nenner in (1.3) für $\tilde{\mathbb{P}}^X$ -f.a. $x \in \mathcal{X}$ größer als null ist.

Beispiel 1.31. Setze $\Theta = \{0, 1\}$, $L(\vartheta, r) = |\vartheta - r|$ (0-1-Verlust) und betrachte eine a-priori-Verteilung π mit $\pi(\{0\}) =: \pi_0$ und $\pi(\{1\}) =: \pi_1 = 1 - \pi_0$. Die Wahrscheinlichkeitsmaße \mathbb{P}_0 und \mathbb{P}_1 mögen Dichten p_0 und p_1 bzgl. einem Maß μ besitzen (z.B. $\mu = \mathbb{P}_0 + \mathbb{P}_1$). Dann ist die a-posteriori-Verteilung durch die Zähldichte

$$f_{T|X=x}(i) = \frac{\pi_i p_i(x)}{\pi_0 p_0(x) + \pi_1 p_1(x)}, \quad i = 0, 1 \quad (\tilde{\mathbb{P}}^X\text{-f.ü.})$$

gegeben. Damit ist das a-posteriori-Risiko eines Schätzers $\hat{\vartheta}: \mathcal{X} \rightarrow \{0, 1\}$ gegeben durch

$$R_\pi(\hat{\vartheta}|x) = \frac{\hat{\vartheta}(x)\pi_0 p_0(x) + (1 - \hat{\vartheta}(x))\pi_1 p_1(x)}{\pi_0 p_0(x) + \pi_1 p_1(x)}.$$

Satz 1.32. *Es gelten die Bedingungen der vorangegangenen Definition. Für das Bayesrisiko eines Schätzers $\hat{\rho}$ gilt*

$$R_\pi(\hat{\rho}) = \int R_\pi(\hat{\rho}|x) f_X(x) \mu(dx).$$

Minimiert $\hat{\rho}(x)$ für $\tilde{\mathbb{P}}^X$ -f.a. das a-posteriori-Risiko $\min_{t \in \text{ran}(\rho)} R_\pi(t|x)$, dann ist $\hat{\rho}$ Bayesschätzer.

Beweis. Aus (1.3) folgt $f_{T|x=x}(\vartheta) f_X(x) = f_{X|T=\vartheta}(x) f_T(\vartheta)$. Der Satz von Fubini ergibt

$$\begin{aligned} R_\pi(\hat{\rho}) &= \int_{\Theta} \int_{\mathcal{X}} L(\vartheta, \hat{\rho}(x)) \mathbb{P}_{\vartheta}(dx) \pi(d\vartheta) \\ &= \int_{\Theta} \int_{\mathcal{X}} L(\vartheta, \hat{\rho}(x)) f_{T|x=x}(\vartheta) f_X(x) \mu(dx) \nu(d\vartheta) = \int_{\mathcal{X}} R_\pi(\hat{\rho}|x) \mu(dx). \quad \square \end{aligned}$$

Korollar 1.33. *Unter quadratischem Verlust ist der Bayesschätzer gegeben durch*

$$\hat{\rho}(x) = \int_{\Theta} \rho(\vartheta) f_{T|X=x}(\vartheta) \nu(d\vartheta) =: \mathbb{E}[\rho(\vartheta)|X = x].$$

Der Bayesschätzer bzgl. absolutem Verlust ist gegeben durch den Median der a-posteriori-Verteilung. Für den 0-1-Verlust ist der Bayesschätzer der Modus der a-posteriori-Verteilung.

Beweis. Übung \square .

Methode 4: Bayesschätzer. Durch die Wahl einer Verlustfunktion und einer a-priori-Verteilung im statistischen Modell erhalten wir nach Berechnung der a-posteriori-Verteilung und durch das vorangegangene Korollar einen expliziten Bayesschätzer.

Beispiel 1.34. Sei $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ eine mathematische Stichprobe mit bekanntem $\sigma^2 > 0$ und a-priori-Verteilung $\mu \sim \mathcal{N}(a, b^2)$. Mittels Bayesformel kann die a-posteriori-Verteilung für eine Realisierung $x = (x_1, \dots, x_n)$ berechnet werden:

$$\begin{aligned} f_{T|X=x}(\mu) &\sim f_{X|T=\mu}(x) f_T(\mu) \\ &\sim \exp\left(-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}\right) \exp\left(-\frac{(\mu - a)^2}{2b^2}\right) \\ &\sim \exp\left(-\frac{\mu^2 - 2\mu\bar{x}_n}{2\sigma^2/n} - \frac{\mu^2 - 2a\mu}{2b^2}\right) \\ &\sim \exp\left(-\frac{(b^2 + \sigma^2/n)\mu^2 - 2\mu(b^2\bar{x}_n + a\sigma^2/n)}{2b^2\sigma^2/n}\right) \\ &\sim \exp\left(-\frac{1}{2}\left(\frac{n}{\sigma^2} + \frac{1}{b^2}\right)\left(\mu - \frac{b^2\bar{x}_n}{b^2 + \sigma^2/n} - \frac{a\sigma^2/n}{b^2 + \sigma^2/n}\right)^2\right). \end{aligned}$$

Gegeben der Beobachtung X ist ϑ also a-posteriori verteilt gemäß

$$\mathcal{N}\left(\frac{b^2}{b^2 + \frac{\sigma^2}{n}}\bar{X}_n - \frac{\frac{\sigma^2}{n}}{b^2 + \frac{\sigma^2}{n}}a, \left(\frac{n}{\sigma^2} + \frac{1}{b^2}\right)^{-1}\right).$$

Der Bayesschätzer bzgl. quadratischem Verlust, gegeben durch den a-posteriori Mittelwert, ist damit

$$\hat{\vartheta}_n = \frac{b^2}{b^2 + \frac{\sigma^2}{n}}\bar{X}_n - \frac{\frac{\sigma^2}{n}}{b^2 + \frac{\sigma^2}{n}}a.$$

Bemerkung 1.35. Erhalten wir bei Wahl einer Klasse von a-priori-Verteilungen für ein statistisches Modell dieselbe Klasse (i.A. mit anderen Parametern) als a-posteriori-Verteilung zurück, so nennt man die entsprechenden Verteilungsklassen konjugiert. Im obigen Beispiel haben wir gesehen, dass die Normalverteilungen zur den Normalverteilungen konjugiert sind (genauer müsste man sagen, dass für unbekanntem Mittelwert in der Normalverteilung a-priori Normalverteilungen konjugiert sind). Als weiteres Beispiel sind die Beta-Verteilungen zur Binomialverteilung konjugiert sind (siehe Übung \square). In diesen (Einzel-)Fällen ist es besonders einfach, die Bayesschätzer zu konstruieren. Für komplexere Modelle werden häufig computer-intensive Methoden wie MCMC (Markov Chain Monte Carlo) verwendet, um die a-posteriori-Verteilung zu berechnen (Problem: i.A. hochdimensionale Integration).

Lemma 1.36. *Unter den Bedingungen der vorangegangenen Definition gilt für jeden Schätzer $\hat{\rho}$*

$$\sup_{\vartheta \in \Theta} R(\vartheta, \hat{\rho}) = \sup_{\pi} R_{\pi}(\hat{\rho}),$$

wobei sich das zweite Supremum über alle a-priori-Verteilungen π erstreckt. Insbesondere ist das Risiko eines Bayesschätzers stets kleiner oder gleich dem Minimaxrisiko.

Beweis. Natürlich gilt $R_{\pi}(\hat{\rho}) = \int_{\Theta} R(\vartheta, \hat{\rho})\pi(d\vartheta) \leq \sup_{\vartheta \in \Theta} R(\vartheta, \hat{\rho})$. Durch Betrachtung der a-priori-Verteilung δ_{ϑ} folgt daher die Behauptung. \square

Durch dieses Lemma können wir untere Schranken für das Minimaxrisiko durch das Risiko von Bayesschätzern abschätzen. Mögliche Anwendungen illustriert folgender Satz.

Satz 1.37. *Sei X_1, \dots, X_n eine $\mathcal{N}(\mu, \sigma^2)$ -verteilte mathematische Stichprobe mit unbekanntem $\mu \in \mathbb{R}$ und bekanntem $\sigma^2 > 0$. Bezüglich quadratischem Risiko ist das arithmetische Mittel \bar{X}_n ein Minimalexschätzer von μ .*

Beweis. Wir betrachten a-priori-Verteilungen $\mu \sim \pi = \mathcal{N}(0, b^2)$. Nach Beispiel 1.34 ist die a-posteriori-Verteilung

$$\mathcal{N}\left(\frac{b^2\bar{X}_n}{b^2 + \frac{\sigma^2}{n}}, \left(\frac{n}{\sigma^2} + b^{-2}\right)^{-1}\right),$$

der Bayesschätzer bzgl. quadratischem Risiko ist gegeben durch den a-posteriori-Erwartungswert $\hat{\mu}_n = b^2\bar{X}_n/(b^2 + \sigma^2 n^{-1})$ und dessen a-posteriori-Risiko ist gegeben durch die Varianz der a-posteriori-Verteilung. Ist f_X die Randdichte von X von $\tilde{\mathbb{P}}$, folgt aus Satz 1.32

$$\begin{aligned} R_{\pi}(\hat{\mu}_n) &= \int_{\mathbb{R}^n} \text{Var}_{T|X=x}(\mu) f_X(x) dx \\ &= \int_{\mathbb{R}^n} (n\sigma^{-2} + b^{-2})^{-1} f_X(x) dx = (n\sigma^{-2} + b^{-2})^{-1}. \end{aligned}$$

Somit können wir das Minimaxrisiko nach unten abschätzen:

$$\begin{aligned} \inf_{\tilde{\mu}} \sup_{\mu \in \mathbb{R}} R(\mu, \tilde{\mu}) &= \inf_{\tilde{\mu}} \sup_{\pi} R_{\pi}(\tilde{\mu}) \geq \inf_{\tilde{\mu}} \sup_{b>0} R_{\mathcal{N}(0, b^2)}(\tilde{\mu}) \\ &\geq \sup_{b>0} \inf_{\tilde{\mu}} R_{\mathcal{N}(0, b^2)}(\tilde{\mu}) = \sup_{b>0} (n\sigma^2 + b^{-2})^{-1} = \frac{\sigma^2}{n}, \end{aligned}$$

wie behauptet, da $R(\mu, \bar{X}_n) = \sigma^2/n$. \square

1.3 Ergänzungen: Quantile

Definition. Sei \mathbb{P} ein Wahrscheinlichkeitsmaß auf $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ mit Verteilungsfunktion $F(x) = \mathbb{P}((-\infty, x])$. Für $\alpha \in (0, 1)$ ist das α -Quantil $q_\alpha \in \mathbb{R}$ von \mathbb{P} definiert durch

$$\mathbb{P}((-\infty, q_\alpha)) \leq \alpha \leq \mathbb{P}((-\infty, q_\alpha]).$$

Die Quantilfunktion ist definiert als verallgemeinertes Inverses von F :

$$F^{-1}(p) := \inf\{x \in \mathbb{R} : F(x) \geq p\}, \quad p \in [0, 1].$$

α -Quantile sind nicht eindeutig, falls F auf dem Niveau α irgendwo konstant ist. Es gilt aber

Lemma. $F^{-1}(\alpha)$ ist ein α -Quantil.

Beweis. Aufgrund der Rechtsstetigkeit von F gilt $F(F^{-1}(\alpha)) \geq \alpha$. Für alle $x < F^{-1}(\alpha)$ gilt $F(x) < \alpha$ und wegen der linken Grenzwerte von F

$$\alpha \geq \lim_{r \uparrow F^{-1}(\alpha)} F(x) = \lim_{r \uparrow F^{-1}(\alpha)} \mathbb{P}((-\infty, r]) = \mathbb{P}((-\infty, r)). \quad \square$$

Das verallgemeinerte Inverse hat folgende **Eigenschaften**:

- (i) $F^{-1}(p) \leq x \Leftrightarrow p \leq F(x)$;
- (ii) $F \circ F^{-1}(p) \geq p$ und Gleichheit gilt genau dann, wenn $p \in \text{ran } F$. Die Gleichheit kann nur dann nicht gelten, wenn F unstetig bei $F^{-1}(p)$ ist;
- (iii) $F^{-1} \circ F(x) \leq x$, wobei Gleichheit genau dann nicht gilt wenn x im Inneren oder am rechten Rand einer "Ebene" (kein Anstieg) von F liegt.

Damit gilt $F \circ F^{-1}(p) = p$ auf $(0, 1)$ genau dann, wenn F stetig ist (d.h. $\text{ran } F = [0, 1]$) und $F^{-1} \circ F(x) = x$ gilt auf \mathbb{R} genau dann, wenn F strikt monoton wachsend ist. Folglich ist F^{-1} ein echtes Inverses genau dann, wenn F stetig und streng monoton wachsend ist.

Satz. Ist $U \sim \text{Uni}([0, 1])$, dann besitzt die Zufallsvariable $F^{-1}(U)$ die Verteilungsfunktion F (Quantilstransformation). Besitzt X die Verteilungsfunktion F , dann gilt $F(X) \sim \text{Uni}([0, 1])$ genau dann, wenn F stetig ist.

Beweis. Aus (i) folgt $\mathbb{P}(F^{-1}(U) \leq x) = \mathbb{P}(U \leq F(x)) = F(x)$ für alle $x \in \mathbb{R}$. Andererseits gilt für $p \in (0, 1)$ wegen (i) und (ii)

$$\mathbb{P}(F(X) \leq p) = \mathbb{P}(X \leq F^{-1}(p)) = F(F^{-1}(p)) = p \iff p \in \text{ran } F. \quad \square$$

Schließlich wollen wir noch den **QQ-Plot** (Quantil-Quantil-Plot) verstehen: Die empirische Verteilungsfunktion einer mathematischen Stichprobe X_1, \dots, X_n ist gegeben durch $F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_i \leq x\}}$. Die Verteilungsfunktion der Standardnormalverteilung ist $\Phi(x) = \int_{-\infty}^x (2\pi)^{-1/2} e^{-y^2/2} dy$. Für große n approximiert F_n die wahre Verteilungsfunktion F , da nach dem starken Gesetz der großen Zahlen $F_n(x) \rightarrow \mathbb{E}[\mathbb{1}_{\{X_1 \leq x\}}] = F(x)$ \mathbb{P} -f.s. für alle $x \in \mathbb{R}$ gilt (tatsächlich gilt diese Konvergenz sogar gleichmäßig auf \mathbb{R} nach dem Satz von Borel-Cantelli). Falls $X_i \sim \mathcal{N}(\mu, \sigma^2)$, so gilt $F(x) = \Phi(\frac{x-m}{\sigma})$. Für die Quantilfunktion gilt also

$$F^{-1}(\Phi(x)) = \Phi^{-1}(\Phi(x)) \cdot \sigma + m = \sigma \cdot x + m,$$

d.h. $F^{-1} \circ \Phi$ ist eine Gerade. Im QQ-Plot wird F_n^{-1} (die empirischen Quantile) gegen Φ^{-1} aufgetragen und unter einer $\mathcal{N}(\mu, \sigma^2)$ -Annahme sollten die Werte in etwa auf einer Geraden liegen.

2 Lineares Modell

2.1 Regression und kleinste Quadrate

Regression ist eine Methode um den Zusammenhang zwischen einer *Zielgröße* (*Response-Variable*) Y und einem Vektor von erklärenden Variablen (*Kovariablen*, *Regressoren*) $X = (x_1, \dots, x_k)$ zu analysieren. Beginnen wir mit dem *einfachen linearen Modell*

$$Y_i = ax_i + b + \varepsilon_i, \quad i = 1, \dots, n,$$

mit Zufallsvariablen $\varepsilon_1, \dots, \varepsilon_n$, die zentriert sind ($\mathbb{E}_i[\varepsilon_i] = 0$) und endliche Varianz $\text{Var}(\varepsilon_i) = \sigma^2 > 0$ haben. Die Parameter $a, b \in \mathbb{R}, \sigma > 0$ sind unbekannt. Gesucht ist eine *Regressionsgerade* der Form $y = ax + b$, die die Beobachtungen möglichst gut erklärt. Der Parameter σ ist typischerweise nicht das Ziel der statistischen Inferenz und somit ein *Störparameter*.

Beispiel 2.1. Y_i ist das Wachstum von Deutschlands Bruttoinlandsproduktes im Jahr i . Die Kovariable x_i ist die Veränderung der Arbeitslosenquote im Vergleich zum Vorjahr. Unter Verwendung der Daten von 1992 bis 2012 aus den "World Development Indicators" der Weltbank erhalten als Regressionsgrade erhalten wir $y = -1,080 \cdot x + 1,338$. Betrachten wir alle sechs Gründungsmitglieder der EU im gleichen Zeitraum ergibt ganz ähnlich $y = -1,075 \cdot x + 1,819$. Der lineare Zusammenhang beider Größen ist als *Okuns Gesetz* bekannt.

Um die Situation weiter zu vereinfachen nehmen wir zunächst an, dass $\varepsilon_1, \dots, \varepsilon_n$ unabhängig und $\mathcal{N}(0, \sigma^2)$ -verteilt sind. Nun können wir den Maximum-Likelihood-Schätzer bestimmen: Der Beobachtungsvektor ist verteilt gemäß der Lebesgue-dichte

$$\begin{aligned} L(a, b, \sigma; y) &= \prod_{i=1}^n (2\pi\sigma^2)^{-1/2} \exp\left(-\frac{(y_i - ax_i - b)^2}{2\sigma^2}\right) \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - ax_i - b)^2\right), \quad y \in \mathbb{R}^n. \end{aligned}$$

Somit ist die Loglikelihoodfunktion

$$l(a, b, \sigma; y) := \log L(a, b, \sigma; y) = -\frac{n}{2}(\log \sigma^2 + \log(2\pi)) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - ax_i - b)^2.$$

Das Maximieren der Likelihood über a, b ist also äquivalent zum Minimieren der Summe der quadrierten Residuen (RSS: residual sum of squares). Auch wenn die Fehler nicht normalverteilt sind, kann diese Methode gute Ergebnisse erzielen.

Methode 5: Methode der kleinsten Quadrate. Im einfachen linearen Modell sind die Kleinste-Quadrate-Schätzer \hat{a}, \hat{b} durch Minimierung der Summe quadratischen Abstände

$$(\hat{a}, \hat{b}) := \arg \min_{a, b} \sum_{i=1}^n (Y_i - ax_i - b)^2$$

gegeben.

Satz 2.2. *Im einfachen linearen Modell mit unabhängigen und $\mathcal{N}(0, \sigma^2)$ -verteilten Fehlern, ist der Maximum-Likelihood-Schätzer gleich dem Kleinste-Quadrate-Schätzer und es gilt*

$$\hat{a} = \frac{\sum_{i=1}^n (x_i - \bar{x}_n)(Y_i - \bar{Y}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \quad \text{und} \quad \hat{b} = \bar{Y}_n - \hat{a}\bar{x}_n,$$

wobei $\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$ und $\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i$.

Beweis. Es bleibt festzustellen, dass wir durch Differentiation folgende Normalgleichungen erhalten:

$$0 = \sum_{i=1}^n x_i(Y_i - ax_i - b) \quad \text{und} \quad 0 = \sum_{i=1}^n (Y_i - ax_i - b),$$

die leicht gelöst werden können. □

Bemerkung 2.3. Bei der Wahl anderer Fehlerverteilungen ergibt das Maximum-Likelihood-Prinzip andere (nicht weniger sinnvolle) Schätzer (Übung □), die aber im Allgemeinen nicht in geschlossener Form darstellbar sind. Populäre nicht gaußsche Fehlerverteilungen sind Laplace- und Exponential-Verteilungen.

Haben wir $k \geq 2$ Kovariablen und n Beobachtungen Y_i , führt das zur *multiplen linearen Regression*

$$Y_i = \beta_0 + \sum_{j=1}^k \beta_j x_{i,j} + \varepsilon_i, \quad i = 1, \dots, n,$$

wobei die Fehlerterme (ε_i) iid. und zentriert sind mit $0 < \text{Var}(\varepsilon_i) =: \sigma^2 < \infty$. In Vektorschreibweise erhalten wir

$$\begin{aligned} Y &= (Y_1, \dots, Y_n)^\top \in \mathbb{R}^n && \text{Response-Vektor,} \\ X &:= \begin{pmatrix} 1 & x_{1,1} & \cdots & x_{1,k} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n,1} & \cdots & x_{n,k} \end{pmatrix} \in \mathbb{R}^{n \times (k+1)} && \text{Design-Matrix,} \\ \varepsilon &:= (\varepsilon_1, \dots, \varepsilon_n)^\top \in \mathbb{R}^n && \text{Vektor der Fehlerterme,} \\ \beta &:= (\beta_0, \dots, \beta_k)^\top \in \mathbb{R}^{k+1} && \text{Parametervektor,} \end{aligned}$$

so dass das multiple Regressionsmodell in der Form

$$Y = X\beta + \varepsilon$$

geschrieben werden kann. Der kleinste-Quadrate-Schätzer löst folglich das Minimierungsproblem

$$\min_b |Xb - Y|^2.$$

Beispiel 2.4. Im “crime”-Datensatz von Agresti and Finlay (1997, Kap. 9) stehen für die 51 Staaten der USA die beiden Responsevariablen

- Anzahl der Gewaltverbrechen pro 100.000 Einwohnern (crime),
- Morde pro 1.000.000 Einwohner (murder),

und folgende Kovariablen zur Verfügung:

- Prozentualer Anteil der Bevölkerung die in Ballungs-/ Großstadtgebieten leben (pctmetro),
- Prozentualer Anteil der weißen Bevölkerung (pctwhite),
- Prozentualer Anteil der Bevölkerung mit einem High-School-Abschluss (pcths),
- Prozentualer Anteil der Bevölkerung der unter der Armutsgrenze leben (poverty) und
- Prozentualer Anteil der Bevölkerung mit alleinerziehenden Eltern (single).

Bemerkung 2.5. Wechselwirkungen zwischen zwei Kovariablen x_i und x_j werden durch Interaktionsterme $x_i \cdot x_j$ modelliert. Kategorielle Kovariablen sollten durch eine Menge von sogenannten “Dummy-Indikatoren” kodiert werden, um nicht implizit eine (inadäquate) Metrisierung auf dem diskreten Wertebereich solcher Kovariablen zu induzieren. Eine kategorielle Kovariable mit ℓ möglichen Ausprägungen wird dabei durch $(\ell - 1)$ Indikatoren (d.h. $\{0, 1\}$ -wertige Variablen) repräsentiert. Der j -te Dummy-Indikator kodiert dabei das Ereignis, dass die Kategorie $(j + 1)$ bei der zugehörigen Kovariablen vorliegt, $j = 1, \dots, \ell - 1$. Sind also alle $(\ell - 1)$ Indikatoren gleich Null, so entspricht dies der (Referenz-) Kategorie 1 der zugehörigen kategoriellen Kovariable (vgl. Varianzanalyse).

Dies führt uns zur allgemeinen Definition des *linearen Modells*:

Definition 2.6. Ein lineares Modell mit n reellwertigen Beobachtungen $Y = (Y_1, \dots, Y_n)^\top$ und k -dimensionalem Parameter $\beta \in \mathbb{R}^k$, $k < n$, besteht aus einer reellen Matrix $X \in \mathbb{R}^{n \times k}$ von vollem Rang k , der Designmatrix, und einem Zufallsvektor $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^\top$, den Fehler- oder Störgrößen, mit $\mathbb{E}[\varepsilon_i] = 0$, $\text{Cov}(\varepsilon_i, \varepsilon_j) = \Sigma_{i,j}$ für eine Kovarianzmatrix $\Sigma > 0$. Beobachtet wird eine Realisierung von

$$Y = X\beta + \varepsilon.$$

Der (gewichtete) Kleinste-Quadrate-Schätzer $\hat{\beta}$ von β minimiert den gewichteten Euklidischen Abstand zwischen Beobachtungen und Modellvorhersage:

$$|\Sigma^{-1/2}(X\hat{\beta} - Y)|^2 = \inf_{b \in \mathbb{R}^k} |\Sigma^{-1/2}(Xb - Y)|^2.$$

Im gewöhnlichen Fall $\Sigma = \sigma^2 E_n$ mit Fehlerniveau $\sigma > 0$ erhalten wir den gewöhnlichen Kleinste-Quadrate-Schätzer (OLS: ordinary least squares)

$$|X\hat{\beta} - Y|^2 = \inf_{b \in \mathbb{R}^k} |Xb - Y|^2,$$

der unabhängig von der Kenntniss von σ^2 ist.

Bemerkung 2.7. Wir schreiben $\Sigma > 0$, falls Σ eine symmetrische, strikt positiv-definite Matrix ist. Dann ist Σ diagonalisierbar mit $\Sigma = TDT^\top$, $D = \text{diag}(\lambda_1, \dots, \lambda_n)$ Diagonalmatrix und T orthogonale Matrix, und wir setzen $\Sigma^{-1/2} := TD^{-1/2}T^\top$ mit $D^{1/2} := \text{diag}(\lambda_1^{-1/2}, \dots, \lambda_n^{-1/2})$. Wie erwartet, gilt $(\Sigma^{-1/2})^2 = \Sigma^{-1}$ und somit $|\Sigma^{-1/2}v|^2 = \langle \Sigma^{-1}v, v \rangle$.

Zusätzlich zur einfachen und multiplen Regression umfasst das lineare Modell weitere Beispiele.

Beispiel 2.8 (Polynomiale Regression). Wir beobachten

$$Y_i = a_0 + a_1x_i + a_2x_i^2 + \dots + a_{k-1}x_i^{k-1} + \varepsilon_i, \quad i = 1, \dots, n.$$

Damit ergibt sich als Parameter $\beta = (a_0, \dots, a_{k-1})^\top$ und eine Designmatrix vom Vandermonde-Typ

$$X = \begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^{k-1} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^{k-1} \end{pmatrix}.$$

Die Matrix hat vollen Rang, sofern k der Designpunkte (x_i) verschieden sind.

Lemma 2.9. Setze $X_\Sigma := \Sigma^{-1/2}X$. Mit Π_{X_Σ} werde die Orthogonalprojektion von \mathbb{R}^n auf den Bildraum $\text{ran}(X_\Sigma)$ bezeichnet. Dann gilt

$$\Pi_\Sigma = X_\Sigma(X_\Sigma^\top X_\Sigma)^{-1}X_\Sigma^\top$$

und für den Kleinste-Quadrate-Schätzer

$$\hat{\beta} = (X^\top \Sigma^{-1}X)^{-1}X^\top \Sigma^{-1}Y.$$

Insbesondere existiert der Kleinste-Quadrate-Schätzer, ist eindeutig und erwartungstreu.

Beweis. Zunächst beachte, dass $X_\Sigma^\top X_\Sigma = X^\top \Sigma^{-1} X$ invertierbar ist wegen der Invertierbarkeit von Σ und der Rangbedingung an X :

$$X^\top \Sigma^{-1} X v = 0 \Rightarrow v^\top X^\top \Sigma^{-1} X v = 0 \Rightarrow |\Sigma^{-1/2} X v| = 0 \Rightarrow |X v| = 0 \Rightarrow v = 0.$$

Setze $P_{X_\Sigma} := X_\Sigma (X_\Sigma^\top X_\Sigma)^{-1} X_\Sigma^\top$ und $w = P_{X_\Sigma} v$ für ein $v \in \mathbb{R}^n$. Dann folgt $w \in \text{ran}(X_\Sigma)$ und im Fall $v = X_\Sigma u$ durch Einsetzen $w = P_{X_\Sigma} X_\Sigma u = v$, so dass P_{X_Σ} eine Projektion auf $\text{ran}(X_\Sigma)$ ist. Da P_{X_Σ} selbstadjungiert (symmetrisch) ist, handelt es sich um die Orthogonalprojektion Π_{X_Σ} :

$$\forall u \in \mathbb{R}^n, \forall w \in \text{ran } X_\Sigma : \langle u - P_{X_\Sigma} u, w \rangle = \langle u, w \rangle - \langle u, P_{X_\Sigma} w \rangle = 0.$$

Aus der Eigenschaft $\hat{\beta} = \arg \min_b |\Sigma^{-1/2}(Y - Xb)|^2$ folgt, dass $\hat{\beta}$ die beste Approximation von $\Sigma^{-1/2} Y$ durch $X_\Sigma b$ liefert. Diese ist durch die Orthogonalprojektionseigenschaft $\Pi_{X_\Sigma} \Sigma^{-1/2} Y = X_\Sigma \hat{\beta}$ bestimmt. Es folgt

$$X_\Sigma^\top \Pi_{X_\Sigma} \Sigma^{-1/2} Y = (X_\Sigma^\top X_\Sigma) \hat{\beta} \Rightarrow (X_\Sigma^\top X_\Sigma)^{-1} X_\Sigma^\top \Sigma^{-1/2} Y = \hat{\beta}.$$

Schließlich folgt aus der Linearität des Erwartungswertes und $\mathbb{E}[\varepsilon] = 0$:

$$\mathbb{E}[\hat{\beta}] = \mathbb{E}[(X_\Sigma^\top X_\Sigma)^{-1} X_\Sigma^\top \Sigma^{-1/2} (X \beta + \varepsilon)] = \beta + 0 = \beta. \quad \square$$

Bemerkung 2.10.

- Im gewöhnlichen linearen Modell bzw. der multiplen linearen Regression gilt $\hat{\beta} = (X^\top X)^{-1} X^\top Y$ und ist somit unabhängig vom unbekanntem Parameter $\sigma > 0$.
- $X_\Sigma^\dagger := (X_\Sigma^\top X_\Sigma)^{-1} X_\Sigma^\top$ heißt auch Moore-Penrose-(Pseudo-)Inverse von X_Σ , so dass $\hat{\beta} = X_\Sigma^\dagger \Sigma^{-1/2} Y$ bzw. $\hat{\beta} = X^\dagger Y$ im gewöhnlichen linearen Modell gilt.

Wir kommen zum zentralen Satz in der Regressionsanalyse:

Satz 2.11 (Gauß-Markov). *Ist der Parameter $\rho = \langle \beta, v \rangle$ für ein $v \in \mathbb{R}^k$ im linearen Modell zu schätzen, so ist $\hat{\rho} = \langle \hat{\beta}, v \rangle$ ein (in den Daten Y) linearer erwartungstreuer Schätzer, der unter allen linearen erwartungstreuen Schätzern minimale Varianz besitzt, nämlich $\text{Var}(\hat{\rho}) = |X_\Sigma (X_\Sigma^\top X_\Sigma)^{-1} v|^2$.*

Beweis. Die Linearität ist klar und aus dem vorangegangenen Lemma folgt, dass $\hat{\rho}$ erwartungstreu ist. Sei nun $\tilde{\rho} = \langle Y, w \rangle$ ein beliebiger linearer erwartungstreuer Schätzer von ρ . Dies impliziert für alle $\beta \in \mathbb{R}^k$

$$\mathbb{E}[\langle Y, w \rangle] = \rho \Rightarrow \langle X \beta, w \rangle = \langle \beta, v \rangle \Rightarrow \langle X^\top w - v, \beta \rangle = 0$$

und somit $v = X^\top w = X_\Sigma^\top \Sigma^{1/2} w$. Nach Pythagoras erhalten wir

$$\begin{aligned} \text{Var}(\tilde{\rho}) &= \mathbb{E}[\langle \varepsilon, w \rangle^2] = \mathbb{E}[w^\top \varepsilon \varepsilon^\top w] \\ &= w^\top \Sigma w = |\Sigma^{1/2} w|^2 = |\Pi_{X_\Sigma}(\Sigma^{1/2} w)|^2 + |(E_n - \Pi_\Sigma)(\Sigma^{1/2} w)|^2. \end{aligned}$$

Damit gilt $\text{Var}(\tilde{\rho}) \geq |\Pi_{X_\Sigma}(\Sigma^{1/2} w)|^2 = |X_\Sigma (X_\Sigma^\top X_\Sigma)^{-1} X^\top w| = |X_\Sigma (X_\Sigma^\top X_\Sigma)^{-1} v| = \text{Var}(\hat{\rho})$. \square

Bemerkung 2.12. Man sagt, dass der Schätzer $\hat{\rho}$ im Satz von Gauß-Markov best linearer erwartungstreuer Schätzer (blue: best linear unbiased estimator) ist. Eingeschränkt auf lineare Schätzer ist der Kleinste-Quadrate-Schätzer damit minimax. Ob es einen besseren nichtlinearen Schätzer geben kann, werden wir in Kapitel 3 beantworten.

Im gewöhnlichen linearen Modell ist die optimale Varianz insbesondere $\sigma^2 |X (X^\top X)^{-1} v|^2$. In diesem Spezialfall ist es auch von Interesse das Rauschniveau σ^2 zu schätzen. Dies ermöglicht es insbesondere Tests und Konfidenzbereiche zu konstruieren.

Lemma 2.13. *Im gewöhnlichen linearen Modell mit $\sigma > 0$ und Kleinste-Quadrate-Schätzer $\widehat{\beta}$ gilt $X\widehat{\beta} = \Pi_X Y$ und $R := Y - X\widehat{\beta}$ bezeichne den Vektor der Residuen. Die geeignet normalisierte Stichprobenvarianz*

$$\widehat{\sigma}^2 := \frac{|R|^2}{n-k} = \frac{|Y - X\widehat{\beta}|^2}{n-k}$$

ist erwartungstreu Schätzer von σ^2 .

Beweis. $X\widehat{\beta} = \Pi_X Y$ folgt aus Lemma 2.9. Einsetzen zeigt $\mathbb{E}[|Y - X\widehat{\beta}|^2] = \mathbb{E}[|Y - \Pi_X Y|^2] = \mathbb{E}[|(E_n - \Pi_X)\varepsilon|^2]$. Ist nun e_1, \dots, e_{n-k} eine Orthonormalbasis vom $(n-k)$ -dimensionalen Bild $\text{ran}(E_n - \Pi_X) \subseteq \mathbb{R}^n$, so folgt

$$\mathbb{E}[|(E_n - \Pi_X)\varepsilon|^2] = \sum_{i=1}^{n-k} \mathbb{E}[\langle \varepsilon, e_i \rangle^2] = \sigma^2(n-k),$$

was die Behauptung impliziert. □

Beachte, dass der Maximum-Likelihood-Schätzer von σ^2 gegeben ist durch $\widehat{\sigma}_{ML}^2 = n^{-1}|R|^2 \neq \widehat{\sigma}^2$ (Übung □). Der erwartungstreu Schätzer $\widehat{\sigma}^2$ wird in der Praxis bevorzugt, hat jedoch größere Varianz als $\widehat{\sigma}_{ML}^2$.

Bevor wir uns mit statistischer Inferenz, also der Konstruktion von Tests und Konfidenzintervallen, im linearen Modell beschäftigen, soll der Bayesianische Ansatz auf das Regressionsproblem angewendet werden.

Satz 2.14. *Im gewöhnlichen linearen Modell $Y = X\beta + \varepsilon$ mit $\varepsilon \sim \mathcal{N}(0, \sigma^2 E_n)$ und bekanntem $\sigma > 0$ genüge $\beta \in \mathbb{R}^k$ der a-priori-Verteilung*

$$\beta \sim \mathcal{N}(m, \sigma^2 M)$$

mit Parametern $m \in \mathbb{R}^k$ und symmetrisch positiv definiten Matrix $M \in \mathbb{R}^{k \times k}$. Dann ist die a-posteriori-Verteilung von β gegeben einer Realisierung $y \in \mathbb{R}^n$ gegeben durch

$$\beta|Y=y \sim \mathcal{N}(\mu_y, \Sigma_y) \quad \text{mit} \quad \Sigma_y = \sigma^2(X^\top X + M^{-1})^{-1}, \mu_y = \Sigma_y(\sigma^{-2}X^\top y + \sigma^{-2}M^{-1}m).$$

Insbesondere ist der Bayesschätzer bzgl. quadratischem Verlust gegeben durch $\widehat{\beta}_{Bayes} = (X^\top X + M^{-1})^{-1}(X^\top Y + M^{-1}m)$.

Beweis. Für die a-posteriori-Dichte an der Stelle $t \in \mathbb{R}^k$ gilt

$$\begin{aligned} f_{\beta|Y=y}(t) &\sim \exp\left(-\frac{1}{2\sigma^2}(y - Xt)^\top(y - Xt)\right) \exp\left(-\frac{1}{2\sigma^2}(t - m)^\top M^{-1}(t - m)\right) \\ &\sim \exp\left(\frac{1}{\sigma^2}t^\top X^\top y - \frac{1}{2\sigma^2}t^\top X^\top X t - \frac{1}{2\sigma^2}t^\top M^{-1}t + \frac{1}{\sigma^2}t^\top M^{-1}m\right) \\ &= \exp\left(\frac{1}{\sigma^2}t^\top(X^\top y + M^{-1}m) - \frac{1}{2\sigma^2}t^\top(X^\top X + M^{-1})t\right). \end{aligned}$$

Daher ist β gegeben $Y = y$ normalverteilt mit Kovarianzmatrix $\Sigma_y = (\sigma^{-2}X^\top X + \sigma^{-2}M^{-1})^{-1}$ und Mittelwert $\mu_y = \Sigma_y(X^\top y + M^{-1}m)/\sigma^2$. □

Es ist erneut bemerkenswert, dass der Bayesschätzer $\widehat{\beta}_{Bayes}$ nicht von σ^2 abhängt.

Bemerkung 2.15. Indem wir auch den Parameter σ^2 mit einer a-priori-Verteilung versehen, erhalten wir ein (mehrstufiges) Bayesmodell. Da wir besonders an konjugierten Verteilungsklassen interessiert sind, wird hierzu oft die inverse Gamma-Verteilung verwendet: Ist $Z \sim \Gamma(a, b)$ so ist $1/Z \sim IG(a, b)$ invers Gamma-verteilt mit Parametern $a, b > 0$ und Lebesgue-dichte

$$f_{a,b}(x) = \frac{b^a}{\Gamma(a)} x^{-(a-1)} e^{-a/x} \mathbb{1}_{(0,\infty)}(x), \quad x \in \mathbb{R}.$$

Das Bayesmodell ist also gegeben durch

$$Y|\beta, \sigma^2 \sim \mathcal{N}(X\beta, \sigma^2 E_n), \quad \beta|\sigma^2 \sim \mathcal{N}(m, \sigma^2 M), \quad \sigma \sim IG(a, b).$$

Die gemeinsame Verteilung von $(\beta, \sigma^2) \sim NIG(m, M, a, b)$ wird Normal-inverse Gammaverteilung genannt und besitzt die Dichte

$$\begin{aligned} f(\beta, \sigma^2) &= \frac{1}{(2\pi\sigma^2)^{k/2}|M|^{1/2}} \exp\left(-\frac{1}{2\sigma^2}(\beta - m)^\top M^{-1}(\beta - m)\right) \frac{b^a}{\Gamma(a)(\sigma^2)^{a+1}} e^{-a/\sigma^2} \\ &\sim \frac{1}{(\sigma^2)^{k/2+a+1}} \exp\left(\frac{1}{2\sigma^2}((\beta - m)^\top M^{-1}(\beta - m) + b)\right), \quad \beta \in \mathbb{R}^k, \sigma^2 > 0. \end{aligned}$$

In diesem Modell ist die a-posteriori-Verteilung von σ^2 gegeben β und Y gegeben durch $\sigma^2|\beta, Y \sim IG(a', b')$ mit $a' = a + \frac{n}{2} + \frac{k}{2}$ und

$$b' = b + \frac{1}{2}(Y - X\beta)^\top (Y - X\beta) + \frac{1}{2}(\beta - m)^\top M^{-1}(\beta - m).$$

Die a-posteriori-Verteilung von (β, σ^2) gegeben Y ist $(\beta, \sigma^2)|Y \sim NIG(\tilde{m}, \tilde{M}, \tilde{a}, \tilde{b})$ mit Parametern

$$\begin{aligned} \tilde{M} &= (X^\top X + M^{-1})^{-1}, \quad \tilde{m} = \tilde{M}(M^{-1}m + X^\top y), \\ \tilde{a} &= a + \frac{n}{2}, \quad \tilde{b} = b + \frac{1}{2}\left(Y^\top Y + m^\top M^{-1}m - \tilde{m}^\top \tilde{M}^{-1}\tilde{m}\right), \end{aligned}$$

siehe Fahrmeir et al. (2009, Kap. 3.5).

Korollar 2.16. *Unter den Voraussetzungen des vorangegangenen Satzes mit $m = 0$ und $M = \tau^2 E_k, \tau > 0$, gilt für den Bayeschätzer unter quadratischem Verlust*

$$\hat{\beta}_{Bayes} = \arg \min_{\beta \in \mathbb{R}^k} |Y - X\beta|^2 + \frac{1}{\tau^2} |\beta|^2.$$

Beweis. Im Spezialfall $m = 0$ und $M = \tau^2 E_k$ folgt aus obigem Satz $\hat{\beta}_{Bayes} = (X^\top X + \tau^{-2} E_k)^{-1} X^\top y$. Andererseits gilt

$$\begin{aligned} &\arg \min_{\beta} \left((Y^\top - \beta^\top X^\top)(Y - X\beta) + \frac{1}{\tau^2} \beta^\top \beta \right) \\ &= \arg \min_{\beta} \left(-2Y^\top X\beta + \beta^\top (X^\top X + \frac{1}{\tau^2} E_k) \beta \right). \end{aligned}$$

Null setzen des Differenzials der Funktion $\beta \mapsto -2Y^\top X\beta + \beta^\top (X^\top X + \frac{1}{\tau^2} E_k) \beta$ liefert $0 = -2Y^\top X + 2\beta^\top (X^\top X + \frac{1}{\tau^2} E_k)$, so dass aus der positiv Definitheit und Symmetrie von $X^\top X + \frac{1}{\tau^2} E_k$ die Behauptung folgt. \square

Der Bayesansatz führt uns also zu einer neuen Schätzmethode im linearen Modell:

Methode 6: Ridge-Regression. Im linearen Modell $Y = X\beta + \varepsilon$ ist der Ridge-Regressionsschätzer oder Schrumpfungsschätzer (engl.: Shrinkage) mit Schrumpfungskoeffizient $\lambda \geq 0$ definiert als

$$\hat{\beta}_{ridge} = \arg \min_{\beta \in \mathbb{R}^k} |Y - X\beta|^2 + \lambda |\beta|^2.$$

Durch Einführung des Strafterms (engl.: *penalty*) $\lambda|\beta|^2$ wird die Varianz auf Kosten eines Bias verringert. Dies ist insbesondere sinnvoll, wenn einige (wenige) Koeffizienten von β groß sind und die übrigen klein und liefert in diesen Fällen gute Schätzergebnisse auch wenn die Parameterdimension in einer ähnlichen Größenordnung liegt wie die Anzahl der Beobachtungen ($n \sim p$). Dies wird im nächsten Beispiel illustriert. Die richtige Wahl des Strumpfungparameters λ ist allerdings ein schwieriges Problem.

Beispiel 2.17. Betrachten wir das Modell $Y_i = x_i^\top \beta + \varepsilon_i$ mit Kovariablenvektor $x_i \in \mathbb{R}^p$, Parameter $\beta \in \mathbb{R}^p$ und $\varepsilon_i \stackrel{iid.}{\sim} \mathcal{N}(0, 1)$ mit $i = 1, \dots, n$. Wir wählen $n = 50$ und $p = 30$ wobei 10 Koeffizienten groß sind (zwischen 0,5 und 1) und 20 klein (zwischen 0 und 0,3) und bestimmen den mittleren Quadratischen Fehler aus 200 Simulationen für verschiedene Werte von $\lambda \in [0, 20]$ (Übung \square).

2.2 Inferenz unter Normalverteilungsannahme

Im Folgenden werden wir das gewöhnliche lineare Modell unter der Normalverteilungsannahme $(\varepsilon_i) \sim \mathcal{N}(0, \sigma^2 E_n)$ betrachten.

Beispiel 2.18. Sind die Messfehler $(\varepsilon_i) \sim \mathcal{N}(0, \sigma^2 E_n)$ gemeinsam normalverteilt und $\rho = \langle v, \beta \rangle$ für $v \in \mathbb{R}^k$, so gilt

$$\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 (X^\top X)^{-1}) \quad \text{und} \quad \hat{\rho} = \langle v, \hat{\beta} \rangle \sim \mathcal{N}(\rho, \sigma^2 v^\top (X^\top X)^{-1} v).$$

Ist $\sigma > 0$ bekannt, so ist ein Konfidenzintervall zum Niveau 95% für ρ gegeben durch

$$I_{0,95}(\rho) := [\hat{\rho} - 1,96\sigma\sqrt{v^\top (X^\top X)^{-1} v}, \hat{\rho} + 1,96\sigma\sqrt{v^\top (X^\top X)^{-1} v}].$$

Dabei ist der Wert 1,96 gerade das 0,975-Quantil bzw. 0,025 Fraktile der Standardnormalverteilung. Analog (Korrespondenzsatz) wird der zweiseitige Gauß-Test der Hypothese $H_0 : \rho = \rho_0$ gegen $H_1 : \rho \neq \rho_0$ zum Niveau $\alpha \in (0, 1)$ konstruiert: Wähle die Teststatistik $|\hat{\rho} - \rho_0|$ und den kritischen Wert $q_{1-\alpha/2}\sigma\sqrt{v^\top (X^\top X)^{-1} v}$ mit dem $(1 - \alpha/2)$ -Quantil von $\mathcal{N}(0, 1)$.

Ist σ unbekannt, so ist eine Idee, einfach σ durch den Schätzer $\hat{\sigma}$ in obiger Formel zu ersetzen. Allerdings wird dann das vorgegebene Niveau nur noch asymptotisch erreicht für einen konsistenten Schätzer (Slutsky-Lemma). Im vorliegenden Fall können wir aber sogar die Verteilung für endliche Stichprobenumfänge exakt bestimmen.

Definition 2.19. Die t-Verteilung $t(n)$ (oder Student-t-Verteilung) mit $n \in \mathbb{N}$ Freiheitsgraden auf $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ ist gegeben durch die Lebesgue-dichte

$$t_n(x) = \frac{\Gamma(\frac{n+1}{2})}{\Gamma(\frac{n}{2})\sqrt{\pi n}} \left(1 + \frac{x^2}{n}\right)^{-(n+1)/2}, \quad x \in \mathbb{R}.$$

Die F-Verteilung $F(m, n)$ (oder Fisher-Verteilung) mit $(m, n) \in \mathbb{N}^2$ Freiheitsgraden auf $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ ist gegeben durch die Lebesgue-dichte

$$f_{m,n}(x) = \frac{m^{m/2} n^{n/2}}{B(\frac{m}{2}, \frac{n}{2})} \frac{x^{m/2-1}}{(mx+n)^{(m+n)/2}} \mathbb{1}_{\mathbb{R}^+}(x), \quad x \in \mathbb{R}.$$

Dabei bezeichnet $\Gamma(p) = \int_0^\infty t^{p-1} e^{-t} dt$ die Gamma-Funktion und $B(p, q) = \frac{\Gamma(p)\Gamma(q)}{\Gamma(p+q)}$ die Beta-Funktion.

Erinnerung: Für $X_1, \dots, X_m \sim \mathcal{N}(0, 1)$ ist $X := \sum_{i=1}^m X_i^2 \sim \chi^2(m)$ verteilt mit Lebesgue-dichte $f_X(x) = (2^{m/2} \Gamma(\frac{m}{2}))^{-1} x^{m/2-1} e^{-x/2} \mathbb{1}_{\mathbb{R}^+}(x)$.

Lemma 2.20. Es seien $X_1, \dots, X_m, Y_1, \dots, Y_n$ unabhängige $\mathcal{N}(0, 1)$ -verteilte Zufallsvariablen. Dann gilt

$$T_n := \frac{X_1}{\sqrt{\frac{1}{n} \sum_{j=1}^n Y_j^2}} \sim t(n) \quad \text{und} \quad F_{m,n} := \frac{\frac{1}{m} \sum_{i=1}^m X_i^2}{\frac{1}{n} \sum_{j=1}^n Y_j^2} \sim F(m, n).$$

Beweis. Es gilt $T_n^2 = F_{1,n}$, so dass mittels Dichtetransformation $f_{|T_n|}(x) = f_{F_{1,n}}(x^2)2x, x \geq 0$, gilt. Da T_n symmetrisch (wie $-T_n$) verteilt ist, folgt $f_{T_n} = F_{F_{1,n}}(x^2)|x|, x \in \mathbb{R}$, und Einsetzen zeigt die Behauptung für T_n , sofern $F_{1,n}$ $F(1, n)$ -verteilt ist.

Um die Behauptung für $F_{m,n}$ nachzuweisen, benutze, dass $X := \sum_{i=1}^m X_i^2$ $\chi^2(m)$ -verteilt und $Y := \sum_{j=1}^n Y_j^2$ $\chi^2(n)$ -verteilt sind. Wegen Unabhängigkeit von X und Y gilt für $z > 0$ (setze $w = x/y$)

$$\begin{aligned}\mathbb{P}(X/Y \leq z) &= \int \int \mathbb{1}_{\{x/y \leq z\}} f_X(x) f_Y(y) dx dy \\ &= \int \mathbb{1}_{\{w \leq z\}} \left(\int f_X(wy) f_Y(y) y dy \right) dw,\end{aligned}$$

so dass sich die Dichte wie folgt ergibt (setze $w = (z+1)y$)

$$\begin{aligned}f_{X/Y}(z) &= \int f_X(zy) f_Y(y) y dy \\ &= \frac{2^{-(m+n)/2}}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} \int_0^\infty (zy)^{m/2-1} y^{n/2} e^{-(zy+y)/2} dy \\ &= \frac{2^{-(m+n)/2}}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} \int_0^\infty (zw/(z+1))^{m/2-1} (w/(z+1))^{n/2} e^{-w/2} (z+1)^{-1} dw \\ &= \frac{\Gamma(\frac{m+n}{2})}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} z^{m/2-1} (z+1)^{-(m+n)/2}, \quad z > 0.\end{aligned}$$

Dichtetransformation ergibt damit für $F_{m,n} = \frac{n}{m} \frac{X}{Y}$ die Dichte $\frac{m}{n} f_{X/Y}(\frac{m}{n}x) = f_{m,n}(x)$. \square

Bemerkung 2.21. Es gilt $T_n^2 = F_{1,n}$. Für $n = 1$ ist die $t(n)$ -Verteilung gerade die Cauchy-Verteilung und für $n \rightarrow \infty$ konvergiert sie schwach gegen die Standardnormalverteilung. Für jedes $n \in \mathbb{N}$ besitzt $t(n)$ nur Momente bis zur Ordnung $p < n$ (sie ist *heavy-tailed*). Ähnliches gilt für die F-Verteilung, insbesondere konvergiert die Verteilung von $mF_{m,n}$ für $n \rightarrow \infty$ gegen die $\chi^2(m)$ -Verteilung.

Aus diesem Lemma ergeben sich die Standardtests für die Parameter der Normalverteilung, siehe Witting (1985, S. 200-204).

Bevor wir zur Konstruktion von Tests und Konfidenzbändern im linearen Modell kommen noch ein weiteres nützliches Hilfsresultat zur Verteilung quadratischer Formen:

Lemma 2.22. *Seien $X \sim \mathcal{N}(0, E_n)$ und R eine symmetrische, idempotente $(n \times n)$ -Matrix (d.h. $R = R^\top$ und $R^2 = R$) mit $\text{rank}(R) = r \leq n$. Dann gilt*

- (i) $X^\top R X \sim \chi^2(r)$,
- (ii) $X^\top R X$ ist unabhängig von BX für jede Matrix $B \in \mathbb{R}^{p \times n}$ mit $p \leq n$ und $BR = 0$,
- (iii) für jede weitere symmetrische, idempotente Matrix $S \in \mathbb{R}^{n \times n}$ mit $\text{rank}(S) = s \leq n$ und $RS = 0$ sind $X^\top R X$ und $X^\top S X$ unabhängig und

$$\frac{s}{r} \frac{X^\top R X}{X^\top S X} \sim F(r, s).$$

Beweis. (i) Da R symmetrisch und idempotent ist, existiert eine Orthogonalmatrix P mit $R = PD_r P^\top$, wobei $D_r = \begin{pmatrix} E_r & 0 \\ 0 & 0 \end{pmatrix}$. Da P orthogonal ist und X standardnormalverteilt, folgt $W := P^\top X \sim \mathcal{N}(0, E_n)$. Wegen

$$X^\top R X = X^\top R^2 X = (RX)^\top (RX) = (PD_r W)^\top (PD_r W) = W^\top D_r W = \sum_{i=1}^r W_i^2$$

ist $X^\top R X$ $\chi^2(r)$ -verteilt.

(ii) Wir setzen $Y := BX \sim \mathcal{N}(0, B^\top B)$ und $Z := RX \sim \mathcal{N}(0, R)$. Dann gilt

$$\text{Cov}(Y, Z) = B \text{Var}(X) R^\top = BR = 0.$$

Da (X, Y) als Lineartransformation von X gemeinsam normalverteilt ist, folgt aus der Unkorreliertheit bereits die Unabhängigkeit.

(iii) Genau wie in (ii) folgt die Unabhängigkeit von $Y := SX$ und $Z := RX$ und somit auch die Unabhängigkeit von $Y^\top Y = X^\top SX$ und $Z^\top Z = X^\top RX$. Zusammen mit (i) und dem vorangegangenen Lemma folgt die Behauptung. \square

Als Korollar erhalten wir Konfidenzbereiche für die Schätzung von β und linearen Funktionalen im gewöhnlichen linearen Modell unter der Normalverteilungsannahme.

Satz 2.23. *Im gewöhnlichen linearen Modell unter der Normalverteilungsannahme $(\varepsilon_i) \sim \mathcal{N}(0, \sigma^2 E_n)$ für $\sigma > 0$ gelten folgende Konfidenzaussagen für gegebenes Niveau $\alpha \in (0, 1)$:*

(i) *Ist $q_{F(k, n-k); 1-\alpha}$ das $(1-\alpha)$ -Quantil der $F(k, n-k)$ -Verteilung, so ist*

$$C := \{\beta \in \mathbb{R}^k \mid |X(\beta - \hat{\beta})|^2 < k\hat{\sigma}^2 q_{F(k, n-k); 1-\alpha}\}$$

ein Konfidenzellipsoid zum Konfidenzniveau $1-\alpha$ für β .

(ii) *Ist $q_{t(n-k); 1-\alpha/2}$ das $(1-\frac{\alpha}{2})$ -Quantil der $t(n-k)$ -Verteilung, so ist*

$$I := \left[\hat{\rho} - \hat{\sigma} \sqrt{v^\top (X^\top X)^{-1} v} q_{t(n-k); 1-\alpha/2}, \hat{\rho} + \hat{\sigma} \sqrt{v^\top (X^\top X)^{-1} v} q_{t(n-k); 1-\alpha/2} \right]$$

ein Konfidenzintervall zum Konfidenzniveau $1-\alpha$ für $\rho = \langle v, \beta \rangle$.

Beweis. (i) Nach Konstruktion gilt

$$X\hat{\beta} = XX^\dagger Y = \Pi_X Y = X\beta + \Pi_X \varepsilon, \quad \hat{\sigma}^2 = \frac{|(E_n - \Pi_X)\varepsilon|^2}{(n-k)}.$$

Da Π_X und $(E_n - \Pi_X)$ symmetrische, idempotente Matrizen mit Rang k bzw. $(n-k)$ sind (Projektionen auf $\text{ran } X$ bzw. $(\text{ran } X)^\perp$) und es gilt $(E_n + \Pi_X)\Pi_X = 0$, folgt aus Lemma 2.22:

$$\frac{|X(\beta - \hat{\beta})|^2}{k\hat{\sigma}^2} = \frac{(n-k)}{k} \frac{\varepsilon^\top (E_n - \Pi_X)\varepsilon}{\varepsilon^\top \Pi_X \varepsilon} \sim F(k, n-k).$$

Durch die Wahl des Quantiles folgt die Konfidenzaussage $\mathbb{P}_\beta(\beta \in C) = 1-\alpha$.

(ii) Wegen $\hat{\rho} \sim \mathcal{N}(\rho, \sigma^2 v^\top (X^\top X)^{-1} v)$ nach dem Satz von Gauß-Markov, ist

$$\frac{\rho - \hat{\rho}}{\sigma \sqrt{v^\top (X^\top X)^{-1} v}} \sim \mathcal{N}(0, 1).$$

Andererseits sind $\hat{\rho}$ und $\hat{\sigma}^2$ unabhängig und es gilt $\hat{\sigma}^2 = \sigma^2 Z / (n-k)$ für eine Zufallsvariable $Z \sim \chi^2(n-k)$. Damit ist

$$\frac{\rho - \hat{\rho}}{\sqrt{\hat{\sigma}^2 v^\top (X^\top X)^{-1} v}} \sim t(n-k). \quad \square$$

Bemerkung 2.24. Ebenso kann man ein Konfidenzintervall für die Varianz konstruieren (Übung \square).

Zusammen mit dem Korrespondenzsatz liefert dieses Resultat:

Methode 7: t-Test und F-Test. Im gewöhnlichen linearen Modell unter Normalverteilungsannahme $(\varepsilon_i) \sim \mathcal{N}(0, \sigma^2 E_n)$ ist der (zweiseitige) t-Test der Hypothese $H_0 : \rho = \rho_0$ gegen die Alternative $H_1 : \rho \neq \rho_0$ für $\rho_0 = \langle v, \beta_0 \rangle$ zum Niveau $\alpha \in (0, 1)$ gegeben durch

$$\varphi_{\rho_0}(Y) = \mathbb{1}_{\{|T_{n-k}(Y)| > q_{t(n-k); 1-\alpha/2}\}} \quad \text{mit} \quad T_{n-k}(Y) := \frac{\rho_0 - \widehat{\rho}}{\widehat{\sigma} \sqrt{v^\top (X^\top X)^{-1} v}}.$$

Der F-Test der Hypothese $H_0 : \beta = \beta_0$ vs. $H_1 : \beta \neq \beta_0$ zum Niveau $\alpha \in (0, 1)$ ist gegeben durch

$$\varphi_{\beta_0}(Y) = \mathbb{1}_{\{F_{k, n-k}(Y) > q_{F(k, n-k); 1-\alpha}\}} \quad \text{mit} \quad F_{k, n-k}(Y) := \frac{|X(\beta_0 - \widehat{\beta})|^2}{k \widehat{\sigma}^2}.$$

Schließlich wollen wir Hypothesentests noch für den allgemeineren Fall von linearen (bzw. affinen) Hypothesen konstruieren.

Definition 2.25. Im gewöhnlichen linearen Modell ist ein (zweiseitiges) lineares Testproblem gegeben durch

$$H_0 : K\beta = d \quad \text{versus} \quad H_1 : K\beta \neq d$$

für eine (deterministische) Matrix $K \in \mathbb{R}^{r \times k}$ mit vollem Rang $\text{rank}(K) = r \leq k$ und einem Vektor $d \in \mathbb{R}^r$. K wird Kontrastmatrix genannt. Unter der Hypothese H_0 sind also insgesamt $r \leq k$ linear unabhängige Bedingungen an die Parameter des linearen Modells gestellt.

Beispiel 2.26. Test auf Gleichheit zweier Regressionskoeffizienten: Für $2 \leq j < l \leq k$ ist das Testproblem gegeben durch

$$H_0 : \beta_j = \beta_l \quad \text{versus} \quad H_1 : \beta_j \neq \beta_l.$$

Damit ist die Kontrastmatrix $K = (a_{1,i}) \in \mathbb{R}^{1 \times k}$ gegeben durch $a_{1,i} = \mathbb{1}_{\{i=j\}} - \mathbb{1}_{\{i=l\}}$ und $d = 0$.

Weitere Beispiele sind der Globaltest (Übung \square):

$$H_0 : \forall j \in \{1, \dots, k\} : \beta_j = 0 \quad \text{versus} \quad H_1 : \exists j \in \{1, \dots, k\} : \beta_j \neq 0$$

sowie der Test eines Subvektors $\beta^* = (\beta_1^*, \dots, \beta_r^*)^\top$ mit $r \leq k$ (Übung \square):

$$H_0 : \forall j \in \{1, \dots, r\} : \beta_j = \beta_j^* \quad \text{versus} \quad H_1 : \exists j \in \{1, \dots, r\} : \beta_j \neq \beta_j^*.$$

Die Grundidee für das Testen linearer Hypothesen ist, die Residuen $RSS = |Y - X\widehat{\beta}|^2$ des Kleinste-Quadrat-Schätzers mit den Residuen des auf $H_0 : K\beta = d$ eingeschränkten Kleinste-Quadrat-Schätzers $\widehat{\beta}_{H_0}$, d.h.

$$RSS_{H_0} := |Y - X\widehat{\beta}_{H_0}|^2 \quad \text{mit} \quad |Y - X\widehat{\beta}_{H_0}|^2 = \min_{\beta \in \mathbb{R}^k : K\beta = d} |Y - X\beta|^2,$$

zu vergleichen. Ist die Abweichung (relativ zu RSS) zu groß, spricht dies gegen die Hypothese.

Satz 2.27. Im gewöhnlichen linearen Modell unter Normalverteilungsannahme $(\varepsilon_j) \sim \mathcal{N}(0, \sigma^2 E_n)$ ist die lineare Hypothese

$$H_0 : K\beta = d \quad \text{versus} \quad H_1 : K\beta \neq d$$

mit Kontrastmatrix $K \in \mathbb{R}^{r \times k}$ und $d \in \mathbb{R}^r$ zu testen. Es gilt

$$(i) \quad \widehat{\beta}_{H_0} = \widehat{\beta} - (X^\top X)^{-1} K^\top (K(X^\top X)^{-1} K^\top)^{-1} (K\widehat{\beta} - d),$$

$$(ii) \quad RSS_{H_0} - RSS = (K\widehat{\beta} - d)^\top (K(X^\top X)^{-1} K^\top)^{-1} (K\widehat{\beta} - d) \quad \text{und} \quad (RSS_{H_0} - RSS)/\sigma^2 \sim \chi^2(r) \quad \text{unter } H_0$$

$$(iii) \quad \text{die Fisher-Statistik } F := \frac{n-k}{r} \frac{RSS_{H_0} - RSS}{RSS} \text{ ist unter } H_0 \text{ gemäß } F(r, n-k) \text{ verteilt.}$$

Beweis. (i) Für jeden Vektor $\gamma \in \mathbb{R}^k$, der die Nebenbedingung $K\gamma=d$ erfüllt, gilt

$$|Y - X\gamma|^2 = |Y - X\hat{\beta} + X(\hat{\beta} - \gamma)|^2 = |Y - X\hat{\beta}|^2 + |X(\hat{\beta} - \gamma)|^2$$

nach Pythagoras, da $Y - X\hat{\beta} = (E_n - \Pi_X)Y \perp \text{ran}(X)$. Außerdem ist

$$|X(\hat{\beta} - \gamma)|^2 = |X(\hat{\beta} - \hat{\beta}_{H_0})|^2 + |X(\hat{\beta}_{H_0} - \gamma)|^2 + 2\langle X(\hat{\beta} - \hat{\beta}_{H_0}), X(\hat{\beta}_{H_0} - \gamma) \rangle.$$

Die Wahl von $\hat{\beta}_{H_0}$ impliziert jedoch

$$\begin{aligned} \langle X(\hat{\beta} - \hat{\beta}_{H_0}), X(\hat{\beta}_{H_0} - \gamma) \rangle &= ((X^\top X)^{-1}K^\top(K(X^\top X)^{-1}K^\top)^{-1}(K\hat{\beta} - d))^\top X^\top X(\hat{\beta}_{H_0} - \gamma) \\ &= (K\hat{\beta} - d)^\top (K(X^\top X)^{-1}K^\top)^{-1}(K\hat{\beta}_{H_0} - K\gamma) = 0, \end{aligned}$$

denn $\hat{\beta}_{H_0}$ erfüllt die Nebenbedingung:

$$K\hat{\beta}_{H_0} = K\hat{\beta} - K(X^\top X)^{-1}K^\top(K(X^\top X)^{-1}K^\top)^{-1}(K\hat{\beta} - d) = d.$$

Insgesamt erhalten wir also

$$|Y - X\gamma|^2 = |Y - X\hat{\beta}|^2 + |X(\hat{\beta} - \hat{\beta}_{H_0})|^2 + |X(\hat{\beta}_{H_0} - \gamma)|^2, \quad (2.1)$$

was offensichtlich für $\gamma = \hat{\beta}_{H_0}$ minimal ist.

(ii) Aus (2.1) mit $\gamma = \hat{\beta}_{H_0}$ folgt durch Einsetzen von $\hat{\beta}_{H_0}$

$$\begin{aligned} RSS_{H_0} - RSS &= |Y - X\hat{\beta}_{H_0}|^2 - |Y - X\hat{\beta}|^2 = |X(\hat{\beta} - \hat{\beta}_{H_0})|^2 \\ &= (\hat{\beta} - \hat{\beta}_{H_0})^\top X^\top X(\hat{\beta} - \hat{\beta}_{H_0}) \\ &= (K\hat{\beta} - d)^\top (K(X^\top X)^{-1}K^\top)^{-1}(K\hat{\beta} - d). \end{aligned}$$

Unter H_0 gilt für die Zufallsvariable $Z := K\hat{\beta}$, dass $\mathbb{E}[Z] = d$ und $\text{Var}(Z) = \sigma^2 K(X^\top X)^{-1}K^\top$. Aus der Normalverteilung von $\hat{\beta}$ folgt daher $(RSS_{H_0} - RSS)/\sigma^2 \sim \chi^2(r)$.

(iii) Da $RSS_{H_0} - RSS$ eine Funktion von $\hat{\beta}$ ist und somit unabhängig von RSS ist (Lemma 2.22), folgt die Verteilungsaussage für F aus der Charakterisierung der $F(r, n-p)$ -Verteilung. \square

Bemerkung 2.28. $W := rF$ heißt auch *Wald-Statistik*. Im Fall $d = 0$ ist $L := \{X\beta | \beta \in \mathbb{R}^k, K\beta = 0\}$ ein linearer Unterraum von $\text{ran } X$ und $X\hat{\beta}_{H_0} = \Pi_L Y$ die Orthogonalprojektion der Beobachtungen Y auf L . In diesem Fall gilt nach Pythagoras

$$RSS_{H_0} = |Y - \Pi_L Y|^2 = |Y - \Pi_X Y + (\Pi_X - \Pi_L)Y|^2 = |Y - \Pi_X Y|^2 + |X\hat{\beta} - X\hat{\beta}_{H_0}|^2,$$

so dass die Fisher-Statistik auch als

$$F = \frac{|X\hat{\beta} - X\hat{\beta}_{H_0}|^2}{r\hat{\sigma}^2}$$

geschrieben werden kann.

Beispiel 2.26 (fortgesetzt). Einsetzen von K und d liefert

$$F = \frac{n-k}{RSS} \frac{(\hat{\beta}_j - \hat{\beta}_l)^2}{K(X^\top X)^{-1}K^\top}.$$

Wegen $\text{Var}(\hat{\beta}_j - \hat{\beta}_l) = \text{Var}(K\hat{\beta}) = \sigma^2 K(X^\top X)^{-1}K^\top$ ist $\widehat{\text{Var}}(\hat{\beta}_j - \hat{\beta}_l) = K(X^\top X)^{-1}K^\top \hat{\sigma}^2$ mit $\hat{\sigma}^2 = RSS/(n-k)$ der natürliche (plug-in) Varianzschätzer. Damit können wir die Test-Statistik F als

$$F = \frac{(\hat{\beta}_j - \hat{\beta}_l)^2}{\widehat{\text{Var}}(\hat{\beta}_j - \hat{\beta}_l)} \stackrel{H_0}{\sim} F(1, n-k)$$

schreiben. Dieser F -Test ist äquivalent zum (zweiseitigen) t -Test mit der Teststatistik

$$T = \frac{\widehat{\beta}_j - \widehat{\beta}_l}{(\widehat{\text{Var}}(\widehat{\beta}_j - \widehat{\beta}_l))^{1/2}} \sim t(n - k).$$

Beispiel 2.29 (Klimaentwicklung). Wir folgen Beispiel 12.24 von Georgii (2007) und betrachten die mittleren Augusttemperaturen von 1799 bis 2008 in Karlsruhe (Quelle: <http://www.klimadiagramme.de/Europa/special101.htm>). Für die Jahre 1854 und 1945 liegen keine Daten vor, so dass wir $n = 208$ Beobachtungen haben. Eine polynomielle Regression in der Zeit t (in Jahrhunderten beginnend bei 1799) mit Graden $d = 1, \dots, 4$ liefert

$$\begin{aligned} p_1(t) &= 18,7 + 0,1t, \\ p_2(t) &= 20,0 - 3,5t + 1,7t^2, \\ p_3(t) &= 19,5 - 0,6t - 1,7t^2 + 1,1t^3, \\ p_4(t) &= 19,4 + 0,5t - 4,1t^2 + 2,9t^3 - 0,4t^4. \end{aligned}$$

Zunächst ist es plausibel, dass die zufälligen Schwankungen unabhängig von einander sind und als näherungsweise normalverteilt angenommen werden können (QQ-Plot). Um statistisch verwertbare Aussagen zu treffen, setzen wir noch das Niveau $\alpha = 0,05$ fest. Der Parametervektor ist $\beta = (\beta_0, \dots, \beta_d)^\top$. Welcher Grad des Regressionspolynoms ist sinnvoll?

Frage 1: Ist der positive Trend von p_1 signifikant? $H_0 : \beta_1 \leq 0$ vs. $H_1 : \beta_1 > 0$. Die zugehörige t -Statistik $T = \frac{\widehat{\beta}_1}{\widehat{\sigma} \sqrt{v^\top (X^\top X)^{-1} v}} \approx 0,62$ liegt deutlich unter dem kritischen Wert $q_{t(n-2), 1-\alpha} \approx 1,65$ (einseitiger T -Test), so dass die Hypothese nicht verworfen werden kann.

Frage 2: Liegt den Beobachtungen ein linearer Zusammenhang zugrunde (im Modell mit $d = 4$)? $H_0 : \beta_2 = \beta_3 = \beta_4 = 0$. Mittels Bemerkung 2.28 berechnen wir die Fisher-Statistik

$$F = \frac{\sum_{k=1}^n (p_4(t_k) - p_1(t_k))^2}{3\widehat{\sigma}^2} \approx 13,68 > 2,65 \approx q_{F(3, n-5), 1-\alpha}.$$

Folglich kann die Hypothese abgelehnt werden und wir schlussfolgern, dass eine Regressionsgerade unzureichend ist.

Frage 3: Benötigen wir ein Polynom vierten Grades? $H_0 : \beta_4 = 0$. Die zugehörige t -Statistik hat den Wert $-0,41$ dessen Absolutbetrag kleiner als das Quantil $q_{t(n-5), 0.975} \approx 1,97$ ist (zweiseitiger t -Test). Diese Nullhypothese kann also akzeptiert werden.

Frage 4: Benötigen wir ein Polynom dritten Grades? $H_0 : \beta_3 = 0$ (im Modell mit $d = 3$). Die zugehörige t -Statistik hat den Wert $2,05$ dessen Absolutbetrag größer als das Quantil $q_{t(n-4), 0.975} \approx 1,97$ ist. Die Hypothese kann also abgelehnt werden und der kubische Anteil im Regressionspolynom ist signifikant, d.h. p_3 ist signifikant besser geeignet die Beobachtungen zu beschreiben als p_2 .

p_3 zeigt einen deutlichen Anstieg der Temperaturen im 19. Jahrhundert. Es sei bemerkt, dass wir hier nur eine Zeitreihe betrachtet haben und somit nicht auf einen allgemeinen Zusammenhang schließen können (Aufgabe der Klimatologen).

2.3 Varianzanalyse

Beispiel 2.30. Um den Einfluss von $k \in \mathbb{N}$ verschiedenen Düngemitteln auf den Ernteertrag zu vergleichen wird jedes Düngemittel $i \in \{1, \dots, k\}$ auf n_i verschiedenen Agrarflächen ausgebracht. Der durch Witterungseinflüsse etc. zufällige Ernteertrag kann mittels $Y_{ij} = \mu_i + \varepsilon_{ij}$ für $j = 1, \dots, n_i$ und $i = 1, \dots, k$ modelliert werden, wobei μ_i der mittlere Ernteertrag von Düngemittel i ist und ε_{ij} unabhängige, zentrierte Störgrößen sind. Wir fragen uns also ob $\mu_1 = \dots = \mu_k$ gilt oder nicht.

Definition 2.31. Das Modell der einfaktoriellen Varianzanalyse (ANOVA1: (one-way) analysis of variance) ist gegeben durch Beobachtungen

$$Y_{ij} = \mu_i + \varepsilon_{ij}, \quad i = 1, \dots, k, j = 1, \dots, n_i,$$

mit iid.-verteilten Störgrößen $\varepsilon_{ij} \sim \mathcal{N}(0, \sigma^2)$. Wir bezeichnen die erste Dimension als den Faktor und den Wert $i = 1, \dots, k$ als die Faktorstufe. Folglich geben $(n_i)_{i=1, \dots, k}$ die Anzahl der unabhängigen Versuchswiederholungen pro Faktor an und $n := \sum_{i=1}^k n_i$ ist der Gesamtstichprobenumfang. Gilt $n_1 = \dots = n_k$, so sprechen wir von balanciertem Design.

Damit ist das ANOVA1-Modell ein Spezialfall des gewöhnlichen linearen Modells der Form

$$\mathbb{R}^n \ni Y := \begin{pmatrix} Y_{11} \\ \vdots \\ Y_{1n_1} \\ \vdots \\ Y_{k1} \\ \vdots \\ Y_{kn_k} \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & 0 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 1 & 0 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 1 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix}}_{=: X \in \mathbb{R}^{n \times k}} \cdot \underbrace{\begin{pmatrix} \mu_1 \\ \vdots \\ \mu_k \end{pmatrix}}_{=: \mu \in \mathbb{R}^k} + \begin{pmatrix} \varepsilon_{11} \\ \vdots \\ \varepsilon_{1n_1} \\ \vdots \\ \varepsilon_{k1} \\ \vdots \\ \varepsilon_{kn_k} \end{pmatrix}.$$

Beachte, dass $\text{rank } X = k$. Die klassische Fragestellung der Varianzanalyse lautet: “Existieren Unterschiede in den Faktorstufen-spezifischen Mittelwerten μ_i ?” oder anders formuliert “Hat der Faktor einen Einfluss auf die Response oder nicht?”. Dies führt auf das *Testproblem*

$$H_0 : \mu_1 = \dots = \mu_k \quad \text{versus} \quad H_1 : \exists i, l \in \{1, \dots, k\} : \mu_i \neq \mu_l.$$

Satz 2.32 (Streuungszerlegung). *Im ANOVA1-Modell definieren wir das i -te Gruppenmittel, $i = 1, \dots, k$, bzw. das Gesamtmittel als*

$$\bar{Y}_{i\bullet} := \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij} \quad \text{bzw.} \quad \bar{Y}_{\bullet\bullet} := \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij}$$

sowie

$$SSB := \sum_{i=1}^k n_i (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})^2 \quad \text{und} \quad SSW := \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet})^2$$

(*SSB*: sum of squares between groups; *SSW*: sum of squares within groups). Dann gilt

$$SST := \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{\bullet\bullet})^2 = SSB + SSW.$$

Beweis. Es gilt

$$\begin{aligned} SST &= \sum_i \sum_j (Y_{ij} - \bar{Y}_{\bullet\bullet})^2 = \sum_i \sum_j (Y_{ij} - \bar{Y}_{i\bullet} + \bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})^2 \\ &= \sum_i \sum_j ((Y_{ij} - \bar{Y}_{i\bullet})^2 + 2(Y_{ij} - \bar{Y}_{i\bullet})(\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet}) + (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})^2), \end{aligned}$$

wobei

$$\begin{aligned} \sum_i \sum_j (Y_{ij} - \bar{Y}_{i\bullet})(\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet}) &= \sum_i (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet}) \sum_j (Y_{ij} - \bar{Y}_{i\bullet}) \\ &= \sum_i (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})(n_i \bar{Y}_{i\bullet} - n_i \bar{Y}_{i\bullet}) = 0. \quad \square \end{aligned}$$

Offenbar spricht es gegen die Nullhypothese, wenn die Streuung zwischen den Gruppen größer ist als die Streuung innerhalb der Gruppen. Dies motiviert sowohl den Namen ANOVA als auch folgende Methode:

	Fg	Quadratsummen	Quadratmittel	F-Statistik
zwischen	$k - 1$	$SSB = \sum_{i=1}^k n_i (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})^2$	$SSB/(k - 1)$	$\frac{n - k}{k - 1} \frac{SSB}{SSW}$
innerhalb	$n - k$	$SSW = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet})^2$	$SSW/(n - k)$	
total	$n - 1$	$SST = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{\bullet\bullet})^2$	$SST/(n - 1)$	

Tabelle 1: ANOVA-Tafel

Methode 8: Einfaktorielle Varianzanalyse (ANOVA1). Im Modell der einfaktoriellen Varianzanalyse testen wir

$$H_0 : \mu_1 = \dots = \mu_k \quad \text{versus} \quad H_1 : \exists i, l \in \{1, \dots, k\} : \mu_i \neq \mu_l$$

zum Niveau $\alpha \in (0, 1)$ durch den F-Test

$$\varphi_\mu(Y) = \mathbb{1}_{\{F(Y) > q_{F(k-1, n-k); 1-\alpha}\}} \quad \text{mit} \quad F(Y) := \frac{n-k}{k-1} \frac{SSB}{SSW},$$

wobei $q_{F(k-1, n-k); 1-\alpha}$ das $(1 - \alpha)$ -Quantil der $F(k - 1, n - k)$ -Verteilung ist.

Satz 2.33. *Im einfaktoriellen Varianzanalysemodell gilt:*

(i) *Der Kleinste-Quadrate-Schätzer von $\mu = (\mu_1, \dots, \mu_k)^\top$ ist gegeben durch $\hat{\mu} = (\bar{Y}_{1\bullet}, \dots, \bar{Y}_{k\bullet})^\top$.*

(ii) *$SSW/\sigma^2 \sim \chi^2(n - k)$ und unter H_0 gilt $SSB/\sigma^2 \sim \chi^2(k - 1)$*

(iii) *SSW und SSB sind unabhängig und somit $F := \frac{n-k}{k-1} \frac{SSB}{SSW} \stackrel{H_0}{\sim} F(k - 1, n - k)$.*

Beweis. (i) Nachrechnen zeigt

$$\hat{\mu} = (X^\top X)^{-1} X^\top Y = \begin{pmatrix} 1/n_1 & & 0 \\ & \ddots & \\ 0 & & 1/n_k \end{pmatrix} \begin{pmatrix} \sum_{j=1}^{n_1} Y_{1j} \\ \vdots \\ \sum_{j=1}^{n_k} Y_{kj} \end{pmatrix} = \begin{pmatrix} \bar{Y}_{1\bullet} \\ \vdots \\ \bar{Y}_{k\bullet} \end{pmatrix}.$$

(ii)+(iii) Wegen $RSS = |Y - X\hat{\mu}|^2 = SSW$ folgt $SSW/\sigma^2 \sim \chi^2(n - k)$ und die Unabhängigkeit von SSW und $\hat{\mu}$ aus Lemma 2.22. Nach dem vorangegangenen Satz gilt weiterhin $SSB = SST - SSW$. Somit folgt die Behauptung aus Satz 2.27, falls $SST = RSS_{H_0}$. Nun gilt

$$RSS_{H_0} = \min_{\mu \in \mathbb{R}^k} \left| Y - X \underbrace{\begin{pmatrix} \mu \\ \vdots \\ \mu \end{pmatrix}}_{\in \mathbb{R}^k} \right|^2 = \min_{\mu \in \mathbb{R}^k} \left| Y - \underbrace{\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}}_{=: X_0 \in \mathbb{R}^{n \times 1}} \mu \right|^2.$$

Dieses Minimierungsproblem wird gelöst durch $\hat{\mu}_{H_0} = (X_0^\top X_0)^{-1} X_0^\top Y = n^{-1} \sum_{i,j} Y_{ij} = \bar{Y}_{\bullet\bullet}$. Damit folgt $RSS_{H_0} = SST$. \square

Bemerkung 2.34. In der Effektdarstellung wird das einfaktorielle Varianzanalysemodell als

$$Y_{ij} = \mu_0 + \alpha_i + \varepsilon_{ij}, \quad i = 1, \dots, k, j = 1, \dots, n_i,$$

geschrieben mit "Intercept" $\mu_0 := \frac{1}{n} \sum_{i=1}^k n_i \mu_i = \mathbb{E}[\bar{Y}_{\bullet\bullet}]$ und $\alpha_i := \mu_i - \mu_0$, den Effekt der Faktorstufe $i = 1, \dots, k$. Insbesondere muss in dieser Darstellung die Nebenbedingung $0 = \sum_{i=1}^k n_i \alpha_i$ oder äquivalent $n_k \alpha_k = -\sum_{i=1}^{k-1} n_i \alpha_i$ beachtet werden, damit die Designmatrix weiter vollen Rang hat. Der Parametervektor ist also gegeben durch $(\mu_0, \alpha_1, \dots, \alpha_{k-1})^\top$. Die F-Statistik um die Globalhypothese $H_0 : \alpha_1 = \dots = \alpha_{k-1} = 0$ zu überprüfen, ist identisch zur Statistik aus Satz 2.33.

Beispiel 2.35 (Zweistichproben t-Test). Soll die Gleichwertigkeit von bspw. zwei Düngemitteln getestet werden, ist $k = 2$ und das Testproblem $H_0 : \mu_1 = \mu_2$ versus $H_1 : \mu_1 \neq \mu_2$. Wegen $n\bar{Y}_{\bullet\bullet} = n_1\bar{Y}_{1\bullet} + n_2\bar{Y}_{2\bullet}$ gilt

$$\begin{aligned} SSB &= n_1(\bar{Y}_{1\bullet} - \bar{Y}_{\bullet\bullet})^2 + n_2(\bar{Y}_{2\bullet} - \bar{Y}_{\bullet\bullet})^2 \\ &= n_1\bar{Y}_{1\bullet}^2 + n_2\bar{Y}_{2\bullet}^2 + n\bar{Y}_{\bullet\bullet}^2 - 2(n_1\bar{Y}_{1\bullet} + n_2\bar{Y}_{2\bullet})\bar{Y}_{\bullet\bullet} \\ &= n_1\bar{Y}_{1\bullet}^2 + n_2\bar{Y}_{2\bullet}^2 - \frac{1}{n}(n_1\bar{Y}_{1\bullet} + n_2\bar{Y}_{2\bullet})^2 = \frac{n_1n_2}{n}(\bar{Y}_{1\bullet} - \bar{Y}_{2\bullet})^2. \end{aligned}$$

Somit ist

$$\varphi = \mathbb{1}_{\{|T| > q_{t(n-2), 1-\alpha/2}\}} \quad \text{mit} \quad T := \frac{\bar{Y}_{1\bullet} - \bar{Y}_{2\bullet}}{\sqrt{(\frac{1}{n_1} + \frac{1}{n_2})SSW/(n-2)}}$$

mit dem $(1 - \alpha/2)$ -Quantil der $t(n - 2)$ -Verteilung $q_{t(n-2), 1-\alpha/2}$ ein Test der Hypothese H_0 zum Niveau $\alpha \in (0, 1)$.

Definition 2.36. Das Modell der zweifaktoriellen Varianzanalyse mit balanciertem Design (ANOVA2) ist gegeben durch Beobachtungen

$$\begin{aligned} Y_{ijk} &= \mu_{ij} + \varepsilon_{ijk}, \quad i = 1, \dots, I, j = 1, \dots, J, k = 1, \dots, K \\ &= \mu_0 + \alpha_i + \beta_j + \gamma_{ij} + \varepsilon_{ijk} \end{aligned}$$

mit $I, J, K \geq 2$, iid.-verteilten Störgrößen $\varepsilon_{ijk} \sim \mathcal{N}(0, \sigma^2)$ und Nebenbedingungen (der Effektdarstellung)

$$\sum_{i=1}^I \alpha_i = \sum_{j=1}^J \beta_j = \sum_{i=1}^I \gamma_{ij} = \sum_{j=1}^J \gamma_{ij} = 0.$$

Wir haben also zwei Faktoren mit Faktorstufen $i = 1, \dots, I$ und $j = 1, \dots, J$. (α_i) bzw. (β_j) heißen Haupteffekte des ersten bzw. zweiten Faktors. (γ_{ij}) heißen Interaktions- bzw. Wechselwirkungseffekte.

Das ANOVA2-Modell ist also ein lineares Modell mit zwei kategoriellen Kovariablen. Die Gesamtanzahl an Beobachtungen ist gegeben durch $n = I \cdot J \cdot K$. Die typische Testprobleme sind

$$H_0 : \forall i : \alpha_i = 0 \quad \text{versus} \quad H_1 : \exists i \in \{1, \dots, I\} : \alpha_i \neq 0, \quad (2.2)$$

$$H_0 : \forall j : \beta_j = 0 \quad \text{versus} \quad H_1 : \exists j \in \{1, \dots, J\} : \beta_j \neq 0, \quad (2.3)$$

$$H_0 : \forall i, j : \gamma_{ij} = 0 \quad \text{versus} \quad H_1 : \exists i \in \{1, \dots, I\}, j \in \{1, \dots, J\} : \gamma_{ij} \neq 0. \quad (2.4)$$

Satz 2.37. Im zweifaktoriellen Varianzanalysemodell mit balanciertem Design gilt:

- (i) Die Kleinsten-Quadrate-Schätzer für μ_0, α_i, β_j und γ_{ij} , $i = 1, \dots, I - 1, j = 1, \dots, J - 1$, sind gegeben durch (\bullet heißt, dass über die jeweilige Koordinate gemittelt wird)

$$\begin{aligned} \hat{\mu}_0 &= \bar{Y}_{\bullet\bullet\bullet}, \quad \hat{\alpha}_i = \bar{Y}_{i\bullet\bullet} - \bar{Y}_{\bullet\bullet\bullet}, \quad \hat{\beta}_j = \bar{Y}_{\bullet j \bullet} - \bar{Y}_{\bullet\bullet\bullet}, \\ \hat{\gamma}_{ij} &= (\bar{Y}_{ij\bullet} - \bar{Y}_{\bullet\bullet\bullet}) - \hat{\alpha}_i - \hat{\beta}_j = \bar{Y}_{ij\bullet} - \bar{Y}_{i\bullet\bullet} - \bar{Y}_{\bullet j \bullet} + \bar{Y}_{\bullet\bullet\bullet}. \end{aligned}$$

- (ii) Definieren wir

$$\begin{aligned} SSW &:= \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (Y_{ijk} - \bar{Y}_{ij\bullet})^2, \\ SSB_1 &:= JK \sum_{i=1}^I (\bar{Y}_{i\bullet\bullet} - \bar{Y}_{\bullet\bullet\bullet})^2, \quad SSB_2 := IK \sum_{j=1}^J (\bar{Y}_{\bullet j \bullet} - \bar{Y}_{\bullet\bullet\bullet})^2, \\ SSB_{12} &:= K \sum_{i=1}^I \sum_{j=1}^J (\bar{Y}_{ij\bullet} - \bar{Y}_{i\bullet\bullet} - \bar{Y}_{\bullet j \bullet} + \bar{Y}_{\bullet\bullet\bullet})^2, \end{aligned}$$

dann können die Hypothesen (2.2), (2.3) bzw. (2.4) mit den F -Statistiken

$$\frac{IJ(K-1)}{I-1} \frac{SSB_1}{SSW} \sim F(I-1, IJ(K-1)), \quad \frac{IJ(K-1)}{J-1} \frac{SSB_2}{SSW} \sim F(J-1, IJ(K-1)) \quad \text{bzw.}$$

$$\frac{IJ(K-1)}{(I-1)(J-1)} \frac{SSB_{12}}{SSW} \sim F((I-1)(J-1), IJ(K-1))$$

getestet werden.

Beweis. Übung \square .

Bemerkung 2.38. Selbstverständlich erhält man analoge Resultate, wenn wir für jede "Zelle" $(i, j) \in \{1, \dots, I\} \times \{1, \dots, J\}$ verschiedene Stichprobenumfänge $n_{ij} \geq 2$ beobachten.

Beispiel 2.39. Ein Bauer möchte wissen ob die Größe seiner geernteten Kohlköpfe sich für zwei verschiedene Kultursorten unterscheidet. Auch der Pflanztag könnte eine Rolle spielen.

3 Exponentialfamilien and verallgemeinerte lineare Modelle

3.1 Die Informationsungleichung

Der Satz von Gauß-Markov hat uns bereits ein Optimalitätsresultat geliefert, dass allerdings auf lineare Schätzer im linearen Modell eingeschränkt ist. Wir suchen nun allgemeiner nach unverzerrten Schätzern deren Schätzwerte möglichst wenig um den korrekten Wert streuen.

Definition 3.1. Sei $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ ein statistisches Modell. Ein erwartungstreuer Schätzer T eines abgeleiteten Parameters $\rho(\vartheta)$ heißt varianzminimierend bzw. (gleichmäßig) bester Schätzer (UMVUE: uniformly minimum variance unbiased estimator), wenn für jeden weiteren erwartungstreuen Schätzer S gilt:

$$\text{Var}_\vartheta(T) \leq \text{Var}_\vartheta(S) \quad \text{für alle } \vartheta \in \Theta.$$

Wir werden zunächst eine untere Schranke für die Varianz beweisen und anschließend untersuchen, für welche Schätzer diese erreicht wird.

Definition 3.2. Ein vom Maß μ dominiertes, statistisches Modell $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ heißt regulär, wenn die folgenden Eigenschaften erfüllt sind:

- (i) Θ ist eine offene Menge in \mathbb{R}^d , $d \geq 1$.
- (ii) Die Likelihood-Funktion $L(\vartheta, x)$ ist auf $\Theta \times \mathcal{X}$ strikt positiv und nach ϑ stetig differenzierbar. Bezeichnen wir den Gradienten in ϑ mit $\nabla_\vartheta = (\frac{\partial}{\partial \vartheta_1}, \dots, \frac{\partial}{\partial \vartheta_d})^\top$, existiert insbesondere die Scorefunktion

$$U_\vartheta(x) := \nabla_\vartheta \log L(\vartheta, x) = \frac{\nabla_\vartheta L(\vartheta, x)}{L(\vartheta, x)}.$$

- (iii) Für jedes $\vartheta \in \Theta$ existiert die Fisher-Information

$$I(\vartheta) := \mathbb{E}_\vartheta \left[U_\vartheta(X) U_\vartheta(X)^\top \right]$$

und ist positiv definit.

- (iv) Es gilt die Vertauschungsrelation

$$\int h(x) \nabla_\vartheta L(\vartheta, x) \mu(dx) = \nabla_\vartheta \int h(x) L(\vartheta, x) \mu(dx) \quad (3.1)$$

für $h(x) = 1$.

Ein Schätzer $T: \mathcal{X} \rightarrow \mathbb{R}$ heißt regulär, falls $\mathbb{E}[|T(X)|^2] < \infty$ und (3.1) auch für $h(x) = T(x)$ gilt.

Bemerkung 3.3.

- (i) Der Satz von Lebesgue liefert eine hinreichende Bedingung für die Vertauschungsrelation (3.1): Sie gilt falls für jedes $\vartheta_0 \in \Theta$ eine Umgebung $V_{\vartheta_0} \subseteq \Theta$ existiert, so dass

$$\int_{\mathcal{X}} \sup_{\vartheta \in V_{\vartheta_0}} \left| \nabla_{\vartheta} L(\vartheta, x) \right| \mu(dx) < \infty.$$

Außerdem kann man (3.1) für jedes gegebene Modell (und jeden Schätzer) explizit nachprüfen.

- (ii) Als Konsequenz von (3.1) ergibt sich

$$\mathbb{E}_{\vartheta}[U_{\vartheta}] = \int \nabla_{\vartheta} L(\vartheta, x) \mu(dx) = \nabla_{\vartheta} \int L(\vartheta, x) \mu(dx) = \nabla_{\vartheta} 1 = 0$$

und damit $\text{Var}_{\vartheta}(U_{\vartheta}) = I(\vartheta)$.

- (iii) Ist $L(\vartheta, x)$ in ϑ zweimal stetig differenzierbar und gilt (3.1) mit $h(x) = 1$ und L ersetzt mit $\frac{\partial L}{\partial \vartheta_i}$ für alle $i \in \{1, \dots, d\}$, dann gilt $I(\vartheta) = -\mathbb{E}_{\vartheta}[H_{U_{\vartheta}(X)}(\vartheta)]$ für die Hesse-Matrix $H_{U_{\vartheta}(X)}$ der Scorefunktion $\vartheta \mapsto U_{\vartheta}(x)$ (Übung \square).
- (iv) Warum heißt $I(\vartheta)$ Information? Erstens: $I(\vartheta) = 0$ gilt auf einer Umgebung $\Theta_0 \subseteq \Theta$ genau dann, wenn $U_{\vartheta}(x) = 0$ für alle $\vartheta \in \Theta_0$ und μ -f.a. $x \in \mathcal{X}$, also wenn $L(\vartheta, x)$ für μ -f.s. konstant ist und somit keine Beobachtung die Parameter in Θ_0 unterscheiden kann (dieser Fall ist daher in der Definition ausgeschlossen). Zweitens, verhält sich die Fisher-Information bei unabhängigen Beobachtungen additiv: Ist $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_{\vartheta})_{\vartheta \in \Theta})$ ein reguläres Modell mit Fisher-Information I , so hat das Produktmodell $(\mathcal{X}^n, \mathcal{F}^{\otimes n}, (\mathbb{P}_{\vartheta}^{\otimes n})_{\vartheta \in \Theta})$ die Fisher-Information $I^{\otimes n} = nI$ (Beweis als Übung \square).

Satz 3.4 (Cramér-Rao-Ungleichung, Informationsschranke). *Gegeben seien ein reguläres statistisches Modell $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_{\vartheta})_{\vartheta \in \Theta})$, eine zu schätzende stetig differenzierbare Funktion $\rho: \Theta \rightarrow \mathbb{R}$, und ein regulärer erwartungstreuer Schätzer T von ρ . Dann gilt*

$$\text{Var}_{\vartheta}(T) \geq \left(\nabla \rho(\vartheta) \right)^{\top} I(\vartheta)^{-1} \nabla \rho(\vartheta) \quad \text{für alle } \vartheta \in \Theta. \quad (3.2)$$

Beweis. Aus der Zentriertheit von U_{ϑ} und der Regularität und Erwartungstreue von T erhalten wir

$$\begin{aligned} \text{Cov}_{\vartheta}(U_{\vartheta}, T) &= \mathbb{E}_{\vartheta}[TU_{\vartheta}] = \int_{\mathcal{X}} T(x) \nabla_{\vartheta} L(\vartheta, x) \mu(dx) \\ &= \nabla \int_{\mathcal{X}} T(x) L(\vartheta, x) \mu(dx) = \nabla \mathbb{E}_{\vartheta}[T] = \nabla \rho \end{aligned}$$

für alle $\vartheta \in \Theta$. Für jeden Vektor $e \in \mathbb{R}^d$ ergibt die Cauchy-Schwarz-Ungleichung somit

$$\langle e, \nabla \rho \rangle^2 = \text{Cov}_{\vartheta}(\langle e, U_{\vartheta} \rangle, T)^2 \leq \text{Var}_{\vartheta}(\langle e, U_{\vartheta} \rangle) \text{Var}_{\vartheta}(T) = \langle I(\vartheta)e, e \rangle \text{Var}_{\vartheta}(T),$$

also

$$\text{Var}_{\vartheta}(T) \geq \frac{\langle \nabla \rho, e \rangle^2}{\langle I(\vartheta)e, e \rangle}.$$

Maximieren über $e \in \mathbb{R}^d$ ergibt mit $e = I(\vartheta)^{-1} \nabla \rho(\vartheta)$ die Behauptung. \square

Definition 3.5. Ein regulärer erwartungstreuer Schätzer für den Gleichheit in (3.2) gilt, heißt Cramér-Rao-effizient.

Im Folgenden beschränken wir uns auf einparametrische ($d = 1$) Modelle.

Satz 3.6. *Unter den Bedingungen von Satz 3.4 mit $\Theta \subseteq \mathbb{R}$ erreicht der Schätzer T die untere Schranke für alle $\vartheta \in \Theta$ genau dann, wenn μ -f.ü. gilt*

$$T - \rho(\vartheta) = \rho'(\vartheta)I(\vartheta)^{-1}U_\vartheta \quad \text{für alle } \vartheta \in \Theta.$$

Falls $\rho' \neq 0$ ist dies äquivalent zu

$$L(\vartheta, x) = \exp\left(\eta(\vartheta)T(x) - \zeta(\vartheta)\right)c(x),$$

wobei $\eta: \Theta \rightarrow \mathbb{R}$ eine Stammfunktion von I/ρ' , $c: \mathcal{X} \rightarrow (0, \infty)$ messbar und $\zeta(\vartheta) = \log \int c(x) \exp(\eta(\vartheta)T(x))\mu(dx)$ eine Normierungsfunktion sind.

Beweis. Definieren wir $v(\vartheta) := \rho'(\vartheta)I^{-1}(\vartheta)$ (konstant in x) erhalten wir wegen $\text{Cov}_\vartheta(U_\vartheta, T) = \rho'(\vartheta)$

$$\begin{aligned} 0 &\leq \text{Var}_\vartheta(T - v(\vartheta)U_\vartheta) \\ &= \text{Var}_\vartheta(T) + v(\vartheta)^2 \text{Var}_\vartheta(U_\vartheta) - 2v(\vartheta) \text{Cov}_\vartheta(U_\vartheta, T) = \text{Var}_\vartheta(T) - \rho'(\vartheta)^2 I^{-1}(\vartheta), \end{aligned}$$

also wieder die Informationsungleichung. Gleichheit gilt genau dann, wenn $T - v(\vartheta)U_\vartheta$ \mathbb{P}_ϑ -f.s. konstant also gleich seinem Erwartungswert $\rho(\vartheta)$ ist. Da \mathbb{P}_ϑ eine strikt positive μ -Dichte hat gilt $\mu(T - \rho(\vartheta) \neq v(\vartheta)U_\vartheta) = 0$. Wenn dies nun für alle $\vartheta \in \Theta$ gilt, so folgt sogar

$$\mu(T - \rho(\vartheta) \neq v(\vartheta)U_\vartheta \text{ für ein } \vartheta \in \Theta) = 0,$$

denn aus Stetigkeitsgründen kann man sich auf rationale ϑ beschränken und die abzählbare Vereinigung von Nullmengen ist wieder eine Nullmenge. Die explizite Form der Likelihood-Funktion folgt durch unbestimmte Integration bzgl. ϑ . \square

Dieser Satz führt uns in natürlicher Weise auf eine wichtige Klasse von statistischen Modellen:

Definition 3.7. Es sei $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\vartheta)_{\vartheta \in \Theta})$ ein von μ dominiertes statistisches Modell mit $\Theta \subseteq \mathbb{R}$ offen. Dann heißt $(\mathbb{P}_\vartheta)_{\vartheta \in \Theta}$ (einparametrische) Exponentialfamilie in $\eta(\vartheta)$ und T , wenn messbare Funktionen $\eta: \Theta \rightarrow \mathbb{R}, T: \mathcal{X} \rightarrow \mathbb{R}$ und $c: \mathcal{X} \rightarrow (0, \infty)$ existieren, so dass

$$\frac{d\mathbb{P}_\vartheta}{d\mu}(x) = c(x) \exp(\eta(\vartheta)T(x) - \zeta(\vartheta)), \quad x \in \mathcal{X}, \vartheta \in \Theta,$$

wobei $\zeta(\vartheta) := \log \int c(x) \exp(\eta(\vartheta)T(x))\mu(dx)$. Dabei wird angenommen, dass T nicht μ -f.s. konstant ist. $\eta(\vartheta)$ heißt natürlicher Parameter der Exponentialfamilie und

$$\Xi := \left\{ \eta \in \mathbb{R} : \int_{\mathcal{X}} c(x) e^{\eta T(x)} \mu(dx) \in (0, \infty) \right\}$$

heißt natürlicher Parameterraum. Ist die Exponentialfamilie durch $\eta \in \Xi$ parametrisiert, dann wird sie als natürliche Exponentialfamilie bezeichnet.

Bemerkung 3.8.

(i) Die Darstellung ist nicht eindeutig, mit $a \neq 0$ erhält man beispielsweise eine Exponentialfamilie in $\tilde{\eta}(\vartheta) = a\eta(\vartheta)$ und $\tilde{T}(x) = T(x)/a$. Außerdem kann die Funktion c in das dominierenden Maß absorbiert werden: $\tilde{\mu}(dx) := c(x)\mu(dx)$.

(ii) Die Identifizierbarkeitsforderung $\mathbb{P}_\vartheta \neq \mathbb{P}_{\vartheta'}$ für alle $\vartheta \neq \vartheta'$ ist äquivalent zur Injektivität von η .

Beispiel 3.9.

(i) $(\mathcal{N}(\mu, \sigma^2))_{\mu \in \mathbb{R}}$ mit $\sigma > 0$ bekannt ist eine Exponentialfamilie in $\eta(\mu) = \mu/\sigma^2$ und $T(x) = x$:

$$L(\vartheta, x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/(2\sigma^2)} = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x^2-2\mu x+\mu^2)/(2\sigma^2)}, \quad x \in \mathbb{R}.$$

(ii) Die Familie der Poissonverteilungen $(\text{Poiss}(\lambda))_{\lambda > 0}$ mit Intensitätsparameter λ bildet eine Exponentialfamilie mit natürlichem Parameter $\eta(\lambda) = \log \lambda$ und $T(x) = x$:

$$L(\lambda, x) = \frac{\lambda^x e^{-\lambda}}{x!} = \frac{1}{x!} e^{x \log \lambda - \lambda}, \quad x \in \mathbb{Z}_+.$$

Lemma 3.10. *Ist ein statistisches Modell durch eine Exponentialfamilie in $\eta: \Theta \rightarrow \mathbb{R}$ und $T: \mathcal{X} \rightarrow \mathbb{R}$ mit differenzierbarem η gegeben, so ist dieses regulär. Ferner gilt*

(i) *Jede Statistik $S: \mathcal{X} \rightarrow \mathbb{R}$ mit existierendem Erwartungswert ist regulär. $\rho(\vartheta) := \mathbb{E}_\vartheta[T]$ ist stetig differenzierbar mit $\rho'(\vartheta) = \eta'(\vartheta) \text{Var}_\vartheta(T) \neq 0$, $\vartheta \in \Theta$.*

(ii) *Die Normierungsfunktion ζ ist auf $\Theta \subseteq \mathbb{R}$ stetig differenzierbar mit $\zeta'(\vartheta) = \eta'(\vartheta) \mathbb{E}_\vartheta[T]$ für $\vartheta \in \Theta$. Die Scorefunktion ist $U_\vartheta = \eta'(\vartheta)T - \zeta'(\vartheta)$.*

(iii) *Für die Fisher-Information gilt $I(\vartheta) = \eta'(\vartheta)\zeta''(\vartheta) = \eta'(\vartheta)\rho'(\vartheta)$ für alle $\vartheta \in \Theta$.*

Beweis. O.B.d.A. ist $\eta(\vartheta) = \vartheta$ und somit $\eta' = 1$ für alle $\vartheta \in \Theta$. Der allgemeine Fall ergibt sich durch Reparametrisierung und Anwendung der Kettenregel.

Schritt 1: Sei S eine beliebige reelle Statistik mit $S \in \mathcal{L}^1(\mathbb{P}_\vartheta)$ für alle $\vartheta \in \Theta$. Dann ist die Funktion

$$u_S(\vartheta) := e^{\zeta(\vartheta)} \mathbb{E}_\vartheta[S] = \int_{\mathcal{X}} S(x) e^{\vartheta T(x)} c(x) \mu(dx)$$

auf Θ wohl definiert. Wir zeigen nun, dass u_S beliebig oft differenzierbar ist.

Ist $\vartheta \in \Theta$ und $t \in \mathbb{R}$ so klein, dass auch $\vartheta \pm t \in \Theta$, so gilt mittels monotoner Konvergenz

$$\begin{aligned} \sum_{k \geq 0} \frac{|t|^k}{k!} \int_{\mathcal{X}} |S(x)| |T(x)|^k e^{\vartheta T(x)} c(x) \mu(dx) &= \int_{\mathcal{X}} |S(x)| e^{\vartheta T(x) + |t|T(x)} c(x) dx \\ &\leq \int_{\mathcal{X}} |S(x)| (e^{(\vartheta+t)T(x)} + e^{(\vartheta-t)T(x)}) c(x) dx < \infty. \end{aligned}$$

Also ist $ST^k \in \mathcal{L}^1(\mathbb{P}_\vartheta)$ für alle $\vartheta \in \Theta$ und insbesondere $T \in \mathcal{L}^2(\mathbb{P}_\vartheta)$ für alle ϑ . Ferner ist die Reihe

$$\sum_{k \geq 0} \frac{t^k}{k!} \int_{\mathcal{X}} S(x) T(x)^k e^{\vartheta T(x)} c(x) \mu(dx)$$

absolut konvergent und Summation und Integration können vertauscht werden. Die Reihe nimmt daher den Wert $u_S(\vartheta + t)$ an. Damit ist u_S sogar analytisch.

Schritt 2: Es folgt $u'_S(\vartheta) = e^{\zeta(\vartheta)} \mathbb{E}_\vartheta[ST]$ und insbesondere $u'_1(\vartheta) = u_1(\vartheta) \mathbb{E}_\vartheta[T]$ sowie $u''_1(\vartheta) = u_1(\vartheta) \mathbb{E}_\vartheta[T^2]$. Für $\zeta(\vartheta) = \log u_1(\vartheta)$ bekommen wir also $\zeta'(\vartheta) = \mathbb{E}_\vartheta[T] =: \rho(\vartheta)$ und

$$\rho'(\vartheta) = \zeta''(\vartheta) = u''_1(\vartheta)/u_1(\vartheta) - (u'_1(\vartheta)/u_1(\vartheta))^2 = \text{Var}_\vartheta(T).$$

Aus der Differenzierbarkeit von ζ folgt

$$U_\vartheta = \frac{\partial}{\partial \vartheta} \log L(\vartheta, x) = T - \zeta'(\vartheta), \quad \vartheta \in \Theta$$

und somit $I(\vartheta) = \text{Var}_\vartheta(U_\vartheta) = \text{Var}_\vartheta(T) > 0$. Weiter können wir schreiben

$$\begin{aligned} \frac{d}{d\vartheta} \mathbb{E}_\vartheta[S] &= (u_S(\vartheta) e^{-\zeta(\vartheta)})' = (u'_S(\vartheta) - u_S(\vartheta) \zeta'(\vartheta)) e^{-\zeta(\vartheta)} \\ &= \mathbb{E}_\vartheta[ST] - \mathbb{E}_\vartheta[S] \zeta'(\vartheta) = \mathbb{E}_\vartheta[SU_\vartheta] \\ &= \int_{\mathcal{X}} S(x) \frac{\partial}{\partial \vartheta} L(\vartheta, x) \mu(dx). \end{aligned}$$

Daher gilt einerseits (3.1) für alle $h \in \mathcal{L}^1(\mathbb{P}_\vartheta)$ und andererseits folgt die Regularität des Modells. \square

Korollar 3.11 (Existenz von besten Schätzern). *Für jedes statistische Modell gegeben durch eine Exponentialfamilie mit differenzierbarem η und $\eta' \neq 0$ ist die zugrunde liegende Statistik T ein bester und Cramér-Rao-effizienter Schätzer für $\rho(\vartheta) := \mathbb{E}_\vartheta[T] = \zeta'(\vartheta)/\eta'(\vartheta)$. In dem Fall gilt*

$$\text{Var}_\vartheta(T) = \rho'(\vartheta)/\eta'(\vartheta) \quad \text{und} \quad I(\vartheta) = \eta'(\vartheta)\rho'(\vartheta) \quad \text{für alle } \vartheta \in \Theta.$$

Für natürliche Exponentialfamilien gilt also insbesondere $\text{Var}_\eta(T) = I(\eta)$.

Beweis. Folgt unmittelbar aus Satz 3.4 und Lemma 3.10. Für natürliche Exponentialfamilien gilt also $\text{Var}_\eta(T) = \rho'(\eta) = I(\eta)$ und die Informationsschranke ist gegeben durch $\rho'(\eta)^2/I(\eta) = I(\eta)$. \square

Beispiel 3.9 (fortgesetzt).

(i) $(\mathcal{N}(\mu, \sigma^2))_{\mu \in \mathbb{R}}$ und bekanntem $\sigma > 0$ ist wie oben gesehen eine Exponentialfamilie in $\eta(\mu) = \mu/\sigma^2$, $T(x) = x$ und mit $\zeta(\mu) = \mu^2/(2\sigma^2)$. Somit ist $\rho(\mu) = \mathbb{E}_\mu[T] = \mu$ und $\text{Var}_\mu(T) = \sigma^2$. Da T nicht von $\sigma > 0$ abhängt, ist T sogar bester Schätzer für den Erwartungswert für alle Normalverteilungen.

(ii) Für die Exponentialfamilie $(\text{Poiss}(\lambda))_{\lambda > 0}$ in $\eta(\lambda) = \log \lambda$ und $T(x) = x$ gilt $\zeta(\lambda) = \lambda$. Wegen $\rho(\lambda) = \mathbb{E}_\lambda[T] = \lambda$ und $\text{Var}_\lambda(T) = \lambda$ ist T bester Schätzer für λ .

Lemma 3.12. *Ist $(\mathbb{P}_\vartheta)_{\vartheta \in \Theta}$ auf $(\mathcal{X}, \mathcal{F})$ eine Exponentialfamilie in $\eta: \Theta \rightarrow \mathbb{R}$ und $T: \mathcal{X} \rightarrow \mathbb{R}$ so ist $(\mathbb{P}_\vartheta^{\otimes n})_{\vartheta \in \Theta}$ eine Exponentialfamilie auf $(\mathcal{X}^n, \mathcal{F}^{\otimes n})$ mit zugrundeliegender Statistik $T_n = \frac{1}{n} \sum_{i=1}^n T \circ X_i$. Ist η differenzierbar mit $\eta' \neq 0$, folgt insbesondere, dass T_n ein bester Schätzer für $\rho(\vartheta) = \mathbb{E}_\vartheta[T]$ ist.*

Beweis. Übung \square .

Abschließend klären wir noch die Frage was das Maximum-Likelihood-Prinzip für natürliche Exponentialfamilien ergibt.

Lemma 3.13. *Ist $(\mathbb{P}_\vartheta)_{\vartheta \in \Theta}$ auf $(\mathcal{X}, \mathcal{F})$ eine natürliche Exponentialfamilie in $\eta \in \Xi$ und $T: \mathcal{X} \rightarrow \mathbb{R}$, dann ist T auf dem Ereignis $\{T(X) \in \text{ran}(\zeta')\}$ der eindeutige Maximum-Likelihood-Schätzer des Parameters $\rho(\eta) := \mathbb{E}_\eta[T]$. Ferner ist $\zeta': \Theta \rightarrow \mathbb{R}$ invertierbar und der eindeutige Maximum-Likelihood-Schätzer des natürlichen Parameters η ist gegeben durch*

$$\hat{\eta} = (\zeta')^{-1}(T(X)).$$

Beweis. Um die Maximalstelle der Likelihood-Funktion zu finden, setzen wir die Scorefunktion gleich null. Auf $\{T(X) \in \text{ran}(\zeta')\}$ gilt

$$\partial_\eta \log L(\eta, x) = U_\eta(x) = 0 \quad \Leftrightarrow \quad T(x) = \zeta'(\eta).$$

Da $\partial_\eta^2 \log L(\eta, x) = -\zeta''(\eta) = -\text{Var}_\eta(T) < 0$, ist $\eta \mapsto -\log L(\eta, x)$ konvex und somit T der eindeutige Maximum-Likelihood-Schätzer des Parameters $\rho(\eta) = \zeta'(\eta)$. Aus $\zeta'' > 0$ folgt außerdem, dass ζ' invertierbar ist, so dass der Maximum-Likelihood-Schätzer des natürlichen Parameters gegeben ist durch $(\zeta')^{-1} \circ T$. \square

3.2 Verallgemeinerte Lineare Modelle

Mit Hilfe von Exponentialfamilien wollen wir nun lineare Modelle verallgemeinern. Wie in Beispiel 3.9 gesehen bildet $(\mathcal{N}(\mu, \sigma^2))_{\mu \in \mathbb{R}}$ eine Exponentialfamilie mit natürlichem Parameter $\eta(\mu) = \mu/\sigma^2$ und Statistik $T(x) = x$, die ein effizienter Schätzer des Parameters $\rho(\mu) = \mathbb{E}_\mu[T] = \mu$ ist. Im gewöhnlichen linearen Modell sind nun die Beobachtungen gegeben durch

$$\mathbb{R}^n \ni Y = X\beta + \varepsilon,$$

mit Parametervektor $\beta \in \mathbb{R}^k$, Designmatrix $X \in \mathbb{R}^{n \times k}$ und iid. Fehlervariablen $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ mit Varianz $\sigma > 0$. Schreiben wir die Designmatrix als

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \quad \text{mit Zeilenvektoren } x_1, \dots, x_n \in \mathbb{R}^k,$$

ist Beobachtung Y_i gemäß $\mathcal{N}(x_i\beta, \sigma^2)$ verteilt, folgt also einer Exponentialfamilie mit $\eta_i(\beta) = x_i\beta/\sigma^2$ und $\rho_i(\beta) = x_i\beta$, $i = 1, \dots, n$. Lassen wir nun andere Exponentialfamilien zu, können wir sowohl Situationen modellieren in den der Zusammenhang zwischen $\mathbb{E}[Y_i]$ und den Kovariablen (codiert in der Designmatrix X) nichtlinear ist als auch diskrete Beobachtungen Y_i zulassen.

Definition 3.14. Auf einem Produktmodell $(\mathcal{X}^n, \mathcal{F}^{\otimes n})$ liegt ein verallgemeinertes lineares Modell (GLM: generalized linear model) mit n unabhängigen Beobachtungen Y_1, \dots, Y_n vor, falls die Randverteilungen von Y_i durch natürliche Exponentialfamilien gegeben sind mit Dichten

$$\frac{d\mathbb{P}_{\eta_i}^{Y_i}}{d\mu}(y_i) = \exp\left(\frac{\eta_i y_i - \zeta(\eta_i)}{\varphi}\right) c(y_i, \varphi), \quad i = 1, \dots, n,$$

bzgl. einem dominierenden Maß μ , mit unbekanntem Dispersionsparameter $\varphi > 0$,

$$\eta_i \in \Xi = \left\{ \eta \in \mathbb{R} : \int_{\mathcal{X}} e^{\eta y / \varphi} c(y, \varphi) \mu(dy) \in (0, \infty) \right\} \subseteq \mathbb{R}$$

für alle i und bekannten Funktionen $\zeta: \Xi \rightarrow \mathbb{R}$ und $c: \mathcal{X} \rightarrow \mathbb{R}_+$ mit $\zeta''(\eta) > 0$ für alle inneren Punkte $\eta \in \Xi^\circ$. Setze $\rho(\eta_i) := \mathbb{E}_{\eta_i}[Y_i]$. Für einen unbekanntem Parametervektor $\beta \in \mathbb{R}^k$, eine Designmatrix $X \in \mathbb{R}^{n \times k}$ und eine bijektive, stetig differenzierbare Funktion $g: \mathbb{R} \rightarrow \mathbb{R}$ gelte weiter

$$\begin{pmatrix} g(\rho(\eta_1)) \\ \vdots \\ g(\rho(\eta_n)) \end{pmatrix} = X\beta.$$

g heißt Linkfunktion. Falls $\rho = g^{-1}$, gilt $(\eta_1, \dots, \eta_n)^\top = X\beta$ und g heißt kanonische Linkfunktion (oder kanonischer Link).

Während β der interessierende Parameter ist, wird φ als Störparameter angesehen. Für fixiertes φ ist Y_i also gemäß einer natürlichen Exponentialfamilie in $T(y) = y/\varphi$ verteilt. Aus den Eigenschaften natürlicher Exponentialfamilien folgt

$$\mathbb{E}_{\beta, \varphi}[Y_i] = \zeta'(\eta_i) \quad \text{und} \quad \text{Var}_{\beta, \varphi}(Y_i) = \varphi \zeta''(\eta_i), \quad i = 1, \dots, n.$$

Beispiel 3.15. Das gewöhnliche lineare Modell ist ein GLM mit kanonischer Linkfunktion $g(x) = x$, $\zeta(\eta) = \eta^2/2$ und Dispersionsparameter $\varphi = \sigma^2$. Lassen wir allgemeinere Linkfunktionen zu erhalten wir *nicht-lineare Regressionsmodelle* (mit normalverteilten Fehlern) gegeben durch Beobachtungen $Y_i \sim \mathcal{N}(g^{-1}((X\beta)_i), \varphi)$.

Der Dispersionsparameter wird dazu verwendet eine Unterschätzung der (empirisch beobachteten) Varianz durch das Modell auszugleichen (siehe Übung \square).

Um den unbekanntem Parametervektor β in einem verallgemeinerten linearen Modell zu schätzen, verwenden wir den Maximum-Likelihood-Ansatz. Da ζ' streng monoton wachsend und die Linkfunktion g invertierbar sind, existiert die Funktion $\psi := (g \circ \rho)^{-1}$. Ist $x_i \in \mathbb{R}^k$ wieder die i -te Zeile von X , kann Loglikelihood-Funktion geschrieben werden als

$$\log L(\beta, \varphi; y) = \sum_{i=1}^n \left(\frac{\psi(x_i\beta)y_i - \zeta(\psi(x_i\beta))}{\varphi} + \log(c(y_i, \varphi)) \right).$$

Als notwendige Bedingung an einen Maximum-Likelihood-Schätzer $\hat{\beta}$ erhalten wir durch Ableiten

$$\nabla_{\beta} \log L(\hat{\beta}, \varphi; y) = \frac{1}{\varphi} \sum_{i=1}^n (y_i - \rho(\psi(x_i\hat{\beta}))) \psi'(x_i\hat{\beta}) x_i^\top = 0. \quad (3.3)$$

Lemma 3.16. *In einem verallgemeinerten linearen Modell mit kanonischer Linkfunktion ist die Fisher-Information gegeben durch*

$$I(\beta) = \frac{1}{\varphi} \sum_{i=1}^n \zeta''(x_i \beta) x_i^\top x_i \in \mathbb{R}^{k \times k}.$$

Ist $I(\beta)$ positiv definit für alle β und existiert eine Lösung $\hat{\beta}$ von (3.3), so ist $\hat{\beta}$ der eindeutige Maximum-Likelihood-Schätzer von β .

Beweis. Aus Lemma 3.10 folgt, dass die Fisher-Information im natürlichen Parameter $(\eta_1, \dots, \eta_n)^\top$ gegeben ist durch $\frac{1}{\varphi} \sum_{i=1}^n \zeta''(\eta_i)$. Die Reparametrisierung $\eta_i = x_i \beta$ zusammen mit der Kettenregel ergibt die Darstellung von $I(\beta)$.

Der kanonische Link ist gegeben durch $g = \rho^{-1}$, so dass ψ in (3.3) die Identität ist. Wegen $\rho = \zeta'$, gilt also

$$\frac{\partial^2 \log L(\beta, \varphi; y)}{\partial \beta \partial \beta^\top} = -\frac{1}{\varphi} \sum_{i=1}^n \zeta''(x_i \beta) x_i^\top x_i = -I(\beta).$$

Da $I(\beta) > 0$, ist $\beta \mapsto -\log L(\beta, \varphi; y)$ streng konvex und somit $\hat{\beta}$ der eindeutige Maximum-Likelihood-Schätzer. \square

Bemerkung 3.17.

- (i) Typischerweise besitzt $\hat{\beta}$ keine geschlossene Form mehr und muss durch numerische Verfahren bestimmt werden. *Fishers Scoring-Methode* verwendet hierfür das iterative Verfahren

$$\hat{\beta}^{(t+1)} = \hat{\beta}^{(t)} + I(\hat{\beta}^{(t)})^{-1} \nabla_{\beta} \log L(\hat{\beta}^{(t)}, \varphi; y), \quad t = 0, 1, \dots$$

(Beachte, dass sich der unbekannte Dispersionsparameter φ gerade rauskürzt). Für den kanonischen Link ist dieses Verfahren äquivalent zur *Newton-Raphson-Methode*.

- (ii) Ist g nicht der kanonische Link ist eine Lösung von (3.3) nicht notwendigerweise ein Maximum-Likelihood-Schätzer.

Zwei wichtige Beispielklassen für verallgemeinerte lineare Modelle sind die *Poisson-Regression* und die *logistische Regression*, die abschließend eingeführt werden.

Die Poisson-Regression modelliert unabhängige Poisson-verteilte Beobachtungen, deren Intensitätsparameter von Kovariablen abhängen. Sie eignet sich also für Beobachtungen die Zählstruktur haben. Wir hatten bereits gesehen dass die Familie $(\text{Poiss}(\lambda))_{\lambda > 0}$ eine Exponentialfamilie in $\eta(\lambda) = \log \lambda$ und $T(x) = x$ ist: Bezüglich des Zählmaßes ist die Likelihood-Funktion gegeben durch

$$L(\lambda, x) = \frac{\lambda^x e^{-\lambda}}{x!} = \frac{1}{x!} e^{x \log \lambda - \lambda}, \quad x \in \mathbb{Z}_+,$$

und es gilt $\rho(\lambda) = \mathbb{E}_\lambda[T] = \lambda$.

Definition 3.18. Ein verallgemeinertes lineares Modell auf $(\mathbb{Z}_+^n, \mathcal{P}(\mathbb{Z}_+^n))$ heißt Poisson-Regression, falls die unabhängigen Beobachtungen Y_i *Poiss* (λ_i) -verteilt sind, wobei $\lambda_i = e^{\eta_i}$ mit natürlichen Parametern $\eta_i \in \mathbb{R}, i = 1, \dots, n$, und

$$\eta = (\eta_1, \dots, \eta_n)^\top = X\beta$$

mit unbekanntem $\beta \in \mathbb{R}^k$ und Designmatrix $X \in \mathbb{R}^{n \times k}$.

Bemerkung 3.19. Wir verwenden hier also den kanonischen Link $g(\lambda) = \log \lambda$. In der Praxis wird oft das erweiterte Modell $Y_i \sim \text{Poiss}(\lambda_i \cdot s_i)$ verwendet für einen so genannten *Zählrahmen* $s_i > 0, i = 1, \dots, n$. Dann gilt $\mathbb{E}_\beta[Y_i] = \exp(x_i \beta + \log(s_i))$ mit den Zeilen x_i von X . Der Term $\log(s_i)$ wird als *Offset* bezeichnet, da er jeder Beobachtung einen individuellen "Intercept" zuweist.

Beispiel 3.20. In einem großen Krankenhaus wird die Anzahl der Beschwerden über $n = 44$ Notfallärzten untersucht (Daten aus Le (2003)). Der Zählrahmen pro Arzt ist die Anzahl an Patientenbesuchen, die vier zu berücksichtigenden Kovariablen sind Vergütung (in \$/h), Erfahrung (in h), Geschlecht und Facharztausbildung (ja/nein).

Lemma 3.21. Die Familie der Bernoulli-Verteilungen $(\text{Bernoulli}(p))_{p \in (0,1)}$ bildet eine Exponentialfamilie in $\eta(p) = \log \frac{p}{1-p}$ und $T(x) = x$.

Beweis. Die vom Zählmaß dominierte Familie besitzt die Likelihood-Funktion

$$L(p, x) = p^x (1-p)^{1-x} = (1-p) \left(\frac{p}{1-p} \right)^x = \exp \left(x \log \left(\frac{p}{1-p} \right) + \log(1-p) \right), \quad x \in \{0, 1\}. \quad \square$$

Definition 3.22. Ein verallgemeinertes lineares Modell auf $(\{0, 1\}^n, \mathcal{P}(\{0, 1\}^n))$ heißt logistische Regression, falls die unabhängigen Beobachtungen Y_i Bernoulli (p_i) -verteilt sind, $i = 1, \dots, n$, mit natürlichem Parameterraum \mathbb{R} , der kanonischen Link-Funktion $g: (0, 1) \rightarrow \mathbb{R}, g(p) = \log(p/(1-p))$ und

$$\eta = (g(p_1), \dots, g(p_n))^\top = X\beta$$

mit unbekanntem $\beta \in \mathbb{R}^k$ und Designmatrix $X \in \mathbb{R}^{n \times k}$. Die Funktion g heißt Logit-Funktion und ihre Umkehrfunktion $g^{-1}: \mathbb{R} \rightarrow (0, 1), g^{-1}(x) = (1 + e^{-x})^{-1}$ heißt logistische Funktion.

Bemerkung 3.23. Es gilt also $\mathbb{E}[Y_i] = g^{-1}(\eta_i) = e^{\eta_i} / (1 + e^{\eta_i})$, wobei die Funktion g^{-1} gerade die Verteilungsfunktion der standardisierten logistischen Verteilung ist (welche im Allgemeinen einen Mittelwerts- und einen Streuungsparameter besitzt). Das motiviert ein populäres Beispiel für eine nicht kanonische Linkfunktion: die Probit-Funktion $g(\lambda) = \Phi^{-1}(\lambda)$ mit der Verteilungsfunktion der Standardnormalverteilung Φ .

Da wir hier ein Modell gefunden haben um $\{0, 1\}$ -wertige Zufallsvariablen durch Kovariablen zu erklären, werden wir die logistische Regression im nächsten Kapitel zur Klassifikation verwenden.

3.3 Ergänzung: Numerische Bestimmung des Maximum-Likelihood-Schätzers

Das vermutlich grundlegendste numerische Verfahren zur Bestimmung von Nullstellen ist das Newton-Verfahren oder Newton-Raphson-Verfahren:

- Ziel: Finde $x^* \in \mathbb{R}: f(x^*) = 0$ für eine Funktion $f: \mathbb{R} \rightarrow \mathbb{R}$.
- Verfahren:
 - (i) Wähle einen Startpunkt $x_0 \in \mathbb{R}$ (der möglichst nahe an x^* liegen sollte).
 - (ii) Approximiere x^* mit der rekursiven Vorschrift

$$x_{n+1} := x_n - \frac{f(x_n)}{f'(x_n)} \quad \text{falls } f'(x_n) \neq 0$$

- Abbruchkriterien: $|f(x_n)| < \varepsilon$ oder $|x_{n+1} - x_n| < \varepsilon$ für ein $\varepsilon > 0$.

Geometrisch ist x_{n+1} genau die Nullstelle der Tangente $y = f(x_n) + f'(x_n)(x - x_n)$ an f im Punkt $(x_n, f(x_n))$. Im allgemeineren Fall $f: \mathbb{R}^k \rightarrow \mathbb{R}^k$ erhalten wir die Rekursionsvorschrift

$$J_f(x_n)(x_{n+1} - x_n) = -f(x_n) \quad \iff \quad x_{n+1} = x_n - J_f(x_n)^{-1} f(x_n)$$

mit der Jacobi-Matrix $J_f(x) = \left(\frac{\partial f_i}{\partial x_j} \right)_{i,j=1,\dots,k} \in \mathbb{R}^{k \times k}$ falls diese positiv definit ist.

Das Newton-Verfahren soll nun verwendet werden um den Maximum-Likelihood-Schätzer $\hat{\beta}$ in einem verallgemeinerten linearen Modell $(\mathcal{X}^n, \mathcal{F}^{\otimes n}, (\mathbb{P}_{\beta, \varphi}^{\otimes n})_{\beta \in \mathbb{R}^k, \varphi > 0})$ mit kanonischem Link zu bestimmen. Zur Erinnerung ist in diesem Fall die Likelihood-Funktion gegeben durch

$$L(\beta, \varphi; y) = \prod_{i=1}^n \exp \left(\frac{(x_i \beta) y_i - \zeta(x_i \beta)}{\varphi} \right) c(y_i, \varphi)$$

mit n Zeilenvektoren $x_i \in \mathbb{R}^k$. Setzen wir also

$$f(\beta) = \nabla_{\beta} \log L(\beta, \varphi; y) = \frac{1}{\varphi} \sum_{i=1}^n (y_i - \zeta'(x_i \beta)) x_i^{\top},$$

dann ist die Jacobi-Matrix gleich der Hesse-Matrix der Loglikelihood-Funktion

$$\begin{aligned} J_f(\beta) &= \left(\frac{\partial \beta_i \log L(\beta, \varphi; y)}{\partial \beta_j} \right)_{l,j=1,\dots,k} = -\frac{1}{\varphi} \left(\sum_{i=1}^n \zeta''(x_i \beta) x_{i,l} x_{i,j} \right)_{l,j=1,\dots,k} \\ &= -\frac{1}{\varphi} \sum_{i=1}^n \zeta''(x_i \beta) x_i^{\top} x_i. \end{aligned}$$

Da diese nicht mehr von y abhängt, erhalten wir

$$J_f(\beta) = \mathbb{E}_{\beta} [H_{\log L(\cdot, \varphi, Y)}(\beta)] = -I(\beta).$$

Einsetzen in obige Iterationsvorschrift liefert *Fishers Scoring-Methode*:

$$\begin{aligned} \widehat{\beta}^{(t+1)} &= \widehat{\beta}^{(t)} + I(\beta)^{-1} \nabla_{\beta} \log L(\widehat{\beta}^{(t)}, \varphi; y) \\ &= \widehat{\beta}^{(t)} - \left(\sum_{i=1}^n \zeta''(x_i \widehat{\beta}^{(t)}) x_i^{\top} x_i \right)^{-1} \sum_{i=1}^n (Y_i - \zeta'(x_i \widehat{\beta}^{(t)})) x_i^{\top}, \quad t \in \mathbb{N}, \end{aligned}$$

wobei wir $\widehat{\beta}^{(0)} = 0$ setzen.

4 Klassifikation

Während im linearen Modell die Zielvariable quantitativ ist, gibt es viele Situationen in denen die Daten *qualitativ* bzw. *kategorisch* sind. Das Grundprinzip der *Klassifikation* ist anhand einer sogenannten *Trainingsmenge* $(x_1, Y_1), \dots, (x_n, Y_n)$ zu lernen, die Klassen zu unterscheiden, um anschließend vorherzusagen, zu welcher Klasse Beobachtungen zu neuen x_{n+1}, \dots, x_{n+m} gehören (*statistisches Lernen*). Anders ausgedrückt, sollen x_{n+1}, \dots, x_{n+m} *klassifiziert* werden.

Beispiel 4.1. Auf Grundlage vom monatlichen Kontostand der Kreditkarte und dem Jahreseinkommen soll vorhergesagt werden ob jemand zahlungsunfähig wird oder nicht. Als Trainingsdatensatz haben wir Daten $(x_{i,1}, x_{i,2}, Y_i)$ für $i = 1, \dots, n$ Personen gegeben, wobei $x_{i,1}$ bzw. $x_{i,2}$ der monatliche Kontostand der Kreditkarte bzw. das Jahreseinkommen von Person i sind und Y_i die Frage ‘‘Zahlungsunfähig?’’ mit ‘‘Ja’’ oder ‘‘Nein’’ beantwortet (simulierter *default*-Datensatz aus James et al. (2013)). Etwa 3% der Personen sind zahlungsunfähig. Beachte, dass in realen Anwendungen die Beziehung zwischen Ko- und Zielvariablen typischerweise nicht so eindeutig sind.

4.1 Logistische Regression

Stammen die Zielvariablen nur aus zwei verschiedenen Klassen (die mit 0 und 1 codiert werden), bietet sich die logistische Regression aus Kapitel 3 als Modell an. Zur Erinnerung heißt ein statisches Experiment $(\{0, 1\}^n, \mathcal{P}(\{0, 1\}^n), (\mathbb{P}_{\beta}^{\otimes n})_{\beta \in \mathbb{R}^k})$ *multiple logistische Regression* mit Kovariablen $x_i = (1, x_{i,1}, \dots, x_{i,k-1}) \in \mathbb{R}^k$ (Zeilenvektor mit Absolutglied), falls Y_i *Bernoulli*(p_i)-verteilt ist, wobei $p_i = p(x_i, \beta)$ gegeben ist durch

$$\log \left(\frac{p_i}{1 - p_i} \right) = \beta_0 + \sum_{j=1}^{k-1} \beta_j x_{i,j} \quad \text{für } i = 1, \dots, n.$$

Äquivalent gilt

$$p(x_i, \beta) = \frac{e^{x_i \beta}}{1 + e^{x_i \beta}}.$$

Die Wahrscheinlichkeit, dass Y_i zur Klasse 1 gehört, wird also durch die $k - 1$ Kovariablen erklärt.

Methode 9: Klassifikation mittels logistischer Regression. Nach Schätzung des Parametervektors $\hat{\beta}$ auf der Trainingsmenge $(x_i, Y_i)_{i=1, \dots, n}$ können wir für eine jede neue Kovariablenrealisierung $x_{n+1} = (1, x_{n+1,1}, \dots, x_{n+1,k-1})$ (Zeilenvektor) einen zugehörigen Wert

$$\hat{p}_{n+1} = p(x_{n+1}, \hat{\beta}) = \frac{e^{x_{n+1}\hat{\beta}}}{1 + e^{x_{n+1}\hat{\beta}}}$$

“vorhersagen” und x_{n+1} der Klasse 1 zuordnen, falls $\hat{p}_{n+1} \geq \tau$ für einen Schwellenwert $\tau \in [0, 1]$. Andernfalls klassifizieren wir x_{n+1} mit 0.

Der Maximum-Likelihood-Ansatz führt auf die Maximierung der Loglikelihood-Funktion

$$\begin{aligned} \ell(\beta, y) &:= \log L(\beta, y) = \sum_{i=1}^n (y_i \log p(x_i, \beta) + (1 - y_i)(1 - p(x_i, \beta))) \\ &= \sum_{i=1}^n (y_i(x_i\beta) - \log(1 + e^{x_i\beta})). \end{aligned}$$

Nullsetzen des Gradienten führt auf k Gleichungen, die nichtlinear in β sind. Um den Maximum-Likelihood-Schätzer numerisch zu bestimmen, verwenden wir wieder das Newton-Verfahren. Dieses führt uns auf die *iterativ neugewichteten Kleinste-Quadrate-Methode* (IRLS: iteratively reweighted least squares):

Lemma 4.2. *In der logistischen Regression mit Designmatrix X von vollem Rang ist der $(t+1)$ ste Iterationsschritt von Fishers Scoring-Methode gegeben durch*

$$\hat{\beta}^{(t+1)} = (X^\top W_{\hat{\beta}^{(t)}} X)^{-1} X^\top W_{\hat{\beta}^{(t)}} Z_{\hat{\beta}^{(t)}} = \arg \min_b |W_{\hat{\beta}^{(t)}}^{1/2} (Z_{\hat{\beta}^{(t)}} - Xb)|^2$$

mit adjustiertem Responsevektor $Z_\beta = X\beta + W_\beta^{-1}(Y - p_\beta)$, wobei

$$\begin{aligned} p_\beta &= (p(x_1, \beta), \dots, p(x_n, \beta))^\top \in \mathbb{R}^n \quad \text{und} \\ W_\beta &= \text{diag}(p(x_1, \beta)(1 - p(x_1, \beta)), \dots, p(x_n, \beta)(1 - p(x_n, \beta))) \in \mathbb{R}^{n \times n}. \end{aligned}$$

Beweis. Für die logistische Funktion $g(x) = e^x/(1 + e^x)$ gilt $g'(x) = g(x)(1 - g(x))$. Aus der expliziten Form der Loglikelihood-Funktion $\ell(\beta)$ folgt damit, dass Scorefunktion und Hesse-Matrix gegeben sind durch

$$\nabla_\beta \ell(\beta, y) = X^\top (y - p_\beta) \quad \text{bzw.} \quad H_{\ell(\cdot, y)}(\beta) = -X^\top W_\beta X.$$

Somit ist der Iterationsschritt von Fishers Scoring-Methode

$$\begin{aligned} \hat{\beta}^{(t+1)} &= \hat{\beta}^{(t)} + (X^\top W_{\hat{\beta}^{(t)}} X)^{-1} X^\top (Y - p_{\hat{\beta}^{(t)}}) \\ &= (X^\top W_{\hat{\beta}^{(t)}} X)^{-1} X^\top W_{\hat{\beta}^{(t)}} (X\hat{\beta}^{(t)} + W^{-1}(Y - p_{\hat{\beta}^{(t)}})) \\ &= (X^\top W_{\hat{\beta}^{(t)}} X)^{-1} X^\top W_{\hat{\beta}^{(t)}} Z_{\hat{\beta}^{(t)}}. \end{aligned}$$

Wie in Kapitel 2 gesehen, ist dies gerade die Lösung des gewichteten Kleinste-Quadrate-Problems. \square

Bemerkung 4.3.

- (i) Der Maximum-Likelihood-Schätzer $\hat{\beta}$ erfüllt $\nabla_\beta \ell(\hat{\beta}, y) = \sum_{i=1}^n x_i^\top (y - p(x_i, \hat{\beta})) = 0$. Da der erste Koeffizient von x_i gleich 1 ist folgt $\sum_i y_i = \sum_i p(x_i, \hat{\beta})$, d.h. die erwartete Anzahl der Beobachtungen in Klasse eins stimmt mit der beobachteten Anzahl überein.

- (ii) Insbesondere zeigt dieses Lemma, dass der Maximum-Likelihood-Schätzer $\hat{\beta}$ die Lösung eines gewichteten Kleinste-Quadrate-Problems mit Responsevektor $Z_{\hat{\beta}} = X\hat{\beta} + W_{\hat{\beta}}^{-1}(Y - p_{\hat{\beta}})$ und Gewichten $w_i = \hat{p}_i(1 - \hat{p}_i)$, wobei beides wieder von $\hat{\beta}$ abhängt. Die gewichteten Quadratsummen der Residuen sind dann

$$\sum_{i=1}^n \frac{(Y_i - \hat{p}_i)^2}{\hat{p}_i(1 - \hat{p}_i)}$$

und messen die Abweichung der Daten von der Modellvorhersage.

Beispiel 4.4. Wir betrachten wieder den Datensatz aus Beispiel 4.1, wobei wir eine zusätzliche Kovariable “Student” mit Werten “Ja” oder “Nein” zur Verfügung haben. Eine Logistische Regression, die nur “Student” und einen Intercept verwendet führt zu einem positiven Koeffizienten der Dummy-Variable, d.h. die Zahlungsunfähigkeitswahrscheinlichkeit ist für Studenten höher als für Nicht-Studenten. Verwenden wir alle drei Kovariablen erhalten wir jedoch einen negativen Zusammenhang! Wie kann man diesen so genannten *Konfundierungseffekt* erklären?

Die logistische Regression kann auch auf mehr als zwei Klassen ausgeweitet werden, indem wir statt der Bernoulli-Verteilung die Multinomialverteilung verwenden. Häufig wird jedoch die Methode des nächsten Abschnittes dieser Variante vorgezogen. Inferenz für die logistische Regression beruht auf asymptotischen Überlegungen auf die wir in dieser Vorlesung nicht weiter eingehen werden.

4.2 Bayesklassifikation¹

Die logistische Regression modelliert die Wahrscheinlichkeit $\mathbb{P}(Y = 1)$ unter Verwendung des Regressorvektors x für zwei Klassen 0 und 1. Verstehen wir den Kovariablenvektor als Zufallsvariable X , wird also die bedingte Wahrscheinlichkeit $\mathbb{P}(Y = 1|X = x)$ der Klasse 1 gegeben einer Kovariablenrealisierung $X = x$ modelliert. Stattdessen wird nun ein Bayesianischer Ansatz verfolgt.

Gegeben sei das zufällige Paar (X, Y) , welches Werte in $\mathbb{R}^d \times \{1, \dots, K\}$ annimmt. Hierbei bezeichnet Y die *Klassifizierung* von X . Das heißt, dass die Verteilung \mathbb{P}_X von X durch die bedingte Verteilung $\mathbb{P}_{X|Y}$ festgelegt wird. Dieser Zusammenhang wird später durch die Bayesformel genauer erläutert. Wie in der logistischen Klassifizierung möchten wir einen *Klassifikator* \mathcal{C} konstruieren, der einer Realisierung $X = x$ eine Klasse $\mathcal{C}(x) \in \{1, \dots, K\}$ zuordnet. Formal haben wir es also mit einer *deterministischen* Abbildung

$$\mathcal{C}: \mathbb{R}^d \mapsto \{1, \dots, K\}$$

zu tun. In der Praxis wird \mathcal{C} normalerweise über ein *Trainingsample* $\mathcal{X}_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$ konstruiert, worauf wir später genauer eingehen werden. Eine wesentliche Frage bezüglich eines Klassifikators ist seine Qualität hinsichtlich einer korrekten Klassifizierung. Diese können wir über die Wahrscheinlichkeit einer fehlerhaften Klassifizierung $\mathbb{P}(\mathcal{C}(X) \neq Y)$ beschreiben, und bezeichnen sie allgemein mit

$$\mathcal{R}(\mathcal{C}) = \mathbb{P}(\mathcal{C}(X) \neq Y).$$

$\mathcal{R}(\mathcal{C})$ entspricht also dem 0-1-Risiko. Ein optimaler Klassifizierer \mathcal{C}^{opt} wäre daher (gegeben Messbarkeit)

$$\mathcal{C}^{\text{opt}} = \arg \min_{\mathcal{C}} \mathcal{R}(\mathcal{C}).$$

¹Vielen Dank an Moritz Jirak für die Ausarbeitung dieses Abschnitts

Wie können wir diesen konstruieren? Zunächst ist es günstig das Risiko $\mathcal{R}(\mathcal{C})$ umzuformen:

$$\begin{aligned}\mathcal{R}(\mathcal{C}) &= \mathbb{P}(\mathcal{C}(X) \neq Y) = \int \mathbb{P}(\mathcal{C}(x) \neq Y | X = x) \mathbb{P}_X(dx) \\ &= \int (1 - \mathbb{P}(\mathcal{C}(x) = Y | X = x)) \mathbb{P}_X(dx) \\ &= 1 - \int \mathbb{P}(\mathcal{C}(x) = Y | X = x) \mathbb{P}_X(dx).\end{aligned}$$

Wir sehen also, dass $\mathcal{R}(\mathcal{C})$ klein ist, wenn die bedingte Wahrscheinlichkeit $\mathbb{P}(\mathcal{C}(x) = Y | X = x)$ möglichst groß ist. Bedingen auf Y liefert weiter

$$\mathcal{R}(\mathcal{C}) = 1 - \int \left(\sum_{k=1}^K \mathbb{P}(\mathcal{C}(x) = k | Y = k, X = x) \mathbb{P}(Y = k | X = x) \right) \mathbb{P}_X(dx).$$

Nun benutzen wir die Tatsache, dass \mathcal{C} deterministisch ist (eine leichte Verallgemeinerung ist Unabhängigkeit von Y). Dadurch erhalten wir

$$\mathbb{P}(\mathcal{C}(x) = k | Y = k, X = x) = \mathbb{P}(\mathcal{C}(x) = k) = \mathbb{1}_{\{\mathcal{C}(x)=k\}},$$

und somit

$$\mathcal{R}(\mathcal{C}) = 1 - \int \left(\sum_{k=1}^K \mathbb{1}_{\{\mathcal{C}(x)=k\}} \mathbb{P}(Y = k | X = x) \right) \mathbb{P}_X(dx).$$

Wir haben es nun mit einer überraschend einfachen Optimierung zu tun. Um $\mathcal{R}(\mathcal{C})$ zu minimieren, genügt es den Ausdruck

$$\mathcal{A}(x) := \sum_{k=1}^K \alpha(k, x) \mathbb{P}(Y = k | X = x) \quad \text{mit} \quad \alpha(k, x) = \mathbb{1}_{\{\mathcal{C}(x)=k\}},$$

für jedes $x \in \mathbb{R}^d$ zu maximieren. Es gilt nun allerdings $\alpha(k, x) \in \{0, 1\}$. Die Größe $\mathcal{A}(x)$ ist folglich genau dann maximal, wenn wir das meiste Gewicht (und somit $\alpha(k, x) = 1$) auf

$$\max_{1 \leq k \leq K} \mathbb{P}(Y = k | X = x)$$

legen (dies kann leicht bewiesen werden). $\mathcal{A}(x)$ ist also genau dann maximal, wenn wir für jedes $x \in \mathbb{R}^d$

$$\alpha(k, x) = \begin{cases} 1, & \text{für } k = k^*, \\ 0, & \text{sonst,} \end{cases} \quad \text{wobei} \quad k^* = \arg \max_{k=1, \dots, K} \mathbb{P}(Y = k | X = x).$$

Dieses k^* liefert uns automatisch die optimale Klassifikation:

$$\mathcal{C}^{\text{opt}}(x) = \arg \max_{k=1, \dots, K} \mathbb{P}(Y = k | X = x). \quad (4.1)$$

Theorem 4.5. *Der deterministische Klassifikator, welcher das Risiko $\mathcal{R}(\mathcal{C})$ minimiert, ist durch die Klassifikation in (4.1) gegeben und wird Bayesklassifikator genannt.*

Bemerkung 4.6.

- (i) Die Optimalität der Bayesklassifikation hängt essentiell mit der Definition des Risikos $\mathcal{R}(\mathcal{C})$ zusammen.
- (ii) Falls $K = 2$, erhalten wir eine sehr einfache Klassifizierungsregel: Wenn $\mathbb{P}(Y = 1 | X = x) \geq 1/2$, dann wählen wir Klasse $k = 1$, ansonsten Klasse $k = 2$.

Viele Klassifikationsalgorithmen versuchen die Bayesklassifikation zu imitieren. Ein allgemeiner Zugang ist die bedingte Wahrscheinlichkeit $\mathbb{P}(Y = k | X = x)$ anhand eines Trainingsamples zu schätzen. Ein bekannter Repräsentant dieser Art ist das KNN-Verfahren (K-nearest neighbour, Übung \square). Ein anderes ist die lineare Diskriminanzanalyse, die im nächsten Abschnitt behandelt wird.

4.3 Lineare Diskriminanzanalyse²

Das Problem der Bayesklassifikation ist, dass es nicht so einfach ist, gute Schätzer für die bedingten Wahrscheinlichkeiten $\mathbb{P}(Y = k|X = x)$ zu konstruieren. Allerdings kann die Bayesformel hier helfen. Wir modellieren die Verteilung von X für jede Klasse $k \in \{1, \dots, K\}$ mit $K \geq 2$ (also gegeben Y) durch eine Dichte

$$f_k(x) = \mathbb{P}(X = dx|Y = k)$$

und wählen a-priori-Wahrscheinlichkeiten der Klassen $\pi_k = \mathbb{P}(Y = k) \in [0, 1]$ für $k = 1, \dots, K$ mit $\sum_k \pi_k = 1$. Die Bayesformel liefert dann die a-posteriori-Zähldichte von Y

$$p_k(x) = \mathbb{P}(Y = k|X = dx) = \frac{\mathbb{P}(X = dx|Y = k)\mathbb{P}(Y = k)}{\mathbb{P}(X = x)} = \frac{\pi_k f_k(x)}{\sum_{l=1}^K \pi_l f_l(x)}.$$

Die Idee der *linearen Diskriminanzanalyse (LDA)* ist nun, $f_k(x)$ als Gaußdichte (univariat) zu modellieren, also

$$f_k(x) = \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{(x - \mu_k)^2}{2\sigma_k^2}\right),$$

wobei μ_k und σ_k der Mittelwert und die Varianz der k -ten Klasse sind. Der Einfachheit halber sei $\sigma_1^2 = \dots = \sigma_K^2$ in der folgenden Diskussion. Dann erhalten wir

$$p_k(x) = \frac{\pi_k \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(-\frac{(x - \mu_k)^2}{2\sigma_k^2}\right)}{\sum_{l=1}^K \pi_l \frac{1}{\sqrt{2\pi}\sigma_l} \exp\left(-\frac{(x - \mu_l)^2}{2\sigma_l^2}\right)}. \quad (4.2)$$

Durch umformen erhalten wir, dass $p_k(x)$ genau dann maximal ist, wenn $\delta_k(x)$ maximal ist, gegeben durch

$$\delta_k(x) = x \frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k). \quad (4.3)$$

Tatsächlich sind die Werte π_k , μ_k und σ^2 aber unbekannt, und müssen geschätzt werden.

Methode 10: Lineare Diskriminanzanalyse. Wir definieren

$$\hat{\pi}_k = \frac{n_k}{n}, \quad \hat{\mu}_k = \frac{1}{|n_k|} \sum_{j: y_j=k} x_j \quad \text{und} \quad \hat{\sigma}^2 = \frac{1}{n - K} \sum_{k=1}^K \sum_{j: y_j=k} (x_j - \hat{\mu}_k)^2,$$

wobei n die Gesamtanzahl des Trainingsamples \mathcal{X}_n und n_k die Anzahl des Trainingsamples in der k -ten Klasse sind. Dann ist der Klassifizierer gegeben durch

$$\mathcal{C}(x) = \arg \max_{k=1, \dots, K} \hat{\delta}_k(x) \quad \text{mit} \quad \hat{\delta}_k(x) = x \frac{\hat{\mu}_k}{\hat{\sigma}^2} - \frac{\hat{\mu}_k^2}{2\hat{\sigma}^2} + \log(\hat{\pi}_k).$$

Im multivariaten Fall erhalten wir analog die Klassifizierungsregel

$$\delta_k(x) = x^\top \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^\top \Sigma^{-1} \mu_k + \log(\pi_k), \quad (4.4)$$

wobei Σ die d -dimensionale Kovarianzmatrix von $X \in \mathbb{R}^d$ ist, und $\mu_k \in \mathbb{R}^d$ der Vektor der komponentenweisen Erwartungswerte. Dabei können Σ , μ_1, \dots, μ_K sowie π_1, \dots, π_K wieder über Plug-in bzw. relative Häufigkeiten geschätzt werden.

Bemerkung 4.7. Eine weitere Verallgemeinerung stellt die *quadratische Diskriminanzanalyse (QDA)* dar, wo jede Klasse k eigene, im allgemeinen unterschiedliche Kovarianzmatrizen Σ_k besitzen. Dies führt zu einer quadratischen Klassifizierungsregel.

²Vielen Dank an Moritz Jirak für die Ausarbeitung dieses Abschnitts

Obwohl die Motivation für die logistische Klassifikation und LDA unterschiedlich ist, gibt es einen engen Zusammenhang. Betrachten wir den Fall $K = 2$. Dann gilt $p_2(x) = 1 - p_1(x)$ und eine kurze Rechnung ergibt für die LDA

$$\log\left(\frac{p_1(x)}{1 - p_1(x)}\right) = \log\left(\frac{p_1(x)}{p_2(x)}\right) = c_1 + c_2x,$$

wobei die Konstanten c_1, c_2 von μ_1, μ_2 und σ^2 abhängen. Im Fall der logistischen Klassifizierung haben wir:

$$\log\left(\frac{p_1(x)}{1 - p_1(x)}\right) = \log\left(\frac{p_1(x)}{p_2(x)}\right) = \beta_1 + \beta_2x.$$

Der Unterschied liegt also nur in der Art und Weise, wie die Konstanten geschätzt werden! In der Praxis führt das oft zu sehr ähnlichen Ergebnissen, aber nicht immer.

5 Ausblick

Im letzten Teil der Vorlesung werden (voraussichtlich) noch folgende Themen behandelt:

- (i) Modellwahl und statistisches Lernen
 - (a) Variablenselektion (C_p , AIC, BIC, R^2)
 - (b) Lasso
 - (c) Dimensionsreduktion / Hauptkomponentenanalyse
- (ii) Resampling
 - (a) Bootstrap
 - (b) Kreuzvalidierung

Literatur

- Agresti, A. and Finlay, B. (1997). *Statistical Methods for Social Sciences*. Prentice Hall.
- Fahrmeir, L., Kneib, T., and Lang, S. (2009). *Regression: Modelle, Methoden und Anwendungen*. Springer, Berlin.
- Georgii, H.-O. (2007). *Stochastik. Einführung in die Wahrscheinlichkeitstheorie und Statistik*. de Gruyter, Berlin.
- James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013). *An Introduction to Statistical Learning (with Applications in R)*. Springer, New York.
- Le, C. T. (2003). *Introductory biostatistics*. John Wiley & Sons.
- Witting, H. (1985). *Mathematische Statistik I*. Teubner.