

Local Case-Control Sampling

basierend auf Fithian und Hastie 2014

Nils Kober

29. Juni 2021

- 1 Einführung und Motivation
- 2 Local Case-Control Sampling im Detail
- 3 Konsistenz des LCC Schätzers
- 4 Asymptotische Verteilung
- 5 Zusammenfassung

- Für große Datensätze ist klassische logistische Regression (zu) rechenaufwendig.
- Ziel des Subsamplings: Verringerung des Rechenaufwands unter Beibehaltung einer möglichst hohen Effizienz des Schätzers (d.h. geringe Varianz)
- → Ermöglicht weitere statistische Verfahren (Bagging, Boosting, ...)

- Uniform sampling:
 - Jeder Datensatz hat gleiche Akzeptanzwahrscheinlichkeit.
 - Nicht optimiert für unausgewogene Datenmengen
- Case-Control (CC) sampling:
 - Akzeptanzwahrscheinlichkeit von (x, y) hängt nur von y ab.
 - Nutzt *marginal imbalance* aus, aber keine *conditional imbalance*.
 - Nicht konsistent, falls das Modell falsch spezifiziert ist.
- Local Case-Control (LCC) Sampling
 - Akzeptanzwahrscheinlichkeit hängt von x und y ab.
 - Behält besonders die Datensätze, bei denen das Label $Y = y$ gegeben die Features $X = x$ unerwartet ist.
 - Nutzt beide Arten der *imbalance* aus.
 - Auch konsistent, wenn das Modell falsch spezifiziert ist.

- Wahrscheinlichkeitsräume

$$\mathcal{X} \subset \mathbb{R}^p, \quad \mathcal{Y} = \{0, 1\} \quad (1)$$

mit gemeinsamem Wahrscheinlichkeitsmaß $\mathbb{P}^{\mathcal{X}, \mathcal{Y}}$

- iid Stichprobe

$$(x_1, y_1), \dots, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y} \quad (2)$$

- $\mathbb{P}^{\mathcal{X}}$ sei Lebesgue-stetig.

- Bedingte Wahrscheinlichkeit & log-odds:

$$p(x) = \mathbb{P}(Y = 1|X = x) \quad (3)$$

$$f(x) = \log \left(\frac{p(x)}{1 - p(x)} \right) \quad (4)$$

- Logistische Regressionsmodell (augmented notation):

$$f_{\theta}(x) = \theta^{\top} x \text{ mit } \theta \in \mathbb{R}^P \quad (5)$$

- Modell muss nicht korrekt spezifiziert sein (außer in Theorem 8).
- Zusätzliche Voraussetzung:

$$\nexists v \in \mathbb{R}^P : \mathbb{E} \left| v^{\top} X \right| = 0 \quad (6)$$

Local Case-Control Sampling (LCC)

Erstschätzung (pilot estimate):

$$\tilde{p}_\lambda(x) = \frac{\exp(\lambda^\top x)}{1 + \exp(\lambda^\top x)} \quad (7)$$

Akzeptanzwahrscheinlichkeit:

$$a_\lambda(x, y) = |y - \tilde{p}_\lambda(x)| = \begin{cases} 1 - \tilde{p}_\lambda(x), & y = 1 \\ \tilde{p}_\lambda(x), & y = 0 \end{cases} \quad (8)$$

$$= y(1 - \tilde{p}_\lambda(x)) + (1 - y)\tilde{p}_\lambda(x) \quad (9)$$

Algorithmus:

- 1 Generiere $z_i \sim \text{Ber}(a_\lambda(x_i, y_i))$
- 2 $S \leftarrow \{(x_i, y_i) | z_i = 1\}$
- 3 Verwende logistische Regression auf S und erhalte Schätzung $\hat{\theta}_S$
- 4 Passe Schätzung an: $\hat{\theta} \leftarrow \hat{\theta}_S + \lambda$

Sei $\mathbb{P}_\lambda^{X,Y}$ die Verteilung von X und Y auf dem Subsample S .

Dann ist $\mathbb{P}_\lambda^{X,Y}$ stetig bzgl. $\mathbb{P}^{X,Y}$ mit Dichte

$$\frac{d\mathbb{P}_\lambda^{X,Y}}{d\mathbb{P}^{X,Y}} = \frac{a_\lambda(x, y)}{\bar{a}(\lambda)}, \quad (10)$$

wobei $\bar{a}(\lambda) = \mathbb{E}[a_\lambda(X, Y)]$

Beweis: im Vortrag

Warum $\hat{\theta} \leftarrow \hat{\theta}_S + \lambda$?

Sei

$$g(x) = \log \frac{\mathbb{P}_\lambda(Y = 1|X = x)}{\mathbb{P}_\lambda(Y = 0|X = x)} \quad (11)$$

$$= \log \frac{\mathbb{P}(Y = 1|X = x, Z = 1)}{\mathbb{P}(Y = 0|X = x, Z = 1)} \quad (12)$$

$$f(x) = \log \frac{\mathbb{P}(Y = 1|X = x)}{\mathbb{P}(Y = 0|X = x)} \quad (13)$$

Dann gilt:

$$f(x) = g(x) + \lambda^\top x \quad (14)$$

Beweis: im Vortrag

Die Erstschätzung λ sollte konsistent sein, denn nur so ist $\hat{\theta}$ konsistent.
Möglichkeiten:

- Generierung aus den Daten (x_i, y_i) bzw. einer Teilmenge davon
→ z.B. mit weighted case-control sampling
- Verwendung der Schätzung zu einem früheren Zeitpunkt (bei Zeitreihen)

Verlustfunktion (negative log-likelihood):

$$\rho(\theta; x, y) = -y\theta^\top x + \log(1 + \exp(\theta^\top x)) \quad (15)$$

Risikofunktion für \mathbb{P} :

$$R(\theta) = \int \rho(\theta; x, y) d\mathbb{P}^{X,Y}(x, y) \quad \text{in gesamter Stichprobe} \quad (16)$$

Risikofunktion für \mathbb{P}_λ :

$$R_\lambda(\theta) = \int \rho(\theta; x, y) d\mathbb{P}_\lambda^{X,Y}(x, y) \quad \text{in Subsample } S \quad (17)$$

$$Q_\lambda(\theta) = R_\lambda(\theta - \lambda) \quad \text{in gesamter Stichprobe} \quad (18)$$

Erwartete Akzeptanzwahrscheinlichkeit:

$$\bar{a}(\lambda) = \mathbb{E} [a_\lambda(X, Y)] \quad (19)$$

Empirisches Risiko:

$$\hat{R}_\lambda^{(0)}(\theta) = \left(\sum_{i=1}^n z_i \right)^{-1} \sum_{i=1}^n z_i \rho(\theta; x_i, y_i) \quad \text{in Subsample } S \quad (20)$$

$$\hat{R}_\lambda(\theta) = \frac{1}{n\bar{a}(\lambda)} \sum_{i=1}^n z_i \rho(\theta; x_i, y_i) \quad \text{in Subsample } S \quad (21)$$

$$\hat{Q}_\lambda(\theta) = \hat{R}_\lambda(\theta - \lambda) \quad \text{in gesamter Stichprobe} \quad (22)$$

$$\hat{\theta} = \underset{\theta \in \mathbb{R}^p}{\operatorname{argmin}} \hat{Q}_\lambda(\theta) \quad (\text{LCC Schätzung}) \quad (23)$$

$$\bar{\theta}(\lambda) = \underset{\theta \in \mathbb{R}^p}{\operatorname{argmin}} Q_\lambda(\theta) \quad (24)$$

$$\theta^* = \underset{\theta \in \mathbb{R}^p}{\operatorname{argmin}} R(\theta) = \bar{\theta}(0) \quad (25)$$

Die folgenden Resultate gelten unter gewissen Voraussetzungen:

- Der LCC Schätzer $\hat{\theta}_n$ ist konsistent, falls die Erstschätzung λ_n konsistent ist.
- Der LCC Schätzer (genauer: $\sqrt{n}(\hat{\theta}_n - \bar{\theta}(\lambda_n))$) ist asymptotisch normalverteilt.
- Falls das logistische Regressionsmodell korrekt spezifiziert ist mit Parameter θ_0 , ist $\sqrt{n}(\hat{\theta}_n - \theta_0)$ asymptotisch normalverteilt mit der doppelten asymptotischen Kovarianzmatrix des MLE.

$R_\lambda(\theta)$ und $Q_\lambda(\theta)$ sind konvex in θ , da der Integrand $\rho(\theta; X, Y)$ bzw. $\rho(\theta - \lambda; X, Y)$ konvex ist (folgt mit Monotonie und Linearität des Integrals).

Für die strikte Konvexität genügt eine der folgenden Bedingungen (wird in Fithian und Hastie 2014 nicht vorausgesetzt):

- $\rho(x)$ hat höchstens abzählbar viele Unstetigkeitsstellen.
- $\mathbb{E} \|X\|^2 < \infty$

Definition

Die Verteilung der Daten / die Klassen (X, Y) heißt nicht-separabel, falls kein $v \in \mathbb{R}^P$ existiert, so dass

$$\mathbb{P}(Y = 0, v^T X > 0) = \mathbb{P}(Y = 1, v^T X < 0) = 0. \quad (26)$$

⇒ Es gibt keine Hyperebene, welche die Klassen $Y = 1$ und $Y = 0$ fast sicher trennt.

Lemma (Lemma 1 in Fithian und Hastie 2014)

Angenommen die Daten seien nicht-separabel. Dann hat $R_\lambda(\theta)$ für alle $\lambda \in \mathbb{R}^p$ einen eindeutigen (globalen) Minimierer.

Beweisidee:

- $R_\lambda(\theta)$ ist strikt konvex in θ , hat also maximal einen Minimierer.
- Für die Existenz genügt es zu zeigen, dass $R_\lambda(\theta) \rightarrow \infty$ für $\|\theta\| \rightarrow \infty$ (da $R_\lambda(\theta)$ konvex und \mathbb{R}^p σ -kompakt).

Proposition (Proposition 2 in Fithian und Hastie 2014)

Angenommen $\mathbb{E} \|X\| < \infty$ und die Daten seien nicht separabel. Sei $\theta^* = \operatorname{argmin}_{\theta \in \mathbb{R}^p} R(\theta) = \bar{\theta}(0)$. Dann gilt:

$$\theta^* = \operatorname{argmin}_{\theta \in \mathbb{R}^p} Q_{\theta^*}(\theta) = \bar{\theta}(\theta^*). \quad (27)$$

Beweis: im Vortrag

$$Q_\lambda(\theta) = \int \rho(\theta - \lambda; x, y) d\mathbb{P}_\lambda^{X, Y}(x, y) \quad (28)$$

$$\rho(\theta - \lambda; x, y) = y(\theta - \lambda)^\top x - \log(1 + \exp((\theta - \lambda)^\top x)) \quad (29)$$

Differentiationslemma siehe Klenke 2013, Satz 6.28.

Proposition (Proposition 3 in Fithian und Hastie 2014)

Sei $\mathbb{E} \|X\| < \infty$ und gelte $\lambda_n \xrightarrow{\mathbb{P}} \lambda_\infty$. Dann gilt für alle $\theta \in \mathbb{R}^p$

$$\hat{Q}_{\lambda_n}(\theta) \xrightarrow{\mathbb{P}} Q_{\lambda_\infty}(\theta). \quad (30)$$

Beweis: siehe Fithian und Hastie 2014

Proposition (Proposition 4 in Fithian und Hastie 2014)

Sei $\mathbb{E} \|X\| < \infty$ und gelte $\lambda_n \xrightarrow{\mathbb{P}} \lambda_\infty$. Weiter sei $\Theta \subset \mathbb{R}^p$ kompakt. Dann gilt

$$\sup_{\theta \in \Theta} \left| \widehat{Q}_{\lambda_n}(\theta) - Q_{\lambda_\infty}(\theta) \right| \xrightarrow{\mathbb{P}} 0. \quad (31)$$

Beweisidee:

- 1 Zeige, dass $F_n(\theta) = \widehat{Q}_{\lambda_n}(\theta) - Q_{\lambda_\infty}(\theta)$ mit Wahrscheinlichkeit, die gegen 1 konvergiert, gleichgradig Lipschitz-stetig ist (d.h. Lipschitz-Konstante unabhängig von n).
- 2 Mithilfe einer endlichen offenen Überdeckung von Θ kann gleichmäßige Konvergenz in Wahrscheinlichkeit gefolgert werden.

Alternativer Beweis: Satz von Arzelá-Ascoli (benötigt etwas Vorarbeit)

Theorem (Theorem 5 in Fithian und Hastie 2014)

Angenommen $\mathbb{E} \|X\| < \infty$ und die Daten seien nicht separabel. Weiter gelte $\lambda_n \xrightarrow{\mathbb{P}} \theta^*$. Dann ist der LCC-Schätzer konsistent, d.h.

$$\hat{\theta}_n \xrightarrow{\mathbb{P}} \theta^*. \quad (32)$$

Beweis: im Vortrag

$$\varepsilon := \inf_{\theta \in \partial\Theta} Q_{\theta^*}(\theta) - Q_{\theta^*}(\theta^*) > 0 \quad (\text{I})$$

$$\mathbb{P} \left(\sup_{\theta \in \Theta} \left| \widehat{Q}_{\lambda_n}(\theta) - Q_{\theta^*}(\theta^*) \right| \geq \frac{\varepsilon}{4} \right) \rightarrow 0 \quad \text{für } n \rightarrow \infty \quad (\text{II})$$

$$Q_{\lambda_n}(\mu_n) - \inf_{\theta \in \partial\Theta} \widehat{Q}_{\lambda_n}(\theta) < \frac{\varepsilon}{4} \quad (\text{III})$$

Ziel: $\sqrt{n}(\hat{\theta}_n - \bar{\theta}(\lambda_n))$ ist asymptotisch normalverteilt.

Voraussetzungen:

- Erstschätzungen $(\lambda_n)_{n \in \mathbb{N}}$ sind unabhängig von den Daten $((x_i, y_i)_{i \in \mathbb{N}})$
→ z.B. Erstschätzungen durch Daten von einem früheren Zeitpunkt oder einer getrennten Stichprobe
- $E \|X\|^2 < \infty$

$$J(\theta, \lambda) = \text{Var}_{\mathbb{P}_\lambda} [\nabla_\theta - \rho(\theta - \lambda; X, Y)] \quad (33)$$

$$H(\theta, \lambda) = -\nabla_\theta^2 Q_\lambda(\theta) \quad (34)$$

$$= -\bar{a}(\lambda)^{-1} \int \frac{e^{(\theta-\lambda)^\top x}}{(1 + e^{(\theta-\lambda)^\top x})^2} \left(\frac{e^{\lambda^\top x} + e^{f(x)}}{(1 + e^{\lambda^\top x})(1 + e^{f(x)})} \right) x x^\top d\mathbb{P}^X(x) \quad (35)$$

$H(\theta, \lambda)$ ist positiv definit, da per Voraussetzung $\mathbb{E} |v^\top X| \neq 0$ für alle $v \in \mathbb{R}^p$.

Theorem (Theorem 6 in Fithian und Hastie 2014)

Angenommen $\mathbb{E} \|X\|^2 < \infty$. Sei $(\lambda_n)_{n \in \mathbb{N}}$ unabhängig von $((x_i, y_i)_{i \in \mathbb{N}})$. Falls $\lambda_n \xrightarrow{\mathbb{P}} \theta^*$, dann gilt

$$\sqrt{n}(\hat{\theta}_n - \bar{\theta}(\lambda_n)) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \bar{a}(\theta^*)^{-1} \Sigma), \quad (36)$$

wobei $\Sigma = H(\theta^*, \theta^*)^{-1} J(\theta^*, \theta^*) H(\theta^*, \theta^*)^{-1}$.

Beweis: Siehe Fithian und Hastie 2014

Theorem (Theorem 8 in Fithian und Hastie 2014)

Angenommen das logistische Regressionsmodell sei korrekt mit Parameter $\theta_0 \in \mathbb{R}^p$, d.h. $f(x) = \log \frac{p(x)}{1-p(x)} = \theta_0^\top x$. Sei Σ_{full} die asymptotische Varianz des ML-Schätzers für die gesamte Stichprobe (genauer:

$\sqrt{n}(\beta_n - \theta_0) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \Sigma_{full})$ für MLE β_n). Gelte $\mathbb{E} \|X\|^2 < \infty$, sei $(\lambda_n)_{n \in \mathbb{N}}$ unabhängig von $((x_i, y_i)_{i \in \mathbb{N}})$. Falls $\lambda_n \xrightarrow{\mathbb{P}} \theta_0$, dann

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{\mathcal{D}} \mathcal{N}(0, 2\Sigma_{full}). \quad (37)$$

Beweis: im Vortrag

Theorem 6:

$$\sqrt{n}(\hat{\theta}_n - \bar{\theta}(\lambda_n)) \xrightarrow{\mathcal{D}} \bar{a}(\theta^*)^{-1} H(\theta^*, \theta^*)^{-1} J(\theta^*, \theta^*) H(\theta^*, \theta^*)^{-1} \quad (\text{I})$$

Unter Regularitätsvoraussetzungen (die hier erfüllt sind) ist ein MLE asymptotisch effizient (Czado und Schmidt 2011, 4.26. & 4.27.)

$$\Sigma_{full} = I(\theta_0)^{-1}. \quad (\text{II})$$

Es gilt für die Fisher-Information von \mathbb{P}_λ

$$I_\lambda(\theta) = -\nabla_\theta^2 Q_\lambda(\theta). \quad (\text{III})$$

$$J(\theta, \lambda) = \text{Var}_{\mathbb{P}_\lambda} [\nabla_\theta - \rho(\theta - \lambda; X, Y)] \quad (\text{IV})$$

$$\begin{aligned} H(\theta, \lambda) &= -\nabla_\theta^2 Q_\lambda(\theta) \\ &= -\bar{a}(\lambda)^{-1} \int \frac{e^{(\theta-\lambda)^\top x}}{(1 + e^{(\theta-\lambda)^\top x})^2} \\ &\quad \left(\frac{e^{\lambda^\top x} + e^{f(x)}}{(1 + e^{\lambda^\top x})(1 + e^{f(x)})} \right) \mathbf{x}\mathbf{x}^\top d\mathbb{P}^X(x) \end{aligned} \quad (\text{V})$$

$$H(\theta_0, 0) = \int \frac{\exp(\theta_0^\top x)}{(1 + \exp(\theta_0^\top x))^2} \mathbf{x}\mathbf{x}^\top d\mathbb{P}^X(x) \quad (\text{VI})$$

Die asymptotische Kovarianzmatrix in Theorem 8 ist unabhängig von der Akzeptanzwahrscheinlichkeit $a_\lambda(x, y)$ bzw. $\bar{a}(\lambda)$.

Erwartete Subsample-Größe:

$$n\bar{a}(\lambda) = n\mathbb{E}a_\lambda(X, Y) = n\mathbb{E}[Y(1 - \tilde{p}_\lambda(X)) + (1 - Y)\tilde{p}_\lambda(X)] \quad (38)$$

$$= n\mathbb{E}[p(X)(1 - \tilde{p}_\lambda(X)) + (1 - p(X))\tilde{p}_\lambda(X)] \quad (39)$$

→ Subsample klein bei starker *imbalance*

→ Günstiges Verhältnis aus Subsample-Größe und Varianz bei starker *imbalance*

- Der LCC Schätzer nutzt *marginal* und *conditional imbalance* in den Daten aus.
- Falls die Erstschätzung konsistent ist, so auch die LCC Schätzung (sogar wenn Modell nicht korrekt spezifiziert ist).
- Die LCC Schätzungen sind in einem gewissen Sinne asymptotisch normalverteilt.
- Im korrekt spezifizierten Fall ist die Effizienz des LCC Schätzers halb so groß wie die des MLE bei potentiell deutlich kleinerem Subsample und Rechenaufwand.

Danke für eure Aufmerksamkeit!



Czado, Claudia und Thorsten Schmidt (2011). *Mathematische Statistik*. Berlin, Heidelberg: Springer Berlin Heidelberg. ISBN: 978-3-642-17260-1 978-3-642-17261-8. DOI: 10.1007/978-3-642-17261-8. URL: <http://link.springer.com/10.1007/978-3-642-17261-8> (besucht am 21.06.2021).



Fithian, William und Trevor Hastie (Okt. 2014). „Local case-control sampling: Efficient subsampling in imbalanced data sets“. In: *The Annals of Statistics* 42.5, S. 1693–1724. ISSN: 0090-5364. DOI: 10.1214/14-AOS1220. URL: <http://projecteuclid.org/euclid.aos/1410440622>.



Klenke, Achim (2013). *Wahrscheinlichkeitstheorie*. Berlin, Heidelberg: Springer Berlin Heidelberg. ISBN: 978-3-642-36017-6 978-3-642-36018-3. DOI: 10.1007/978-3-642-36018-3. URL: <http://link.springer.com/10.1007/978-3-642-36018-3> (besucht am 27.06.2021).