

## Dichte von $P_x^{x,y}$

geg: • Mit  $\pi_{x,y}(x,y)$  bzw.  $\pi_{x,y}^\lambda(x,y)$  bezeichne Dichte von  $P^{x,y}$  bzw.  $P_x^{x,y}$   
bzgl.  $\lambda \otimes \nu$ , wobei  $\lambda$  Lebesgue-Maß,  $\nu$  Zählmaß

$$\text{Z: } \frac{dP_x^{x,y}}{dP^{x,y}} = \frac{a_x(x,y)}{\bar{a}(x)}$$

Beweis:

$$\pi_{x,y}^\lambda(x,y) = \pi_{x,y|z=1}(x,y) \quad = a_x(x,y)$$

Bayes  $\rightarrow$

$$= \frac{P(z=1 | x=x, y=y)}{P(z=1)} \cdot \pi_{x,y}(x,y)$$
$$= \int a_x(x,y) dP^{x,y}(x,y)$$
$$= \bar{a}(x)$$

$$= \frac{a_x(x,y)}{\bar{a}(x)} \cdot \pi_{x,y}(x,y)$$

$$\Rightarrow \frac{dP_x^{x,y}}{dP^{x,y}} = \frac{a_x(x,y)}{\bar{a}(x)}$$



## Offset bei LCC:

geg:  $\cdot g(x)$  bedingte log-odds Funktion von  $P_\lambda$   
 $\cdot f(x)$  bedingte log-odds Funktion von  $P$  } gegeben  $X=x$

$$\exists: g(x) = f(x) - \lambda^T x$$

Beweis:

$$g(x) = \log \frac{P_\lambda(Y=1|X=x)}{P_\lambda(Y=0|X=x)}$$

$$= \log \frac{\frac{\pi_{x,Y}^\lambda(1,x)}{\pi_x^\lambda(x)}}{\frac{\pi_{x,Y}^\lambda(0,x)}{\pi_x^\lambda(x)}}$$

$$= \log \frac{\pi_{x,Y}^\lambda(1,x) \cdot \frac{a_\lambda(1,x)}{\bar{a}(x)}}{\pi_{x,Y}^\lambda(0,x) \cdot \frac{a_\lambda(0,x)}{\bar{a}(x)}}$$

$$= \log \frac{\pi_{x,Y}^\lambda(1,x)}{\pi_{x,Y}^\lambda(0,x)} + \log \frac{a_\lambda(1,x)}{a_\lambda(0,x)}$$
$$= \frac{P(Y=1|X=x)}{P(Y=0|X=x)} = \frac{p(x)}{1-p(x)}$$
$$= \frac{1 - \tilde{p}_\lambda(x)}{\tilde{p}_\lambda(x)}$$

$$= f(x) - \lambda^T x$$



## Proposition 2

- geg:  $\cdot \mathbb{E}\|x\| < \infty$   
 $\cdot$  Daten nicht separabel  
 $\cdot \theta^* = \bar{\theta}(\theta) = \underset{\theta \in \mathbb{R}^p}{\operatorname{argmin}} R(\theta)$

$\hat{?}: \theta^* = \bar{\theta}(\theta^*) = \underset{\theta \in \mathbb{R}^p}{\operatorname{argmin}} Q_{\theta^*}(\theta)$  ← gilt per Def.

Beweis:

Da  $Q_{\theta^*}(\theta)$  strikt konvex ist, gilt

$$\theta = \bar{\theta}(\gamma) \Leftrightarrow \nabla_{\theta} Q_{\gamma}(\theta) = 0.$$

Betrachte also

$$\nabla_{\theta} Q_{\lambda}(\theta) = \int \nabla_{\theta} p(\theta - \lambda; x, y) dP_{\lambda}^{x, y}(x, y)$$

$$= \int yx - x \cdot \frac{\exp((\theta - \lambda)^T x)}{1 + \exp((\theta - \lambda)^T x)} dP_{\lambda}^{x, y}(x, y)$$

Für  $\theta = \lambda = \theta^*$  gilt:

$$\nabla_{\theta} Q_{\theta^*}(\theta^*) = \int yx - x \cdot \frac{\exp(\underbrace{(\theta^* - \theta^*)^T x}_{=0})}{1 + \exp(\underbrace{(\theta^* - \theta^*)^T x}_{=0})} dP_{\theta^*}^{x, y}(x, y)$$

$$= \int (y - \frac{1}{2})x \cdot \underbrace{\frac{\alpha_{\theta^*}(x, y)}{\bar{\alpha}(\theta^*)}}_{= \frac{1}{2}} dP^{x, y}(x, y)$$

$$= \frac{1}{\bar{\alpha}(\theta^*)} \mathbb{E} \left[ (y - \frac{1}{2})x \left( y(1 - p^*(x)) + (1 - y)p^*(x) \right) \right]$$

$$= \frac{1}{\bar{\alpha}(\theta^*)} \mathbb{E} \left[ x \left( \underbrace{y^2}_{=y} (1 - p^*(x)) + \underbrace{y(1-y)}_{=0} p^*(x) \right) \right]$$

$$- \frac{1}{2} y(1 - p^*(x)) - \frac{1}{2} (1 - y)p^*(x) \Big] = -\frac{1}{2} \alpha_{\theta^*}(x, y)$$

$$= \frac{1}{\bar{\alpha}(\theta^*)} \mathbb{E} \left[ \frac{1}{2} x y (1 - p^*(x)) - \frac{1}{2} x (1 - y) p^*(x) \right]$$

$$= \frac{1}{2} \cdot \frac{1}{\bar{\alpha}(\theta^*)} \mathbb{E} \left[ x \left( \underbrace{p(x)(1 - p^*(x)) - (1 - p(x))p^*(x)}_{= p(x) - p^*(x)} \right) \right]$$

$$= \frac{1}{2} \frac{1}{\bar{\alpha}(\theta^*)} \mathbb{E} \left[ x(y - p^*(x)) \right] = 0$$

$= \nabla_{\theta} R(\theta^*) = 0$





Mit der speziellen Wahl  $\Theta = \overline{B_\delta(\theta^*)}$  für  $\delta > 0$  bel., folgt

$$P(|\hat{\Theta}_n - \theta^*| > \delta) = P(\hat{\Theta}_n \notin \Theta) \xrightarrow{n \rightarrow \infty} 0,$$

also  $\hat{\Theta}_n \xrightarrow{P} \theta^*$ .



## Theorem 8

- geg: • Logistisches Regressionsmodell korrekt spezifiziert mit wahren Parameter  $\theta_0$
- $\Sigma_{\text{full}}$ : asymptotische Kovarianzmatrix des MLE auf dem gesamten Datensatz
  - $E\|X\|^2 < \infty$
  - $\lambda_n \xrightarrow{P} \theta_0$  unabh. von den Daten  $(x_i, y_i)_{i \in n}$

$$\tilde{z}: \sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} \mathcal{N}(0, 2\Sigma_{\text{full}})$$

Beweis:

Da das Modell für  $P$  korrekt spezifiziert ist, ist es das auch für  $P_\lambda$ .

(Denn die Parameter unterscheiden sich unter um  $\lambda$ .)

Also gilt  $\bar{\theta}(\lambda) = \underset{\theta \in \mathbb{R}^p}{\text{argmin}} Q_\lambda(\theta) = \theta_0$ .  $\rightarrow$  Da  $\nabla_\theta Q_\lambda(\theta) = 0$  für  $\theta = \theta_0$ .

Da das Modell korrekt spezifiziert ist, gilt  $J(\theta_0, \lambda) = H(\theta_0, \lambda)$ .

Nach Theorem 6 ergibt sich also:

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \rightarrow \mathcal{N}(0, \underbrace{\bar{\alpha}(\theta_0)^{-1}}_{=: \Sigma_{\text{loc}}} \underbrace{H(\theta_0, \theta_0)^{-1}})$$

Bemerke zunächst:

$$\begin{aligned} H(\theta_0, \theta) &= \underbrace{\bar{\alpha}(\theta)^{-1}}_{=2} \int \frac{\exp(\theta_0^T x)}{(1 + \exp(\theta_0^T x))^2} \cdot \frac{\exp(\theta^T x) + \exp(\theta_0^T x)}{(1 + \exp(\theta^T x))(1 + \exp(\theta_0^T x))} x x^T dP^*(x) \\ &= \int \frac{\exp(\theta_0^T x)}{(1 + \exp(\theta_0^T x))^2} x x^T dP^*(x) \end{aligned}$$

Für  $\theta = \theta_0$ ,  $\lambda = \theta_0$ ,  $f(x) = \theta_0^T x$  vereinfacht sich  $H(\theta, \lambda)$  zu:

$$\begin{aligned} H(\theta_0, \theta_0) &= \bar{\alpha}(\theta_0)^{-1} \int \frac{\exp(\theta_0^T x)}{(1 + \exp(\theta_0^T x))^2} \cdot \frac{\exp(\theta_0^T x) + \exp(\theta_0^T x)}{(1 + \exp(\theta_0^T x))(1 + \exp(\theta_0^T x))} x x^T dP^*(x) \\ &= \frac{1}{4} \int \frac{\exp(\theta_0^T x)}{(1 + \exp(\theta_0^T x))^2} x x^T dP^*(x) \end{aligned}$$

$$= \bar{\alpha}(\theta_0)^{-1} \frac{1}{2} \int \frac{\exp(\theta_0^T x)}{(1 + \exp(\theta_0^T x))^2} x x^T dP^*(x)$$

$$\stackrel{(I)}{=} \bar{\alpha}(\theta_0)^{-1} \frac{1}{2} \underbrace{H(\theta_0, 0)}_{= \sum_{full}^{-1}}$$

denn  $Q_0(\theta_0) = R_0(\theta_0) = R(\theta_0)$

Disitribut für normales  
logistic regression

$$= (2 \bar{\alpha}(\theta_0) \sum_{full})^{-1}$$

$$\Rightarrow \sum_{lcc} = \bar{\alpha}(\theta_0)^{-1} \cdot H(\theta_0, \theta_0)^{-1} = \bar{\alpha}(\theta_0)^{-1} \cdot 2 \bar{\alpha}(\theta_0) \sum_{full} = 2 \sum_{full}$$



$$\xi: \mathcal{J}(\theta_0, \lambda) = H(\theta_0, \lambda)$$

Für eine bel. 2-fach diffbare Funktion  $\varphi: \mathbb{R}^p \rightarrow \mathbb{R}$  gilt:

$$\frac{\partial}{\partial \theta_i} \log \varphi(\theta) \stackrel{\text{Kettenregel}}{=} \frac{1}{\varphi(\theta)} \cdot \frac{\partial}{\partial \theta_i} \varphi(\theta)$$

$$\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log \varphi(\theta) = \frac{\partial}{\partial \theta_i} \left[ \frac{1}{\varphi(\theta)} \cdot \frac{\partial}{\partial \theta_j} \varphi(\theta) \right]$$

$$= \frac{1}{\varphi(\theta)} \cdot \frac{\partial}{\partial \theta_i \partial \theta_j} \varphi(\theta) - \frac{1}{\varphi(\theta)^2} \cdot \frac{\partial}{\partial \theta_i} \varphi(\theta) \cdot \frac{\partial}{\partial \theta_j} \varphi(\theta)$$

$$\Rightarrow \nabla_{\theta}^2 \log \varphi(\theta) = \frac{1}{\varphi(\theta)} \cdot \nabla_{\theta}^2 \varphi(\theta) - \frac{1}{\varphi(\theta)^2} \cdot \nabla_{\theta} \varphi(\theta) \cdot (\nabla_{\theta} \varphi(\theta))^T$$

Betrachte nun  $p_{x,y}(\theta) = \gamma p_{\theta}(x) + (1-\gamma)(1-p_{\theta}(x))$

(Bedingte Dichte von  $Y$  geg.  $X=x$ )

$\Rightarrow$  obiges gilt für festes  $x, y$

$$\Rightarrow -\rho(\theta; x, y) = \log p_{x,y}(\theta) \quad (\text{log-likelihood Funktion})$$

Im Folgenden seien alle Erwartungswerte nur über  $X$  und  $Y$  gebildet, nicht über  $\lambda$ .

Wenn nun  $E_{\lambda} \left[ \frac{1}{p_{x,y}(\theta_0 - \lambda)} \nabla_{\theta}^2 p_{x,y}(\theta_0 - \lambda) \right] = 0$  ( $E[\dots]$  bzgl.  $P_{\lambda}^{x,y}$ ), folgt:

$$H(\theta_0, \lambda) = -\nabla_{\theta}^2 Q_{\lambda}(\theta_0, \lambda)$$

$$= \nabla_{\theta}^2 E_{\lambda} [-\rho(\theta_0 - \lambda; x, y)]$$

$$= E_{\lambda} [\nabla_{\theta}^2 \log p_{x,y}(\theta_0 - \lambda)]$$

$$= E_{\lambda} \left[ \underbrace{\frac{1}{p_{x,y}(\theta_0 - \lambda)} \cdot \nabla_{\theta}^2 p_{x,y}(\theta_0 - \lambda)}_{=0} + E_{\lambda} \left[ \frac{1}{p_{x,y}(\theta_0 - \lambda)^2} \nabla_{\theta} p_{x,y}(\theta_0 - \lambda) (\nabla_{\theta} p_{x,y}(\theta_0 - \lambda))^T \right] \right]$$

Hier muss eigentlich ein  $-$  statt dem  $+$  stehen  
 $\hookrightarrow H(\theta, \lambda)$  um Faktor  $-1$  falsch definiert?



$$= \mathbb{E}_\lambda \left[ \left( \nabla_\theta \log \varphi_{x,y}(\theta_0 - \lambda) / \left( \nabla_\theta \log \varphi_{x,y}(\theta_0 - \lambda) \right)' \right) \right]$$

$$= \text{Var}_\lambda \left[ \nabla_\theta \log \varphi_{x,y}(\theta_0 - \lambda) \right]$$

$$= \text{Var}_\lambda \left[ \nabla_\theta (-\rho(\theta_0 - \lambda; x, y)) \right]$$

$$= \mathcal{J}(\theta, \lambda)$$

Also bleibt zu zeigen:  $\mathbb{E}_\lambda \left[ \frac{1}{\varphi_{x,y}(\theta_0 - \lambda)} \nabla_\theta^2 \varphi_{x,y}(\theta_0 - \lambda) \right] = 0$

Mit Turmeigenschaft gilt:

$$\mathbb{E}_\lambda \left[ \frac{1}{\varphi_{x,y}(\theta_0 - \lambda)} \nabla_\theta^2 \varphi_{x,y}(\theta_0 - \lambda) \right] = \mathbb{E} \left[ \underbrace{\mathbb{E}_\lambda \left[ \frac{1}{\varphi_{x,y}(\theta_0 - \lambda)} \nabla_\theta^2 \varphi_{x,y}(\theta_0 - \lambda) \mid X \right]}_{(*)} \right]$$

Es genügt zu zeigen, dass der bedingte Erwartungswert (\*) gleich 0 ist.

$$\mathbb{E}_\lambda \left[ \frac{1}{\varphi_{x,y}(\theta_0 - \lambda)} \nabla_\theta^2 \varphi_{x,y}(\theta_0 - \lambda) \mid X \right] = \sum_{y=0}^1 \frac{1}{\varphi_{x,y}(\theta_0 - \lambda)} \nabla_\theta^2 \varphi_{x,y}(\theta_0 - \lambda) \cdot \varphi_{x,y}(\theta_0 - \lambda)$$

$\varphi_{x,y}(\theta_0 - \lambda)$  ist bedingte  
Zahldichte von  $Y$  geg.  $X$   
auf Subsample

$$= \nabla_\theta^2 \underbrace{\sum_{y=0}^1 \varphi_{x,y}(\theta_0 - \lambda)}_{=1 \text{ (da Dichte)}}$$

$$= 0$$

z:  $J(\theta_0, \lambda) = I_\lambda(\theta_0)$  (Fisher-Information)

Beweis:

Per Def. gilt:

$$\begin{aligned} J(\theta_0, \lambda) &= \text{Var}_\lambda(\nabla_{\theta_0} - \rho(\theta_0 - \lambda; x, Y)) \\ &= \mathbb{E}_\lambda[(\nabla_{\theta_0} - \rho(\theta_0 - \lambda; x, Y))(\nabla_{\theta_0} - \rho(\theta_0 - \lambda; x, Y))^T] - \underbrace{(\mathbb{E}_\lambda[\nabla_{\theta_0} - \rho(\theta_0 - \lambda; x, Y)])^2}_{=0} \\ &= \mathbb{E}_\lambda[(\nabla_{\theta_0} \log \varphi_{x, Y}(\theta_0 - \lambda))(\nabla_{\theta_0} \log \varphi_{x, Y}(\theta_0 - \lambda))^T] \\ &= I_\lambda(\theta_0) \end{aligned}$$

(da  $\varphi_{x, Y}(\theta_0 - \lambda)$  bedingte Zehldichte von  $Y$  geg.  $X$  ist und Multiplikation mit Dichte von  $X$  fällt weg, da  $\log(a \cdot b) = \log(a) + \log(b)$  und Dichte von  $X$  nicht von  $\theta$  abhängt.)

Also betrachte:

$$\begin{aligned} \mathbb{E}_\lambda[\nabla_{\theta_0} - \rho(\theta_0 - \lambda; x, Y) | X] &= \mathbb{E}_\lambda[\nabla_{\theta_0} \log \varphi_{x, Y}(\theta_0 - \lambda)] \\ &= \mathbb{E}_\lambda\left[\frac{1}{\varphi_{x, Y}(\theta_0 - \lambda)} \cdot \nabla_{\theta_0} \varphi_{x, Y}(\theta_0 - \lambda)\right] \\ &= \sum_{y=0}^1 \frac{1}{\varphi_{x, Y}(\theta_0 - \lambda)} \cdot \nabla_{\theta_0} \varphi_{x, Y}(\theta_0 - \lambda) \cdot \varphi_{x, Y}(\theta_0 - \lambda) \\ &= \nabla_{\theta_0} \underbrace{\sum_{y=0}^1 \varphi_{x, Y}(\theta_0 - \lambda)}_{=1} \\ &= 0 \end{aligned}$$