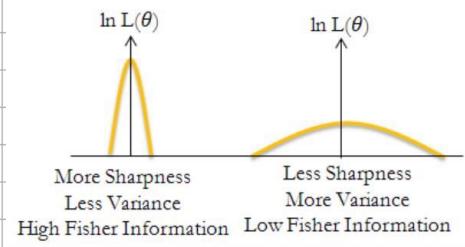


Fisher Information

$$X \sim f_\theta$$

$$I(\theta) = E_\theta [\dot{\ell}_{x_i}(\theta) \dot{\ell}_{x_i}(\theta)^T] = -E_\theta [\ddot{\ell}_{x_i}(\theta)]$$

where $\dot{\ell}_{x_i}(\theta) = \log f_\theta(x_i)$



Theorem: Cramér - Rao

Let $X_1, \dots, X_n \stackrel{iid}{\sim} f_\theta$ where

- Domain of f_θ $\mathcal{D}(f_\theta) = \{x \mid f_\theta(x) > 0\}$ does not depend on θ
- f_θ is 3 times continuously differentiable w.r.t. θ

Let $\hat{\theta} = g(X_1, \dots, X_n)$ be an unbiased estimator of θ_0

$$\text{Var}(\hat{\theta}) \geq [n I(\theta_0)]^{-1}$$

where $I(\theta_0) = E_{\theta_0} [\dot{\ell}_{x_i}(\theta_0) \dot{\ell}_{x_i}(\theta_0)^T] = -E_{\theta_0} [\ddot{\ell}_{x_i}(\theta_0)]$ and $\dot{\ell}_{x_i}(\theta) = \log f_\theta(x_i)$

Proof:

$$\begin{aligned} \text{cov}(\hat{\theta}, \dot{\ell}_{\hat{x}}(\theta)) &= \frac{\text{Cov}(\hat{\theta}, \dot{\ell}_{\hat{x}}(\theta))}{\sqrt{\text{Var}(\hat{\theta}) \text{Var}(\dot{\ell}_{\hat{x}}(\theta))}} \leq 1 \\ \Leftrightarrow \frac{\text{Cov}(\hat{\theta}, \dot{\ell}_{\hat{x}}(\theta))^2}{\text{Var}(\dot{\ell}_{\hat{x}}(\theta))} &\leq \text{Var}(\hat{\theta}) \end{aligned}$$

$$\begin{aligned} \text{Var}(\dot{\ell}_{\hat{x}}(\theta)) &= \text{Var}\left(\frac{\partial}{\partial \theta} \log \prod_{i=1}^n f_\theta(x_i)\right) \\ &= \text{Var}\left(\sum_{i=1}^n \frac{\partial}{\partial \theta} \log f_\theta(x_i)\right) \\ &= \sum_{i=1}^n \text{Var}\left(\frac{\partial}{\partial \theta} \log f_\theta(x_i)\right) \\ &= n \cdot \text{Var}\left(\frac{\partial}{\partial \theta} \log f_\theta(x_i)\right) \\ &= n \cdot \text{Var}(\dot{\ell}_{x_i}(\theta)) \\ &= n \cdot (E[\dot{\ell}_{x_i}(\theta)^2] - \underbrace{E[\dot{\ell}_{x_i}(\theta)]^2}_0) \\ &= n \cdot I(\theta) \end{aligned}$$

$$\begin{aligned} E[\dot{\ell}_{x_i}(\theta)] &= E\left[\frac{\partial}{\partial \theta} \log f_\theta(x_i)\right] \\ &= \int \frac{\partial}{\partial \theta} (\log f_\theta(x_i)) \cdot f_\theta(x_i) dx_i \\ &= \int \frac{\frac{\partial}{\partial \theta} f_\theta(x_i)}{f_\theta(x_i)} f_\theta(x_i) dx_i \\ &= \frac{\partial}{\partial \theta} \int f_\theta(x_i) dx_i \\ &= \frac{\partial}{\partial \theta} 1 \\ &= 0 \end{aligned}$$

$$\begin{aligned} \text{Cov}(\hat{\theta}, \dot{\ell}_{\hat{x}}(\theta)) &= E[\hat{\theta} \cdot \dot{\ell}_{\hat{x}}(\theta)] - E[\hat{\theta}] \cdot \overbrace{E[\dot{\ell}_{\hat{x}}(\theta)]}^0 \\ &= \int \dots \int \hat{\theta} \cdot \frac{\partial}{\partial \theta} (\log f_\theta(x_1, \dots, x_n)) \cdot f_\theta(x_1, \dots, x_n) dx_1 \dots dx_n \\ &= \int \dots \int \hat{\theta} \cdot \frac{\frac{\partial}{\partial \theta} f_\theta(x_1, \dots, x_n)}{f_\theta(x_1, \dots, x_n)} \cdot f_\theta(x_1, \dots, x_n) dx_1 \dots dx_n \\ &= \frac{\partial}{\partial \theta} \int \dots \int \hat{\theta} \cdot f_\theta(x_1, \dots, x_n) dx_1 \dots dx_n \\ &= \frac{\partial}{\partial \theta} E[\hat{\theta}] \\ &= \frac{\partial}{\partial \theta} \theta \quad \Rightarrow \hat{\theta} \text{ unbiased} \\ &= 1 \end{aligned}$$

\nwarrow smoothness of density function

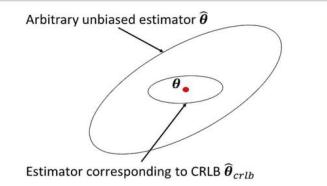
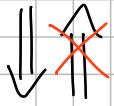


Illustration of the Cramér-Rao bound: there is no unbiased estimator which is able to estimate the (2-dimensional) parameter with less variance than the Cramér-Rao bound, illustrated as standard deviation ellipse. \square

Stochastic Convergence

a sequence of R.V. $\{X_n\}$ converges in probability to a R.V. X , shown by $X_n \xrightarrow{P} X$,
if $\lim_{n \rightarrow \infty} P(|X_n - X| \geq \varepsilon) = 0$, for all $\varepsilon > 0$



a sequence of R.V. $\{X_n\}$ converges in distribution to a R.V. X , shown by $X_n \xrightarrow{D} X$,
if $\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$, for all x at which $F_X(x)$ is continuous

Consistency

an estimator $\hat{\theta}_n$ of θ_0 is consistent, if $\hat{\theta}_n \xrightarrow{P} \theta_0$, i.e. $\hat{\theta}_n$ "converges in probability" to the true value

↳ e.g. an unbiased $\hat{\theta}_n$ with $\text{Var}(\hat{\theta}_n) \rightarrow 0$

Efficiency

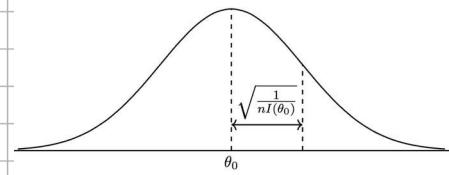
an estimator $\hat{\theta}_n$ of θ_0 is efficient, if $\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{D} N(0, I(\theta_0)^{-1})$



Theorem: Asymptotic efficiency of MLE

Let $X_1, \dots, X_n \stackrel{iid}{\sim} f_{\theta_0}$ where

- Domain of f_{θ} $\mathcal{D}(f_{\theta}) = \{x \mid f_{\theta}(x) > 0\}$ does not depend on θ
- f_{θ} is 3 times continuously differentiable w.r.t. θ



Then $\hat{\theta}_{MLE}$ is efficient

Proof:

$$0 = \dot{\ell}_{x_1, \dots, x_n}(\hat{\theta}_n^{MLE}) = \dot{\ell}_{x_1, \dots, x_n}(\theta_0) + (\hat{\theta}_n^{MLE} - \theta_0) \ddot{\ell}_{x_1, \dots, x_n}(\theta_0) + \text{"small"}$$

$$\Leftrightarrow \sqrt{n}(\hat{\theta}_n^{MLE} - \theta_0) \stackrel{(B)}{=} -\frac{\dot{\ell}_{x_1, \dots, x_n}(\theta_0) \frac{1}{\sqrt{n}}}{\ddot{\ell}_{x_1, \dots, x_n}(\theta_0) \frac{1}{n}} \quad (A)$$

$$(A) \quad -\frac{\ddot{\ell}_{x_1, \dots, x_n}(\theta_0)}{n} = -\frac{\sum_{i=1}^n \ddot{\ell}_{x_i}(\theta_0)}{n} \xrightarrow[n \rightarrow \infty]{P} -E[\ddot{\ell}_{x_i}(\theta_0)] =: I(\theta_0)$$

$$(B) \quad \frac{\dot{\ell}_{x_1, \dots, x_n}(\theta_0)}{\sqrt{n}} = \sqrt{n} \frac{\sum_{i=1}^n \dot{\ell}_{x_i}(\theta_0)}{n} \xrightarrow[n \rightarrow \infty]{CLT} \mathcal{N}(0, \text{Var}[\dot{\ell}_{x_i}(\theta_0)]) = \mathcal{N}(0, I(\theta_0))$$

$$(A) + (B): \sqrt{n}(\hat{\theta}_n^{MLE} - \theta_0) \xrightarrow{D} \frac{\mathcal{N}(0, I(\theta_0))}{I(\theta_0)} = \mathcal{N}(0, I(\theta_0)^{-1})$$

i.e. $\hat{\theta}_n^{MLE} \sim \mathcal{N}(\theta_0, (n I(\theta_0))^{-1})$

□

Confidence Intervals

use efficiency of MLE to construct approximate CI:

$$\hat{\theta}_{MLE} \sim N(\theta_0, (n I(\theta_0))^{-1})$$

$$\Rightarrow P(-z_{1-\frac{\alpha}{2}} < \sqrt{n I(\theta_0)} (\hat{\theta} - \theta_0) < z_{1-\frac{\alpha}{2}}) \approx 1 - \alpha$$

$$= P(-\hat{\theta} - \sqrt{\frac{1}{n I(\theta_0)}} z_{1-\frac{\alpha}{2}} < -\theta_0 < \hat{\theta} - \sqrt{\frac{1}{n I(\theta_0)}} z_{1-\frac{\alpha}{2}})$$

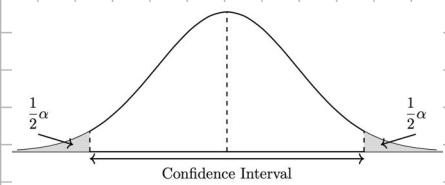
$$= P(\hat{\theta} + \sqrt{\frac{1}{n I(\theta_0)}} z_{1-\frac{\alpha}{2}} > \theta_0 > \hat{\theta} - \sqrt{\frac{1}{n I(\theta_0)}} z_{1-\frac{\alpha}{2}})$$

$$= P(\hat{\theta} - \sqrt{\frac{1}{n I(\theta_0)}} z_{1-\frac{\alpha}{2}} < \theta_0 < \hat{\theta} + \sqrt{\frac{1}{n I(\theta_0)}} z_{1-\frac{\alpha}{2}})$$

$$\Rightarrow C(\vec{x}) = \left(\hat{\theta}_n^{MLE}(\vec{x}) - \sqrt{\frac{1}{n I(\theta_0)}} z_{1-\frac{\alpha}{2}}, \hat{\theta}_n^{MLE}(\vec{x}) + \sqrt{\frac{1}{n I(\theta_0)}} z_{1-\frac{\alpha}{2}} \right)$$

$$P(\theta_0 \in C(\vec{x})) = 1 - \alpha$$

replace θ_0 in $C(\vec{x})$ by $\hat{\theta}_n^{MLE}$ to get a practical CI



Hypothesis Testing

Let Θ_0 be accepted body of knowledge

We observe data $X_1, \dots, X_n \stackrel{iid}{\sim} f_\theta$

$H_0: \theta = \theta_0$ status quo / null hypothesis

$H_1: \theta \neq \theta_0$ alternative hypothesis

Do the data give sufficient evidence to reject H_0 ?

Consider a test statistic T

$$\textcircled{A} \quad T = \hat{\theta}^{\text{MLE}}(\vec{X})$$

$t = \hat{\theta}(\vec{X})$ observed value at estimator

$$p\text{-value}(\vec{X}) = P(|\hat{\theta}(\vec{X}) - \theta_0| > |t - \theta_0| \mid H_0 \text{ true})$$

$$= Z \cdot P(\hat{\theta}(\vec{X}) - \theta_0 > |t - \theta_0| \mid H_0 \text{ true})$$

$$= 2 \cdot P\left(\frac{\hat{\theta}(\vec{X}) - \theta_0}{\sqrt{1/(nI(\theta_0))}} > \frac{|t - \theta_0|}{\sqrt{1/(nI(\theta_0))}} \mid H_0 \text{ true}\right)$$

$$= Z \cdot \left(1 - \Phi\left(\frac{|t - \theta_0|}{\sqrt{1/(nI(\theta_0))}}\right)\right) \quad \Phi(\cdot) \text{ is cdf of } N(0, 1)$$

asymptotic distribution of $\hat{\theta} - \theta_0$ is $N(0, I(\theta_0)^{-1})$

We reject H_0 if $p\text{-value} < \alpha$ where $\alpha = \text{significance level}$
 → e.g. $\alpha = 0.05$ in science, $\alpha \leq 0.001$ in physics

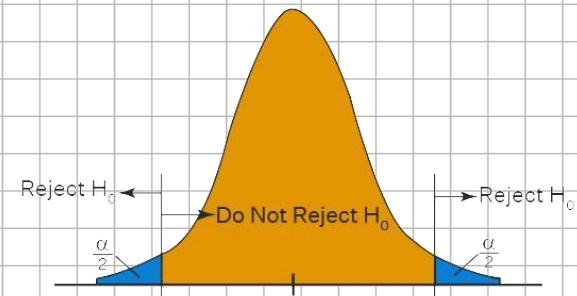
Interpretation p-value: "How likely is it to observe the data that we did observe (or more extreme) if the null hypothesis is true."

$$\textcircled{B} \quad T = -2 \log \frac{f_{\theta_0}(\vec{X})}{\sup_{\theta} f_{\theta}(\vec{X})}$$

likelihood ratio (LR) $\xrightarrow{\text{close to 1}} H_0 \text{ pretty good}$
 $\xrightarrow{\text{close to 0}} H_0 \text{ pretty bad}$

$\in (0, \infty)$

$\xrightarrow{\text{close to 0}} H_0 \text{ pretty good}$
 $\xrightarrow{\text{close to } \infty} H_0 \text{ pretty bad}$



Wilks Theorem

If $X_1, \dots, X_n \stackrel{iid}{\sim} f_{\theta}$ with usual conditions

$$T = -2 \log \frac{f_{\theta_0}(\vec{x})}{\sup_{\theta} f_{\theta}(\vec{x})} \Big| H_0 \xrightarrow{D} \chi^2_{df} \quad df = \# \text{ parameters kept fixed in } H_0$$

Proof:

$$T = -2(\ell_{\vec{x}}(\theta_0) - \ell_{\vec{x}}(\hat{\theta}))$$

$$\ell_{\vec{x}}(\theta_0) = \ell_{\vec{x}}(\hat{\theta}) + \overset{\circ}{\ell'_{\vec{x}}(\hat{\theta})} (\theta_0 - \hat{\theta}) + \frac{1}{2} \ddot{\ell}_{\vec{x}}(\hat{\theta})(\theta_0 - \hat{\theta})^2 + \text{"small"}$$

$$-2(\ell_{\vec{x}}(\theta_0) - \ell_{\vec{x}}(\hat{\theta})) = T \approx -\ddot{\ell}_{\vec{x}}(\hat{\theta})(\theta_0 - \hat{\theta})^2$$

$$= \frac{-\ddot{\ell}_{\vec{x}}(\hat{\theta})}{n I(\theta_0)} (\theta_0 - \hat{\theta})^2 \cdot n I(\theta_0)$$

$$= \frac{\sum_{i=1}^n (-\ddot{\ell}_{x_i}(\hat{\theta}))}{n I(\theta_0)} \cdot \underbrace{\left(\sqrt{n I(\theta_0)} \cdot (\theta_0 - \hat{\theta}) \right)^2}_{N(0, 1)}$$

$$\frac{E[-\ddot{\ell}_{x_i}(\hat{\theta})]}{I(\theta_0)}$$

||

$$\chi^2_1$$

□

Inverting the LRT

Rejection region: $R(\theta_0) = \{\vec{x} \mid T(\vec{x}) \geq K\}$

→ We reject H_0 if we observe $T(\vec{x}) \geq K$

$$P(\vec{x} \in R(\theta_0) \mid H_0 \text{ true}) = P(T(\vec{x}) \geq K \mid H_0 \text{ true}) = \alpha$$

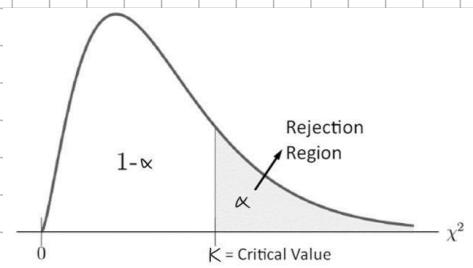
Using the asymptotic dist'r of $T(\vec{x}) \mid H_0 \text{ true}$ we can find K s.t. above equality holds:

$$K = \chi^2_{1,1-\alpha} = F_{\chi^2}^{-1}(1-\alpha)$$

$$P(T(\vec{x}) \geq \chi^2_{1,1-\alpha} \mid H_0 \text{ true}) = \alpha$$

$$\Rightarrow P(\vec{x} \notin R(\theta_0) \mid H_0 \text{ true}) = 1 - \alpha$$

"we accept H_0 "



Hypothesis Testing

$$P(\vec{x} \notin R(\theta_0) \mid H_0 \text{ true}) = 1 - \alpha$$

$$\underbrace{P(T(\vec{x}) \leq \chi^2_{1,1-\alpha})}_{\text{if we "invert" this by isolating } \theta_0 \text{ on one side}} = 1 - \alpha$$

Confidence Interval

$$P(\theta_0 \in C(\vec{x})) = 1 - \alpha$$

$$C(\vec{x}) = \{\theta \mid T_\theta(\vec{x}) \leq \chi^2_{1,1-\alpha}\}$$

we can construct

Linear Regression

Data: $\{(\vec{x}_i, Y_i)\}_{i=1}^n$, $\vec{x}_i \in \mathbb{R}^p$

$$Y_i = \vec{x}_i^\top \vec{\beta} + \varepsilon_i, \quad \vec{\beta} \in \mathbb{R}^p, \quad \varepsilon_i \sim N(0, \sigma^2)$$

response
 error term

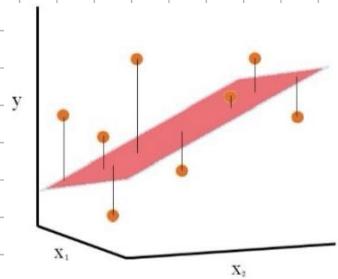
Estimation

$$\ell_p(\vec{\beta}) = \sum_{i=1}^n \log f(Y_i | \vec{\beta}, \vec{x}_i) = \log f(\vec{y} | \vec{\beta}, \vec{x}_1, \dots, \vec{x}_n)$$

$$= \log \left((2\pi)^{-\frac{n}{2}} |\sigma^2 I|^{-\frac{1}{2}} \exp(-\frac{1}{2} (\vec{y} - \underline{\vec{x}} \vec{\beta})^\top (\sigma^2 I)^{-1} (\vec{y} - \underline{\vec{x}} \vec{\beta})) \right)$$

$$= C - \frac{1}{2} \log \sigma^{2n} - \frac{1}{2\sigma^2} (\vec{y} - \underline{\vec{x}} \vec{\beta})^\top (\vec{y} - \underline{\vec{x}} \vec{\beta})$$

$$= C - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \vec{x}_i^\top \vec{\beta})^2$$



$$\underline{X} = \begin{bmatrix} -\vec{x}_1^\top - \\ \vdots \\ -\vec{x}_n^\top - \end{bmatrix}^T \underbrace{\vec{x}}_p$$

$$\dot{\ell}_p(\vec{\beta}) = -\frac{1}{2\sigma^2} \sum_{i=1}^n 2(Y_i - \vec{x}_i^\top \vec{\beta}) \cdot (-\vec{x}_i^\top) = \frac{1}{\sigma^2} \sum_{i=1}^n (Y_i - \vec{x}_i^\top \vec{\beta}) \cdot \vec{x}_i^\top = \frac{1}{\sigma^2} (\vec{y} - \underline{\vec{x}} \vec{\beta})^\top \underline{X} = \frac{1}{\sigma^2} (\vec{y}^\top - \vec{\beta}^\top \underline{X}^\top) \underline{X}$$

$$= \frac{1}{\sigma^2} (\vec{y}^\top \underline{X} - \vec{\beta}^\top \underline{X}^\top \underline{X}) \stackrel{!}{=} \vec{0}_p^\top, \quad \vec{\beta}^\top \underline{X}^\top \underline{X} = \left[\sum_{i=1}^p \beta_i (\underline{X}^\top \underline{X})_{ii} \quad \dots \quad \sum_{i=1}^p \beta_i (\underline{X}^\top \underline{X})_{ip} \right]$$

$$\Rightarrow \vec{\beta}^\top \underline{X}^\top \underline{X} = \vec{y}^\top \underline{X} \quad \Leftrightarrow \underline{X}^\top \underline{X} \vec{\beta} = \underline{X}^\top \vec{y} \quad \Leftrightarrow \vec{\beta} = (\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \vec{y}$$

directly normal (D)
(not just asymptotically)

$$\ddot{\ell}_p(\vec{\beta}) = -\frac{1}{\sigma^2} \frac{d}{d\vec{\beta}} (\vec{\beta}^\top \underline{X}^\top \underline{X}) = -\frac{1}{\sigma^2} \underline{X}^\top \underline{X} \stackrel{\text{Löwner}}{<} 0$$

$$\Rightarrow \hat{\vec{\beta}} = (\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \vec{y} \text{ is uncorrelated} \Rightarrow \text{MLE}$$

Properties:

$$\begin{aligned} \text{expected v: } E[\hat{\vec{\beta}}] &= E[(\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \vec{y}] \\ &= (\underline{X}^\top \underline{X})^{-1} \underline{X}^\top E[\vec{y}] \\ &= (\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \underline{X} \vec{\beta} \\ &= \vec{\beta} \quad \rightarrow \text{unbiased} \end{aligned}$$

$$\begin{aligned} \text{variance: } \text{Var}[\hat{\vec{\beta}}] &= \text{Var}[(\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \vec{y}] \\ &= (\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \text{Var}[\vec{y}] ((\underline{X}^\top \underline{X})^{-1} \underline{X}^\top)^\top \\ &= (\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \sigma^2 \underline{I} \underline{X} (\underline{X}^\top \underline{X})^{-1} \\ &= \sigma^2 (\underline{X}^\top \underline{X})^{-1} \end{aligned}$$

$$\begin{aligned} \text{CRLB: } (\underline{I}(\vec{\beta}))^{-1} &= (-E[\ddot{\ell}_p(\vec{\beta})])^{-1} \\ &= (-E[-\frac{1}{\sigma^2} \underline{X}^\top \underline{X}])^{-1} \\ &= \sigma^2 (\underline{X}^\top \underline{X})^{-1} \end{aligned}$$

Inverse of 2x2 Matrix

If $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ then

$$A^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

Inverse of A Determinant of A Adjoint of A

Note: A^{-1} exists only when $ad - bc \neq 0$

achieves CRLB \rightarrow efficient

Confidence Interval

$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ is a linear combination of normally distributed R.V.

\Rightarrow so it is also normally distributed with mean and variance as calculated above (exact!)

$$\hat{\beta} \sim N(\beta, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})$$

\rightarrow we could create a joint confidence region for $\hat{\beta}$, but we stay modest and do univariate confidence for each of the β_j :

$$\hat{\beta}_j \sim N(\beta_j, \sigma^2 ((\mathbf{X}^\top \mathbf{X})^{-1})_{jj})$$

a $(1-\alpha) \cdot 100\%$ CI for β_j :

$$(\hat{\beta}_j \pm z_{1-\alpha} \cdot \sigma \cdot \sqrt{(\mathbf{X}^\top \mathbf{X})^{-1}_{jj}})$$

Testing particular covariate/predictor x_j significant or not, i.e. $\beta_j = 0$?

$$H_0: \beta_j = 0$$

$$H_1: \beta_j \neq 0 \quad \text{Covariate } j \text{ is interesting}$$

Test-statistic:

$$\hat{\beta}_j | H_0 \sim N(0, \sigma_j^2) \text{ where } \sigma_j^2 = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}_{jj}$$

we estimate β_j as b_j

$$\text{P-value} = P(|\hat{\beta}_j| > |b_j| | H_0) = 2 \cdot P\left(\frac{\hat{\beta}_j}{\sigma_j} > \frac{|b_j|}{\sigma_j} | H_0\right) = 2 \cdot \left(1 - \Phi\left(\frac{|b_j|}{\sigma_j}\right)\right)$$

reject H_0 if P-value $< \alpha$

but: if we estimate

Model Evaluation does adding x_{p+1}, \dots, x_p significantly improve our ability to predict y ?

consider 2 competing models (nested, i.e. predictors in small model all contained in bigger model)

$$\underline{\text{Model 0}}: x_i^{(0)} = (x_{i1}, \dots, x_{ip_0}) \in \mathbb{R}^{p_0}$$

$$\underline{\text{Model 1}}: x_i^{(1)} = (x_{i1}, \dots, x_{ip_1}) \in \mathbb{R}^{p_1}$$

$H_0: \text{model 0 is correct}$

$H_1: \text{model 1 is correct}$

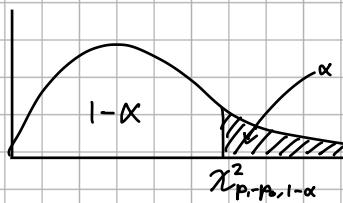
Test-statistic:

$$\begin{aligned} \Lambda &= -2 \log \frac{L_0(\beta_0)}{L_1(\beta_1)} \quad \text{where} \\ &= -2 \left(\ell_{\mathbf{y}}(\hat{\beta}_0) - \ell_{\mathbf{y}}(\hat{\beta}_1) \right) \end{aligned}$$

$$\Lambda | H_0 \sim \chi^2_{p_1 - p_0} \quad (\text{exact!})$$

reject H_0 if $\Lambda > \chi^2_{p_1 - p_0, 1-\alpha}$

$$\begin{aligned} L_0(\beta_0) &= \max_{\beta_0} f(\mathbf{y} | \text{model 0}), \\ L_1(\beta_1) &= \max_{\beta_1} f(\mathbf{y} | \text{model 1}) \end{aligned}$$



this is where the field of data science was in the 70s/80s
there was one big question: how to compare non-nested models?
-> next page

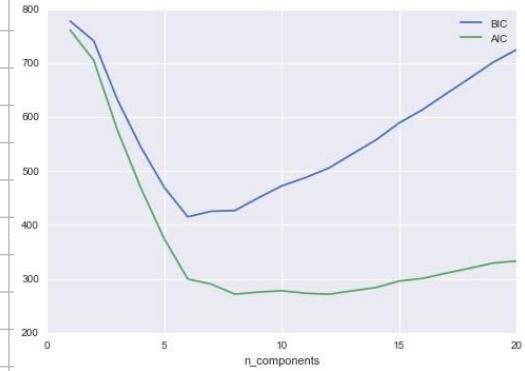
Residual Sum of Squares

$$RSS = \hat{\varepsilon}^T \hat{\varepsilon} = (\vec{y} - X\hat{\beta})^T (\vec{y} - X\hat{\beta}) = \sum_{i=1}^n (y_i - \vec{x}_i^T \hat{\beta})^2$$

Akaike Information Criterion

$$\begin{aligned}
 KL(M_0 | M) &= \int_{-\infty}^{\infty} f_{M_0}(y) \log \left(\frac{f_{M_0}(y)}{f_M(y)} \right) dy \\
 &= E_{Y \sim M_0} [\ell_Y(M_0) - \ell_Y(M)] \\
 &= C - E_{Y \sim M_0} \ell_Y(M) \\
 &\approx \dots \approx C - \underbrace{\frac{1}{n} \sum_{i=1}^n \ell_{Y_i}(\hat{\beta})}_{\ell_Y(\hat{\beta})} + \underbrace{\frac{p}{n}}_{\text{bias}} / \times 2n
 \end{aligned}$$

→ minimize



$$\begin{aligned}
 AIC(M) &= -2 \ell_Y(\vec{\beta}) + 2p \rightarrow \text{prediction} \\
 BIC(M) &= -2 \ell_Y(\vec{\beta}) + p \log(n) \rightarrow \text{"truth"}
 \end{aligned}$$

ANOVA

Mixed Effects Models

Consider \underline{X} to be the covariates describing the quantities of interest, \underline{Z} to be the covariates describing the nuisance effects

Random effects model:

$$\vec{Y} = \underbrace{\underline{X}\vec{\beta}}_{\text{fixed effects}} + \underbrace{\underline{Z}\vec{\gamma}}_{\gamma_j \sim N(0, \sigma_\gamma^2)} + \underbrace{\vec{\epsilon}}_{\sim N(0, \sigma^2 I)}$$

Why is γ random?

- Philosophical: γ is a random draw out of the population of classes
- Pragmatic: if γ is random, then γ is not a parameter, σ_γ^2 is!

Estimation: REML (Restricted Maximum Likelihood)

Way to estimate the variance parameters σ_γ^2 and σ^2

projection matrix into the residual space:

$$\underline{P} = \underline{I} - \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T$$

note that $\underline{P}\vec{Y} = \vec{Y} - \underbrace{\underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T \vec{Y}}_{\vec{\beta}} = \hat{\vec{\epsilon}}$

$$\begin{aligned}
 \mathbb{E}(\underline{P}\vec{Y}) &= \vec{0} \\
 \mathbb{V}(\underline{P}\vec{Y}) &= \underline{P} \mathbb{V}(\vec{Y}) \underline{P}^T \\
 &= \underline{P} (\mathbb{V}(\underline{X}\vec{\beta} + \underline{Z}\vec{\gamma} + \vec{\epsilon})) \underline{P}^T \\
 &= \underline{P} \left(\underbrace{\underline{Z} \mathbb{V}(\vec{\gamma}) \underline{Z}^T}_{\text{function of } \sigma_\gamma^2} + \sigma^2 \underline{I} \right) \underline{P}^T \\
 &\quad \underbrace{\text{known function of } \sigma_\gamma^2 \text{ and } \sigma^2}_{\text{known function of } \sigma_\gamma^2 \text{ and } \sigma^2}
 \end{aligned}$$

If I maximize the likelihood for $\underline{P}\vec{Y}$ with respect to σ_γ^2 and σ^2 , then estimates are unbiased.

⇒ use those estimates to estimate $\vec{\beta}$

REML

problem with ML:

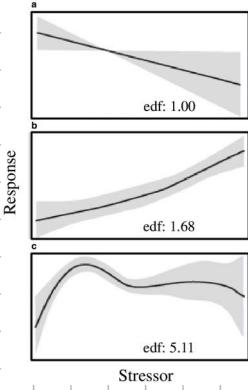
doesn't take into account the degrees of freedom that we lose when we estimate the mean

REML basically corrects for bias in variance components of the ML estimate

Additive Models

non-linearity: $Y = f(X) + \varepsilon$, $\varepsilon \sim N(0, \sigma^2)$

how to find f ?



Option 1: we have a known form of functional space; $f(X) = f(X, \theta)$

find f by estimating θ ; $\hat{\theta} = \operatorname{argmax}_{\theta} -\frac{1}{2} \sum_{i=1}^n (Y_i - f(X_i, \theta))^2 \rightarrow \hat{f} = f(\cdot, \hat{\theta})$

Option 2: data driven non-linear functions

find the true function that belongs to general function class, say $C = \{f \mid f: \mathbb{R} \rightarrow \mathbb{R} \text{ continuous}\}$

construct a sequence of functional spaces $S_1 \subset S_2 \subset S_3 \subset \dots$, such that $\overline{\bigcup_{i=1}^{\infty} S_i} = C \Rightarrow \text{SIEVE}$

IDEA: estimate f inside S_m where n sample size and $m(n)$ increasing sequence with $\lim_{n \rightarrow \infty} m(n) = \infty$

How to choose S_m ?

finite dimensional space with basis $\{b_1, \dots, b_m\}$ where $b_j: \mathbb{R} \rightarrow \mathbb{R}$

Options

(A) $b_j(x) = x^j$, note $S_m = \{\text{all polynomials of order } m\}$ and $\overline{\bigcup_{i=1}^{\infty} S_i} = C$

disadvantage: tends to be quite "unstable" in practice

→ any choice made to fit data at one point affects the fit everywhere else (*NON-LOCAL*)
→ *drawback*

(B) LOCAL BASIS, e.g. cubic spline

3 constraints at knot points ξ : continuous, continuous 1st & 2nd derivative

independent parameters: $(n_\xi + 1) \cdot 4 - n_\xi \cdot 3 = n_\xi + 4$

n_ξ knots $\Rightarrow S_{n_\xi+4}$

by increasing n_ξ , we obtain a "covering" of C

Estimation

given some space S_m with basis $B_m = \{b_1, \dots, b_m\}$ and data $(x_1, Y_1), \dots, (x_n, Y_n)$ MLE of $f \in S_m$

$$\hat{f} = \operatorname{argmax}_{f \in S_m} -\frac{1}{2} \sum_{i=1}^n (Y_i - f(x_i))^2$$

$$\rightarrow \hat{\beta}_{\text{MLE}} = \operatorname{argmax}_{\beta \in \mathbb{R}^m} -\frac{1}{2} \sum_{i=1}^n (Y_i - \sum_{j=1}^m \beta_j b_j(x_i))^2 \quad \rightarrow \text{linear model}$$

$$= (\underline{X}^T \underline{X})^{-1} \underline{X}^T \vec{Y}, \quad \vec{Y} = [Y_1, \dots, Y_n]^T, \quad \underline{X} = \begin{bmatrix} b_1(x_1) & \dots & b_m(x_1) \\ \vdots & & \vdots \\ b_1(x_n) & \dots & b_m(x_n) \end{bmatrix}$$

this may be a problematic solution, e.g. if m is too close to $n \rightarrow \hat{f}$ becomes too "wiggly"
→ consequence of overfitting

solution: regularization / penalty on wigginess W

$$W(f) = \int_L^U (f''(x))^2 dx$$

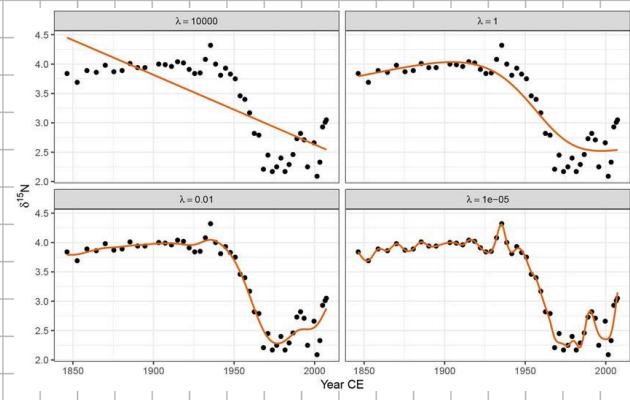
$$\stackrel{f \in S_m}{=} \int_L^U \left(\sum_{j=1}^m \beta_j b_j''(x) \right)^2 dx$$

$$= \int_L^U \sum_{j=1}^m \sum_{k=1}^m \beta_j \beta_k b_j(x) b_k''(x) dx$$

$$= \sum_{j=1}^m \sum_{k=1}^m \beta_j \beta_k \int_L^U b_j''(x) b_k''(x) dx$$

$$= \sum_{j=1}^m \sum_{k=1}^m \beta_j \beta_k W_{jk}^* \quad \rightarrow W_{jk}^* \text{ is known}$$

$$= \vec{\beta}^T \underline{W}^* \vec{\beta}$$



penalized maximum likelihood:

$$\hat{\vec{\beta}} = \underset{\vec{\beta}}{\operatorname{argmax}} \left(-\frac{1}{2} (\vec{y} - \underline{X} \vec{\beta})^T (\vec{y} - \underline{X} \vec{\beta}) - \frac{\lambda}{2} \vec{\beta}^T \underline{W}^* \vec{\beta} \right)$$

if λ is $\begin{cases} = 0 & \text{wiggly solution} \\ \rightarrow \infty & \text{linear solution} \end{cases}$

$$\frac{\partial \mathcal{L}^P}{\partial \vec{\beta}} = \underline{X}^T \vec{y} - \underline{X}^T \underline{X} \vec{\beta} - \lambda \underline{W}^* \vec{\beta} = 0$$

$$\Rightarrow (\underline{X}^T \underline{X} + \lambda \underline{W}^*) \vec{\beta} = \underline{X}^T \vec{y}$$

$$\rightarrow \hat{\vec{\beta}}_{\text{PML}} = (\underline{X}^T \underline{X} + \lambda \underline{W}^*)^{-1} \underline{X}^T \vec{y}$$

how to choose λ : cross-validation

note on wigginess:

among all functions $\in C$ that fit the data "in the same way", cubic splines is the least wiggly

Logistic Regression

consider data (Y_i, \vec{x}_i) where

$$\begin{cases} Y_i \sim \text{Bern}(\pi_i) \\ \vec{x}_i \in \mathbb{R}^p \end{cases}$$

How to model $\pi_i(\vec{x}_i)$?

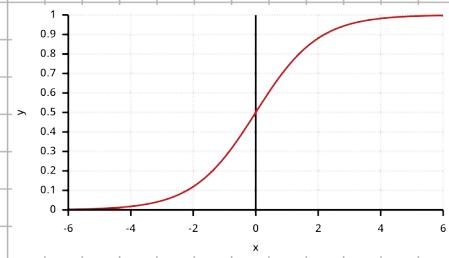
$$\eta = \vec{x}_i^T \vec{\beta} \quad \text{Linear predictor}$$

$$\pi_i = g^{-1}(\eta) \quad \text{for some suitable function } g^{-1}: \mathbb{R} \rightarrow [0, 1] \text{ that fulfills}$$

inverse link function

$$(i) \text{ logistic function } g^{-1}(\eta) = \frac{e^\eta}{1 + e^\eta} \quad \checkmark$$

$$(ii) \text{ cdf of } N(0, 1) \quad g^{-1}(\eta) = \Phi_{0,1}(\eta)$$



- * one-to-one
- * monotone increasing

Logistic Regression

$$\begin{cases} Y_i \sim \text{Bern}(\pi_i) \\ \pi_i = \frac{e^{\eta_i}}{1 + e^{\eta_i}} \\ \eta_i = \vec{x}_i^T \vec{\beta} \end{cases} \implies \eta_i = \underbrace{\log\left(\frac{\pi_i}{1 - \pi_i}\right)}_{\text{log-odds} \in \mathbb{R}} \quad \text{logit function}$$

$\vec{x}_i^T \vec{\beta}$ has interpretation in terms of log-odds:

β_j = increase in log-odds if x_j increases by 1 and rest stayed the same

e^{β_j} = multiplicative increase in odds if x_j increases by 1 and rest stayed the same

"chances of Y are increased by a factor of e^{β_j} if $x_j \rightarrow x_j + 1$ "

Estimation

$$\begin{aligned} \ell(\vec{\beta}) &= \log \prod_{i=1}^n \pi_i^{Y_i} (1 - \pi_i)^{1 - Y_i} \\ &= \sum_{i=1}^n \left(Y_i \log\left(\frac{e^{\eta_i}}{1 + e^{\eta_i}}\right) + (1 - Y_i) \log\left(\frac{1 - e^{\eta_i}}{1 + e^{\eta_i}}\right) \right) \\ &= \sum_{i=1}^n \left(Y_i (\vec{x}_i^T \vec{\beta} - \log(1 + e^{\eta_i})) + (1 - Y_i) (-\log(1 + e^{\eta_i})) \right) \\ &= \sum_{i=1}^n \left(Y_i \vec{x}_i^T \vec{\beta} - Y_i \log(1 + e^{\eta_i}) + Y_i \log(1 + e^{\eta_i}) - \log(1 + e^{\eta_i}) \right) \\ &= \sum_{i=1}^n (Y_i \vec{x}_i^T \vec{\beta} - \log(1 + e^{\vec{x}_i^T \vec{\beta}})) \end{aligned}$$

$$\frac{\partial}{\partial \beta_j} \ell(\vec{\beta}) = \sum_{i=1}^n \left(Y_i x_{ij} - \frac{1}{1 + e^{\vec{x}_i^T \vec{\beta}}} \cdot e^{\vec{x}_i^T \vec{\beta}} \cdot x_{ij} \right)$$

$$= \sum_{i=1}^n (Y_i - \pi_i) x_{ij}$$

$$\begin{cases} \sum_{i=1}^n (Y_i - \pi_i) x_{1i} = 0 \\ \vdots \\ \sum_{i=1}^n (Y_i - \pi_i) x_{pi} = 0 \end{cases}$$

find $\hat{\vec{\beta}}$ numerically using

- (A) Newton-Raphson
- (B) Fisher Scoring

Properties of $\hat{\beta}$

asymptotically (under certain constraints on x):

$$\lim_{n \rightarrow \infty} \sqrt{n}(\hat{\beta} - \vec{\beta}_0) \sim N(0, \underline{I}^{-1}), \quad \underline{I} = \lim_{n \rightarrow \infty} -E\left[\frac{\partial^2 \ell}{\partial \beta^2}\right] / n$$

so approximately:

$$\hat{\beta} \sim N(\vec{\beta}_0, \frac{1}{n} \underline{I}^{-1})$$

→ we can now construct CI's and test hypotheses $H_0: \beta_j = 0$

other things:

Ⓐ testing of nested models $\begin{cases} H_0: M = M_0 \\ H_1: M = M_1 \end{cases}$ where $M_0 \subset M_1$

(likelihood ratio test): $\Lambda = -2(\ell(\hat{\beta}_{M_0}) - \ell(\hat{\beta}_{M_1}))$ where $\Lambda | H_0 \sim \chi^2_{df(M_1) - df(M_0)}$

Ⓑ AIC(M) = $-2\ell(\hat{\beta}_M) + 2df(M)$

Ⓒ goodness-of-fit test $\begin{cases} H_0: M \text{ is true/correct} \\ H_1: M \text{ is not true} \end{cases}$

test-statistic: $T = -2\ell(\hat{\beta}_M)$ deviance where $T | H_0 \sim \chi^2_{n - df(M)}$

reject H_0 if $T > \chi^2_{n - df(M), 1-\alpha}$

Generalized Linear Models

Let $Y \sim \text{Exponential Family}$, i.e.

$$f_Y(y) = \exp\left(\frac{y\theta - b(\theta)}{a(\varphi)} + c(y, \varphi)\right)$$

for parameters

θ : location parameter
 φ : dispersion parameter

for some functions a, b, c

expected value and variance of Y :

$$\mathbb{E}[\dot{\ell}_Y(\theta)] \quad \text{trick...}$$

$$= \mathbb{E}\left[\frac{\partial}{\partial\theta}\left(\frac{y\theta - b(\theta)}{a(\varphi)} + c(y, \varphi)\right)\right]$$

$$= \mathbb{E}\left[\frac{y - \frac{\partial b(\theta)}{\partial\theta}}{a(\varphi)}\right]$$

$$= \frac{\mathbb{E}[Y] - b'(\theta)}{a(\varphi)}$$

$$\rightarrow \mathbb{E}[Y] = b'(\theta) \Rightarrow \theta = b^{-1}(\mathbb{E}[Y]) = b^{-1}(\mu)$$

$$\mathbb{E}[\ddot{\ell}_Y(\theta)^2] + \mathbb{E}[\dot{\ell}_Y(\theta)]$$

$$= \mathbb{E}\left[\left(\frac{y - b(\theta)}{a(\varphi)}\right)^2\right] - \frac{b'(\theta)}{a(\varphi)}$$

$$= \mathbb{E}[(Y - \mathbb{E}[Y])^2] - \frac{b'(\theta)}{a(\varphi)^2}$$

$$\rightarrow \text{Var}[Y] = b''(\theta) a(\varphi)$$

Linear Model and Link Function

consider (\vec{x}, Y) where $Y \sim \text{Exponential Family}$ and $\vec{x} \in \mathbb{R}^p$

linear predictors for Y :

$$\eta = \vec{x}^\top \vec{\beta}, \quad \vec{\beta} \in \mathbb{R}^p$$

the aim of the link function $g: M \rightarrow \mathbb{R}$ where $EY \in M \subset \mathbb{R}$ is to map EY into the space of linear predictors $\mathbb{R} \ni \vec{x}^\top \vec{\beta}$:

$$E[Y] = g^{-1}(\vec{x}^\top \vec{\beta}) \Leftrightarrow g(M) = \eta$$

constraints on g :

* 1-to-1

* continuous

* strictly increasing

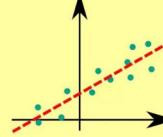
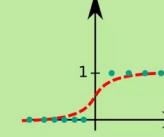
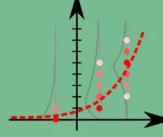
→ computationally more efficient
 but not always best choice from interpretation point of view

canonical / link function: $g = b^{-1}$ then $\eta = \vec{x}^\top \vec{\beta} = g(E[Y]) = \theta$

Generalized Linear Models (GLM) extend the ordinary linear regression and allow the response variable y to have an error distribution other than the normal distribution.

GLMs are:

- a) Easy to understand
- b) Simple to fit and interpret in any statistical package
- c) Sufficient in a lot of practical applications

LINEAR REGRESSION	LOGISTIC REGRESSION	POISSON REGRESSION
<ul style="list-style-type: none"> ① Econometric modelling ② Marketing Mix Model ③ Customer Lifetime Value  <p>Continuous \Rightarrow Continuous</p> $y = \alpha_0 + \sum_{i=1}^N \alpha_i x_i$ <p>1 unit increase in x increases y by α</p>	<ul style="list-style-type: none"> ① Customer Choice Model ② Click-through Rate ③ Conversion Rate ④ Credit Scoring  <p>Continuous \Rightarrow True/False</p> $y = \frac{1}{1 + e^{-z}}$ $z = \alpha_0 + \sum_{i=1}^N \alpha_i x_i$ $\text{glm}(y \sim x1 + x2, \text{data}, \text{family}=binomial())$ <p>1 unit increase in x increases log odds by α</p>	<ul style="list-style-type: none"> ① Number of orders in lifetime ② Number of visits per user  <p>Continuous $\Rightarrow 0, 1, 2, \dots$</p> $y \sim \text{Poisson}(\lambda)$ $\ln \lambda = \alpha_0 + \sum_{i=1}^N \alpha_i x_i$ $\text{glm}(y \sim x1 + x2, \text{data}, \text{family}=poisson())$ <p>1 unit increase in x multiplies y by e^α</p>

Estimation

$$\ell_{\vec{Y}}(\vec{\beta}) = \sum_{i=1}^n \ell_{Y_i}(\vec{\beta}) \\ = \sum_{i=1}^n \left(\frac{Y_i \theta_i - b(\theta_i)}{a_i(\varphi)} + c(Y_i, \varphi) \right)$$

$$\begin{aligned} \frac{d}{d\beta_k} \ell_{\vec{Y}}(\vec{\beta}) &= \sum_{i=1}^n \frac{d}{d\beta_k} \ell_{Y_i}(\vec{\beta}) \\ &= \sum_{i=1}^n \frac{d\ell_{Y_i}}{d\theta_i} \frac{d\theta_i}{d\mu_i} \frac{d\mu_i}{d\eta_i} \frac{d\eta_i}{d\beta_k} \\ &= \sum_{i=1}^n \frac{Y_i - b(\theta_i)}{a_i(\varphi)} \frac{1}{b'(\theta_i)} \frac{d\mu_i}{d\eta_i} \frac{d\eta_i}{d\beta_k} \\ &= \sum_{i=1}^n \frac{(Y_i - M_i) \cdot d\eta_i/d\mu_i}{(d\eta_i/d\mu_i)^2 b(\theta_i) a_i(\varphi)} \frac{d\eta_i}{d\beta_k} \end{aligned}$$

Iterative Re-weighted Least Squares. Aiming to maximize the log, we set the derivative of the log-likelihood to zero. In class we used the usual Taylor method. Here we consider another technique, called *iterative reweighted least squares*.

- The problem of the above equation for the derivative of the log-likelihood is that it is not linear in β . Write $\mu_i = \mu_i(\beta_0) + (\mu_i(\beta) - \mu_i(\beta_0))$ to derive the linear approximation of the numerator using the Taylor expansion of η_i around $\mu_i(\beta_0)$.
- Write now the full derivative of the likelihood in the format

$$\frac{d\ell_Y}{d\beta} \approx \sum_{i,j} \frac{(Y_i^* - x_i^t \beta)}{\left(\frac{d\eta_i}{d\mu_i}\right)^2 V(\mu_i)} \frac{d\eta_i}{d\beta}.$$

What is the *adjusted dependent variable* Y^* .

- Show that we can further approximate the denominator of the above expression to get,

$$\frac{d\ell_Y}{d\beta} \approx X^t W (Y^* - X\beta),$$

where X is the matrix of dyadic covariates $\{x_i\}_{i,j}$ and W is the diagonal matrix with $1/w_i$ on the diagonal, where $w_i = \left(\frac{d\eta_i}{d\mu_i}(\beta_0)\right)^2 V(\mu_i(\beta_0))$. The resulting expression is now linear in β .

- What is the expression of the term β that maximizes the approximation of the likelihood?

Aim: write this as a linear function in β .

- write μ_i as approx linear in η_i :

$$\mu_i(\eta_i) \approx \mu_i(\eta_{i0}) + (\eta_i - \eta_{i0}) \frac{\partial \mu_i}{\partial \eta_i}(\eta_{i0})$$

$$(ii) \quad \frac{\partial \ell_Y}{\partial \beta_k} \approx \sum_{i=1}^n \frac{\left(Y_i - \left[\underbrace{\mu_i(\eta_{i0})}_{\text{approx}} + (\eta_i - \eta_{i0}) \frac{\partial \mu_i}{\partial \eta_i}(\eta_{i0}) \right] \right) \frac{\partial \mu_i}{\partial \eta_i}(\eta_{i0})}{\left(\frac{\partial \eta_i}{\partial \mu_i} \right)^2 b(\theta_i) a(\varphi)}$$

$$Y_i^* = \frac{\partial \eta_i}{\partial \mu_i}(\eta_{i0}) = x_{ik}^t \beta$$

$$= \sum_{i=1}^n \frac{\left(Y_i - \mu_{i0} \frac{\partial \mu_i}{\partial \eta_i}(\eta_{i0}) + \eta_{i0} \right) - \eta_i}{\left(\frac{\partial \eta_i}{\partial \mu_i} \right)^2 b(\theta_i) a(\varphi)} x_{ik}$$

$$= \sum_{i=1}^n \frac{Y_i^* - x_{ik}^t \beta}{\left(\frac{\partial \eta_i}{\partial \mu_i} \right)^2 b(\theta_i) a(\varphi)} \cdot x_{ik}$$

almost linear
but still depend on β

$$\approx \sum_{i=1}^n \frac{Y_i^* - x_{ik}^t \beta}{\left(\frac{\partial \eta_i}{\partial \mu_i} \left[\mu_{i0} \right] \right)^2 b(\theta_i(\mu_{i0})) a(\varphi)} \cdot x_{ik}$$

$$\rightarrow w_i(\mu_{i0})$$

$$= \sum_{i=1}^n \frac{Y_i^* - x_{ik}^t \beta}{w_i(\mu_{i0})} \cdot x_{ik}$$

$$\frac{\partial \ell}{\partial \beta} \approx X^t W (Y^* - X\beta) = 0$$

$\xrightarrow{\text{nxm}} \text{nxm} \quad \text{n} \times 1$
 L ↴ diagonal matrix with $\frac{1}{w_i}$ on diagonal

$$X^t W Y^* - X^t W X \beta = 0$$

$$\beta = (X^t W X)^{-1} X^t W Y^*$$

Poisson Regression

$Y_i = \# \text{ particles that decayed in time interval } i$

$$f_Y(y) = e^{-\lambda} \frac{\lambda^y}{y!} \quad \text{where } y \in \mathbb{N}_0$$

$$= \exp(\log(e^{-\lambda}) \frac{\lambda^y}{y!})$$

$$= \exp(y \log(\lambda) - \log(y!) - \lambda)$$

$$= \exp(y \log(\lambda) - \lambda - \log(y!))$$

$$\stackrel{!}{=} \exp\left(\frac{y\theta - b(\theta)}{a(\epsilon)} + c(y, \epsilon)\right)$$

$$\rightarrow \theta = \log(\lambda), \quad b(\theta) = \exp(\theta), \quad a(\epsilon) = 1, \quad c(y, \epsilon) = -\log(y!)$$

link

$$\theta \xrightarrow{e^\theta} \lambda \xrightarrow{g} \eta = \vec{x}^\top \vec{\beta}$$

$$M = \{\lambda = EY\} = (0, \infty) = \text{set of all Poisson rates}$$

$$g: (0, \infty) \rightarrow \mathbb{R}$$

canonical link:

$$g(\lambda) = b^{-1}(\lambda) = \exp'(\lambda) = \log(\lambda) = \eta = \vec{x}^\top \vec{\beta}$$

offset

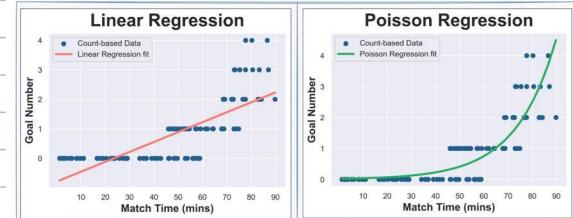
particles that decay depends on length of time interval:

$$Y \sim \text{Poisson}(\text{length} \cdot e^{\vec{x}^\top \vec{\beta}})$$

$$= \text{Poisson}(e^{\vec{x}^\top \vec{\beta} + \log(\text{length})})$$

no additional covariate (and ideal!), we do not estimate the effect of length

in Poisson Regression offset are quite common and need to be included "manually"



✗ Linear Regression

✓ Poisson Regression

Unsuitable for modeling count data

May return negative output for some inputs

Errors must follow a symmetric distribution

Specifically designed to model count data

All outputs are non-negative

Works well even if errors are asymmetrically distributed