$$\hat{\theta}_N = \underset{\theta \in \Theta \subseteq \mathbb{R}^L}{\text{argmax}} \quad E\left[\ell_i(y_i \mid x_i ; \theta)\right]$$

Average log-likelyhood

$$E_N[\cdot] = \sum_{i=1}^{N}[\cdot]$$

in this case
loglikelyhood function

estimator is implicit
function of the sample size

• $\theta$ is an $L \times 1$ vector of unknown parameters

returns a
scalar

$$E_{\theta_0}\left[\ell_i(\theta_0 ; y_i)\right]$$

$\hookrightarrow$ is the true loglikelyhood function

$$\hookrightarrow = \int_S \ell_i(y_i \mid x_i ; \theta) \, P(y_i \mid x_i ; \theta_0) \, dy$$

Example

$$\{x_i\}_{i=1}^{N} \quad iid \sim \exp(\lambda)$$

scalar, i
can draw on
a cartesian
axis

$$f_x(x_i, \lambda) = \lambda \exp(-\lambda x_i)$$

$$\downarrow$$

$$\ell_x(x_i, \lambda) = \log \lambda - \lambda x_i$$

loglikelyhood for each individual
observation

$$E_N\left[\ell_x(x_i, \lambda)\right] = \frac{1}{N} \sum \left[\log \lambda - \lambda x_i\right] = \log \lambda - \lambda \boxed{\frac{1}{N} \sum x_i}$$

$$E_{\theta_0}[x_i] = \frac{1}{\lambda_0}$$

$$= \log \lambda - \lambda \bar{x}$$

$$\hat{\theta} = \underset{\lambda > 0}{\text{argmax}}\left[\log \lambda - \lambda \bar{x}\right]$$

$$\frac{\delta E_N[x, \lambda]}{\delta \lambda} = \frac{1}{\lambda} - \bar{x} = 0$$

$$\hat{\lambda} = \frac{1}{\bar{x}} = \frac{N}{\sum x_i}$$

$$E_{\theta}\left[\ell(x_i, \lambda)\right] = \log \lambda - \frac{\lambda}{\lambda_0}$$

does not depend out anymore because we took expectation

(we realize we do not know) this quantity

↓

but at least i can study if theoretically the model is sound

$$\underset{\lambda > 0}{\text{argmax}}\left[\log \lambda - \frac{\lambda}{\lambda_0}\right]$$

true pop. exp.

$E_{\theta_0}[\theta; \theta_0]$

no more know the data

↓

$$\frac{d\, E_{\theta_0}\left[\ell(\lambda, \lambda_0)\right]}{d\lambda} = \frac{1}{\lambda} - \frac{1}{\lambda_0} = 0 \implies \boxed{\lambda = \lambda_0}$$

Exactly what we wanted

↓

the solution is <u>unique</u> !

$$E_N\left[\ell(y_i | x_i, \theta)\right] \xrightarrow{\ P\ } \underset{\theta_0}{E}\left[\ell(y_i | x_i; \theta]\right]$$

(sample average)          Convergence pointwise, for each value of θ          (population count)

⇓

$$\hat{\theta}_N \rightarrow \theta_0$$

$$E_N\left[\ell_i(y_i; \theta)\right] \xrightarrow[\text{pointwise}]{?} E_{\theta_0}\left[\ell_i(y_i; \theta)\right]$$

ASYMPTOTIC NORMALITY

$$\sqrt{N}\left(\hat{\theta}-\theta_0\right) \xrightarrow{d} N_L\left(0; \Sigma_{\theta_0}\right)$$

(LxL)

↳ multidimensional

$$\Sigma_{\theta_0} = -E_{\theta_0}\left[\frac{\partial^2(\ell(\theta))}{\partial\theta\partial\theta'}\Big|_{\theta=\theta_0}\right]^{-1}$$

inverse of the covariance is sort of a measure of the precision

$$E_N\left[\frac{\partial\ell(\theta)}{\partial\theta}\Big|_{\theta=\hat{\theta}}\right] = 0$$

by construction the empirical score evaluated in $\hat{\theta}$ is zero

the more curve is the function in $\theta_0$ the more we are precise with the estimation, the smaller is the curvature, the less the data are informative

EXPANSION OF THE SCORE AROUND $\theta_0$

$$E_N\left[\frac{\partial\ell(\theta)}{\partial\theta}\Big|_{\theta=\theta_0}\right] + E_N\left[\frac{\partial^2\ell(\theta)}{\partial\theta\partial\theta}\Big|_{\theta=\bar{\theta}}\right]\left(\hat{\theta}-\theta_0\right)$$

$= 0$  |  TAYLOR  $\hat{\theta}_N < \bar{\theta} < \theta_0$

$$f(x) = f(x_0) + f'(x_0)(x-x_0)$$

less curvature, less precision, more variance and uncertainty

less is to have a precise estimation and to not care about the other terms of the taylor polynomial

Now apply the CLT

MEAN VALUE THEOREM

$$\sqrt{N}\left(\hat{\theta}-\theta_0\right) = -E_N\left[\frac{\partial^2\ell(\theta)}{\partial\theta\partial\theta}\Big|_{\theta=\bar{\theta}}\right]^{-1}\sqrt{N}\;E_N\left[\frac{\partial\ell(\theta)}{\partial\theta}\Big|_{\theta=\theta_0}\right]$$

$$\sqrt{N}\;E_N\left[\frac{\partial\ell(\theta)}{\partial\theta}\Big|_{\theta=\theta_0}\right] \xrightarrow{d} N_L\left(0; -E_{\theta_0}\left[\frac{\partial^2\ell(\theta)}{\partial\theta\partial\theta'}\Big|_{\theta=\theta_0}\right]\right)$$

L×1

WHY THE MEAN OF THIS IS ZERO?

$$E_{\theta_0}\left(\frac{\partial\ell(\theta)}{\partial\theta}\Big|_{\theta=\theta_0}\right)$$

$$Var\left(\frac{\partial\ell(\theta)}{\partial\theta}\Big|_{\theta=\theta_0}\right) = E_{\theta_0}\left[\frac{\partial\ell(\theta)}{\partial\theta}\cdot\frac{\partial\ell(\theta)}{\partial\theta'}\Big|_{\theta=\theta_0}\right]$$

outer product

$$= -E_{\theta_0}\left[\frac{\delta^2 \ell(\theta)}{\delta\theta\delta\theta'}\Big|_{\theta=\theta_0}\right] \text{ Fisher's information}$$

$L \times L$

- plim $E_N\left[\frac{\delta^2 \ell(\theta)}{\delta\theta\delta\theta'}\Big|_{\theta=\bar\theta}\right] = E_{\theta_0}\left[\frac{\delta^2 \ell(\theta)}{\delta\theta\delta\theta'}\Big|_{\theta=\bar\theta}\right]$

$N\to\infty$

true expected Hessian
evaluated in $\bar\theta$

$$E_N[\cdot] \xrightarrow{P} E_{\theta_0}[\cdot]$$

WLLN
or
ULLN

if $\hat\theta \xrightarrow{P} \theta_0$, $\bar\theta = w\hat\theta + (1-w)\theta_0$

$w(0,1)$

$$\bar\theta \longrightarrow \theta_0$$

we call the expected Hessian $H$ and the expected
gradient $g$

$$\sqrt{N}(\hat\theta-\theta_0) \xrightarrow{d} -H^{-1}\boxed{\sqrt{N}\,\bar g} \to N(0; -H)$$

$$\sim N(0; H^{-1}\text{\st{XX}}H^{-1})$$

Outer product of the gradient estimator to estimate the inverse
of the Hessian

$$H = E_{\theta_0}\left[\frac{\delta^2 \ell(\theta)}{\delta\theta\delta\theta'}\Big|_{\theta=\bar\theta}\right]$$

$$g = E_{\theta_0}\left[\frac{\delta \ell(\theta)}{\delta\theta}\Big|_{\theta=\theta_0}\right]$$

$I \longrightarrow$ Fisher information

measure of the information we have about a parameter

it is the variance of the score

$$I = Var\left(\frac{\partial \ell_i(y_i;\theta)}{\partial \theta}\bigg|_{\theta=\theta_0}\right)$$

$$I = -E\left[\frac{\partial \ell_i^2(y_i;\theta)}{\partial\theta\partial\theta'}\bigg|_{\theta=\theta_0}\right]$$

measure of curvature of the function of models

the more the function is curve, less is the uncertainty about the parameter

If we have a sample of 100 obs. the loglikelihood for that sample will be a function of $\theta$

we want to maximise it, we do the first derivative and we pose it = 0

if $k>1$ the derivative is a vector (gradient) and the variance is a matrix (hessian)