

Statistical Data Analysis - Revision Lec. 2

Goodness of fit / significance test
(\rightarrow p-value)

$p_H = \text{Prob} \left(\begin{array}{l} \text{data in region} \\ \text{of equal or} \\ \text{lesser compat.} \end{array} \middle| H \right)$
w/ H \leftarrow more compat. w/ some H'

~~$\neq P(H)$~~

data values where p-value

$P_H \leq \alpha$ = a critical region for a test of H of size α .

If $P_H \leq \alpha \Rightarrow$ reject H .

Parameter estimation

given $x \sim f(x|\theta)$

construct $\hat{\theta}(\vec{x}) = \text{estimator for } \theta$

$\hat{\theta}$ is a r.v.

$$\hat{\theta} \sim g(\hat{\theta}|\theta)$$

bias $b = E[\hat{\theta}] - \theta$

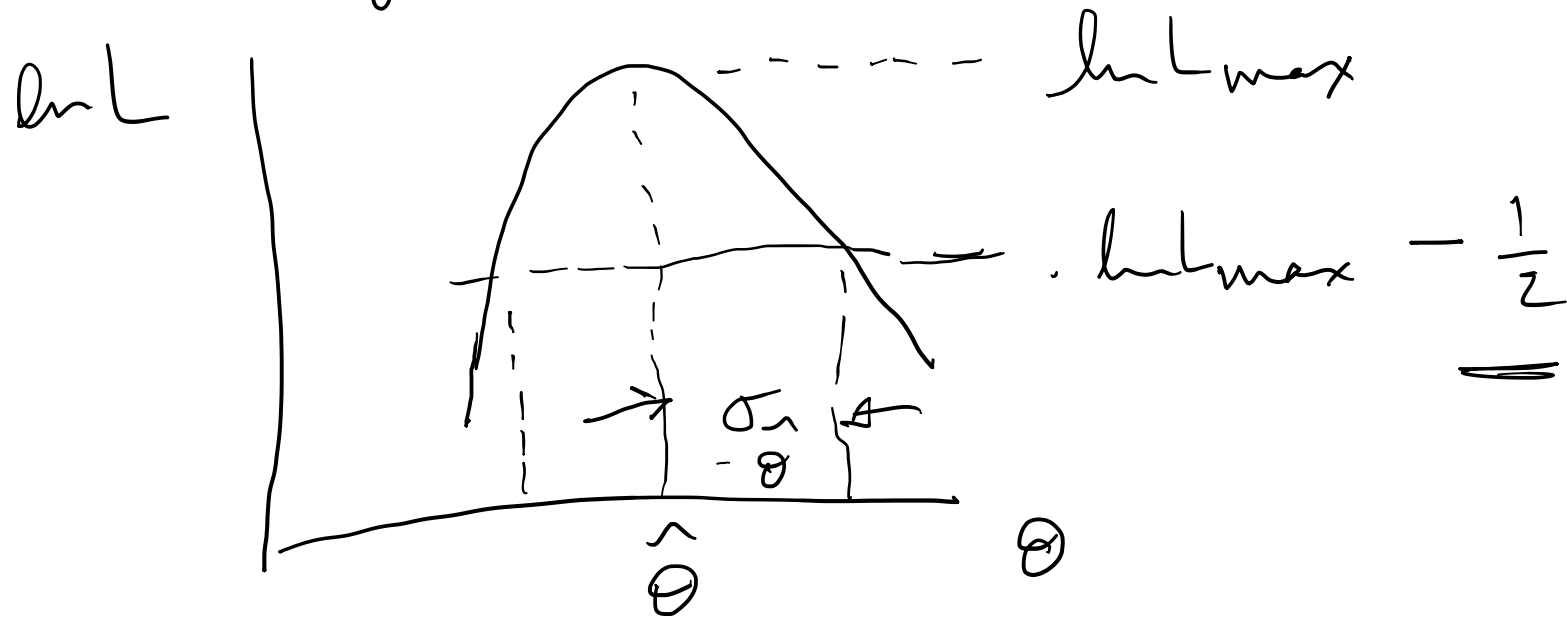
variance $V[\hat{\theta}]$

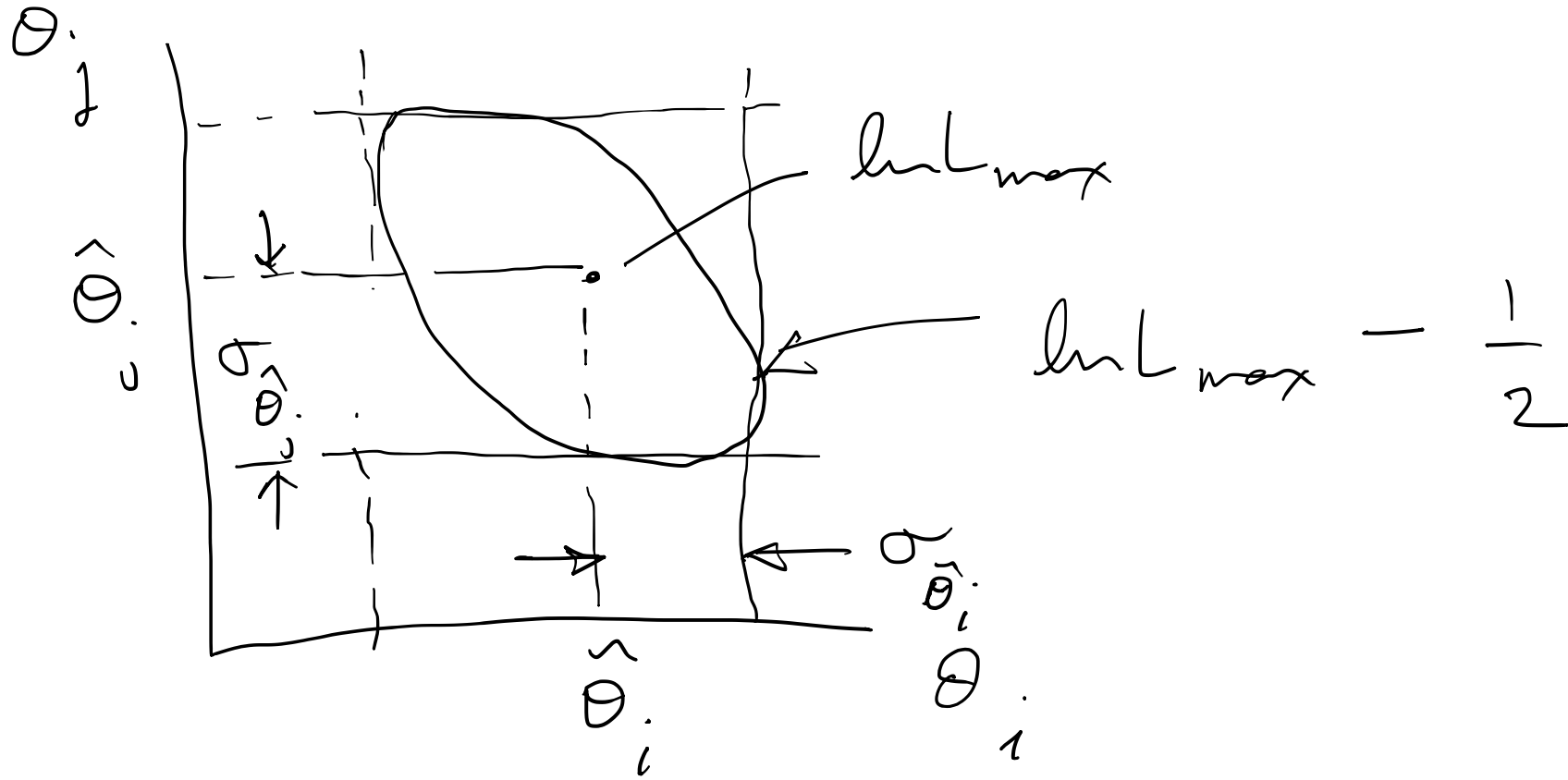
} competing measures
of quality

Likelihood $L(\theta) = P(x | \theta)$
↑ fixed.

Often use $\ln L(\theta)$

Method of max. likelihood: $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} L(\theta)$





Information inequality after b small

$$V[\hat{\theta}] \geq \text{MVB} = - \frac{\left(1 + \frac{\partial b}{\partial \theta}\right)^2}{E\left[\frac{\partial^2 \ln L}{\partial \theta^2}\right]}$$

↑
approx.
equality

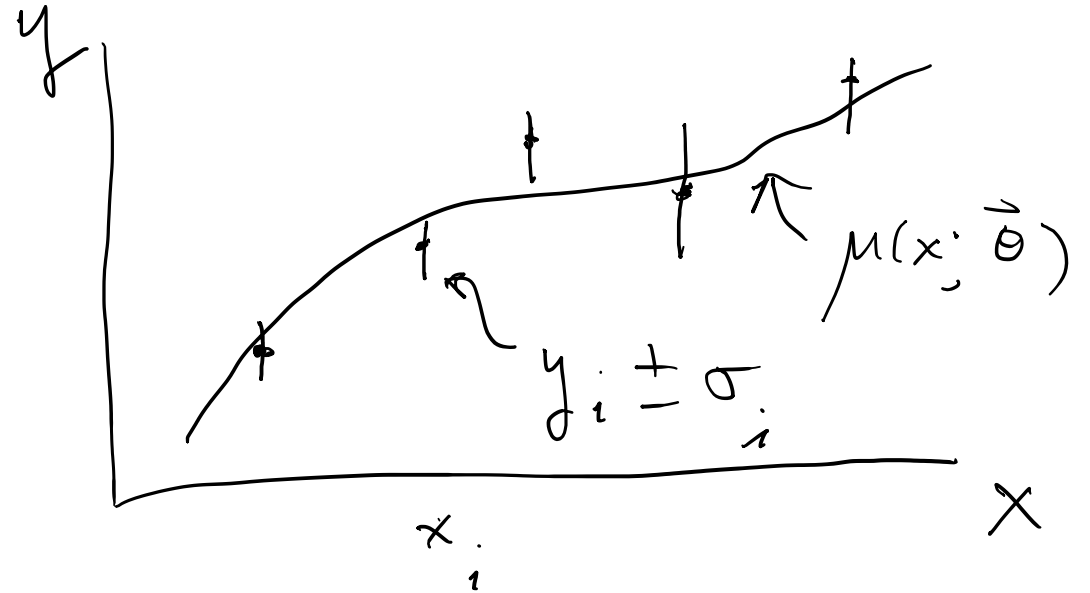
$$V[\hat{\theta}] \approx - \frac{1}{E\left[\frac{\partial^2 \ln L}{\partial \theta^2}\right]}$$

$$\hat{V}[\hat{\theta}] \approx - \frac{1}{\left(\frac{\partial^2 \ln L}{\partial \theta^2}\right)_{\hat{\theta}}}$$

Least Squares

$$y_i \sim \text{Gauss}(\mu(x_i; \vec{\theta}), \sigma_i)$$

↑
indep.



$$-2 \ln L(\vec{\theta}) = \underbrace{\sum_{i=1}^N \frac{(y_i - \mu(x_i; \vec{\theta}))^2}{\sigma_i^2}}_{\chi^2(\vec{\theta})} + C$$

Extended ML

$$n \sim \text{Poisson}(\nu)$$

$$x_1, \dots, x_n$$

$$x \sim f(x; \vec{\theta})$$

} indep.

$$L(\nu, \vec{\theta}) = \frac{\nu^n}{n!} e^{-\nu}$$

$$\prod_{i=1}^n f(x_i; \vec{\theta})$$



extended

Interval estimation

Freq. int est: Test each value of θ (size α)

Set of θ values not reject \equiv

confidence int for θ @ $CL = 1 - \alpha$

$$\uparrow P_{\theta} > \alpha$$

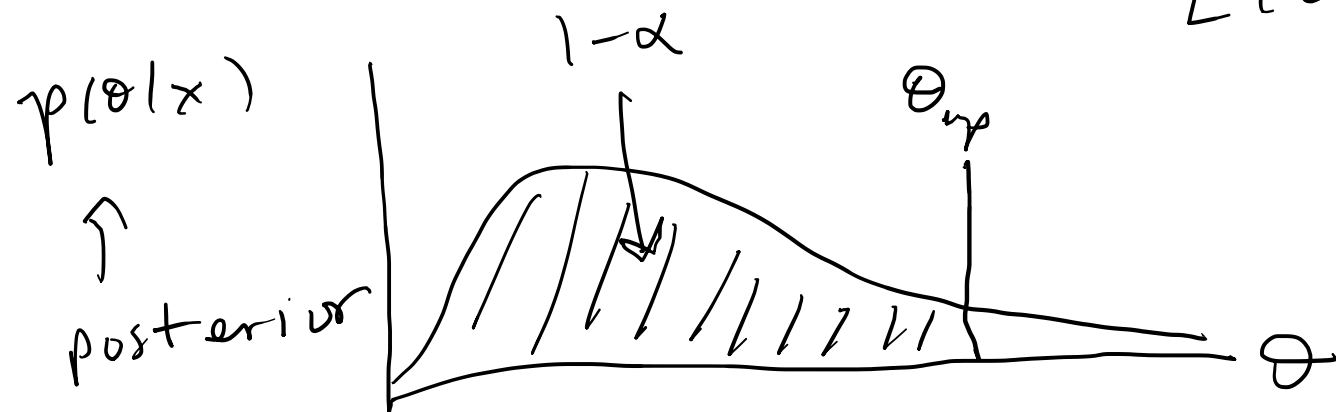
To find boundary of conf. int

set $p_{\theta} = \alpha$ & solve for θ .

Bayesian interval est.

$$p(\theta|x) \propto P(x|\theta) \pi(\theta)$$

\uparrow
 $L(\theta)$



$$\int_{-\infty}^{\theta_{yp}} p(\theta|x) d\theta = 1 - \alpha \quad \leftarrow \text{solve for } \theta_{yp}$$

More on intervals (multi param.)

Works' then

$$\vec{\theta} = (\theta_1, \dots, \theta_n) \rightarrow t_{\vec{\theta}} = -2 \ln \frac{L(\vec{\theta})}{L(\hat{\vec{\theta}})}$$

in "asymptotic limit", $f(t_{\vec{\theta}} | \vec{\theta})$ is χ_n^2 ↑ d.o.f.
↑ large sample, $\hat{\vec{\theta}} \sim \text{Gauss}$

⇒ multiparam cont. region

$$\ln L(\vec{\theta}) = \ln L_{\max} - \frac{1}{2} F_{\chi_n^2}^{-1}(1-\alpha)$$

↖ = # of param.

Special case $n=1$, $1-\alpha = 68.3\%$

$$F_{x_1}^{-1}(0.683) = 1$$

$$\ln L(\theta) = \ln L_{\max} - \frac{1}{2}$$

$[\hat{\theta} - \sigma_{\hat{\theta}}, \hat{\theta} + \sigma_{\hat{\theta}}]$ is a 68.3% CL
conf. int.

Bayes factors

$$B_{ij} = \frac{L_i}{L_j} = \frac{\int p(x|H_i, \theta) \pi(\theta) d\theta}{\int p(x|H_j, \gamma) \pi(\gamma) d\gamma}$$

marginal likelihoods

$$B_{ji} = \frac{1}{B_{ij}}$$

hard to compute
use e.g. harmonic mean
importance sampling

2015 Exam #4

$$n \sim \text{Poisson}(\mu + \nu) \quad \mu + \nu \geq 0$$

a) suppose ν known

$$L(\mu) = P(n|\mu) = \frac{(\mu + \nu)^n}{n!} e^{-(\mu + \nu)}$$

$$\ln L(\mu) = n \ln(\mu + \nu) - (\mu + \nu) + C$$

↑
doesn't depend
on params.

$$b) \quad \frac{d \ln L}{d \mu} = \frac{n}{\mu + \nu} - 1 \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \hat{\mu} = n - \nu$$

$$E[\hat{\mu}] = E[n - \nu] = E[n] - \nu = \mu + \nu - \nu = \mu$$

$$\Rightarrow b = E[\hat{\mu}] - \mu = 0 \quad \swarrow \text{for Poisson}$$

$$V[\hat{\mu}] = V[n - \nu] = V[n] = E[n] = \mu + \nu$$

$$\sigma_{\hat{\mu}} = \sqrt{\mu + \nu}$$

$$c) \text{ Test } \mu = 0 \quad w/ \quad t_0 = -2 \ln \lambda(0)$$

$$\lambda(\mu) = \frac{L(\mu)}{L(\hat{\mu})}$$

$\hat{\mu} = n - \nu$

For $\mu \neq 0$: $t_0 = -2 \left[\ln L(0) - \ln L(\hat{\mu}) \right]$

$$= -2 \left[n \ln \nu - \nu - n \ln (n - \nu + \nu) + \underline{n - \nu + \nu} \right]$$

$$= -2 \left[n \ln \frac{\nu}{n} + n - \nu \right] \quad \checkmark$$

$$n=0, \quad L(n=0 | \mu) = \frac{(\mu+\nu)^0}{0!} e^{-(\mu+\nu)}$$

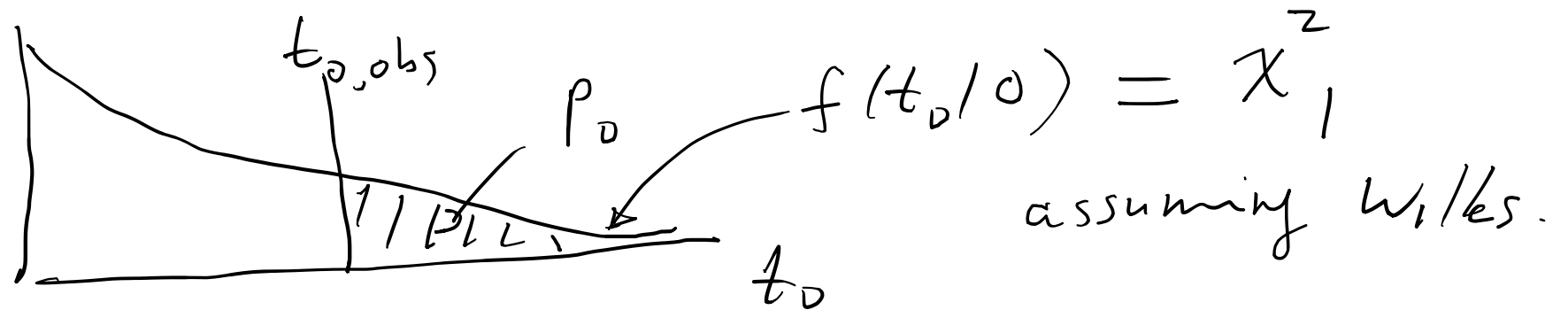
$$\ln L(n=0 | \mu) = -(\mu+\nu)$$

$$\begin{aligned} \underline{t_0(n=0)} &= -2 \left[\ln L(n=0 | \mu=0) - \ln L(n=0 | \mu=\hat{\mu} = -\nu) \right] \\ &= -2 \left[-(0+\nu) + \cancel{(-\nu+\nu)} \right] \\ &= \underline{2\nu} \quad \checkmark \end{aligned}$$

$e - ii)$

$$p_0 = P(t_0 \geq t_{0,obs} \mid \mu = 0)$$

$$= 1 - \int_0^{t_{0,obs}} f(t_0 \mid \mu = 0) dt_0$$



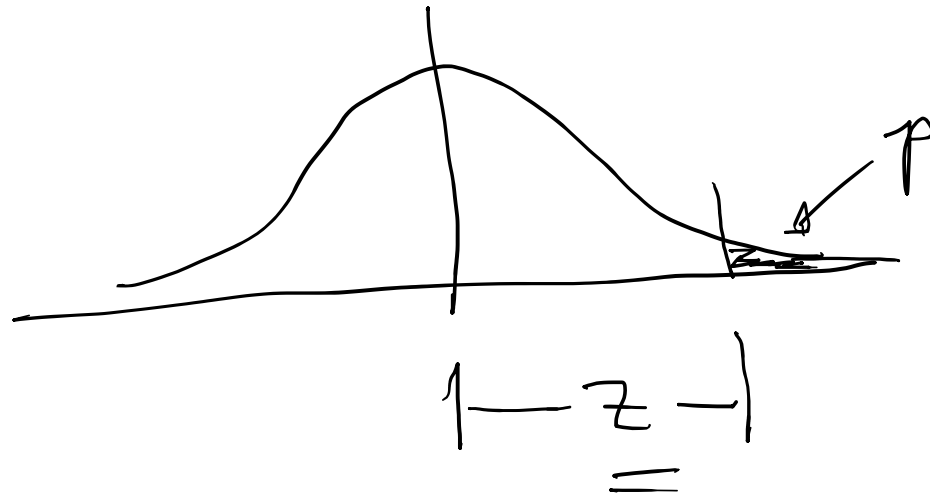
$$p_0 = 1 - F_{\chi^2_1}(t_0) \quad \leftarrow \text{obs}$$

Significance

$$p = 1 - \Phi(z)$$

$$z = \Phi^{-1}(1-p)$$

$$= \Phi^{-1}(F_{X_1}(t_0))$$



d) $\nu = \text{nuisance param.}$

$y \sim \text{Gauss}(\nu, \sigma)$ \leftarrow indep of n
 \uparrow known

$$L(\mu, \nu) = P(n, y_n | \mu, \nu)$$

$$= \frac{(\mu + \nu)^n}{n!} e^{-(\mu + \nu)} \times \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(y - \nu)^2}{2\sigma^2}}$$

$$\ln L(\mu, \nu) = n \ln(\mu + \nu) - (\mu + \nu) - \frac{1}{2} \frac{(y - \nu)^2}{\sigma^2} + C$$

$$\frac{\partial \ln L}{\partial \mu} = \frac{n}{\mu + \nu} - 1 \stackrel{\text{set}}{=} 0$$

$$\frac{\partial \ln L}{\partial \nu} = \left(\frac{n}{\mu + \nu} - 1 \right) + \frac{y - \nu}{\sigma^2}$$

$$\stackrel{\text{set}}{=} 0$$

} solve for $\mu \neq \nu$

$$\left. \begin{aligned} \hat{\mu} &= n - y \\ \hat{\nu} &= y \end{aligned} \right)$$

$$\begin{aligned}V[\hat{\mu}] &= V[n-y] = V[n] + V[y] \\ &= \mu + \nu + \sigma^2\end{aligned}$$

$$V[\hat{v}] = V[y] = \sigma^2$$

$$\begin{aligned}\text{cov}[\hat{\mu}, \hat{v}] &\Rightarrow \text{cov}[n-y, y] = \text{cov}[n, y] - \text{cov}[y, y] \\ &= -\sigma^2\end{aligned}$$

$$\rho_{\hat{\mu}, \hat{v}} = \frac{\text{cov}[\hat{\mu}, \hat{v}]}{\sigma_{\hat{\mu}} \sigma_{\hat{v}}}$$

$$= \frac{\sigma}{\sqrt{\mu + v + \sigma^2}}$$

