

Multi-parameter problems so far

STAT 206
16 Feb 21

θ has been a $(k \times 1)$ or $(1 \times k)$ vector with $k=1$ (i.e., a scalar) Do the likelihood

Lecture

and Bayesian stories still work with $k > 1$?

let's see Ex. NB 10 core study Ill-fitting but interesting model for practice with $k=2$

lec. notes p. 23

Gaussian sampling model The likelihood story: marginal PDF

$(Y_i | \mu, \sigma, N, B) \sim N(\mu, \sigma^2)$
 $(i=1, \dots, n)$ ($k=2$)
 $(\mu \in \mathbb{R})$
 $(\sigma > 0)$
 $Y = (Y_1, \dots, Y_n)$

$p(Y_i | \mu, \sigma, N, B) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left[-\frac{1}{2} \left(\frac{Y_i - \mu}{\sigma}\right)^2\right]$
 $Y = (Y_1, \dots, Y_n)$

joint PDF $p(Y | \mu, \sigma, N, B) = \prod_{i=1}^n p(Y_i | \mu, \sigma, N, B)$
 $= \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} \exp\left[-\frac{1}{2} \left(\frac{Y_i - \mu}{\sigma}\right)^2\right] = \sigma^{-n} (2\pi)^{-\frac{n}{2}}$

$\theta = (\mu, \sigma) (k=2)$ $\exp\left[-\frac{1}{2} \sum_{i=1}^n \left(\frac{Y_i - \mu}{\sigma}\right)^2\right]$

③ likelihood $l(\mu, \sigma | \mathbf{y} \sim NB) = \frac{1}{\sigma^n} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2\right]$ (c-1)

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

④ log likelihood

$$\begin{aligned}
 \text{ll}(\mu, \sigma | \mathbf{y} \sim NB) &= -n \log \sigma - \frac{1}{2} \sum_{i=1}^n \left(\frac{y_i - \mu}{\sigma}\right)^2 \\
 &= -n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2 \\
 &= -n \log \sigma - \frac{1}{2\sigma^2} \left(\sum_{i=1}^n y_i^2 - 2\mu \sum_{i=1}^n y_i + n\mu^2 \right)
 \end{aligned}$$

④.1 Visualize

likelihood & log likelihood functions with the NB10 data (see R demo)

To do this visualization, we have to define a grid for μ and a grid for σ ; (chicken & egg problem); let's get the MLEs first, postponing ④.1

⑤ MLEs with $k=2$ there are now 2 partial derivatives (with respect to μ and σ), so we now have a system of 2 equations

in 2 unknowns to solve:

$$\frac{\partial}{\partial \mu} \ell(\mu, \sigma | \mathcal{Z}, N, \mathcal{B}) = 0 \quad (3)$$

$$\frac{\partial}{\partial \sigma} \ell(\mu, \sigma | \mathcal{Z}, N, \mathcal{B}) = 0$$

$$\frac{\partial}{\partial \mu} \ell(\mu, \sigma | \mathcal{Z}, N, \mathcal{B}) =$$

$$-\frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \mu) = 0 \quad \text{iff } \mu = \hat{\mu}_{MLE} = \bar{y} \quad \checkmark$$

$$\left(\sum_{i=1}^n y_i \right) - n\mu = 0$$

$$\frac{\partial}{\partial \sigma} \ell(\mu, \sigma | \mathcal{Z}, N, \mathcal{B}) = \frac{n}{\sigma} - \frac{\sum_{i=1}^n 1}{\sigma^3}$$

$$= 0 \quad \text{iff (simplifying)} \quad \sigma^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \mu)^2$$

$$\hat{\sigma}_{MLE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}$$

1st equation $\hat{\mu}_{MLE} = \bar{y}$

$$\text{so } \hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$$

2nd equation $\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \mu_{MLE})^2$

This is

slightly different from the usual estimator of variance in the Gaussian model,

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2, \text{ which many people prefer}$$

because in repeated sampling s^2 is frequentist

unbiased for σ^2 : $E_{\mathcal{F}}(s^2) = \sigma^2$

This means that MLEs can be biased (4)
 although the general result (under mild regularity conditions) is that the bias (if any) of an MLE goes to 0 quickly:

$$\text{bias}(\hat{\sigma}_{MLE}^2) = \underline{\underline{O\left(\frac{1}{n}\right)}} \quad \left\{ \begin{array}{l} \text{here we have} \\ n \hat{\sigma}_{MLE}^2 = (n-1)s^2 \end{array} \right.$$

$$\text{so } \hat{\sigma}_{MLE}^2 = \underline{\underline{\left(\frac{n-1}{n}\right) s^2}} \text{ and } E_F(\hat{\sigma}_{MLE}^2) = \frac{n-1}{n} \sigma^2$$

$$\text{and } \text{bias}(\hat{\sigma}_{MLE}^2) = E_F(\hat{\sigma}_{MLE}^2) - \sigma^2 = -\frac{\sigma^2}{n} = \underline{\underline{O\left(\frac{1}{n}\right)}}, \text{ matching the general pattern}$$

with the NR10 data $\bar{y} = 404.6$

$$\hat{\sigma}_{MLE} = 6.4$$

How wide should we make the μ and σ

grids? A₁ Trial and error A₂ Standard error calculations

⑥ Estimated standard errors for the ⑤

MLEs

Q: what's the analogue of observed information when $k > 1$?

A: The second

partial derivatives of the log likelihood function now are stored in a $(k \times k)$ matrix

called the Hessian of $-l(\mu, \sigma | \mathcal{Z}, N, \mathcal{B})$:

So $H \triangleq \begin{pmatrix} \frac{d^2}{d\mu^2} l(-) & \frac{d^2}{d\mu d\sigma} l(-) \\ \frac{d^2}{d\sigma d\mu} l(-) & \frac{d^2}{d\sigma^2} l(-) \end{pmatrix}$

$\frac{d^2}{d\mu^2} l(-) = -\frac{n}{\sigma^2}$

recall that $\frac{d}{d\sigma} l(-) = -\sum_{i=1}^n \frac{y_i}{\sigma^2} + \frac{n\mu}{\sigma^2}$

recall that $\frac{d}{d\sigma} l(-) = \frac{n}{\sigma} - \frac{\sum_{i=1}^n (y_i - \mu)^2}{\sigma^3}$

SA $\frac{d^2}{d\sigma^2} l(-) = -\frac{n}{\sigma^2} + 3 \frac{\sum_{i=1}^n (y_i - \mu)^2}{\sigma^4}$

and the mixed partial becomes $\left(\frac{2\mu n - 2S}{\sigma^3} \right) (*)$

To get the Hessian in Wolfrum, let's $\textcircled{6}$ rewrite $-\ell(\mu, \sigma | \mathcal{Z}, N, \mathcal{B}) = \left[-n \log \sigma - \frac{\sum_{i=1}^n (y_i - \mu)^2}{2\sigma^2} \right]$

$$\begin{aligned}
 \text{as } -\ell(\mu, \sigma | \mathcal{Z}, N, \mathcal{B}) &= -n \log \sigma + \frac{1}{2\sigma^2} \left(\sum_{i=1}^n y_i^2 - 2\mu \sum_{i=1}^n y_i + n\mu^2 \right) \\
 &= -n \log \sigma + \frac{1}{2\sigma^2} (SS - 2\bar{y}\mu + n\mu^2)
 \end{aligned}$$

of $-\ell(\mu, \sigma)$ $\textcircled{7}$ Then the Hessian turns out to be

$$H = \begin{bmatrix} -\frac{n}{\sigma^2} & -\frac{2\mu n - 2rS}{\sigma^3} \\ -\frac{2\mu n - 2rS}{\sigma^3} & \frac{3\mu^2 n + rS^2 - 2\mu rS}{\sigma^4} - \frac{n}{\sigma^2} \end{bmatrix}$$

The observed information matrix \vec{I} is now obtained by evaluating H at $(\mu, \sigma) = (\hat{\mu}_{MLE}, \hat{\sigma}_{MLE})$; the result is after simplifying

$$\vec{I} = \begin{bmatrix} 1/\sigma^2 & 0 \\ 0 & 2/n \end{bmatrix} \text{ and the } (k \times 1) \text{ generalization of } \vec{V}(\vec{\theta}_{MLE}) = \vec{I}^{-1}$$

$= O(n)$

when $k=1$ is the covariance matrix

$$\vec{V}(\vec{\theta}_{MLE}) = (\vec{I})^{-1} \text{ (the inverse of the information matrix)}$$

Here this is $\vec{V}(\vec{\mu}_{MLE}, \vec{\sigma}_{MLE}) = \begin{bmatrix} 1/\sigma^2 & 0 \\ 0 & 2/n \end{bmatrix}$

(STAT 131)

Recall that a covariance matrix looks

like this

$$\vec{V}(\vec{\mu}_{MLE}, \vec{\sigma}_{MLE}) = \begin{bmatrix} \vec{\mu} & \vec{\sigma} \\ \vec{\sigma} & \vec{\sigma} \end{bmatrix} = \begin{bmatrix} \vec{V}(\vec{\mu}) & \vec{C}(\vec{\mu}, \vec{\sigma}) \\ \vec{C}(\vec{\mu}, \vec{\sigma}) & \vec{V}(\vec{\sigma}) \end{bmatrix}$$

in which

$\vec{C}(\vec{\mu}, \vec{\sigma}) =$ the covariance between $\vec{\mu}_{MLE}$ and $\vec{\sigma}_{MLE}$, here $= 0$ (!)

NB 10

$\left\{ \begin{matrix} y_1 \\ \vdots \\ y_n \end{matrix} \right\} \leftarrow 375 \leftarrow \text{micrograms below } 105 \text{ } (8)^7$
 $n=100$
 $\leftarrow 437$

A: How choose good sampling model?

$(Y_i | \boxed{?}) \stackrel{i.i.d.}{\sim} \boxed{?}$
 $(i=1, \dots, n)$

A: frequently used up through ~ 1985

histogram, normal qq plot, we cheat (d.m. 2)

$S = \sum_{i=1}^n Y_i$

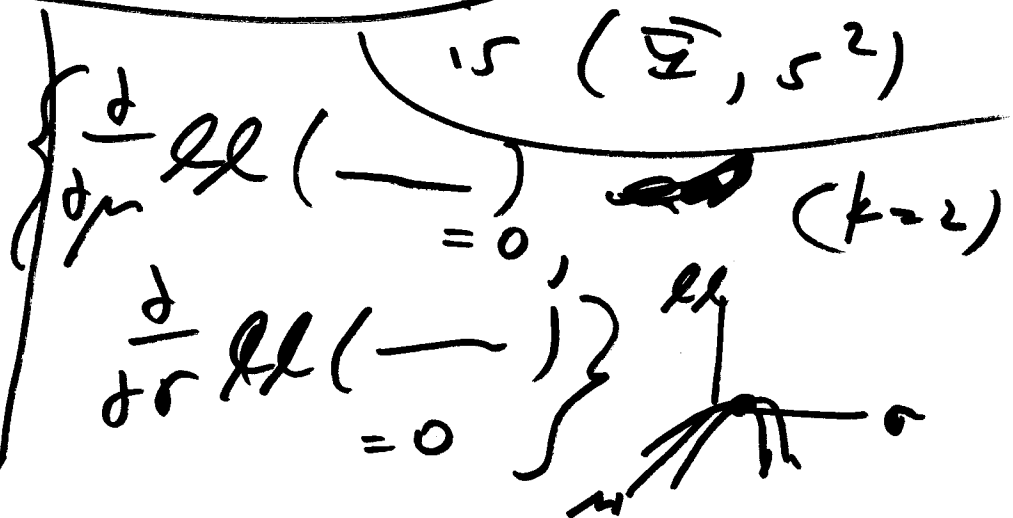
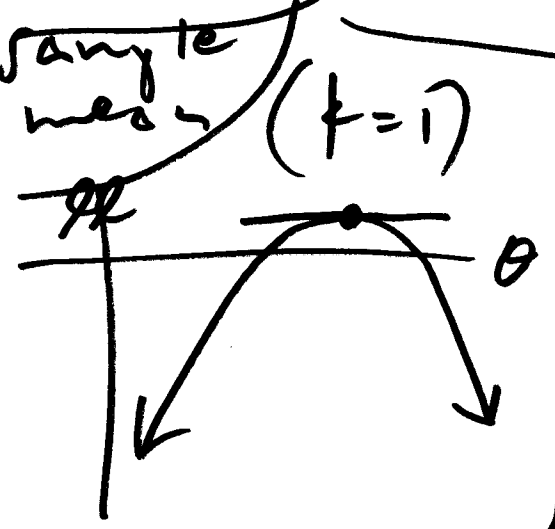
$SS = \sum_{i=1}^n Y_i^2$

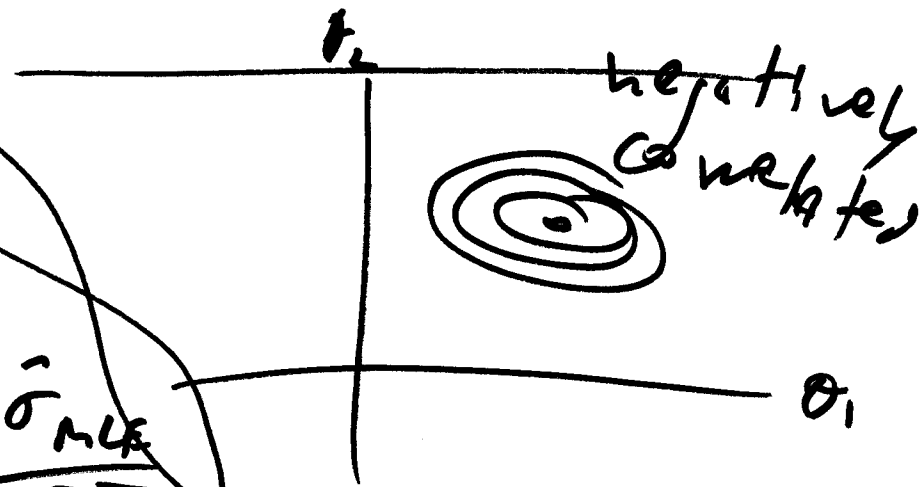
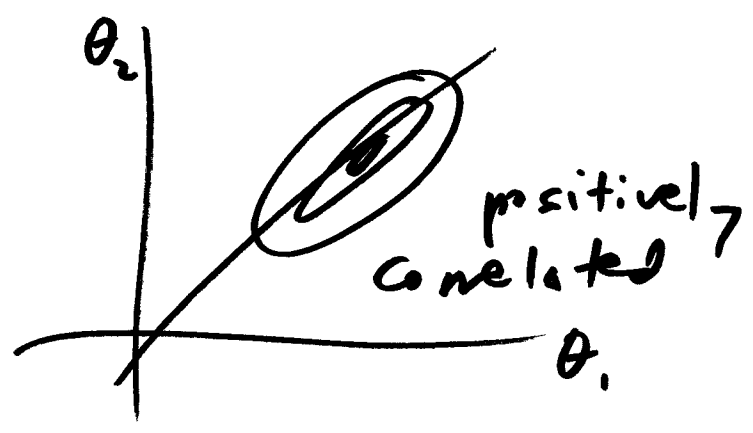
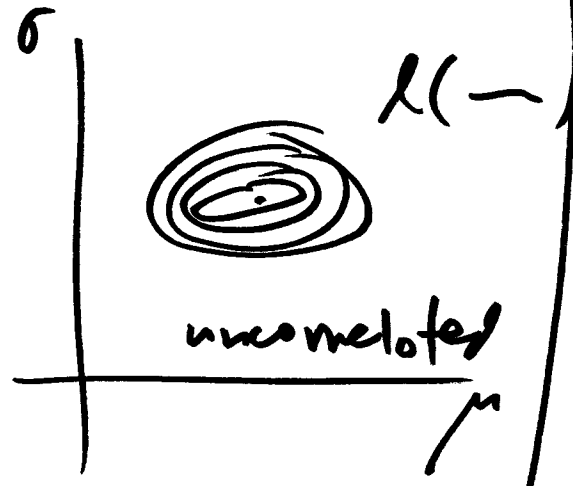
(S, SS) form a minimal suff. stat. ($k=2$)

$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$
 $= \frac{S}{n}$

$s^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$
 sample variance

equiv. set of m.s.s.





$SE(\hat{\mu}_{MLE}) =$

$\sqrt{\hat{V}(\hat{\mu}_{MLE})} = \frac{\hat{\sigma}_{MLE}}{\sqrt{n}}$

approx.

100(1- α)% CI for μ :

$\hat{\mu}_{MLE} \pm \left[\Phi^{-1}\left(1 - \frac{\alpha}{2}\right) \right] \left[SE(\hat{\mu}_{MLE}) \right]$

$SE(\hat{\sigma}_{MLE})$

$= \frac{\hat{\sigma}_{MLE}}{\sqrt{2n}}$

100(1- α)% CI for σ :

$\hat{\sigma}_{MLE} \pm \left[\Phi^{-1}\left(1 - \frac{\alpha}{2}\right) \right] \left[SE(\hat{\sigma}_{MLE}) \right]$

for $k \geq 1$ under mild regularity conditions ⁽¹⁰⁾

$\hat{\theta}_{MLE}$
in repeated sampling

$$N_k \left(\theta, \frac{1}{k} \hat{\Sigma}_{MLE} \right)$$

