

Ratios

one population of size N and two measurements resulting N data

$$X_1, X_2, \dots, X_N \text{ and } Y_1, Y_2, \dots, Y_N.$$

Population parameters:

Population covariance

$$\sigma_{xy} = \frac{1}{N} \sum_{i=1}^N (X_i - \mu_x)(Y_i - \mu_y).$$

two possible goals: ① estimate R via S.R.S
② Often, X values are known, and the ratio estimate can be used to give an estimate of $\mu_y = \frac{\mu_x}{R} \bar{y}$

If X_1, X_2, \dots, X_n and Y_1, Y_2, \dots, Y_n are simple random samples

$$\text{cov}(\bar{X}, \bar{Y}) = \frac{\sigma_{xy}}{N} \left(\frac{N-n}{N-1} \right)$$

$$R = \frac{\frac{1}{n} \sum_{i=1}^n Y_i}{\frac{1}{n} \sum_{i=1}^n X_i} = \frac{\bar{Y}}{\bar{X}} = g(\bar{X}, \bar{Y})$$

We need to estimate $\text{Var}(R)$

$$\text{Var}(R) = \text{Var}(g(\bar{X}, \bar{Y}))$$

$$E(R) = \frac{1}{n} \left(\frac{n-1}{N-1} \right)$$

In the one variable case:

$$Z = g(x) \rightarrow \text{Var}(Z)$$

$$Z = g(x) = g(\mu_x) + g(x - \mu_x) g'(\mu_x) + \frac{1}{2} (x - \mu_x)^2 g''(\mu_x) \quad \text{— Taylor expansion}$$

$$E(Z) = g(\mu_x) + \frac{\sigma_x^2}{2} g''(\mu_x)$$

In the two variables case:

$$Z = g(X, Y) \quad \text{we expand } g(X, Y) \text{ around } \underbrace{g(\mu_x, \mu_y)}_{\vec{\mu}}$$

$$R = g(\bar{X}, \bar{Y}) = g(\vec{\mu}) + (\bar{X} - \mu_x) \frac{\partial g(\vec{\mu})}{\partial x} + (\bar{Y} - \mu_y) \frac{\partial g(\vec{\mu})}{\partial y} + \frac{1}{2} (\bar{X} - \mu_x)^2 \frac{\partial^2 g(\vec{\mu})}{\partial x^2} + \frac{1}{2} (\bar{Y} - \mu_y)^2 \frac{\partial^2 g(\vec{\mu})}{\partial y^2} + (\bar{X} - \mu_x)(\bar{Y} - \mu_y) \frac{\partial^2 g(\vec{\mu})}{\partial x \partial y}$$

$$E[Z] = g(\mu_x, \mu_y) + \frac{\sigma_x^2}{2} \frac{\partial^2 g(\vec{\mu})}{\partial x^2} + \frac{\sigma_y^2}{2} \frac{\partial^2 g(\vec{\mu})}{\partial y^2} + \sigma_{xy} \frac{\partial^2 g(\vec{\mu})}{\partial x \partial y}$$

$$E[R] = \frac{\mu_Y}{\mu_X} + \frac{\sigma_x^2}{2} \frac{2\bar{Y}\bar{Y}}{\mu_X^3} + \frac{\sigma_y^2}{2} \cdot 0 + \sigma_{xy} \frac{1}{\mu_X^2}$$

$$= \frac{\mu_Y}{\mu_X} + \sigma_x^2 \frac{\mu_Y}{\mu_X^3} - \sigma_{xy} \frac{1}{\mu_X^2} = \frac{\mu_Y}{\mu_X} + \frac{\sigma_x^2 \mu_Y - \sigma_{xy} \mu_X}{\mu_X^3}$$

$$(1) E(R) = r + \frac{1}{n} \left(1 - \frac{n-1}{N-1}\right) \cdot \frac{r\sigma_x^2 + \cancel{\sigma_y^2} - \rho\sigma_x\sigma_y}{\mu_x^2}$$

$$(2) \text{Var}[R] = \frac{1}{n} \left(1 - \frac{n-1}{N-1}\right) \cdot \frac{r^2\sigma_x^2 + \sigma_y^2 - 2\rho\sigma_x\sigma_y}{\mu_x^2}$$

(1) comes from a 1st order Taylor expansion of $f(x, y) = \frac{y}{x}$ around the point (μ_x, μ_y) and the squaring.

(2) comes from a 2nd order Taylor expansion of $f(x, y) = \frac{y}{x}$ around (μ_x, μ_y) .

$$f_x(\mu_x, \mu_y) = -\frac{\mu_y}{\mu_x^2}, \quad f_y(\mu_x, \mu_y) = \frac{1}{\mu_x}, \quad f_{xx}(\mu_x, \mu_y) = \frac{2\mu_y}{\mu_x^3}, \quad f_{xy}(\mu_x, \mu_y) = -\frac{1}{\mu_x^2}$$

$$f_{yy}(\mu_x, \mu_y) = 0$$

$$\bar{\sigma}_x^2 = \frac{\sigma_x^2}{n} \left(1 - \frac{n-1}{N-1}\right) \Rightarrow \text{Page. 21. } \Rightarrow (1)$$

$$\bar{\sigma}_{\bar{x}\bar{y}} = \frac{\sigma_{xy}}{n} \left(1 - \frac{n-1}{N-1}\right)$$

$$\text{Var}(R) \equiv (R - r)^2 = \frac{\mu_y^2}{\mu_x^4} \bar{\sigma}_x^2 + \frac{1}{\mu_x^2} \bar{\sigma}_y^2 - \frac{2\mu_y}{\mu_x^3} \bar{\sigma}_{\bar{x}\bar{y}} \Rightarrow (2)$$

- C.I.

$$\bar{S}_R^2 = \frac{1}{n} \left(1 - \frac{n-1}{N-1}\right) \cdot \frac{R^2 S_x^2 + S_y^2 - 2RS_{xy}}{\bar{x}^2}, \text{ where } S_x^2 = \frac{1}{n-1} \sum (X_i - \bar{x})^2, S_y^2 = \frac{1}{n-1} \sum (Y_i - \bar{y})^2$$

$$\frac{R - r}{\sqrt{\text{Var}(R)}} \sim N(0, 1), \quad R \pm z\left(\frac{\alpha}{2}\right) S_R, \quad 100(1-\alpha)\% \text{ C.I. for } r.$$

$$\bar{Y}_R = \mu_x R = \frac{\mu_x}{\bar{x}} \bar{y} \quad \text{Not an unbiased estimator}$$

$$\text{Var}(\bar{Y}_R) = \mu_x^2 \text{Var}(R) = \frac{1}{n} \left(1 - \frac{n-1}{N-1}\right) (r^2 \sigma_x^2 + \sigma_y^2 - 2r\sigma_{xy}) \quad (\times)$$

$$\text{Var}(\bar{y}) = \frac{\sigma_y^2}{n} \left(1 - \frac{n-1}{N-1}\right) = \frac{1}{n} (r^2 \sigma_x^2 + \sigma_y^2 - 2r\sigma_{xy}) \left(1 - \frac{n-1}{N-1}\right) \text{ ignore finite population correction.}$$

when $r^2 \sigma_x^2 < 2r\sigma_{xy}$, $\text{Var}(\bar{Y}_R) < \text{Var}(\bar{y})$.

$$r\sigma_x^2 < 2\rho\sigma_x\sigma_y \Rightarrow \rho > \frac{1}{2} r \cdot \frac{\sigma_x}{\sigma_y} = \frac{1}{2} \frac{C_x}{C_y} = \frac{1}{2} \frac{C_x}{C_y} \Rightarrow \bar{Y}_R \text{ has smaller standard error than } \bar{y}$$

$$C_x, C_y : \text{coefficients of variation}$$

Ex. X_i : # beds in hospital i . Y_i : # discharges in hospital i in 1968/01

$$\begin{array}{lll} M_x = 274.8 & M_y = 814.6 & C_x = 20.78 \\ n = 64 & \sigma_x = 213.2 & C_y = 0.72 \\ N = 393 & \sigma_y = 589.7 & \frac{C_x}{2} \frac{C_y}{C_y} = 0.54 < \rho \\ r = 2.96 & \rho = 0.91 & \end{array}$$

$$\sigma_{\bar{Y}_R} = \frac{\sigma_y}{\sqrt{n}} \sqrt{1 - \frac{1}{N-1}} = \sqrt{\text{Var}(\bar{Y}_R)} = 81.46 \quad \text{page 22.} = 30 \quad (\text{book P.227})$$

$$\sigma_{\bar{Y}} = \frac{\sigma_y}{\sqrt{n}} \sqrt{1 - \frac{1}{N-1}} = 66.3$$

Stratified Sampling

The population is partitioned into L subpopulations. N_1, N_2, \dots, N_L .

$$N_1 + N_2 + \dots + N_L = N$$

$$w_l = \frac{N_l}{N} \text{ weight.}$$

Population measurements and parameters

$X_{11}, X_{12}, \dots, X_{N_11}$. Measurement of individuals in subgroup 1.

$$\mu_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} X_{1i} \quad \sigma_1^2 = \frac{1}{N_1} \sum_{i=1}^{N_1} (X_{1i} - \mu_1)^2$$

$X_{12}, X_{22}, \dots, X_{N_22}$

$$\mu_2 = \frac{1}{N_2} \sum_{i=1}^{N_2} X_{2i} \quad \sigma_2^2 = \frac{1}{N_2} \sum_{i=1}^{N_2} (X_{2i} - \mu_2)^2$$

\vdots
 $X_{1L}, X_{2L}, \dots, X_{NL}$

$$\mu = \frac{1}{N} \sum_{i=1}^L \sum_{e=1}^{N_e} X_{ie} = \frac{1}{N} \sum_{e=1}^L N_e \mu_e = \sum_{e=1}^L \mu_e w_e$$

We want to estimate μ and stratify the population in a way that minimizes the standard error of the estimate.

We select a s.r.s from each strata. n_e individuals from N_e .

$$X_{1,1}, X_{2,1}, \dots, X_{N_e,1}$$

$\{X_{11}, X_{21}, \dots, X_{N_11}\}$ across
 $\{X_{12}, X_{22}, \dots, X_{N_22}\}$ independent, not independent within strata

$$\{X_{1L}, X_{2L}, \dots, X_{NL}\}$$

Our estimate for μ is $\bar{X} = \sum_{e=1}^L \bar{X}_e w_e$. $\bar{X}_e = \frac{1}{n_e} \sum_{i=1}^{n_e} X_{ie}$ unbiased estimator for $\mu_e = E(\bar{X}_e)$

$$E[\bar{X}] = \mu.$$

$$\text{Var}(\bar{X}_S) = \text{Var}\left(\sum_{e=1}^L \bar{X}_e w_e\right) \quad \text{since } \bar{X}_1, \bar{X}_2, \dots, \bar{X}_L \text{ are iid.} \quad | 24.$$

$$= \sum_{e=1}^L \text{Var}(w_e \bar{X}_e) = w_e^2 \sum_{e=1}^L \text{Var}(\bar{X}_e) = \sum_{e=1}^L w_e^2 \frac{\sigma_e^2}{n_e} \left(1 - \frac{n_e - 1}{N_e}\right).$$

$$\text{Var}(\bar{X}_S) \approx \sum_{e=1}^L \frac{\sigma_e^2}{n_e} \quad \text{if } \sigma_e^2 \text{ is the same for each strata.}$$

$$\text{Var}(\bar{X}_S) = \sum_{e=1}^L \frac{\sigma_e^2}{n_e} w_e^2 \quad \text{if } \sigma_1^2, \sigma_2^2, \dots, \sigma_L^2 \text{ are not the same.}$$

- choose n_1, n_2, \dots, n_L to minimize $\text{Var}(\bar{X}_S)$.

$$\min f(n_1, n_2, \dots, n_L) = \sum_{e=1}^L \frac{\sigma_e^2 w_e^2}{n_e} \quad \text{subject to } g(n_1, n_2, \dots, n_L) = n_1 + n_2 + \dots + n_L = n$$

The method of Lagrange multipliers:

extreme values $\nabla f = \lambda \nabla g$
occur where

$$\nabla f = \left(-\frac{\sigma_1^2 w_1^2}{n_1^2}, -\frac{\sigma_2^2 w_2^2}{n_2^2}, \dots, -\frac{\sigma_L^2 w_L^2}{n_L^2}\right)$$

$$\nabla g = (1, 1, \dots, 1)$$

$$\frac{\sigma_1^2 w_1^2}{n_1^2} = \dots = \frac{\sigma_L^2 w_L^2}{n_L^2} = \frac{\sigma_e^2 w_e^2}{n_e^2} \Rightarrow \text{if } n_e = c_e n$$

$$\lambda = \frac{n}{\sum_{k=1}^L \sigma_k^2 w_k} \Rightarrow n_e = \frac{\sigma_e w_e}{\sum_{k=1}^L \sigma_k^2 w_k} n$$

| Stratum | N_e | w_e | M_e | σ_e | optimal sampling allocation |
|---------|-------|-------|--------|------------|-----------------------------|
| A | 98 | 1249 | 182.9 | 103.4 | $n_1 = 0.106n$ |
| B | 98 | 259 | 526.9 | 204.8 | $n_2 = 0.210n$ |
| C | 98 | 249 | 456.3 | 243.5 | $n_3 = 0.250n$ |
| D | 99 | 251 | 1591.2 | 419.2 | $n_4 = 0.434n$ |

If you don't know $\sigma_1^2, \dots, \sigma_L^2$, one method is to use proportional sampling

$$n_e = w_e n = \frac{N_e}{N} n$$

Parameter estimation

[25]

It can sometimes be done by expressing the parameters as functions of moments

$$\textcircled{1} \quad X \sim P(\lambda)$$

$$\mu_1 = E(X) = \lambda \quad \lambda = f(\mu_1) = \mu_1$$

$$\textcircled{2} \quad X \sim N(\mu, \sigma^2)$$

$$\mu_1 = E[X], \quad \mu_2 = E[X^2] \quad \mu = \frac{\mu_1}{f(\mu_1)}, \quad \sigma^2 = \mu_2 - \mu_1^2$$

$$\textcircled{3} \quad X \sim T(\alpha, \lambda)$$

$$\mu_1 = E[X] = \frac{\alpha}{\lambda}, \quad \mu_2 = E[X^2] = \frac{(\alpha+1)\alpha}{\lambda^2}$$

$$\lambda = f_1(\mu_1, \mu_2) = \frac{\mu_1}{\mu_2 - \mu_1^2}$$

$$\lambda = f_2(\mu_1, \mu_2) = \frac{\mu_1}{\mu_2 - \mu_1^2}$$

Method of Moments

$$\hat{\mu}_1 = \frac{1}{n} \sum_{k=1}^n X_k \quad \text{estimates } \mu_1, \quad \hat{\mu}_2 = \frac{1}{n} \sum_{k=1}^n X_k^2 \quad \text{estimates } \mu_2 \quad \text{in the sense}$$

that.

$$P(|\hat{\mu}_1 - \mu_1| > \varepsilon) \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

$$P(|\hat{\mu}_2 - \mu_2| > \varepsilon) \rightarrow 0$$

Moments can be approximated by sample moments by the law of large numbers

when $P(|\hat{\theta}_j - \theta_j| > \varepsilon) \rightarrow 0$, $\forall \varepsilon > 0$, as $n \rightarrow \infty$. $\hat{\theta}_j$ is called a consistent estimator of θ_j .

If f_j is continuous where $\theta_j = f_j(\mu_1, \mu_2, \dots, \mu_e)$, then $\hat{\theta}_j = f_j(\hat{\mu}_1, \hat{\mu}_2, \dots, \hat{\mu}_e)$ is a consistent estimator of θ_j .

Ex. If $X \sim N(\mu, \sigma^2)$, $\mu_1 = E[X]$, $\mu_2 = E[X^2]$

$$\mu = \mu_1 = f_1(\mu_1, \mu_2)$$

$$\sigma^2 = \mu_2 - \mu_1^2 = f_2(\mu_1, \mu_2) \quad \text{since } f_1, f_2 \text{ are continuous,}$$

$$\hat{\mu} = \hat{\mu}_1 = \frac{1}{n} \sum_{i=1}^n X_i \quad \hat{\sigma}^2 = \hat{\mu}_2 - \hat{\mu}_1^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \left(\frac{1}{n} \sum_{i=1}^n X_i \right)^2 \text{ are consistent estimators.}$$

The standard error $\sqrt{\text{Var}(\text{estimator})}$ is a measure of accuracy of the estimator.

In normal dist. the dist'n of $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i \sim N(\mu, \frac{\sigma^2}{n})$. standard error

Since we don't know σ^2 , we resort to $\frac{\hat{\sigma}^2}{n}$ as the S.E.

The dist'n of $\hat{\sigma}^2$, $\frac{n\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-1}$

$$X = Z_1^2 + Z_2^2 + \dots + Z_n^2 \quad \text{if } Z_i \text{ i.i.d. } \sim N(0, 1)$$

$$\text{Var}(X) = E[X^2] - E[X]^2 = 2n$$

$$\begin{aligned} E[X^2] &= E\left[\left(\sum_{i=1}^n Z_i^2\right)^2\right] = E\left[\sum_{i=1}^n Z_i^4 + \sum_{i \neq j} Z_i^2 Z_j^2\right] \\ &= nE[Z_i^4] + n(n-1). \\ &= 3n + n(n-1) = n^2 + 2n \end{aligned}$$

$$\text{Var}(\hat{\sigma}^2) = \text{Var}\left(\frac{\sigma^2}{n} \cdot \frac{n\hat{\sigma}^2}{\sigma^2}\right) = \frac{\sigma^4}{n^2} \text{Var}\left(\frac{n\hat{\sigma}^2}{\sigma^2}\right) = \frac{\sigma^4}{n^2} 2(n-1)$$

$$\text{Standard Error } \sqrt{\frac{2(n-1)}{n^2}} \sigma^2 \approx \sqrt{\frac{2(n-1)}{n^2}} \hat{\sigma}^2$$

Maximum Likelihood Estimator (MLE)

$$L(\theta_1, \theta_2, \dots, \theta_k) = \log f(x_1, \dots, x_n | \theta_1, \dots, \theta_k)$$

$$\underset{\theta_1, \theta_2, \dots, \theta_k}{\text{argmax}} L(\theta_1, \theta_2, \dots, \theta_k) = \hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k \quad \text{where } x_1, x_2, \dots, x_n \text{ are i.i.d.} \sim N(\mu, \sigma^2)$$

$\nabla L(\theta_1, \theta_2, \dots, \theta_k)$.

$$\nabla L(\mu, \sigma^2) = \left(\sum_{i=1}^n \frac{x_i - \mu}{\sigma}, -\frac{\mu}{\sigma} + \sigma^{-2} \sum_{i=1}^n (x_i - \mu)^2 \right) \stackrel{\text{set}}{=} 0.$$

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x} \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

can be used to find C.I. for μ and σ^2

$$\frac{\sqrt{n}(\bar{x} - \mu)}{S} \sim t_{n-1}$$

to obtain a $100(1-\alpha)\%$ C.I. for μ , find

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$t_{n-1}(\frac{\alpha}{2})$ and then

$$P\left(-t_{n-1}(\frac{\alpha}{2}) \leq \frac{\sqrt{n}(\bar{x} - \mu)}{S} \leq t_{n-1}(\frac{\alpha}{2})\right) = 1 - \alpha.$$

where $t_{n-1}(\frac{\alpha}{2})$ is the value for which

To obtain a $100(1-\alpha)\%$ CI for σ^2

27

$$\frac{n\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-1}$$

$$P\left(\chi^2_{n-1}(1-\frac{\alpha}{2}) \leq \frac{n\hat{\sigma}^2}{\sigma^2} \leq \chi^2_{n-1}(1+\frac{\alpha}{2})\right) = 1-\alpha.$$

$$\left[\frac{n\hat{\sigma}^2}{\chi^2_{n-1}(1-\frac{\alpha}{2})}, \frac{n\hat{\sigma}^2}{\chi^2_{n-1}(1+\frac{\alpha}{2})} \right] \text{ is a } 100(1-\alpha)\% \text{ CI for } \sigma^2$$

— CI for MLE's derived from iid samples when sampling dist'n of the MLE is not known

This method requires large samples since it relies on the CLT

For an iid sample, X_1, X_2, \dots, X_n with density $f(x|\theta)$

$$L(\theta) = \sum_{i=1}^n \log f(X_i|\theta)$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} L(\theta) \quad \text{Fisher Information}$$

$$\text{Define: } I(\theta) = E\left[\left(\frac{d}{d\theta} \log f(X|\theta)\right)^2\right] \text{ where } X \text{ has density } f(x|\theta)$$

CLT for MLE (iid case).

If θ_0 is the true value of the parameter, then

$$\sqrt{n}I(\theta_0)(\hat{\theta} - \theta_0) \sim N(0, 1) \text{ as } n \rightarrow \infty$$

$$P\left(-z\left(\frac{\alpha}{2}\right) \leq \sqrt{n}I(\theta_0)(\hat{\theta} - \theta_0) \leq z\left(\frac{\alpha}{2}\right)\right) \approx 1-\alpha.$$

$$I(\theta) = -E\left[\frac{d^2}{d\theta^2} \log f(X|\theta)\right] \text{ for calculating } I(\theta) \text{ easier.}$$

Bayesian
Approach.
parameter
estimation

prior distribution θ

view the data to update the opinion of the data

→ posterior distribution

$$f_{\theta}(\theta) = 1$$

$$f_{x|\theta}(x|\theta) = \begin{cases} \theta & x=1 \\ 1-\theta & x=0 \end{cases} \quad \theta x + (1-\theta)x - x$$

$$f_{x|\theta}(x, \theta) = f_{x|\theta}(x|\theta) \cdot f(\theta) = f_{x|\theta}(x|\theta)$$

$$\text{Beta dist'n: } \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}$$

$$X \sim P(\lambda) \quad f_{\lambda}(\lambda) = \frac{\lambda^{\lambda}}{\Gamma(\lambda)} e^{-\lambda}$$

$$f_{X|\lambda}(x|\lambda) = \frac{\lambda^x}{x!} e^{-\lambda}$$

We observe $\bar{X} = (X_1, X_2, \dots, X_n)$.

$$f_{\bar{X}|\lambda}(x_1, x_2, \dots, x_n | \lambda) = \frac{\lambda^{\sum x_i}}{\prod_{i=1}^n x_i!} e^{-n\lambda}$$

$$f_{\bar{X}, \lambda}(x_1, x_2, \dots, x_n, \lambda) = f_{\bar{X}|\lambda}(x_1, x_2, \dots, x_n | \lambda) \cdot f_{\lambda}(\lambda)$$

$$f_{\lambda|\bar{X}}(\lambda | x_1, x_2, \dots, x_n) = \frac{f_{\bar{X}, \lambda}(x_1, x_2, \dots, x_n, \lambda)}{f_{\bar{X}}(x_1, x_2, \dots, x_n)} = C \lambda^{\frac{n}{2} \bar{x}_i + d - 1} e^{-(n+u)\lambda}$$

$$C = \frac{(n+u)^{\sum x_i + d}}{\Gamma(\sum x_i + d)}$$

$$\xi = \frac{1}{\sigma^2}$$

Ex. case 1. $f_{\bar{X}}(x|\theta, \xi) = \sqrt{\frac{\xi}{2\pi}} e^{-\frac{\xi(x-\theta)^2}{2}}$ θ unknown, ξ known.

give θ a prior distribution of $N(\theta_0, \xi_{\text{prior}})$

$$f_{\theta}(\theta) = \sqrt{\frac{\xi_{\text{prior}}}{2\pi}} e^{-\frac{(\theta-\theta_0)^2}{2\xi_{\text{prior}}}}$$

Given iid observations x_1, x_2, \dots, x_n . $f_{x|\theta}(x|\theta) = \sqrt{\frac{\xi}{2\pi}} e^{-\frac{\xi}{2} \sum (x_i - \theta)^2}$

We're interested in $f_{\theta|x}(\theta | \bar{x}) \propto f(x|\theta) f(\theta)$

$$= \sqrt{\frac{\xi}{2\pi}} e^{-\frac{\xi \sum (x_i - \theta)^2}{2\xi}} \sqrt{\frac{\xi_{\text{prior}}}{2\pi}} e^{-\frac{\xi_{\text{prior}} (\theta - \theta_0)^2}{2\xi}}$$

$$\propto \exp \left\{ -\frac{\xi}{2} \sum (x_i - \theta)^2 - \frac{\xi_{\text{prior}}}{2} (\theta - \theta_0)^2 \right\} \cdot \frac{\xi \bar{x} + \xi_{\text{prior}} \theta_0}{\xi \bar{x} + \xi_{\text{prior}}} \cdot \frac{(\xi + \xi_{\text{prior}})^{-1}}{\xi + \xi_{\text{prior}}}$$

$$\propto \exp \left\{ -\frac{\xi + \xi_{\text{prior}}}{2} \left(\theta - \frac{\xi \bar{x} + \xi_{\text{prior}} \theta_0}{\xi + \xi_{\text{prior}}} \right)^2 \right\} \sim N \left(\frac{\xi \bar{x} + \xi_{\text{prior}} \theta_0}{\xi + \xi_{\text{prior}}}, \frac{(\xi + \xi_{\text{prior}})^{-1}}{\xi + \xi_{\text{prior}}} \right)$$

case 2: θ Known, ξ Unknown.

We give ξ a prior distribution $f_\xi(\xi)$, we derive its posterior given observed data.

$$f(\xi) = \frac{\lambda^2}{\Gamma(\alpha)} \xi^{\alpha-1} e^{-\lambda\xi} \quad \xi > 0 \quad \text{book 292.}$$

$$\begin{aligned} f(\xi | \vec{x}) &\propto \xi^{\frac{n}{2}-1} e^{-(\frac{1}{2}\sum_i(x_i-\theta)^2)} \xi^{\alpha-1} e^{-\lambda\xi} = \xi^{\frac{n}{2}+\alpha-1} e^{-(\frac{1}{2}\sum_i(x_i-\theta)^2 + \lambda)\xi} \\ &= \frac{(\frac{\lambda^2}{2} + \lambda)^{\frac{n}{2}+\alpha}}{\Gamma(\frac{n}{2}+\alpha)} \xi^{\frac{n}{2}+\alpha-1} e^{-(\frac{\lambda^2}{2} + \lambda)\xi} \end{aligned}$$

Case 3: both unknown. → P.293 book.

$$-1 < X = \cos \theta < 1 \quad f(x|\alpha) = \frac{1+2x}{2} \quad -1 \leq x \leq 1.$$

- Method Moment Estimator for α .

$$\mu = E[X] = \int_{-1}^1 x \cdot \frac{1+2x}{2} dx = \frac{2}{3}$$

$\hat{\alpha} = 3\bar{x}^{-1}$ is the MME for α .

- MLE for α .

$$\ell(\alpha) = \sum_{i=1}^n \ln f(x_i | \alpha) = \sum_{i=1}^n \ln(1+2x_i) - n \ln 2$$

$$\ell'(\alpha) = \sum_{i=1}^n \frac{x_i}{1+2x_i} = 0 \Rightarrow \hat{\alpha}$$

Efficiency

$$\frac{\text{Var}(\hat{\alpha})}{\text{Var}(\hat{\alpha})} = \text{efficiency of } \hat{\alpha} \text{ relative to } \hat{\alpha}$$

$$\text{Var}(\hat{\alpha}) = q \text{Var}(\bar{x}) = \frac{q}{n} \text{Var}(X_1) = \frac{q}{n} \left[E(X_1^2) - \frac{E(X_1)^2}{} \right]$$

$$\sqrt{nI(\alpha)} \cdot (\hat{\alpha} - \alpha) \sim N(0, 1)$$

$$\text{Var}(\hat{\alpha}) \approx \frac{1}{nI(\alpha)}$$

$$I(\alpha) = E\left[\left(\frac{\partial}{\partial \alpha} \ln f(x|\alpha)\right)^2\right]$$

$$= \int_{-1}^1 \left(\frac{x}{1+2x}\right)^2 \frac{1+2x}{2} dx$$

MLE for λ in $P(\lambda)$.

30

$$f(x_1, x_2, \dots, x_n | \lambda) = \frac{\lambda^{\sum x_i}}{\prod_{i=1}^n x_i!} e^{-n\lambda}$$

$$\ell(\lambda) = \log f(x_1, x_2, \dots, x_n | \lambda) = n\bar{x} \log \lambda - n\lambda - \log \prod x_i!$$

$$\ell'(\lambda) = \frac{n\bar{x}}{\lambda} - n = 0 \quad \bar{x} = \hat{\lambda}$$

Cramer-Rao.

given iid sample x_1, x_2, \dots, x_n , $T(x_1, x_2, \dots, x_n)$.

If $\hat{\lambda}$ is an unbiased estimator of λ then.

$\text{Var}(T) \geq \frac{1}{n I(\lambda)}$ provided the density $f(x|\lambda)$ of the x_i 's is smooth.

$$I(\lambda) = E\left[\left(\frac{\partial}{\partial \lambda} \log f(x|\lambda)\right)^2\right] = -E\left[\frac{\partial^2}{\partial \lambda^2} \log f(x|\lambda)\right]$$

$$\text{Var}(\bar{x}) = \frac{1}{n}$$

$$I(\lambda) = -E\left[-\frac{1}{\lambda^2}\right] = \frac{1}{\lambda^2} E[X] = \frac{1}{\lambda}$$

Sufficiency

A statistic $T(x_1, x_2, \dots, x_n)$ is sufficient for θ if the conditional dist'n of x_1, x_2, \dots, x_n given $T(x_1, x_2, \dots, x_n)$ does not depend on θ .

Ex. $f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$. The joint density of an iid sample x_1, x_2, \dots, x_n . $f(x_1, x_2, \dots, x_n | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}^n} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2}$

$$= \frac{1}{\sqrt{2\pi\sigma^2}^n} e^{-\frac{1}{2\sigma^2} (\frac{n}{2} \sum_{i=1}^n x_i^2 - 2\mu \sum_{i=1}^n x_i + n\mu^2)}$$

(*) It suggests $T_1, T_2 = (\sum x_i, \sum x_i^2)$ are sufficient stats for (μ, σ^2) .

$$= g(T_1(x_1, x_2, \dots, x_n), T_2(x_1, x_2, \dots, x_n), \mu, \sigma^2)$$

by Factorization theorem

(*) holds

Factorization Theorem

T be a sufficient stats is that $f(x_1, x_2, \dots, x_n | \theta) = g(T(x_1, x_2, \dots, x_n), \theta) h(x_1, x_2, \dots, x_n)$ for some function g and h .

Ex

x_1, x_2, \dots, x_n iid $P(\lambda)$. λ unknown.

$T(x_1, x_2, \dots, x_n) = \bar{x}$ is MLE for λ . Is T sufficient? Yes

$$\begin{aligned} f(x_1, x_2, \dots, x_n | \lambda) &= \prod_{i=1}^n \left(\frac{\lambda}{x_i!} e^{-\lambda} \right) = \frac{\lambda^{\sum x_i}}{\prod x_i!} e^{-n\lambda} \\ &= \underbrace{(\lambda^{\sum x_i} e^{-\lambda})}_{g(T, \lambda)} \underbrace{\left(\frac{1}{\prod x_i!} \right)^{-1}}_{h(T)}. \end{aligned}$$

If T is a sufficient stats for θ then so is $K(T)$ for K any function.

Exponential families of densities.

These densities have the form

$$\begin{aligned} f(x | \theta) &= \prod_{i=1}^n \exp \left\{ c(\theta) k(x_i) + d(\theta) + s(x_i) \right\} \\ &= \exp \left\{ c(\theta) \sum_{i=1}^n k(x_i) + n d(\theta) + \sum_{i=1}^n s(x_i) \right\} \\ &= \underbrace{\exp \left\{ c(\theta) \sum_{i=1}^n k(x_i) + n d(\theta) \right\}}_{g} \cdot \underbrace{\exp \left\{ \sum_{i=1}^n s(x_i) \right\}}_{h}. \end{aligned}$$

$T(x_1, x_2, \dots, x_n) = \sum_{i=1}^n k(x_i)$ is a sufficient statistic.

$$\text{Ex. } f(x_1, x_2, \dots, x_n | \lambda) = \left(\frac{\lambda^x}{x!} \right)^n e^{-\lambda \sum x_i + (\lambda - 1) \sum \log x_i}$$

$$\begin{aligned} x_i^{x_i} e^{-\lambda x_i} &= e^{-\lambda x_i + (\lambda - 1) \sum \log x_i + \lambda \log \frac{x^x}{x!}} \\ &= e^{\sum x_i \cdot \sum \log x_i} \end{aligned}$$

($\sum x_i \cdot \sum \log x_i$) is sufficient statistic.

Sufficient stats and MLEs

32

$$f(x_1, x_2, \dots, x_n | \theta) = g(\theta, T(x_1, x_2, \dots, x_n)) \cdot h(T(x_1, x_2, \dots, x_n))$$

$$\ell(\theta) = \frac{d}{d\theta} g(\theta) = 0 \quad \text{We need to solve } \frac{d}{d\theta} g(T(x_1, x_2, \dots, x_n), \theta) = 0$$

$$g(T(x_1, x_2, \dots, x_n), \hat{\theta}) = 0 \Rightarrow \hat{\theta} = \hat{g}(T(x_1, x_2, \dots, x_n))$$

Rao-Blackwell Theorem

If $\hat{\theta}$ is an estimator of θ w/ $E[\hat{\theta}^2] < \infty$ for all values of θ and

If T is a sufficient statistic for θ , then

$$\tilde{\theta} = E[\hat{\theta} | T] \text{ satisfies} \\ E[(\tilde{\theta} - \theta)^2] \leq E[(\hat{\theta} - \theta)^2] \text{ w/ equality holds} \Leftrightarrow \tilde{\theta} = \hat{\theta}$$

Proof for sufficient statistic T

From sufficiency

$$\left\{ (x) \alpha + (\theta) \beta + (\lambda) \gamma(\theta) \right\} q_{\theta} \prod_{i=1}^n = (\theta) \alpha +$$

$$\left\{ (\lambda) \alpha \sum_{i=1}^n + (\beta) \gamma_i + (\lambda) \gamma(\theta) \right\} q_{\theta} =$$

$$\left\{ (\lambda) \alpha \sum_{i=1}^n \right\} q_{\theta} = \left\{ (\beta) \gamma_i + (\lambda) \gamma(\theta) \right\} q_{\theta} =$$

sufficient requirement: $(\alpha + \frac{\lambda}{n}) = (\alpha + \lambda) \alpha$

$$(\alpha + \lambda) \alpha + (\beta) \gamma_i + (\lambda) \gamma(\theta) = (\alpha + \lambda) \alpha + (\lambda) \gamma(\theta)$$

$$(\alpha + \lambda) \alpha + (\lambda) \gamma(\theta) =$$

sufficient condition: $(\alpha + \lambda) \alpha =$

Hypothesis testing

H_0 : null hypothesis.

H_1 : alternative hypothesis.

Ind. trials are run and based on some test statistic. H_0 is either accepted or rejected

Simple one:

$$H_0: \theta = \theta_0$$

$$H_1: \theta = \theta_1$$

Composite one:

$$H_0: \theta = \theta_0$$

$$H_1: \theta \geq \theta_1 > \theta_0$$

An iid sample is taken from the unknown dist'n $f(x|\theta)$
 $\vec{x} = (x_1, x_2, \dots, x_n)$.

In the case of a simple hypothesis, you'd look at

$$\frac{f(\vec{x}|\theta_0)}{f(\vec{x}|\theta_1)} \quad \text{Find some acceptance level } x_0 \text{ and}$$

Your test says:

$$\text{if } \frac{f(\vec{x}|\theta_0)}{f(\vec{x}|\theta_1)} \geq x_0 \text{ accept } H_0$$

$$\text{if } \frac{f(\vec{x}|\theta_0)}{f(\vec{x}|\theta_1)} < x_0 \text{ reject } H_0$$

$\alpha = P(\text{reject } H_0 \mid \text{given } H_0 \text{ is true})$ $\hat{=} \text{significance level.}$

Type I error.

$$= P\left(\frac{f(\vec{x}|\theta_0)}{f(\vec{x}|\theta_1)} < x_0 \mid H_0\right)$$

Another error. Type II error, is to accept H_0 when H_1 is true.

$$\beta = P(\text{accept } H_0 \mid H_1)$$

$$1 - \beta = P(\text{reject } H_0 \mid H_1) = \text{power of test.}$$

Ex. x_1, x_2, \dots, x_n iid $N(\theta, \sigma^2)$. θ unknown. σ^2 known.

$$H_0: \theta = \theta_0$$

$$H_1: \theta = \theta_1$$

where $\theta_0 > \theta_1$

$$\frac{f(\bar{x} | \theta_0)}{f(\bar{x} | \theta_1)} = \frac{\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} \sum (x_i - \theta_0)^2}}{\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} \sum (x_i - \theta_1)^2}} = e^{-\frac{1}{2\sigma^2} (\sum (x_i - \theta_0)^2 + \sum (x_i - \theta_1)^2)}$$

$$= \exp \left\{ -\frac{n}{2\sigma^2} (2\bar{x}(\theta_1 - \theta_0) + \theta_0^2 - \theta_1^2) \right\} > x_0$$

$$2\bar{x}(\theta_1 - \theta_0) + \theta_0^2 - \theta_1^2 \leq y_0 = -\frac{2\sigma^2}{n} \log x_0$$

If $\bar{x} \geq \frac{y_0 + \theta_1^2 - \theta_0^2}{2(\theta_1 - \theta_0)}$, accept H_0 .

If $\bar{x} < \dots$ reject H_0 .

$A = \{x_1, x_2, \dots, x_n : \bar{x} \leq \frac{y_0 + \theta_1^2 - \theta_0^2}{2(\theta_1 - \theta_0)}\}$ is called acceptance region.
 $A^c = \left(\frac{y_0 + \theta_1^2 - \theta_0^2}{2(\theta_1 - \theta_0)}, \infty \right)$ is called rejection region.

\bar{x} : test statistic.

Given a significance level α , how can we set y_0 so that

$$P(\text{reject } H_0 | H_0) = \alpha$$

$$P(\bar{x} \leq \frac{y_0 + \theta_1^2 - \theta_0^2}{2(\theta_1 - \theta_0)} | H_0) = \alpha$$

$$\bar{x} \sim N(\theta_0, \frac{\sigma^2}{n}) \quad P(\frac{\bar{x} - \theta_0}{\sigma/\sqrt{n}} \leq \frac{\frac{y_0 + \theta_1^2 - \theta_0^2}{2(\theta_1 - \theta_0)} - \theta_0}{\sigma/\sqrt{n}} | H_0) = \alpha$$

$$= -Z(\alpha)$$

$$\beta = P(\text{accept } H_0 | H_1)$$

$$1 - \beta = P(\text{reject } H_0 | H_1) = P(\bar{x} \leq \frac{y_0 + \theta_1^2 - \theta_0^2}{2(\theta_1 - \theta_0)} | H_1)$$

$$\bar{x} \sim N(\theta_1, \frac{\sigma^2}{n})$$

Neyman-Pearson Lemma

Among all tests w/ significance level α , the likelihood ratio test is ~~most~~ most powerful.

what about composite hypothesis

$$H_0: \mu = \mu_0$$

$$H_1: \mu > \mu_0$$

Consider an iid sample X_1, X_2, \dots, X_n from a $N(\mu, \sigma^2)$ distn.

35

w/ σ^2 known, μ unknown, and test

$$H_0: \mu = \mu_0$$

$$H_1: \mu \neq \mu_0$$

Accept H_0 if $|\bar{X} - \mu_0| \leq x_0$

reject H_0 if $|\bar{X} - \mu_0| > x_0$

$$\alpha = P(\text{reject } H_0 | H_0) = P(|\bar{X} - \mu_0| \geq x_0 | H_0). \quad \bar{X} \sim N(\mu_0, \frac{\sigma^2}{n})$$

$$= P\left(\frac{|\bar{X} - \mu_0|}{\sigma/\sqrt{n}} \geq \frac{x_0}{\sigma/\sqrt{n}} | H_0\right).$$

$$\text{Mo} \in [\bar{X} - z(\frac{\alpha}{2}) \frac{\sigma}{\sqrt{n}}, \bar{X} + z(\frac{\alpha}{2}) \frac{\sigma}{\sqrt{n}}] \\ 1 - \alpha(1-\alpha)\% \text{ C.I. for } \mu_0 = P(|Z| \geq \frac{x_0}{\sigma/\sqrt{n}}) = \alpha. \quad x_0 = z(\frac{\alpha}{2}) \frac{\sigma}{\sqrt{n}}$$

$\beta = P(\text{reject } H_0 | H_1)$. — cannot be computed. Since we don't know the dist'n of \bar{X} under H_1 .

Sometimes, hypotheses have the form

$$H_0: \theta \in W_0$$

$$H_1: \theta \in W_1$$

$$W_0 \cap W_1 = \emptyset$$

Generalized Likelihood Ratio Test

Set $\mathcal{R} = W_0 \cup W_1$.

$$l(\theta) = \log f(X_1, X_2, \dots, X_n | \theta).$$

$$\text{define } \lambda = \frac{\max_{\theta \in W_0} l(\theta)}{\max_{\theta \in W_1} l(\theta)}$$

Then a test would take the form. accept H_0 if $\lambda \geq c$.
reject H_0 if $\lambda \leq c$.

Ex. X_1, X_2, \dots, X_n iid $N(\mu, \sigma^2)$

$$H_0: \mu \in W_0 = \{\mu_0\} \quad (\mu = \mu_0)$$

$$H_1: \mu \in W_1 = \mathbb{R} - \{\mu_0\}. \quad (\mu \neq \mu_0)$$

$$\max_{\theta \in W_0} l(\theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} \sum (X_i - \mu_0)^2}$$

$$\max_{\theta \in W_1} l(\theta) = l(\hat{\theta}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} \sum (X_i - \bar{X})^2}$$

$$\lambda = \exp\left\{ \frac{1}{2\sigma^2} \left[\sum_{i=1}^n (X_i - \bar{X})^2 - \sum_{i=1}^n (X_i - \mu_0)^2 \right] \right\}$$

$$= \exp\left\{ \frac{n}{2\sigma^2} (\bar{X} - \mu_0)^2 \right\}$$

$$-2\log \lambda = \left(\frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \right)^2 \sim \chi^2_1$$

$\lambda > c$ becomes $-2\log \lambda \leq -2\log c$

The test becomes.

$$\left(\frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \right)^2 \leq -2\log c$$

CLT for Generalized Likelihood Ratio Test

If the underlying densities are smooth for all values of the parameter

θ and $\lambda = \frac{\max_{\theta \in W_0} l(\theta)}{\max_{\theta \in \mathcal{R}} l(\theta)}$ and $\dim \mathcal{R} - \dim W_0 = k$, then

$-2\log \lambda$ has an approximate χ^2_k for large n .

Multinomial

$$f(x_1, x_2, \dots, x_m | p_1, \dots, p_m) = \frac{n!}{\prod x_i!} \cdot \prod_{i=1}^m p_i^{x_i}$$

parameters $\vec{p} = (p_1, p_2, \dots, p_m)$
where $p_i > 0$ and $\sum p_i = 1$

P-value := smallest significance level at which null hypothesis is rejected.

- Probability of an outcome for T at least as extreme as the observed one t .

$$P(|T| > t | H_0) = P\text{-value}$$

Poisson Dispersion Test

x_i = # particles in i^{th} cell.

Assume dist'n of particles in each cell is Poisson and we want to test

H_0 : λ is the same in all cells

H_1 : λ_i = parameter in i^{th} cell may vary from cell to cell.

use Generalized Likelihood Ratio Test.

$$\lambda = \frac{\max_{\lambda \in W_0} l(\lambda)}{\max_{\lambda \in \mathcal{R}} l(\lambda_1, \lambda_2, \dots, \lambda_n)}$$

$W_0 = [0, \infty)$, $\dim W_0 = 1$

$\mathcal{R} = (0, \infty)^n$, $\dim \mathcal{R} = n$

$$\prod_{i=1}^n \frac{\lambda^{x_i}}{x_i!} e^{-\lambda}$$

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

37

$$\lambda = \frac{\prod_{i=1}^n \frac{\lambda^{x_i}}{x_i!} e^{-\lambda}}{\max \prod_{i=1}^n \frac{\lambda^{x_i}}{x_i!} e^{-\lambda}}$$

$$= \frac{\prod_{i=1}^n \frac{\lambda^{x_i}}{x_i!} e^{-\lambda}}{\prod_{i=1}^n \frac{\bar{\lambda}^{x_i}}{x_i!} e^{-\bar{\lambda}}} = \prod_{i=1}^n \left(\frac{\lambda}{\bar{\lambda}}\right)^{x_i} e^{\bar{\lambda} - \lambda} = \prod_{i=1}^n \left(\frac{\bar{x}}{x_i}\right)^{x_i} e^{x_i - \bar{x}}$$

$$\frac{d}{d\lambda_i} \left(\log \frac{\lambda_i^{x_i}}{x_i!} e^{-\lambda_i} \right) = \frac{d}{d\lambda_i} (x_i \log \lambda_i - \lambda_i)$$

$$= \frac{x_i}{\lambda_i} - 1 = 0$$

$$x_i = \lambda_i$$

$$\bar{\lambda}_i = x_i$$

$$\left. \begin{array}{l} \text{if} \\ \text{else} \end{array} \right\} = 0$$

$$f(x) = x \log \frac{x}{\bar{x}} = f(\bar{x}) + f'(\bar{x})(x - \bar{x})$$

$$+ f''(\bar{x}) \frac{(x - \bar{x})^2}{2}$$

~~so~~

when the number of observations is large.

$$-\cancel{2 \log \lambda} \rightarrow 2 \log \lambda \sim \chi^2_{n-1}$$

$$-2 \sum_{i=1}^n x_i \log \frac{\bar{x}}{x_i} - 2 \sum_{i=1}^n (x_i - \bar{x})$$

$$= 2 \sum_{i=1}^n x_i \log \frac{x_i}{\bar{x}}$$

$$\approx \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\bar{x}} = \boxed{\frac{n \hat{\sigma}^2}{\bar{x}}}$$

is called dispersion.

Test: If $\frac{n \hat{\sigma}^2}{\bar{x}} > c$ reject.

To obtain a significance level α we find c s.t. $P_0 \left(\frac{n \hat{\sigma}^2}{\bar{x}} > c \right) = \alpha$

Regression

11

Tid random points $(x_i, y_i), i=1, 2, \dots, n$

$$y = ax + bx.$$

$$y = ax + bx + \varepsilon \leftarrow \begin{array}{l} \text{everything we don't know} \\ \text{it's not due to randomness.} \end{array}$$

Goal: Find \hat{a}, \hat{b} that approach a, b .

$X \in \mathbb{R}$: univariate regression

$X \in \mathbb{R}^p$: multivariate regression

$$Y = a + X^T b.$$

Why we do this?

ML: relation between X and Y .

Linear function is the simplest.

Ex. $D_i \approx a + b p_i$
 ↑ ↑
 Demand price

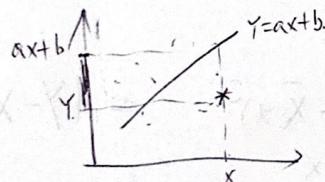
$$PV = nRT \quad (\text{ideal gas law}).$$

Variable transformation

take $\left\{ \begin{array}{l} \text{exponential} \\ \text{square} \\ \log \\ \text{etc.} \end{array} \right\}$

$$\log p + \log V = \log nR + \log T.$$

$$\log p = \log nR + \log T - \log V$$



$$E[(y - (ax + b))^2]$$

function of a and b that measures how close the point (x, y) is to the line $y = ax + b$ in expectation

assumption: $\text{Var}(X) \neq 0$

$$\tilde{x} = \bar{x} - E(x)$$

$$\min_{a, b} E[(y - (ax + b))^2] = E[(y - (a + b\tilde{x}))^2]$$

$$\frac{\partial}{\partial a} = -E[2(y - (\tilde{a} + b\tilde{x}))] = 0$$

$$\tilde{a} = \bar{a}$$

$$\frac{\partial}{\partial b} = -E[2(y - (\tilde{a} + b\tilde{x}))]\tilde{x} = 0$$

$$\left\{ \begin{array}{l} a = E(Y) - bE(X) \\ b = \frac{\text{cov}(X, Y)}{\text{Var}(X)} \end{array} \right\} \leftarrow \begin{array}{l} \text{replace by } \bar{X}, \bar{Y} \\ \text{etc.} \\ \text{Var}(X) \end{array}$$

$$y = a + bx + \varepsilon \Leftrightarrow \varepsilon = y - (a + bx). E(\varepsilon) = E(y) - [a + bE(x)] = 0 \quad [2]$$

$$\text{cov}(X, \varepsilon) = \text{cov}(X, y - (a + bx)) = \text{cov}(X, y) - (a + b\text{cov}(X, X)).$$

$$= E[X(y - (a + bx))] = \text{cov}(X, y) = 0$$

If $X \perp \varepsilon$, then $\text{cov}(X, \varepsilon) = 0$

Assumption $\begin{cases} E(\varepsilon) = 0 \\ \text{cov}(X, \varepsilon) = 0 \end{cases}$

We need more when we do more analysis, e.g. t-test we need ε to be Gaussian.

$$Y_i = a + bx_i + \varepsilon_i \quad \varepsilon_i \text{ i.i.d.}$$

$$E(\varepsilon_i) = 0$$

$$\text{cov}(X_i, \varepsilon_i) = 0$$

$$\hat{b} = \frac{\text{cov}(X, Y)}{\text{var}(X)} = \frac{E(XY) - E(X)E(Y)}{E(X^2) - (E(X))^2} = \frac{\bar{XY} - \bar{X}\bar{Y}}{\bar{X^2} - \bar{X}^2}$$

$$E[(y - (a + bx))^2] \rightarrow \frac{1}{n} \sum_{i=1}^n (y_i - (a + bx_i))^2$$

Multivariate case

$$Y_i = X_i^\top \beta + \varepsilon_i \quad i = 1, 2, \dots, n. \quad \text{cov}(X_i, \varepsilon_i) = 0$$

↑
dep. var.
response var. ↑ covariates
 ↓ ind. Var.

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (Y_i - X_i \beta)^2. \quad = \|Y - X\beta\|_2^2.$$

Matrix notation

$$\begin{pmatrix} Y \\ X_1 \\ X_2 \\ \vdots \\ X_n \end{pmatrix} = \begin{pmatrix} 1 & X_1 & X_1^2 & \dots & X_1^{p-1} \\ 1 & X_2 & X_2^2 & \dots & X_2^{p-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_n & X_n^2 & \dots & X_n^{p-1} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \varepsilon \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$\min \|Y - X\beta\|^2 = \|Y\|^2 + \|X\beta\|^2 - 2X^\top X\beta$$

$$= \|Y\|^2 + \beta^\top X^\top X\beta - 2X^\top X\beta$$

$$\frac{\partial}{\partial \beta} = 0 + 2X^\top X\beta - 2X^\top Y \Rightarrow \beta = (X^\top X)^{-1} X^\top Y$$

if $p > n$. Least squares doesn't exist.

$\hat{X}\hat{\beta}$ $X(X^\top X)^{-1} X^\top Y$ orthogonal projection of y onto the subspace spanned by columns of X

Statistical Inference

13

Assumptions:

- 1) X_i deterministic.
- 2) The model is homoscedastic: $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ are iid.
- 3) $\varepsilon \sim N_p(0, \sigma^2 I_p)$. Unknown $\sigma^2 > 0$.

$$\begin{aligned}\hat{\beta} &= (X^T X)^{-1} X^T Y \quad Y = X\beta + \varepsilon \\ &= (X^T X)^{-1} X^T X \beta + (X^T X)^{-1} X^T \varepsilon \\ &= \beta + (X^T X)^{-1} X^T \varepsilon \\ &= \beta + N_p(0, (X^T X)^{-1} \sigma^2) \\ &= \beta + N_p(0, (X^T X)^{-1} \sigma^2)\end{aligned}$$

1) $LSE = MLE \quad Y \sim N_p(X\beta, \sigma^2 I_p)$

$$P(Y) = \frac{1}{(\sigma^2 \pi)^{n/2}} \exp \left\{ - \frac{\|Y - X\beta\|_2^2}{2\sigma^2} \right\}$$

likelihood
 $\log P(Y) = -\frac{n}{2} \log(\sigma^2 \pi) - \frac{1}{2\sigma^2} \|Y - X\beta\|_2^2$

$$MLE = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \|Y - X\beta\|_2^2 = LSE$$

2) Quadratic Risk of $\hat{\beta}$:

$$\varepsilon \sim N(0, \Sigma)$$

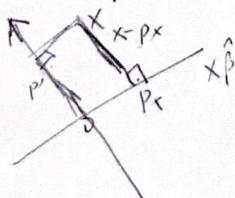
$$E\|\varepsilon\|^2 = \text{Tr}(\Sigma)$$

$$\begin{aligned}E[\|\hat{\beta} - \beta\|_2^2] &= \sum_{i=1}^n E(\hat{\beta}_i - \beta_i)^2 \\ &= \sum_{i=1}^n \text{Var}(\hat{\beta}_i) \\ &= \sigma^2 \text{Tr}(X^T X)^{-1}\end{aligned}$$

trace of a square matrix
 the sum of elements
 on the main diagonal

3) prediction error.

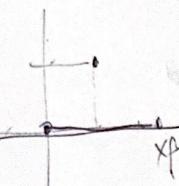
$$Y - X\hat{\beta} = \cancel{IY} - X(X^T X)^{-1} X^T Y = [I - X(X^T X)^{-1} X^T] Y$$



projection onto orthogonal
 of columns span of X

$$E\|y - \hat{\beta}\|^2 = E\|\underbrace{(I - X(X^T X)^{-1} X^T) y}_{P'y}\|^2 = E\|N(0, \sigma^2 P' P)^T\|$$

$$Y = X\beta + \varepsilon$$

$$P'y = P'X\beta + P'\varepsilon$$


$$\begin{aligned} &= \sigma^2 \underbrace{\text{rank}(P')}_n \\ &= \sigma^2(n-p) \end{aligned}$$

L4

4). Unbiased estimator $\hat{\sigma}^2 = \frac{1}{n-p} \|y - \hat{\beta}\|^2$

$b \perp \hat{\sigma}^2$ trans dith projections
 $(n-p) \frac{\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-p}$

Tests

$$\begin{aligned} H_0: \beta_j &= 0 \\ H_1: \beta_j &\neq 0 \end{aligned}$$

$$\hat{\beta} = \beta + N(0, \sigma^2 (X^T X)^{-1})$$

$$\hat{\beta}_j \sim \beta_j N(0, \sigma^2 (X^T X)^{-1}) + \beta_j$$

Test: $\left| \frac{\hat{\beta}_j}{\sigma \sqrt{r_j}} \sim N(0, 1) \right| > q_{\alpha/2}(N(0, 1))$. If σ^2 is known.

If σ^2 is unknown, we have $\hat{\sigma}^2$

$$\frac{\hat{\beta}_j}{\hat{\sigma} \sqrt{r_j}} = \frac{\hat{\beta}_j}{\sigma \sqrt{r_j}} \cdot \frac{1}{\sqrt{\frac{\hat{\sigma}^2}{\sigma^2(n-p)}}} \cdot \sqrt{n-p} = \frac{\hat{\beta}_j / \sigma \sqrt{r_j}}{\sqrt{\frac{\hat{\sigma}^2}{\sigma^2(n-p)}} / (n-p)} \sim \chi^2_{n-p}$$

$\sim t_{n-p}$

$$\left| \frac{\hat{\beta}_j}{\hat{\sigma} \sqrt{r_j}} \right| > q_{\alpha/2}(t_{n-p})$$

| | Value | T | P-value |
|-----------------|-------|---|---------|
| β_0 | | | |
| A | | | |
| $\hat{\beta}_j$ | | | |

$H_0: \beta_j = 0$ for all $j \in S \subset \{1, \dots, p\}$

L5

$H_1: \beta_j \neq 0$ for some $j \in S$

$P_{H_0}(\text{reject } \beta_1 \cup \text{reject } \beta_2 \dots \text{reject } \beta_s) \leq \alpha$
if they disjoint ~~=~~

$$\frac{\hat{\beta}}{\hat{\sigma}\sqrt{r_j}} > q_{\frac{\alpha}{2}}(t_{n-p}) \rightarrow P \geq \alpha$$

Bonferroni correction. (multiple testing simultaneous)

$H_0: G\beta = \lambda$ $H_1: G\beta \neq \lambda$

$$G\hat{\beta} = G\beta + \mathcal{N}(0, \sigma^2(X^T X)^{-1})$$

$$\frac{G\hat{\beta} - \lambda}{\hat{\sigma}} \sim N\left(0, \frac{G(X^T X)^{-1} G^T}{\sigma^2}\right) \quad \text{if } X \sim N(0, \Sigma) \\ X^T \Sigma^{-1} X \sim \chi^2_k$$

$$\sigma^2(G\hat{\beta} - \lambda)^T (G(X^T X)^{-1} G^T)^{-1} (G\hat{\beta} - \lambda) \sim \chi^2_k$$

$$\text{Let } \frac{\sigma^2(G\hat{\beta} - \lambda)^T (G(X^T X)^{-1} G^T)^{-1} (G\hat{\beta} - \lambda)}{k} \sim \frac{\chi^2_k}{k}$$

Bayesian statistics

$$B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

$$\text{Beta distribution } \text{Beta}(a, b). \quad f(x) = \begin{cases} x^{a-1} (1-x)^{b-1} & x \in [0, 1] \\ 0 & \text{otherwise} \end{cases}$$

posterior. $p(\theta | \mathbf{x})$.

Beta is conjugate prior.

$$x_1, x_2, \dots, x_n \stackrel{iid}{\sim} \text{Ber}(p). \quad p(\theta | \mathbf{x}) \propto \frac{1}{p} \cdot p^{a-1} (1-p)^{b-1} \cdot p^{\sum x_i} (1-p)^{n-\sum x_i} \\ = p^{a-1+\sum x_i} \cdot (1-p)^{b-1+n-\sum x_i}$$

Non informative priors

6

- Good candidate: $\pi(\theta) \propto 1$ constant pdf.
- If Θ is bounded, this is uniform prior on Θ .
- If Θ is unbounded, this does not define a proper pdf on Θ .
- If $p \sim U(0,1)$, $p, X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Ber}(p)$.
Posterior: $B(1 + \sum x_i, 1 + n - \sum x_i)$
- If $\pi(\theta) = 1$, given $\theta, X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\theta, 1)$.
Posterior $N(\bar{X}_n, \frac{1}{n})$.

GLMs

$$Y|X \sim N(\mu(x), \sigma^2 I)$$

$$= N(X^T \beta, \sigma^2 I)$$

$$X \rightarrow X^T \beta$$

1. Random component:

$$Y|X \sim N(\mu(x), \sigma^2 I)$$

$\rightarrow Y|X \sim \text{Some distribution from the exponential family}$

$$2. g(\mu(x)) = X^T \beta \quad (\Leftrightarrow \mu(x) = f(X^T \beta))$$

Link function.

log link function

reciprocal link function

exponential family

$$P_\theta(x) = p(\theta, x) = \exp(\theta \cdot x) \cdot f(x) \cdot g(\theta) = \exp\left[\sum_{i=1}^K \eta_i(\theta) T_i(x)\right] \cdot c(\theta, h(x))$$

$$\eta_1(\theta) \cdot T_1(x)$$

$$\eta_2(\theta) \cdot T_2(x)$$

$$\vdots$$

$$\eta_K(\theta) \cdot T_K(x)$$

Ex. $X \sim N(\mu, \sigma^2)$. $\Theta(\mu, \sigma^2)$

$$P_\Theta(X) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(X-\mu)^2}{2\sigma^2}\right\}$$

$$= \exp\left\{-\frac{x^2}{2\sigma^2} - \frac{\mu^2}{2\sigma^2} + \frac{\mu x}{\sigma^2} - \log(\sigma\sqrt{2\pi})\right\}$$

$$\eta_1 = \frac{\mu}{\sigma^2}, \quad \eta_2 = -\frac{1}{2\sigma^2} \quad T_1(x) = x, \quad T_2(x) = x^2.$$

$$B(\Theta) = \frac{\mu^2}{2\sigma^2} + \log(\sigma\sqrt{2\pi}), \quad h(x) = 1.$$

When σ^2 is known.

$$\eta = \frac{\mu}{\sigma^2}. \quad T(x) = x. \quad B(\Theta) = \frac{\mu^2}{2\sigma^2}. \quad h(x) = \frac{e^{-\frac{x^2}{2\sigma^2}}}{\sigma\sqrt{2\pi}}$$

Bernoulli(p). $p^x(1-p)^{1-x}$

Poisson $\frac{\lambda^x}{x!} e^{-\lambda}$.

Gamma(a,b). $\frac{1}{\Gamma(a)b^a} x^{a-1} e^{-\frac{x}{b}}$

a: shape parameter

b: scale parameter

Reparametrize: $\mu = ab$

$$\frac{1}{\Gamma(a)} \left(\frac{a}{\mu}\right)^a x^{a-1} e^{-\frac{ax}{\mu}}$$

Inverse Gamma(α, β): $\frac{\beta^\alpha}{\Gamma(\alpha)} x^{-\alpha-1} e^{-\frac{\beta}{x}}$

Inverse Gaussian(μ, σ^2): $\sqrt{\frac{\sigma^2}{2\pi x^3}} e^{-\frac{\sigma^2(x-\mu)^2}{2\mu^2 x}}$

Chi-square, Beta, Binomial, Negative binomial

7