

# Chapter 4

## Least Squares: Exact Distribution

### 4.1 Maximum Likelihood

There are very few *principles* in statistics, and maximum likelihood is one of them. In this chapter, we first give an introduction of the maximum likelihood estimation. Consider a random sample of  $Z = (z_1, z_2, \dots, z_n)$  drawn from a parametric distribution with density  $f_z(z_i; \theta)$ , where  $z_i$  is either a scalar random variable or a random vector. A parametric distribution is completely characterized by a finite-dimensional parameter  $\theta$ . We know that  $\theta$  belongs to a parameter space  $\Theta$ . We use the data to estimate  $\theta$ .

The log-likelihood of observing the entire sample  $Z$  is

$$L_n(\theta; Z) := \log \left( \prod_{i=1}^n f_z(z_i; \theta) \right) = \sum_{i=1}^n \log f_z(z_i; \theta). \quad (4.1)$$

In reality the sample  $Z$  is given and for each  $\theta \in \Theta$  we can evaluate  $L(\theta; Z)$ . The maximum likelihood estimator is

$$\hat{\theta}_{MLE} := \arg \max_{\theta \in \Theta} L_n(\theta; Z).$$

Why maximizing the log-likelihood function is desirable? An intuitive explanation is that  $\hat{\theta}_{MLE}$  makes observing  $Z$  the “most likely” in the entire parametric space.

A more formal justification requires an explicitly defined distance. Suppose that the true parameter value that generates the data is  $\theta_0$ , so that the true distribution is  $f_z(z_i; \theta_0)$ . Any generic point  $\theta \in \Theta$  produce  $f_z(z_i; \theta)$ . To measure their difference, we introduce the *Kullback-Leibler distance*, or *Kullback-Leibler divergence*, defined as

$$D_f(\theta, \theta_0) = D(f_z(z_i; \theta), f_z(z_i; \theta_0)) := E_{\theta_0} \left[ \log \frac{f_z(z_i; \theta_0)}{f_z(z_i; \theta)} \right].$$

We say it is a “distance” because it is non-negative. To see this, notice that  $-\log(\cdot)$  is strictly convex and then by Jensen’s inequality

$$\begin{aligned} E_{\theta_0} \left[ \log \frac{f_z(z_i; \theta_0)}{f_z(z_i; \theta)} \right] &= E_{\theta_0} \left[ -\log \frac{f_z(z_i; \theta)}{f_z(z_i; \theta_0)} \right] \geq -\log \left( E_{\theta_0} \left[ \frac{f_z(z_i; \theta)}{f_z(z_i; \theta_0)} \right] \right) \\ &= -\log \left( \int \frac{f_z(z_i; \theta)}{f_z(z_i; \theta_0)} f_z(z_i; \theta_0) dz_i \right) = -\log \left( \int f_z(z_i; \theta) dz_i \right) \\ &= -\log 1 = 0, \end{aligned}$$

where  $\int f_z(z_i; \theta) dz_i = 1$  for any PDF. The equality holds if and only if  $f_z(z_i; \theta) = f_z(z_i; \theta_0)$  almost everywhere. Furthermore, if there is a one-to-one mapping between  $\theta$  and  $f_z(z_i; \theta)$  on  $\Theta$  (identification), then  $\theta_0 = \arg \min_{\theta \in \Theta} D_f(\theta, \theta_0)$  is the unique solution.

**Example 4.1.** Consider the Gaussian location model  $z_i \sim N(\mu, 1)$ , where  $\mu$  is the unknown parameter to be estimated. The (averaged) log-likelihood function for  $n$  observations is

$$\ell_n(\mu) = -\frac{1}{2} \log(2\pi) - \frac{1}{2n} \sum_{i=1}^n (z_i - \mu)^2.$$

Here we use the averaged log-likelihood  $\ell_n$ , instead of the (raw) log-likelihood  $L_n$  in (4.1), to make it directly comparable with its population counterpart

$$\begin{aligned} E[\ell_n(\mu)] &= -\frac{1}{2} \log(2\pi) - \frac{1}{2} E[(z_i - \mu)^2] \\ &= -\frac{1}{2} \log(2\pi) - \frac{1}{2} E[((z_i - \mu_0) + (\mu_0 - \mu))^2] \\ &= -\frac{1}{2} \log(2\pi) - \frac{1}{2} E[(z_i - \mu_0)^2] - E[z_i - \mu_0](\mu_0 - \mu) - \frac{1}{2} (\mu_0 - \mu)^2 \\ &= -\frac{1}{2} \log(2\pi) - \frac{1}{2} - \frac{1}{2} (\mu_0 - \mu)^2. \end{aligned}$$

Obviously,  $\ell_n(\mu)$  is maximized at  $\bar{z} = \frac{1}{n} \sum_{i=1}^n z_i$  while  $E[\ell_n(\mu)]$  is maximized at  $\mu = \mu_0$ .

We use the following code to demonstrate the population log-likelihood  $E[\ell_n(\mu)]$  when  $\mu_0 = 2$  and the 3 sample realizations when  $n = 4$ .

```

set.seed(2020-10-7)

mu0 <- 2; gamma0 <- 1

# population likelihood function
L <- function(mu) {
  ell = -0.5 * log(2*pi*gamma0) - 0.5 / gamma0 * ( 1 + (mu - mu0)^2 )
  return(ell) }

# sample likelihood function
Ln <- function(mu) {
  elln = -0.5 * log(2*pi*gamma0) - 0.5 / gamma0 * mean( (z - mu)^2 )
  return(elln) }

mu_base = mu0 + seq(-3, 3, by = 0.01)

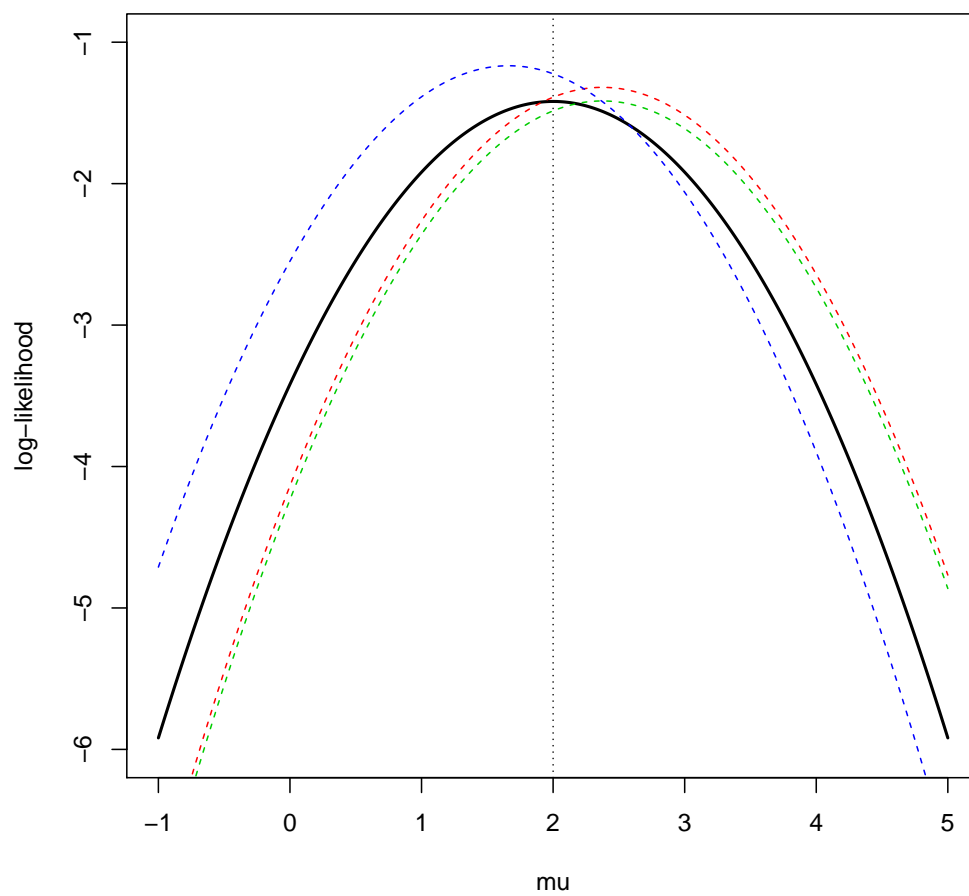
# draw sample log-likelihood graph
n = 4

lnz = matrix(0, length(mu_base), 3)
for (rr in 1:3){
  z <- rnorm(n, mu0, sqrt(gamma0) )
  lnz[,rr] <- plyr::lapply(.data = mu_base, .fun = Ln)
}

matplot(x = mu_base, y = cbind( L(mu_base), lnz),
        type = "l", lty = c(1, rep(2,3)),
        lwd = c(2,rep(1,3)), col = 1:4, ylim = c(-6, -1),

```

```
      xlab = "mu", ylab = "log-likelihood")  
abline(v = mu0, lty = 3)
```



## 4.2 Likelihood Estimation for Regression

Notation:  $y_i$  is a scalar, and  $x_i = (x_{i1}, \dots, x_{iK})'$  is a  $K \times 1$  vector.  $Y$  is an  $n \times 1$  vector, and  $X$  is an  $n \times K$  matrix.

We continue with properties of OLS. Noticing that OLS coincides with the maximum likelihood estimator if the error term follows a normal distribution, we derive its finite-sample exact distribution which can be used for statistical inference. The Gauss-Markov theorem justifies the optimality of OLS under the classical assumptions.

In this chapter we employ the classical statistical framework under restrictive distributional assumption

$$y_i|x_i \sim N(x_i'\beta, \gamma), \quad (4.2)$$

where  $\gamma = \sigma^2$  to ease the differentiation. This assumption is equivalent to  $e_i|x_i = (y_i - x_i'\beta) | x_i \sim N(0, \gamma)$ . Because the distribution of  $e_i$  is invariant to  $x_i$ , the error term  $e_i \sim N(0, \gamma)$  and is statistically independent of  $x_i$ . This is a very strong assumption.

The likelihood of observing a pair  $(y_i, x_i)$  is

$$\begin{aligned} f_{yx}(y_i, x_i) &= f_{y|x}(y_i|x_i) f_x(x) \\ &= \frac{1}{\sqrt{2\pi\gamma}} \exp\left(-\frac{1}{2\gamma} (y_i - x_i'\beta)^2\right) \times f_x(x), \end{aligned}$$

where  $f_{yx}$  is the joint pdf,  $f_{y|x}$  is the conditional pdf and  $f_x$  is the marginal

pdf of  $x$ , and the second equality holds under the assumption (4.2). The likelihood a random sample  $(y_i, x_i)_{i=1}^n$  is

$$\begin{aligned}\prod_{i=1}^n f_{y|x}(y_i, x_i) &= \prod_{i=1}^n f_{y|x}(y_i|x_i) f_x(x) \\ &= \prod_{i=1}^n f_{y|x}(y_i|x_i) \times \prod_{i=1}^n f_x(x) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\gamma}} \exp\left(-\frac{1}{2\gamma} (y_i - x_i'\beta)^2\right) \times \prod_{i=1}^n f_x(x).\end{aligned}$$

The parameters of interest  $(\beta, \gamma)$  are irrelevant to the second term  $\prod_{i=1}^n f_x(x)$  for they appear only in the conditional likelihood

$$\prod_{i=1}^n f_{y|x}(y_i|x_i) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\gamma}} \exp\left(-\frac{1}{2\gamma} (y_i - x_i'\beta)^2\right).$$

We focus on the conditional likelihood. To facilitate derivation, we work with the conditional log-likelihood function

$$L(\beta, \gamma) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \gamma - \frac{1}{2\gamma} \sum_{i=1}^n (y_i - x_i'\beta)^2,$$

for  $\log(\cdot)$  is a monotonic transformation that does not change the maximizer. The maximum likelihood estimator  $\hat{\beta}_{MLE}$  can be found using the FOC:

$$\frac{\partial}{\partial \beta} L(\beta, \gamma) = \frac{1}{2\gamma} \sum_{i=1}^n 2x_i (y_i - x_i'\beta) = \frac{1}{\gamma} \sum_{i=1}^n x_i (y_i - x_i'\beta) = 0.$$

Rearranging the above equation in matrix form  $X'Y = X'X\hat{\beta}_{MLE}$ , we explicitly solve

$$\hat{\beta}_{MLE} = (X'X)^{-1}X'Y.$$

The maximum likelihood estimator (MLE) coincides with the OLS estimator. Similarly, the other FOC with respect to  $\gamma$  gives  $\hat{\gamma}_{MLE} = \hat{e}'\hat{e}/n$ .

### 4.3 Finite Sample Distribution

We can show the finite-sample exact distribution of  $\hat{\beta}$  assuming the error term follows a Gaussian distribution. *Finite sample distribution* means that the distribution holds for any  $n$ ; it is in contrast to *asymptotic distribution*, which is a large sample approximation to the finite sample distribution. We first review some properties of a generic jointly normal random vector.

**Fact 4.1.** Let  $z \sim N(\mu, \Omega)$  be an  $l \times 1$  random vector with a positive definite variance-covariance matrix  $\Omega$ . Let  $A$  be an  $m \times l$  non-random matrix where  $m \leq l$ . Then  $Az \sim N(A\mu, A\Omega A')$ .

**Fact 4.2.** If  $z \sim N(0, 1)$ ,  $w \sim \chi^2(d)$  and  $z$  and  $w$  are independent. Then  $\frac{z}{\sqrt{w/d}} \sim t(d)$ .

The OLS estimator

$$\hat{\beta} = (X'X)^{-1}X'Y = (X'X)^{-1}X'(X'\beta + e) = \beta + (X'X)^{-1}X'e,$$



and its conditional distribution can be written as

$$\begin{aligned}\widehat{\beta}|X &= \beta + (X'X)^{-1} X'e|X \\ &\sim \beta + (X'X)^{-1} X' \cdot N(0_n, \sigma^2 \cdot I_n) \\ &\sim N\left(\beta, \sigma^2 (X'X)^{-1} X'X (X'X)^{-1}\right) \sim N\left(\beta, \sigma^2 (X'X)^{-1}\right)\end{aligned}$$

by Fact 4.1. The  $k$ -th element of the vector coefficient

$$\widehat{\beta}_k|X = \eta'_k \widehat{\beta}|X \sim N\left(\beta_k, \sigma^2 \eta'_k (X'X)^{-1} \eta_k\right) \sim N\left(\beta_k, \sigma^2 (X'X)^{-1}_{kk}\right),$$

where  $\eta_k = (1 \{l = k\})_{l=1, \dots, K}$  is the selector of the  $k$ -th element.

In reality,  $\sigma^2$  is an unknown parameter, and

$$s^2 = \widehat{e}'\widehat{e} / (n - K) = e' M_X e / (n - K)$$

is an unbiased estimator of  $\sigma^2$ . Consider the  $t$ -statistic

$$\begin{aligned}T_k &= \frac{\widehat{\beta}_k - \beta_k}{\sqrt{s^2 [(X'X)^{-1}]_{kk}}} = \frac{\widehat{\beta}_k - \beta_k}{\sqrt{\sigma^2 [(X'X)^{-1}]_{kk}}} \cdot \frac{\sqrt{\sigma^2}}{\sqrt{s^2}} \\ &= \frac{(\widehat{\beta}_k - \beta_k) / \sqrt{\sigma^2 [(X'X)^{-1}]_{kk}}}{\sqrt{\frac{e' M_X e}{\sigma} / (n - K)}}.\end{aligned}$$

The numerator follows a standard normal, and the denominator follows  $\frac{1}{n-K} \chi^2(n-K)$ . Moreover, the numerator and the denominator are sta-

tistically independent (See Section 4.7). As a result, we conclude  $T_k \sim t(n - K)$  by Fact 4.2. This finite sample distribution allows us to conduct statistical inference.

## 4.4 Mean and Variance

Now we relax the normality assumption and statistical independence. Instead, we represent the regression model as  $Y = X\beta + e$  and

$$\begin{aligned} E[e|X] &= 0_n \\ \text{var}[e|X] &= E[ee'|X] = \sigma^2 I_n. \end{aligned}$$

where the first condition is the *mean independence* assumption, and the second condition is the *homoskedasticity* assumption. These assumptions are about the first and second *moments* of  $e_i$  conditional on  $x_i$ . Unlike the normality assumption, they do not restrict the distribution of  $e_i$ .

- Unbiasedness:

$$\begin{aligned} E[\hat{\beta}|X] &= E[(X'X)^{-1}XY|X] = E[(X'X)^{-1}X(X'\beta + e)|X] \\ &= \beta + (X'X)^{-1}XE[e|X] = \beta. \end{aligned}$$

Unbiasedness does not rely on homoskedasticity.

- Variance:

$$\begin{aligned}
\text{var} [\hat{\beta}|X] &= E \left[ \left( \hat{\beta} - E\hat{\beta} \right) \left( \hat{\beta} - E\hat{\beta} \right)' | X \right] \\
&= E \left[ \left( \hat{\beta} - \beta \right) \left( \hat{\beta} - \beta \right)' | X \right] \\
&= E \left[ (X'X)^{-1} X'ee'X (X'X)^{-1} | X \right] \\
&= (X'X)^{-1} X'E[ee'|X] X (X'X)^{-1}
\end{aligned}$$

where the second equality holds as  $E[\hat{\beta}] = E[E[\hat{\beta}|X]] = \beta$ . Under the assumption of homoskedasticity, it can be simplified as

$$\text{var} [\hat{\beta}|X] = (X'X)^{-1} X' \left( \sigma^2 I_n \right) X (X'X)^{-1} = \sigma^2 (X'X)^{-1}.$$

**Example 4.2.** (Heteroskedasticity) If  $e_i = x_i u_i$ , where  $x_i$  is a scalar random variable,  $u_i$  is statistically independent of  $x_i$ ,  $E[u_i] = 0$  and  $E[u_i^2] = \sigma^2$ . Then  $E[e_i|x_i] = 0$  but  $E[e_i^2|x_i] = \sigma^2 x_i^2$  is a function of  $x_i$ . We say  $e_i^2$  is a heteroskedastic error.

```

n = 100; X = rnorm(n)

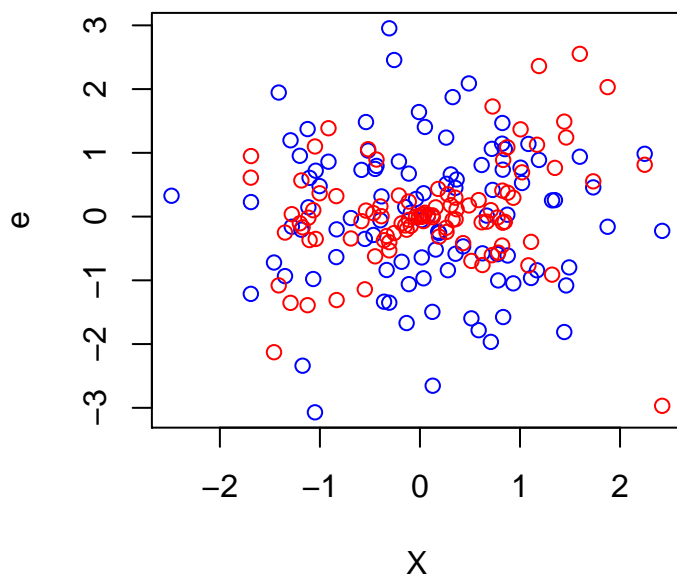
e1 = rnorm(n);

plot( y = e1, x = X, col = "blue", ylab = "e")

e2 = X * rnorm(n);

points( y = e2, x = X, col = "red")

```



It is important to notice that independently and identically distributed sample (iid)  $(y_i, x_i)$  does not imply homoskedasticity. Homoskedasticity or heteroskedasticity is about the relationship between  $(x_i, e_i = y_i - \beta x)$ , whereas iid is about the relationship between  $(y_i, x_i)$  and  $(y_j, x_j)$  for  $i \neq j$ .

## 4.5 Gauss-Markov Theorem

Gauss-Markov theorem is concerned about the optimality of OLS. It justifies OLS as the efficient estimator among all linear unbiased ones. *Efficient* here means that it enjoys the smallest variance in a family of estimators.

We have shown that OLS is unbiased in that  $E[\hat{\beta}] = \beta$ . There are numerous linearly unbiased estimators. For example,  $(Z'X)^{-1} Z'y$  for  $z_i = x_i^2$  is unbiased because  $E[(Z'X)^{-1} Z'y] = E[(Z'X)^{-1} Z'(X\beta + e)] = \beta$ . We cannot say OLS is better than those other unbiased estimators because they are equally good in this aspect. Thus, we move to the second order property of variance: an estimator is better if its variance is smaller.

**Fact 4.3.** *For two generic random vectors  $X$  and  $Y$  of the same size, we say  $X$ 's variance is smaller or equal to  $Y$ 's variance if  $(\Omega_Y - \Omega_X)$  is a positive semi-definite matrix. The comparison is defined this way because for any non-zero constant vector  $c$ , the variance of the linear combination of  $X$*

$$\text{var}(c'X) = c'\Omega_X c \leq c'\Omega_Y c = \text{var}(c'Y)$$

*is no bigger than the same linear combination of  $Y$ .*

Let  $\tilde{\beta} = A'y$  be a generic linear estimator, where  $A$  is any  $n \times K$  functions of  $X$ . As

$$E[A'y|X] = E[A'(X\beta + e)|X] = A'X\beta.$$

So the linearity and unbiasedness of  $\tilde{\beta}$  implies  $A'X = I_n$ . Moreover, the variance

$$\text{var}(A'y|X) = E[(A'y - \beta)(A'y - \beta)'|X] = E[A'ee'A|X] = \sigma^2 A'A.$$

Let  $C = A - X (X'X)^{-1}$ .

$$\begin{aligned} A'A - (X'X)^{-1} &= \left( C + X (X'X)^{-1} \right)' \left( C + X (X'X)^{-1} \right) - (X'X)^{-1} \\ &= C'C + (X'X)^{-1} X'C + C'X (X'X)^{-1} \\ &= C'C, \end{aligned}$$

where the last equality follows as

$$(X'X)^{-1} X'C = (X'X)^{-1} X' \left( A - X (X'X)^{-1} \right) = (X'X)^{-1} - (X'X)^{-1} = 0.$$

Therefore  $A'A - (X'X)^{-1}$  is a positive semi-definite matrix. The variance of any  $\tilde{\beta}$  is no smaller than the OLS estimator  $\hat{\beta}$ . The above derivation shows OLS achieves the smallest variance among all linear unbiased estimators.

Homoskedasticity is a restrictive assumption. Under homoskedasticity,  $\text{var} [\hat{\beta}] = \sigma^2 (X'X)^{-1}$ . Popular estimator of  $\sigma^2$  is the sample mean of the residuals  $\hat{\sigma}^2 = \frac{1}{n} \hat{e}'\hat{e}$  or the unbiased one  $s^2 = \frac{1}{n-K} \hat{e}'\hat{e}$ . Under heteroskedasticity, Gauss-Markov theorem does not apply.

## 4.6 Summary

The linear algebraic properties holds in finite sample no matter the data are taken as fixed numbers or random variables. The exact distribution under the normality assumption of the error term is the classical statistical

results. The Gauss Markov theorem holds under two crucial assumptions: linear CEF and homoskedasticity.

**Historical notes:** MLE was promulgated and popularized by Ronald Fisher (1890–1962). He was a major contributor of the frequentist approach which dominates mathematical statistics today, and he sharply criticized the Bayesian approach. Fisher collected the iris flower dataset of 150 observations in his biological study in 1936, which can be displayed in R by typing `iris`. Fisher invented the many concepts in classical mathematical statistics, such as sufficient statistic, ancillary statistic, completeness, and exponential family, etc.

**Further reading:** Phillips (1983) offers a comprehensive treatment of exact small sample theory in econometrics. After that, theoretical studies in econometrics swiftly shifted to large sample theory, which we will introduce in the next chapter.

## 4.7 Appendix

$Y = (y_1, \dots, y_n)$  consists of  $n$  iid observations. We say  $T(Y)$  is a sufficient statistic for a parameter  $\theta$  if the conditional probability  $f(Y|T(Y))$  does not depend on  $\theta$ . For example, for  $y_i \sim N(\mu, \sigma^2)$  with known  $\sigma^2$  and unknown  $\mu$ , We verify that the sample mean  $\bar{y} = n^{-1} \sum_{i=1}^n y_i$  is a sufficient

statistic for  $\mu$ . Notice that the joint density of  $Y$  is

$$\begin{aligned} f(Y) &= (2\pi\sigma)^{-n/2} \exp\left(-\frac{1}{2} \sum_{i=1}^n (y_i - \mu)^2\right) \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \bar{y})^2\right) \exp\left(-\frac{n}{2\sigma^2} (\bar{y} - \mu)^2\right). \end{aligned}$$

Because  $\bar{y} \sim N(\mu, \sigma^2/n)$ , the marginal density is

$$f(\bar{y}) = (2\pi\sigma^2/n)^{-1/2} \exp\left(-\frac{1}{2\sigma^2/n} (\bar{y} - \mu)^2\right).$$

The conditional density is

$$f(Y|\bar{y}) = \frac{f(Y)}{f(\bar{y})} = \frac{(2\pi\sigma^2)^{-n/2}}{(2\pi\sigma^2/n)^{-1/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \bar{y})^2\right)$$

is independent of  $\mu$ , and thus  $\bar{y}$  is a sufficient statistic for  $\mu$ .

In the meantime, the sample standard deviation  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$  is an *ancillary statistic* for  $\mu$ , because the distribution of  $s^2$  does not depend on  $\mu$ .

*Basu's theorem* says that a *complete* sufficient statistic is statistically independent from any ancillary statistic. For a normal distribution with unknown mean and known variance, the sample mean  $\bar{y}$  is the sufficient statistic and the sample standard deviation  $s^2$  is an ancillary statistic.

A parametric distribution indexed by  $\theta$  is a member of the *exponential*



*family* is its PDF can be written as

$$f(Y|\theta) = h(Y) g(\theta) \exp(\eta(\theta)' T(Y)),$$

where  $g(\theta)$  and  $\eta(\theta)$  are functions depend, only on  $\theta$  and  $h(Y)$  and  $T(Y)$  are functions depend only on  $Y$ . The normal distribution with known  $\sigma^2$  and unknown  $\mu$  belongs to the exponential family in view of the decomposition

$$\begin{aligned} f(Y) &= (\sqrt{2\pi}\sigma)^{-n} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2\right) \\ &= \underbrace{\exp\left(-\sum_{i=1}^n \frac{y_i^2}{2\sigma^2}\right)}_{h(Y)} \cdot \underbrace{(\sqrt{2\pi}\sigma)^{-n} \exp\left(-\frac{n}{2\sigma^2}\mu^2\right)}_{g(\theta)} \cdot \underbrace{\exp\left(\frac{\mu n}{2\sigma^2}\bar{y}\right)}_{\exp(\eta(\theta)' T(Y))}. \end{aligned}$$

The exponential family is a class of distributions with the special functional form which is convenient for deriving sufficient statistics as well as other desirable properties in classical mathematical statistics.

Zhentao Shi. Oct 6.

# Bibliography

Phillips, P. C. (1983). Exact small sample theory in the simultaneous equations model. *Handbook of econometrics 1*, 449–516.