

Statistical models and estimation¹

STA431 Spring 2023

¹See last slide for copyright information.

Overview

- 1 Models
- 2 MOM
- 3 MLE
- 4 Invariance
- 5 Consistency
- 6 Asymptotic Normality

Statistical model

Most good statistical analyses are based on a *model* for the data.

A *statistical model* is a set of assertions that partly specify the probability distribution of the observable data. The specification may be direct or indirect.

- Let x_1, \dots, x_n be a random sample from a normal distribution with expected value μ and variance σ^2 .
- For $i = 1, \dots, n$, let $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \epsilon_i$, where
 - β_0, \dots, β_k are unknown constants.
 - x_{ij} are known constants.
 - $\epsilon_1, \dots, \epsilon_n$ are independent $N(0, \sigma^2)$ random variables, not observable.
 - σ^2 is an unknown constant.
 - y_1, \dots, y_n are observable random variables.

A model is not the same thing as the *truth*.

Statistical models leave something unknown

Otherwise they are probability models

- The unknown part of the model is called the *parameter*.
- Usually, parameters are (vectors of) numbers.
- Usually denoted by θ or $\boldsymbol{\theta}$ or other Greek letters.
- In the non-Bayesian world, parameters are unknown constants.

Parameter Space

The *parameter space* is the set of values that can be taken on by the parameter.

- Let x_1, \dots, x_n be a random sample from a normal distribution with expected value μ and variance σ^2 .

The parameter space is

$$\Theta = \{(\mu, \sigma^2) : -\infty < \mu < \infty, \sigma^2 > 0\}.$$

- For $i = 1, \dots, n$, let $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \epsilon_i$, where
 - β_0, \dots, β_k are unknown constants.
 - x_{ij} are known constants.
 - $\epsilon_1, \dots, \epsilon_n$ are independent $N(0, \sigma^2)$ random variables.
 - σ^2 is an unknown constant.
 - y_1, \dots, y_n are observable random variables.

The parameter space is

$$\Theta = \{(\beta_0, \dots, \beta_k, \sigma^2) : -\infty < \beta_j < \infty, \sigma^2 > 0\}.$$

Parameters need not be numbers

Let X_1, \dots, X_n be a random sample from a continuous distribution with unknown distribution function $F(x)$.

- The parameter is the unknown distribution function $F(x)$.
- The parameter space is a space of distribution functions.
- We may be interested only in a *function* of the parameter, like

$$\mu = \int_{-\infty}^{\infty} x f(x) dx$$

The rest of $F(x)$ is just a nuisance parameter.

General statement of a statistical model

d is for Data

$$d \sim P_{\theta}, \quad \theta \in \Theta$$

- Both d and θ could be vectors
- For example,
 - $d = (\mathbf{y}_1, \dots, \mathbf{y}_n)$ independent multivariate normal.
 - $\theta = (\boldsymbol{\mu}, \boldsymbol{\Sigma})$.
 - P_{θ} is the joint distribution function of $\mathbf{y}_1, \dots, \mathbf{y}_n$, with joint density

$$f(\mathbf{y}_1, \dots, \mathbf{y}_n) = \prod_{i=1}^n f(\mathbf{y}_i; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Estimation

For the model $d \sim P_\theta$, $\theta \in \Theta$

- We don't know θ .
- We never know θ .
- All we can do is guess.
- Estimate θ (or a function of θ) based on the observable data.
- t is an *estimator* of θ (or a function of θ): $t = t(d)$
- t is a *statistic*, a random variable (vector) that can be computed from the data without knowing the values of any unknown parameters.

For example,

- $d = x_1, \dots, x_n \stackrel{i.i.d}{\sim} N(\mu, \sigma^2)$ $t = (\bar{x}, S^2)$.
- For an ordinary multiple regression model, $t = (\hat{\beta}, MSE)$

Parameter estimation

For the model $d \sim P_\theta$, $\theta \in \Theta$

- Estimate θ with $t = t(d)$.
- How do we get a recipe for t ? Guess?
- It's good to be systematic. Lots of methods are available.
- We will consider two: Method of moments and maximum likelihood.

Moments

Based on a random sample like $(x_1, y_1), \dots, (x_n, y_n)$

- Moments are quantities like $E\{x_i\}$, $E\{x_i^2\}$, $E\{x_i y_i\}$, $E\{W_i x_i^2 y_i^3\}$, etc.
- *Central* moments are moments of *centered* random variables:

$$E\{(x_i - \mu_x)^2\}$$

$$E\{(x_i - \mu_x)(y_i - \mu_y)\}$$

$$E\{(x_i - \mu_x)^2 (y_i - \mu_y)^3 (z_i - \mu_z)^2\}$$

- These are all *population* moments.

Population moments and sample moments

Population moment	Sample moment
$E\{x_i\}$	$\frac{1}{n} \sum_{i=1}^n x_i$
$E\{x_i^2\}$	$\frac{1}{n} \sum_{i=1}^n x_i^2$
$E\{x_i y_i\}$	$\frac{1}{n} \sum_{i=1}^n x_i y_i$
$E\{(x_i - \mu_x)^2\}$	$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2$
$E\{(x_i - \mu_x)(y_i - \mu_y)\}$	$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)(y_i - \bar{y}_n)$
$E\{(x_i - \mu_x)(y_i - \mu_y)^2\}$	$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)(y_i - \bar{y}_n)^2$

Estimation by the Method of Moments (MOM)

For the model $d \sim P_\theta$, $\theta \in \Theta$

- Population moments are a function of θ .
- Find θ as a function of the population moments.
- Estimate θ with that function of the *sample* moments.

Symbolically,

- Let m denote a vector of population moments.
- \hat{m} is the corresponding vector of sample moments.
- Find $m = g(\theta)$
- Solve for θ , obtaining $\theta = g^{-1}(m)$.
- Let $\hat{\theta} = g^{-1}(\hat{m})$.

It doesn't matter if you solve first or put hats on first.

Example: $x_1, \dots, x_n \stackrel{i.i.d}{\sim} U(0, \theta)$

$$f(x) = \frac{1}{\theta} \text{ for } 0 < x < \theta$$

First, find the moment (expected value).

$$\begin{aligned} E(x_i) &= \int_0^{\theta} x \frac{1}{\theta} dx \\ &= \frac{1}{\theta} \int_0^{\theta} x dx \\ &= \frac{1}{\theta} \left. \frac{x^2}{2} \right|_0^{\theta} = \frac{1}{2\theta} (\theta^2 - 0) \\ &= \frac{\theta}{2} \end{aligned}$$

So $m = \frac{\theta}{2} \iff \theta = 2m$, and $\hat{\theta} = 2\bar{x}$.

Small numerical example

Let x_1, \dots, x_n be a random sample from a uniform distribution on $(0, \theta)$. Estimate θ by the Method of Moments for the following data. Your answer is a number. Show some work.

4.09 0.13 0.84 3.83 2.13 4.67 4.61 0.40 4.19 0.71

$$\bar{x} = 2.56 \text{ so } \hat{\theta} = 2\bar{x} = 2 * 2.56 = 5.12.$$

Method of moments estimators are not unique

What moments you use are up to you.

$$E(x_i^2) = \frac{1}{\theta} \int_0^\theta x^2 dx = \frac{\theta^2}{3}$$

So set $m = \frac{\theta^2}{3} \Leftrightarrow \theta = \sqrt{3m}$, and

$$\hat{\theta} = \sqrt{\frac{3}{n} \sum_{i=1}^n x_i^2}$$

Compared to $2\bar{x}$.

Compare $\hat{\theta}_1 = 2\bar{x}$ and $\hat{\theta}_2 = \sqrt{\frac{3}{n} \sum_{i=1}^n x_i^2}$

For the numerical example

x	4.09	0.13	0.84	3.83	2.13	4.67	4.61	0.40	4.19
x ²	16.7281	0.0169	0.7056	14.6689	4.5369	21.8089	21.2521	0.16	17.5561

$$\hat{\theta}_1 = 5.12 \quad \hat{\theta}_2 = 5.42$$

Expressions for lower order moments tend to be simpler, and are preferable if only for that reason.

Method of Moments estimator for normal

Let $x_1, \dots, x_n \stackrel{i.i.d}{\sim} N(\mu, \sigma^2)$

From the moment-generating function or a textbook, $E(x_i) = \mu$ and $E(x_i^2) = \sigma^2 + \mu^2$. Solving for the parameters,

$$\begin{aligned}\mu &= E(x_i) \\ \sigma^2 &= E(x_i^2) - (E(x_i))^2\end{aligned}$$

so

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

A regression example

Independently for $i = 1, \dots, n$,

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ where}$$

- $E(x_i) = \mu_x, \text{Var}(x_i) = \sigma_x^2$
 - $E(\epsilon_i) = 0, \text{Var}(\epsilon_i) = \sigma_\epsilon^2$
 - x_i and ϵ_i are independent.
 - The distributions of x_i and ϵ_i are unknown.
 - What's the parameter?
-
- The parameter is $(\beta_0, \beta_1, F_\epsilon(\epsilon), F_x(x))$.
 - We want to estimate β_0 and β_1 , a *function* of the parameter.

Calculate some moments

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

$$E(x_i) = \mu_x$$

$$Var(x_i) = \sigma_x^2$$

$$E(y_i) = \beta_0 + \beta_1 \mu_x$$

$$Cov(x_i, y_i) = \beta_1 \sigma_x^2$$

$$\begin{aligned} Cov(x_i, y_i) &= Cov(x_i, \beta_0 + \beta_1 x_i + \epsilon_i) \\ &= Cov(x_i, \beta_1 x_i) + Cov(x_i, \epsilon_i) \\ &= \beta_1 Cov(x_i, x_i) + 0 \\ &= \beta_1 \sigma_x^2 \end{aligned}$$

Solve for β_0 and β_1

Have $E(x_i) = \mu_x$, $Var(x_i) = \sigma_x^2$, $E(y_i) = \beta_0 + \beta_1\mu_x$, $Cov(x_i, y_i) = \beta_1\sigma_x^2$

Putting hats on first, solve

$$\begin{aligned}\bar{y} &= \hat{\beta}_0 + \hat{\beta}_1\bar{x} \\ \hat{\sigma}_{xy} &= \hat{\beta}_1\hat{\sigma}_x^2\end{aligned}$$

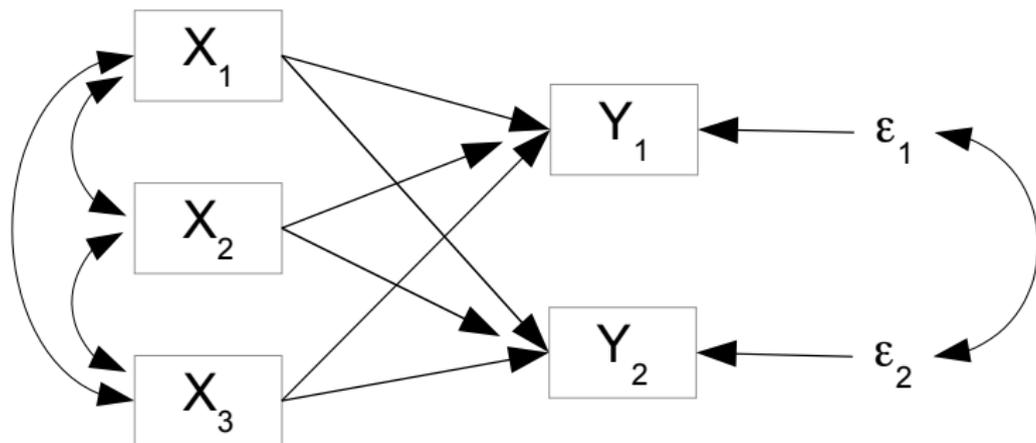
\Rightarrow

$$\begin{aligned}\hat{\beta}_1 &= \frac{\hat{\sigma}_{xy}}{\hat{\sigma}_x^2} = \frac{\sum_{i=1}^n (x_i - \bar{x}_n)(y_i - \bar{y}_n)}{\sum_{i=1}^n (x_i - \bar{x}_n)^2} \text{ and} \\ \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1\bar{x}\end{aligned}$$

These happen to be the same as the least-squares estimates.

Multivariate multiple regression

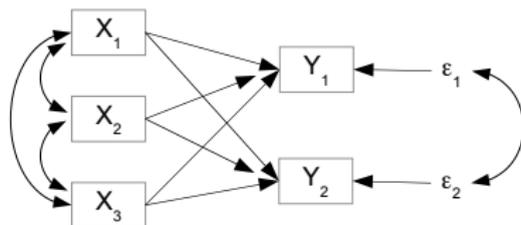
Multivariate means more than one response variable



We will obtain method of moments estimation for this.

One regression equation for each response variable

Give the equations in scalar form.



$$y_{i,1} = \beta_{1,0} + \beta_{1,1}x_{i,1} + \beta_{1,2}x_{i,2} + \beta_{1,3}x_{i,3} + \epsilon_{i,1}$$

$$y_{i,2} = \beta_{2,0} + \beta_{2,1}x_{i,1} + \beta_{2,2}x_{i,2} + \beta_{2,3}x_{i,3} + \epsilon_{i,2}$$

$$\mathbf{y}_i = \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 \mathbf{x}_i + \boldsymbol{\epsilon}_i$$

That's matrix form

In scalar form, had

$$y_{i,1} = \beta_{1,0} + \beta_{1,1}x_{i,1} + \beta_{1,2}x_{i,2} + \beta_{1,3}x_{i,3} + \epsilon_{i,1}$$

$$y_{i,2} = \beta_{2,0} + \beta_{2,1}x_{i,1} + \beta_{2,2}x_{i,2} + \beta_{2,3}x_{i,3} + \epsilon_{i,2}$$

In matrix form,

$$\mathbf{y}_i = \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 \mathbf{x}_i + \boldsymbol{\epsilon}_i$$

$$\begin{pmatrix} y_{i,1} \\ y_{i,2} \end{pmatrix} = \begin{pmatrix} \beta_{1,0} \\ \beta_{2,0} \end{pmatrix} + \begin{pmatrix} \beta_{1,1} & \beta_{1,2} & \beta_{1,3} \\ \beta_{2,1} & \beta_{2,2} & \beta_{2,3} \end{pmatrix} \begin{pmatrix} x_{i,1} \\ x_{i,2} \\ x_{i,3} \end{pmatrix} + \begin{pmatrix} \epsilon_{i,1} \\ \epsilon_{i,2} \end{pmatrix}$$

Note different order from $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \epsilon_i$

Statement of the model in general form

Independently for $i = 1, \dots, n$,

$$\mathbf{y}_i = \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 \mathbf{x}_i + \boldsymbol{\epsilon}_i, \text{ where}$$

- \mathbf{y}_i is an $q \times 1$ random vector of observable response variables, so the regression is multivariate; there are q response variables.
- \mathbf{x}_i is a $p \times 1$ observable random vector; there are p explanatory variables. $E(\mathbf{x}_i) = \boldsymbol{\mu}_x$ and $cov(\mathbf{x}_i) = \boldsymbol{\Phi}_{p \times p}$. The vector $\boldsymbol{\mu}_x$ and the matrix $\boldsymbol{\Phi}$ are unknown parameters.
- $\boldsymbol{\beta}_0$ is a $q \times 1$ vector of unknown constants.
- $\boldsymbol{\beta}_1$ is a $q \times p$ matrix of unknown constants. These are the regression coefficients, with one row for each response variable and one column for each explanatory variable.
- $\boldsymbol{\epsilon}_i$ is a $q \times 1$ unobservable random vector with expected value zero and unknown variance-covariance matrix $cov(\boldsymbol{\epsilon}_i) = \boldsymbol{\Psi}_{q \times q}$.
- $\boldsymbol{\epsilon}_i$ is independent of \mathbf{x}_i .

A Method of Moments estimate of β_1

$$\mathbf{y}_i = \beta_0 + \beta_1 \mathbf{x}_i + \epsilon_i$$

Denote the $p \times q$ matrix of (population) covariances between \mathbf{x}_i and \mathbf{y}_i by

$$\begin{aligned}\Sigma_{xy} &= \text{cov}(\mathbf{x}_i, \mathbf{y}_i) \\ &= \text{cov}(\mathbf{x}_i, \beta_0 + \beta_1 \mathbf{x}_i + \epsilon_i) \\ &= \text{cov}(\mathbf{x}_i, \beta_1 \mathbf{x}_i) + \text{cov}(\mathbf{x}_i, \epsilon_i) \\ &= \text{cov}(\mathbf{x}_i) \beta_1^\top + \mathbf{0} \\ &= \Phi \beta_1^\top\end{aligned}$$

Solve for β_1

In terms of moments of the observable data

$$\begin{aligned}\Phi\beta_1^\top &= \Sigma_{xy} \\ \Rightarrow \Phi^{-1}\Phi\beta_1^\top &= \Phi^{-1}\Sigma_{xy} \\ \Rightarrow \beta_1^\top &= \Phi^{-1}\Sigma_{xy} \\ \Rightarrow \beta_1 &= \Sigma_{xy}^\top(\Phi^{-1})^\top \\ &= \Sigma_{yx}\Phi^{-1} \\ &= \Sigma_{yx}\Sigma_x^{-1},\end{aligned}$$

Where $\Phi = cov(\mathbf{x}_i)$ is written Σ_x .

MOM estimate of β_1 based on $\beta_1 = \Sigma_{yx} \Sigma_x^{-1}$

Just put hats on.

$$\hat{\beta}_1 = \hat{\Sigma}_{yx} \hat{\Sigma}_x^{-1},$$

where

$$\begin{aligned}\hat{\Sigma}_{yx} &= \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top \\ \hat{\Sigma}_x &= \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top\end{aligned}$$

Method of Moments is Least Squares in this case

$$\hat{\beta}_1 = \hat{\Sigma}_{yx} \hat{\Sigma}_x^{-1}$$

- This is $(\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{y}$
- Transposed
- With both x and y variables centered by subtracting off the sample means.

Maximum likelihood estimation

A great idea from R. A. Fisher (1890-1962)

- Given a model and a set of observed data, how should we estimate θ ?
- Find the value of θ that makes the data we observed have the highest probability.
- If the model is continuous, maximize the probability of observing data in a little region surrounding the observed data vector.
- In either case, let $f(\mathbf{d}; \theta)$ denote the joint probability density function or probability mass function evaluated at the observed data vector.
- Maximize $L(\theta) = f(\mathbf{d}; \theta)$ over all $\theta \in \Theta$.
- $L(\theta)$ is called the *likelihood function*.

Maximum likelihood estimation for independent random sampling

$$d_1, \dots, d_n \stackrel{i.i.d.}{\sim} P_\theta, \theta \in \Theta.$$

$$L(\theta) = \prod_{i=1}^n f(d_i; \theta),$$

where $f(d_i; \theta)$ is the density or probability mass function evaluated at d_i .

- Find the value of θ for which $L(\theta)$ is maximum.
- Or equivalently, maximize $\ell(\theta) = \ln L(\theta)$.
- The elementary approach:
 - Take derivatives,
 - Set derivatives to zero,
 - Solve for θ ,
 - Put a hat on the answer.

Example

Maximum likelihood for the univariate normal

Let $x_1, \dots, x_n \stackrel{i.i.d}{\sim} N(\mu, \sigma^2)$.

$$\begin{aligned}\ell(\theta) &= \ln \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2}} \\ &= \ln \left(\sigma^{-n} (2\pi)^{-\frac{n}{2}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2} \right) \\ &= -n \ln \sigma - \frac{n}{2} \ln(2\pi) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\end{aligned}$$

Differentiate with respect to the parameters

$$\ell(\theta) = -n \ln \sigma - \frac{n}{2} \ln(2\pi) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$\begin{aligned} \frac{\partial \ell}{\partial \mu} &= -\frac{1}{2\sigma^2} \sum_{i=1}^n 2(x_i - \mu)(-1) \stackrel{\text{set}}{=} 0 \\ \Rightarrow \mu &= \bar{x} \end{aligned}$$

$$\begin{aligned} \frac{\partial \ell}{\partial \sigma} &= -\frac{n}{\sigma} - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 (-2\sigma^{-3}) \\ &= -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2 \stackrel{\text{set}}{=} 0 \\ \Rightarrow \sigma^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \end{aligned}$$

Substituting

Setting derivatives to zero, we have obtained

$$\mu = \bar{x} \text{ and } \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2, \text{ so}$$

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

Gamma Example

Let x_1, \dots, x_n be a random sample from a Gamma distribution with parameters $\alpha > 0$ and $\beta > 0$

$$f(x; \alpha, \beta) = \frac{1}{\beta^\alpha \Gamma(\alpha)} e^{-x/\beta} x^{\alpha-1}$$

$$\Theta = \{(\alpha, \beta) : \alpha > 0, \beta > 0\}$$

Log Likelihood

$$f(x; \alpha, \beta) = \frac{1}{\beta^\alpha \Gamma(\alpha)} e^{-x/\beta} x^{\alpha-1}$$

$$\begin{aligned} \ell(\alpha, \beta) &= \ln \prod_{i=1}^n \frac{1}{\beta^\alpha \Gamma(\alpha)} e^{-x_i/\beta} x_i^{\alpha-1} \\ &= \ln \left(\beta^{-n\alpha} \Gamma(\alpha)^{-n} \exp\left(-\frac{1}{\beta} \sum_{i=1}^n x_i\right) \left(\prod_{i=1}^n x_i\right)^{\alpha-1} \right) \\ &= -n\alpha \ln \beta - n \ln \Gamma(\alpha) - \frac{1}{\beta} \sum_{i=1}^n x_i + (\alpha - 1) \sum_{i=1}^n \ln x_i \end{aligned}$$

Differentiate with respect to the parameters

$$\ell(\theta) = -n\alpha \ln \beta - n \ln \Gamma(\alpha) - \frac{1}{\beta} \sum_{i=1}^n x_i + (\alpha - 1) \sum_{i=1}^n \ln x_i$$

$$\frac{\partial \ell}{\partial \beta} \stackrel{\text{set}}{=} 0 \Rightarrow \alpha \beta = \bar{x}$$

$$\begin{aligned} \frac{\partial \ell}{\partial \alpha} &= -n \ln \beta - n \frac{\partial}{\partial \alpha} \ln \Gamma(\alpha) + \sum_{i=1}^n \ln x_i \\ &= \sum_{i=1}^n \ln x_i - n \ln \beta - n \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} \stackrel{\text{set}}{=} 0 \end{aligned}$$

Solve for α

$$\sum_{i=1}^n \ln x_i - n \ln \beta - n \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} = 0$$

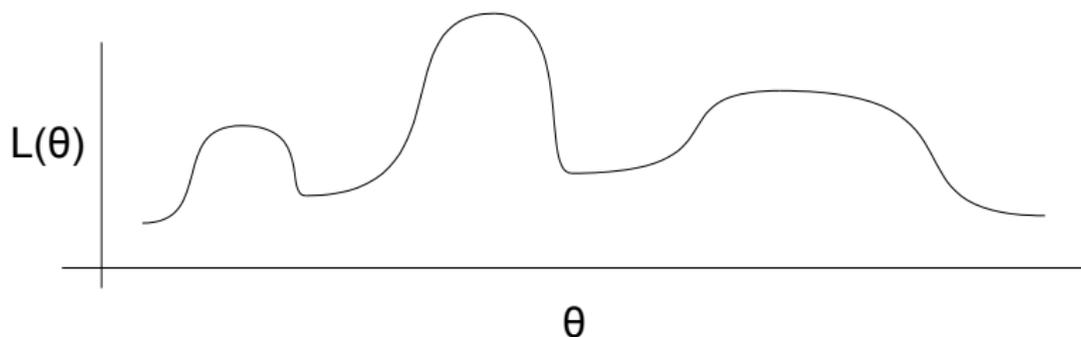
where

$$\Gamma(\alpha) = \int_0^{\infty} e^{-t} t^{\alpha-1} dt.$$

Nobody can do it.

Maximize the likelihood numerically with software

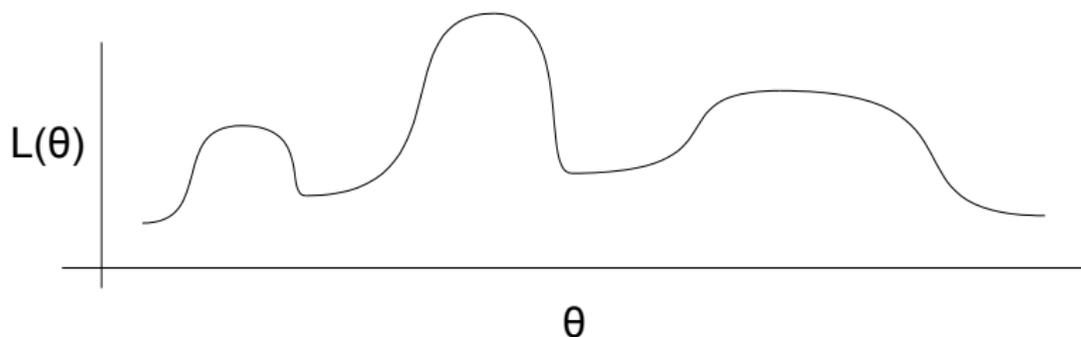
Usually this is in high dimension



- It's like trying to find the top of a mountain by walking uphill blindfolded.
- You might stop at a local maximum.
- The starting place is very important.
- The final answer is a number (or vector of numbers).
- There is no explicit formula for the MLE.

There is a lot of useful theory

Even without an explicit formula for the MLE



- MLE is asymptotically normal.
- Variance of the MLE is deeply related to the curvature of the log likelihood at the MLE.
- The more curvature, the smaller the variance.
- The variance of the MLE can be estimated from the curvature (using the Fisher Information).
- Basis of tests and confidence intervals.

Comparing MOM and MLE

- Sometimes they are identical, sometimes not.
- If the model is right they are usually close for large samples.
- Both are asymptotically normal, centered around the true parameter value(s).
- Estimates of the variance are easy to obtain for both.
- Small variance of an estimator is good.
- As $n \rightarrow \infty$, nothing can beat the MLE.
- Except that the MLE depends on a very specific distribution.
- And sometimes the dependence matters.
- In such cases, MOM may be preferable.

Gamma Example

$$f(x; \alpha, \beta) = \frac{1}{\beta^\alpha \Gamma(\alpha)} e^{-x/\beta} x^{\alpha-1}$$

$$\ell(\alpha, \beta) = -n\alpha \ln \beta - n \ln \Gamma(\alpha) - \frac{1}{\beta} \sum_{i=1}^n x_i + (\alpha - 1) \sum_{i=1}^n \ln x_i$$

R function for the minus log likelihood

```
gmll = function(theta,datta)
{
  aa = theta[1]; bb = theta[2]
  nn = length(datta); sumd = sum(datta)
  sumlogd = sum(log(datta))
  value = nn*aa*log(bb) + nn*lgamma(aa) + sumd/bb - (aa-1)*sumlogd
  return(value)
} # End function gmll
```

Simulated Data

$n = 50$, True $\alpha = 2$, $\beta = 3$

```
> d
```

```
[1] 20.87 13.74 5.13 2.76 4.73 2.66 11.74 0.75 22.07 10.49 7.26 5.82  
[13] 13.08 1.79 4.57 1.40 1.13 6.84 3.21 0.38 11.24 1.72 4.69 1.96  
[25] 7.87 8.49 5.31 3.40 5.24 1.64 7.17 9.60 6.97 10.87 5.23 5.53  
[37] 15.80 6.40 11.25 4.91 12.05 5.44 12.62 1.81 2.70 3.03 4.09 12.29  
[49] 3.23 10.94
```

Where should the numerical search start?

- How about Method of Moments estimates?
- $E(x) = \alpha\beta$, $Var(x) = \alpha\beta^2$
- Solve for α and β , replace population moments by sample moments and put a \sim above the parameters.

$$\tilde{\alpha} = \frac{\bar{x}^2}{s^2} \quad \text{and} \quad \tilde{\beta} = \frac{s^2}{\bar{x}}$$

```
> # MOM for starting values
> momalpha = mean(d)^2/var(d); momalpha
[1] 1.899754
> mombeta = var(d)/mean(d); mombeta
[1] 3.620574
```

Numerical search using the `optim` function

```
> # Error message says: "Bounds can only be used with method
> #                               L-BFGS-B (or Brent)"
> gsearch = optim(par=c(momalpha,mombeta), fn = gml1,
+               method = "L-BFGS-B", lower = c(0,0), hessian=TRUE, datta=d)
```

```
> gsearch
```

```
$par
```

```
[1] 1.805930 3.808674
```

```
$value
```

```
[1] 142.0316
```

```
$counts
```

```
function gradient
```

```
          9          9
```

```
$convergence
```

```
[1] 0
```

```
$message
```

```
[1] "CONVERGENCE: REL_REDUCTION_OF_F <= FACTR*EPSMCH"
```

```
$hessian
```

```
      [,1]      [,2]
```

```
[1,] 36.69402 13.127928
```

```
[2,] 13.12793  6.224773
```

Meaning of the output

Output

Meaning

\$par

[1] 1.805930 3.808674

$\hat{\theta}$

\$value

[1] 142.0316

$-\ell(\hat{\theta})$

\$hessian

	[,1]	[,2]
[1,]	36.69402	13.127928
[2,]	13.12793	6.224773

$$\mathbf{H} = \left[\frac{\partial^2(-\ell)}{\partial \theta_i \partial \theta_j} \right]_{\theta = \hat{\theta}}$$

Second Derivative test

$$\mathbf{H} = \left[\frac{\partial^2(-\ell)}{\partial\theta_i\partial\theta_j} \right]$$

- If the second derivatives are continuous, \mathbf{H} is symmetric.
- If the gradient is zero at a point and $|\mathbf{H}| \neq 0$,
 - If \mathbf{H} is positive definite, local minimum
 - If \mathbf{H} is negative definite, local maximum
 - If \mathbf{H} has both positive and negative eigenvalues, saddle point

MLE

```
> thetahat = gsearch$par
> names(thetahat) = c("alpha-hat","beta-hat"); thetahat
alpha-hat  beta-hat
  1.805930  3.808674

> # Second derivative test
> eigen(gsearch$hessian)$values
[1] 41.569998  1.348796
```

The Invariance principle of maximum likelihood estimation

- The Invariance Principle of maximum likelihood estimation says that *the MLE of a function is that function of the MLE.*²

²Provided the function is one-to-one.

Example of the invariance principle

Let d_1, \dots, d_n be a random sample from a Bernoulli distribution ($1=\text{Yes}$, $0=\text{No}$) with parameter θ , $0 < \theta < 1$.

- The parameter space is $\Theta = (0, 1)$
- MLE is $\hat{\theta} = \bar{d}$, the sample proportion.
- Write the model in terms of the *odds* of $d_i = 1$, a re-parameterization that is often useful in categorical data analysis.
- Denote the odds by $\theta' = \frac{\theta}{1-\theta}$.
- $\theta' = \frac{\theta}{1-\theta} \iff \theta = \frac{\theta'}{1+\theta'}$.
- As θ ranges from zero to one, θ' ranges from zero to infinity.
- So there is a new parameter space: $\theta' \in \Theta' = (0, \infty)$.

MLE of the odds

- $\theta' = \frac{\theta}{1-\theta} \iff \theta = \frac{\theta'}{1+\theta'}$
- Because the re-parameterization is one-to-one, $\hat{\theta}' = \frac{\bar{d}}{1-\bar{d}}$ without any calculation.

Theorem

See text for a proof. The one-to-one part is critical.

Let $g : \Theta \rightarrow \Theta'$ be a one-to-one re-parameterization, with the maximum likelihood estimate $\hat{\theta}$ satisfying $L(\hat{\theta}) > L(\theta)$ for all $\theta \in \Theta$ with $\theta \neq \hat{\theta}$. Then $L'(g(\hat{\theta})) > L'(\theta')$ for all $\theta' \in \Theta'$ with $\theta' \neq g(\hat{\theta})$.

In other words

- The MLE of $g(\theta)$ is $g(\hat{\theta})$.
- $\widehat{g(\theta)} = g(\hat{\theta})$.
- The MLE of θ' is $g(\hat{\theta})$.
- $\hat{\theta}' = g(\hat{\theta})$.

Re-parameterization in general

The parameters of common statistical models are written in a standard way, but other equivalent parameterizations are sometimes useful.

Suppose $x_i \sim N(\mu, \sigma^2)$. Have

$$\hat{\theta} = \left(\bar{x}, \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right)$$

- Write $x_i \sim N(\mu, \sigma)$.
 - $g(\theta) = (\theta_1, \sqrt{\theta_2})$
 - $\hat{\theta}' = \left(\bar{x}, \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \right)$
- Write $x_i \sim N(\mu, \tau)$, where $\tau = 1/\sigma^2$ is called the *precision*.
 - $g(\theta) = (\theta_1, 1/\theta_2)$
 - $\hat{\theta}' = \left(\bar{x}, \frac{n}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)$

Consistency

- The idea is large-sample accuracy.
- As $n \rightarrow \infty$, you get the truth.
- It's a kind of limit, but with probability involved.

The setting

- Let t_1, t_2, \dots be a sequence of random variables.
- Main application: t_n is an estimator of θ based on a sample of size n .
- Think $t_n = \bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i$.

Definition of Convergence in Probability

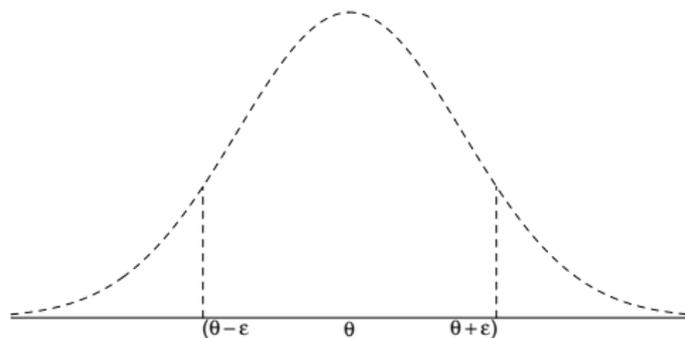
We say that t_n converges *in probability* to the constant θ , and write $t_n \xrightarrow{P} \theta$ if for all $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} P\{|t_n - \theta| < \epsilon\} = 1$$

Convergence in probability to θ means no matter how small the interval around θ , the probability distribution of t_n becomes concentrated in that interval as $n \rightarrow \infty$.

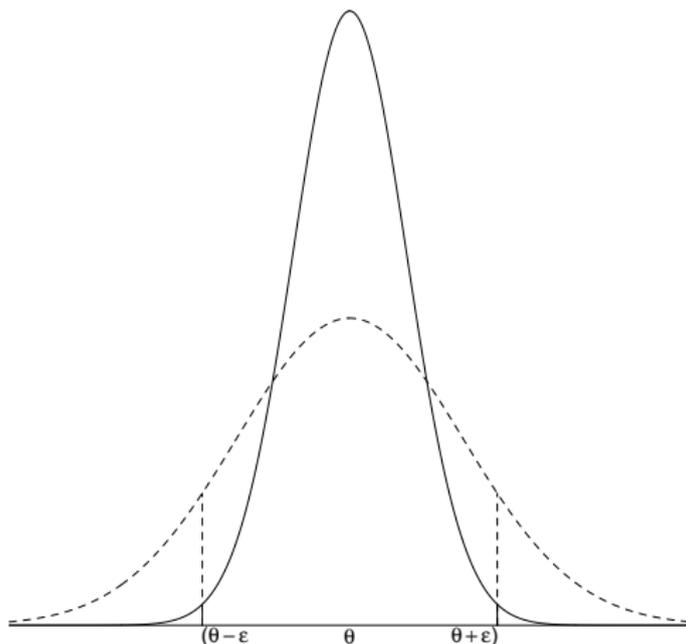
Picture it

$$\begin{aligned} P\{|t_n - t| < \epsilon\} &= P\{-\epsilon < t_n - \theta < \epsilon\} \\ &= P\{\theta - \epsilon < t_n < \theta + \epsilon\} \end{aligned}$$



Picture it

$$\begin{aligned}P\{|t_n - t| < \epsilon\} &= P\{-\epsilon < t_n - \theta < \epsilon\} \\ &= P\{\theta - \epsilon < t_n < \theta + \epsilon\}\end{aligned}$$



The Law of Large Numbers

We will use this a lot

Let x_1, x_2, \dots be independent random variables from a distribution with expected value μ and variance σ^2 . The Law of Large Numbers says

$$\bar{x}_n \xrightarrow{p} \mu$$

Roadmap

Markov's Inequality



Variance Rule



Law of Large Numbers

Markov's Inequality

For $g(x) \geq 0$ and $a \geq 0$,

$$E\{g(x)\} \geq a \Pr\{g(x) \geq a\}$$

To prove, split up the integral.

Variance Rule

- Let t_1, t_2, \dots be a sequence of random variables
- With $E(t_n) = \mu_n$ and $Var(t_n) = \sigma_n^2$
- If $\lim_{n \rightarrow \infty} \mu_n = \theta$ and $\lim_{n \rightarrow \infty} \sigma_n^2 = 0$, then

$$t_n \xrightarrow{p} \theta$$

To prove, let $g(x) = (x - \mu)^2$ and $a = \epsilon^2$ in Markov's inequality.

Proving the Law of Large Numbers

The Variance Rule says

- Let t_1, t_2, \dots be a sequence of random variables
- With $E(t_n) = \mu_n$ and $Var(t_n) = \sigma_n^2$
- If $\lim_{n \rightarrow \infty} \mu_n = \theta$ and $\lim_{n \rightarrow \infty} \sigma_n^2 = 0$, then $t_n \xrightarrow{p} \theta$.

- Let $t_n = \bar{x}_n$ and $\theta = \mu$.
- $E(\bar{x}_n) = \mu$ and $Var(\bar{x}_n) = \frac{\sigma^2}{n} \rightarrow 0$
- Conclude

$$\bar{x}_n \xrightarrow{p} \mu$$

The Change of Variables formula: Let $y = g(x)$

$$E(y) = \int_{-\infty}^{\infty} y f_y(y) dy = \int_{-\infty}^{\infty} g(x) f_x(x) dx$$

Or, for discrete random variables

$$E(y) = \sum_y y p_y(y) = \sum_x g(x) p_x(x)$$

This is actually a big theorem, not a definition.

Applying the change of variables formula

To approximate $E[g(x)]$

Have x_1, \dots, x_n from the distribution of x . Want $E(y)$, where $y = g(x)$.

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n g(x_i) &= \frac{1}{n} \sum_{i=1}^n y_i \xrightarrow{p} E(y) \\ &= E(g(x)) \end{aligned}$$

So for example

$$\frac{1}{n} \sum_{i=1}^n x_i^k \xrightarrow{p} E(x^k)$$

$$\frac{1}{n} \sum_{i=1}^n U_i^2 V_i W_i^3 \xrightarrow{p} E(U^2 V W^3)$$

- That is, sample moments converge in probability to population moments.
- Central sample moments converge to central population moments as well.

Population moments and sample moments

Repeating an earlier slide

Population moment	Sample moment
$E\{x_i\}$	$\frac{1}{n} \sum_{i=1}^n x_i$
$E\{x_i^2\}$	$\frac{1}{n} \sum_{i=1}^n x_i^2$
$E\{x_i y_i\}$	$\frac{1}{n} \sum_{i=1}^n x_i y_i$
$E\{(x_i - \mu_x)^2\}$	$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2$
$E\{(x_i - \mu_x)(y_i - \mu_y)\}$	$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)(y_i - \bar{y}_n)$
$E\{(x_i - \mu_x)(y_i - \mu_y)^2\}$	$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)(y_i - \bar{y}_n)^2$

Convergence in Probability for Random Vectors

Let $\mathbf{t}_1, \mathbf{t}_2, \dots$ be a sequence of k -dimensional random vectors.

We say that \mathbf{t}_n converges in probability to $\boldsymbol{\theta} \in \mathbb{R}^k$, and write

$\mathbf{t}_n \xrightarrow{p} \boldsymbol{\theta}$ if for all $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} P\{\|\mathbf{t}_n - \boldsymbol{\theta}\| < \epsilon\} = 1,$$

where $\|\mathbf{a} - \mathbf{b}\|$ denotes Euclidian distance in \mathbb{R}^k .

Two more Theorems

- The “stack” theorem and continuous mapping.
- Often used together.

The “Stack” Theorem

Because I don't know what to call it.

Let $\mathbf{x}_n \xrightarrow{p} \mathbf{a}$ and $\mathbf{y}_n \xrightarrow{p} \mathbf{b}$. Then the partitioned random vector

$$\begin{pmatrix} \mathbf{x}_n \\ \mathbf{y}_n \end{pmatrix} \xrightarrow{p} \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix}$$

Continuous mapping

One of the Slutsky lemmas: See Appendix A

Let $\mathbf{t}_n \xrightarrow{p} \mathbf{c}$, and let the function $g(\mathbf{x})$ be continuous at $\mathbf{x} = \mathbf{c}$.
Then

$$g(\mathbf{t}_n) \xrightarrow{p} g(\mathbf{c})$$

Note that the function g could be multidimensional, for example mapping \mathbb{R}^5 into \mathbb{R}^2 .

Definition of Consistency

The random vector (of statistics) \mathbf{t}_n is said to be a *consistent* estimator of the parameter vector $\boldsymbol{\theta}$ if

$$\mathbf{t}_n \xrightarrow{p} \boldsymbol{\theta}$$

for all $\boldsymbol{\theta} \in \Theta$.

Consistency of the Sample Variance

This answer gets full marks.

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

By LLN, $\bar{x}_n \xrightarrow{p} \mu$ and $\frac{1}{n} \sum_{i=1}^n x_i^2 \xrightarrow{p} E(x_i^2) = \sigma^2 + \mu^2$.

By continuous mapping,

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 \xrightarrow{p} \sigma^2 + \mu^2 - \mu^2 = \sigma^2$$

Note the silent use of the Stack Theorem.

Method of Moments Estimators are Consistent

For most practical cases

Recall

- Let m denote a vector of population moments.
- \hat{m} is the corresponding vector of sample moments.
- Find $m = g(\theta)$
- Solve for θ , obtaining $\theta = g^{-1}(m)$.
- Let $\hat{\theta}_n = g^{-1}(\hat{m}_n)$.

If g is continuous, so is g^{-1} . Then by continuous mapping,
 $\hat{m} \xrightarrow{P} m \Rightarrow \hat{\theta}_n = g^{-1}(\hat{m}_n) \xrightarrow{P} g^{-1}(m) = \theta$.

Maximum Likelihood Estimators are Consistent

If the model is correct, and given some additional “regularity conditions.”

$$\hat{\boldsymbol{\theta}}_n \xrightarrow{p} \boldsymbol{\theta}$$

Consistency is great but it's not enough.

- It's the least we can ask. Estimators that are *not* consistent are completely unacceptable for most purposes.
- Think of $a_n = 1/n$ as a sequence of degenerate random variables with $P\{a_n = 1/n\} = 1$.
- So, $a_n \xrightarrow{p} 0$.

Suppose

$$t_n \xrightarrow{p} \theta \Rightarrow U_n = t_n + \frac{100,000,000}{n} \xrightarrow{p} \theta.$$

Convergence in Distribution

Sometimes called *Weak Convergence*, or *Convergence in Law*

Denote the cumulative distribution functions of t_1, t_2, \dots by $F_1(x), F_2(x), \dots$ respectively, and denote the cumulative distribution function of t by $F(x)$.

We say that t_n converges *in distribution* to t , and write $t_n \xrightarrow{d} t$ if for every point x at which F is continuous,

$$\lim_{n \rightarrow \infty} F_n(x) = F(x)$$

We will seldom use this definition directly.

Connections among the Modes of Convergence

- $t_n \xrightarrow{p} t \Rightarrow t_n \xrightarrow{d} t.$

- If a is a constant, $t_n \xrightarrow{d} a \Rightarrow t_n \xrightarrow{p} a.$

Univariate Central Limit Theorem

Let x_1, \dots, x_n be a random sample from a distribution with expected value μ and variance σ^2 . Then

$$z_n = \frac{\sqrt{n}(\bar{x}_n - \mu)}{\sigma} \xrightarrow{d} z \sim N(0, 1)$$

Sometimes we say the distribution of the sample mean is approximately normal, or asymptotically normal.

- This is justified by the Central Limit Theorem.
- But it does *not* mean that \bar{x}_n converges in distribution to a normal random variable.
- The Law of Large Numbers says that \bar{x}_n converges in probability to a constant, μ .
- So \bar{x}_n converges to μ in distribution as well.

Why would we say that for large n , the sample mean is approximately $N(\mu, \frac{\sigma^2}{n})$?

Have $z_n = \frac{\sqrt{n}(\bar{x}_n - \mu)}{\sigma} \xrightarrow{d} z \sim N(0, 1)$.

$$\begin{aligned} Pr\{\bar{x}_n \leq x\} &= Pr\left\{\frac{\sqrt{n}(\bar{x}_n - \mu)}{\sigma} \leq \frac{\sqrt{n}(x - \mu)}{\sigma}\right\} \\ &= Pr\left\{z_n \leq \frac{\sqrt{n}(x - \mu)}{\sigma}\right\} \approx \Phi\left(\frac{\sqrt{n}(x - \mu)}{\sigma}\right) \end{aligned}$$

Suppose y is *exactly* $N(\mu, \frac{\sigma^2}{n})$:

$$\begin{aligned} Pr\{y \leq x\} &= Pr\left\{\frac{\sqrt{n}(y - \mu)}{\sigma} \leq \frac{\sqrt{n}(x - \mu)}{\sigma}\right\} \\ &= Pr\left\{z_n \leq \frac{\sqrt{n}(x - \mu)}{\sigma}\right\} = \Phi\left(\frac{\sqrt{n}(x - \mu)}{\sigma}\right) \end{aligned}$$

Asymptotic Normality

- We say \bar{x}_n is asymptotically normal, with asymptotic mean μ and asymptotic variance $\frac{\sigma^2}{n}$.
- Write $\bar{x}_n \overset{\sim}{\sim} N(\mu, \frac{\sigma^2}{n})$
- In tests and confidence intervals, $\frac{\hat{\sigma}^2}{n}$ may be used in place of $\frac{\sigma^2}{n}$, where $\hat{\sigma}^2$ is any consistent estimator of σ^2 .

Asymptotic *Multivariate* Normality

- Multivariate central limit theorem
- Central limit theorem for vectors of MLEs

Multivariate central limit theorem

Let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be i.i.d. p -dimensional random vectors with expected value vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. Then

$$\sqrt{n}(\bar{\mathbf{x}}_n - \boldsymbol{\mu}) \xrightarrow{d} \mathbf{X} \sim N_p(\mathbf{0}, \boldsymbol{\Sigma})$$

- Say $\bar{\mathbf{x}}_n$ is asymptotically $N_p(\boldsymbol{\mu}, \frac{1}{n}\boldsymbol{\Sigma})$.
- Write $\bar{\mathbf{x}}_n \overset{\sim}{\sim} N_p(\boldsymbol{\mu}, \frac{1}{n}\boldsymbol{\Sigma})$.
- The asymptotic covariance matrix of $\bar{\mathbf{x}}_n$ is $\frac{1}{n}\boldsymbol{\Sigma}$.
- $\boldsymbol{\Sigma}$ may be estimated by the sample variance-covariance matrix $\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top$.

Central limit theorem for vectors of MLEs

$$\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)$$

If the model is correct and under some technical conditions that always hold for the models used in this class,

$$\sqrt{n}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}) \xrightarrow{d} \mathbf{t} \sim N_k(\mathbf{0}, \mathcal{I}(\boldsymbol{\theta})^{-1}),$$

where (for the record) $\mathcal{I}(\boldsymbol{\theta})$ is the Fisher information matrix.

$$\mathcal{I}(\boldsymbol{\theta}) = \left[E \left[-\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(y; \boldsymbol{\theta}) \right] \right]$$

See Appendix A.

Asymptotic Multivariate Normality of the MLEs

- Say $\hat{\boldsymbol{\theta}}_n$ is asymptotically $N_k(\boldsymbol{\theta}, \mathbf{V}_n)$, where $\mathbf{V}_n = \frac{1}{n}\mathcal{I}(\boldsymbol{\theta})^{-1}$.
- Write $\hat{\boldsymbol{\theta}}_n \overset{\sim}{\sim} N_k(\boldsymbol{\theta}, \mathbf{V}_n)$.
- For tests and confidence intervals replace \mathbf{V}_n by either
 - $\hat{\mathbf{V}}_n = \frac{1}{n}\mathcal{I}(\hat{\boldsymbol{\theta}})^{-1}$, or
 - $\hat{\mathbf{V}}_n$ = the inverse of the Hessian of the minus log likelihood, evaluated at the MLE.
 - For numerical MLEs, the second choice is usually more convenient.

Back to the Gamma Example

```
gsearch = optim(par=c(momalpha,mombeta), fn = gml1,
               method = "L-BFGS-B", lower = c(0,0), hessian=TRUE, datta=d)
```

Output

Meaning

\$par

[1] 1.805930 3.808674

$$\hat{\boldsymbol{\theta}} = (\hat{\alpha}, \hat{\beta})$$

\$value

[1] 142.0316

$$-\ell(\hat{\boldsymbol{\theta}})$$

\$hessian

	[,1]	[,2]
[1,]	36.69402	13.127928
[2,]	13.12793	6.224773

$$\mathbf{H} = \left[\frac{\partial^2(-\ell)}{\partial \theta_i \partial \theta_j} \right]_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}$$

```
> # Asymptotic variance-covariance matrix is the inverse of the
> # Fisher Information
> Vhat_n = solve(gsearch$hessian); Vhat_n
      [,1]      [,2]
[1,] 0.1110190 -0.2341369
[2,] -0.2341369 0.6544386
> # Confidence interval for alpha (true value is 2)
> thetahat
alpha-hat  beta-hat
 1.805930  3.808674
> se_alphahat = sqrt(Vhat_n[1,1])
> lower95 = thetahat[1] - 1.96*se_alphahat
> upper95 = thetahat[1] + 1.96*se_alphahat
> c(lower95,upper95)
alpha-hat  alpha-hat
 1.152868  2.458992
```

Copyright Information

This slide show was prepared by **Jerry Brunner**, Department of Statistics, University of Toronto. It is licensed under a **Creative Commons Attribution - ShareAlike 3.0 Unported License**. Use any part of it as you like and share the result freely. The $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ source code is available from the course website:
<http://www.utstat.toronto.edu/brunner/oldclass/431s23>