

# Statistical models and estimation<sup>1</sup>

STA431 Spring 2015

---

<sup>1</sup>See last slide for copyright information.

# Overview

**1** Models

**2** MOM

**3** MLE

# 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  $\mu$  and variance  $\sigma^2$ .
- For  $i = 1, \dots, n$ , let  $Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \epsilon_i$ , where
  - $\beta_0, \dots, \beta_k$  are unknown constants.
  - $x_{ij}$  are known constants.
  - $\epsilon_1, \dots, \epsilon_n$  are independent  $N(0, \sigma^2)$  random variables.
  - $\sigma^2$  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 for the data is called the *parameter*.
- Usually, parameters are (vectors of) numbers.
- Usually denoted by  $\theta$  or  $\boldsymbol{\theta}$  or other Greek letters.
- 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  $\mu$  and variance  $\sigma^2$ .

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

$\beta_0, \dots, \beta_k$  are unknown constants.

$x_{ij}$  are known constants.

$\epsilon_1, \dots, \epsilon_n$  are independent  $N(0, \sigma^2)$  random variables.

$\sigma^2$  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  $\theta$  could be vectors
- For example,
  - $D = \mathbf{Y}_1, \dots, \mathbf{Y}_n$  independent multivariate normal.
  - $\theta = (\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .
  - $P_\theta$  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  $\theta$ .
- We never know  $\theta$ .
- All we can do is guess.
- Estimate  $\theta$  (or a function of  $\theta$ ) based on the observable data.
- $T$  is an *estimator* of  $\theta$  (or a function of  $\theta$ ):  $T = T(D)$

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,  $T = (\hat{\beta}, MSE)$

$T$  is a *statistic*, a random variable (vector) that can be computed from the data without knowing the values of any unknown parameters.



# Parameter estimation

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

- Estimate  $\theta$  with  $T = T(D)$ .
- How do we get a recipe for  $T$ ?
- 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  $\theta$ .
- Find  $\theta$  as a function of the population moments.
- Estimate  $\theta$  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  $\theta$ , 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  $m = \frac{\theta}{2} \Leftrightarrow \theta = 2m$ , 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  $\theta$  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 <sup>2</sup>	16.7281	0.0169	0.7056	14.6689	4.5369	21.8089	21.2521	0.16	17.5561

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



# 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

$$\begin{aligned}\hat{\mu} &= \bar{X} \\ \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2\end{aligned}$$

# 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  $\epsilon_i$  are independent.
  - The distributions of  $X_i$  and  $\epsilon_i$  are unknown.
  - What's the parameter?
- 
- The parameter is  $(\beta_0, \beta_1, F_\epsilon(\epsilon), F_x(x))$ .
  - We want to estimate  $\beta_0$  and  $\beta_1$ , a *function* of the parameter.

# Calculate some moments

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

$$E(X_i) = \mu_x$$

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

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

$$Cov(X_i, Y_i) = \beta_1 \sigma_x^2$$

Use the centering rule to get the last one:

$$\begin{aligned} Cov(X_i, Y_i) &= E(\overset{c}{X}_i \overset{c}{Y}_i) \\ &= E\{\overset{c}{X}_i (\beta_1 \overset{c}{X}_i + \epsilon_i)\} \\ &= E\{\beta_1 \overset{c}{X}_i^2 + \overset{c}{X}_i \epsilon_i\} \\ &= \beta_1 E\{\overset{c}{X}_i^2\} + E\{\overset{c}{X}_i\} E\{\epsilon_i\} \\ &= \beta_1 \sigma_x^2 \end{aligned}$$

Solve for  $\beta_0$  and  $\beta_1$

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

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

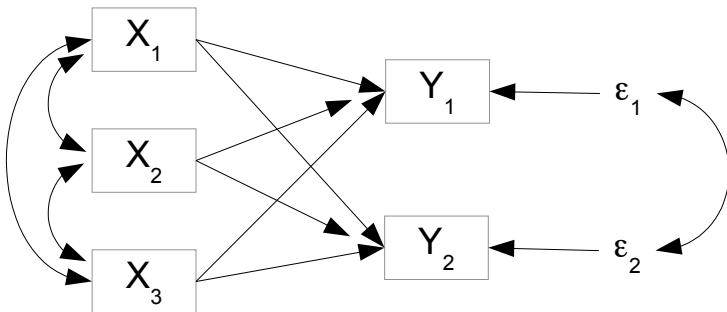
$\Rightarrow$

$$\begin{aligned}\hat{\beta}_1 &= \frac{\hat{\sigma}_{xy}}{\hat{\sigma}_x^2} = \frac{\sum_{i=1}^n (X_i - \bar{X}_n)(Y_i - \bar{Y}_n)}{\sum_{i=1}^n (X_i - \bar{X}_n)^2} \text{ and} \\ \hat{\beta}_0 &= \bar{Y} - \hat{\beta}_1\bar{X}\end{aligned}$$

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

# Multivariate multiple regression

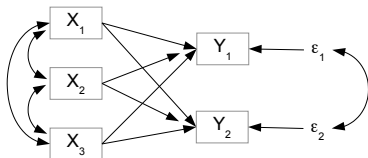
Multivariate means more than one response variable



We will obtain method of moments estimation for this.

# One regression equation for each response variable

Give the equations in scalar form.



$$Y_{i,1} = \beta_{0,1} + \beta_{1,1}X_{i,1} + \beta_{1,2}X_{i,2} + \beta_{1,3}X_{i,3} + \epsilon_{i,1}$$

$$Y_{i,2} = \beta_{0,2} + \beta_{2,1}X_{i,1} + \beta_{2,2}X_{i,2} + \beta_{2,3}X_{i,3} + \epsilon_{i,2}$$

$$\mathbf{Y}_i = \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 \mathbf{X}_i + \boldsymbol{\epsilon}_i$$

In scalar form,

$$\begin{aligned} Y_{i,1} &= \beta_{0,1} + \beta_{1,1}X_{i,1} + \beta_{1,2}X_{i,2} + \beta_{1,3}X_{i,3} + \epsilon_{i,1} \\ Y_{i,2} &= \beta_{0,2} + \beta_{2,1}X_{i,1} + \beta_{2,2}X_{i,2} + \beta_{2,3}X_{i,3} + \epsilon_{i,2} \end{aligned}$$

In matrix form,

$$\begin{aligned} \mathbf{Y}_i &= \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 \mathbf{X}_i + \boldsymbol{\epsilon}_i \\ \begin{pmatrix} Y_{i,1} \\ Y_{i,2} \end{pmatrix} &= \begin{pmatrix} \beta_{0,1} \\ \beta_{0,2} \end{pmatrix} + \begin{pmatrix} \beta_{1,1} & \beta_{1,2} & \beta_{1,3} \\ \beta_{2,1} & \beta_{2,2} & \beta_{2,3} \end{pmatrix} \begin{pmatrix} X_{i,1} \\ X_{i,2} \\ X_{i,3} \end{pmatrix} + \begin{pmatrix} \epsilon_{i,1} \\ \epsilon_{i,2} \end{pmatrix} \end{aligned}$$

# Statement of the model

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  $V(\mathbf{X}_i) = \boldsymbol{\Phi}_{p \times p}$ . The matrix  $\boldsymbol{\Phi}$  is unknown.
- $\boldsymbol{\beta}_0$  is a  $q \times 1$  matrix 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  $q \times 1$  random vector with expected value zero and unknown variance-covariance matrix  $V(\boldsymbol{\epsilon}_i) = \boldsymbol{\Psi}_{q \times q}$ .
- $\boldsymbol{\epsilon}_i$  is independent of  $\mathbf{X}_i$ .



# A Method of Moments estimate of $\beta_1$

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

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

$$\begin{aligned}\Sigma_{xy} &= C(\mathbf{X}_i, \mathbf{Y}_i) \\ &= E\{\overset{c}{\mathbf{X}}_i \overset{c}{\mathbf{Y}}_i^\top\} \\ &= E\{\overset{c}{\mathbf{X}}_i (\beta_1 \overset{c}{\mathbf{X}}_i + \epsilon_i)^\top\} \\ &= E\{\overset{c}{\mathbf{X}}_i (\overset{c}{\mathbf{X}}_i^\top \beta_1^\top + \epsilon_i^\top)\} \\ &= E\{\overset{c}{\mathbf{X}}_i \overset{c}{\mathbf{X}}_i^\top \beta_1^\top + \overset{c}{\mathbf{X}}_i \epsilon_i^\top\} \\ &= E\{\overset{c}{\mathbf{X}}_i \overset{c}{\mathbf{X}}_i^\top\} \beta_1^\top + E\{\overset{c}{\mathbf{X}}_i \epsilon_i^\top\} \\ &= V(\mathbf{X}_i) \beta_1^\top + C(\mathbf{X}_i, \epsilon_i) \\ &= \Phi \beta_1^\top + \mathbf{0} \\ &= \Phi \beta_1^\top\end{aligned}$$

# Solve for $\beta_1$

In terms of moments of the observable data

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

Noting  $\Phi = V(\mathbf{X}_i)$  could be written  $\Sigma_x$ , have  $\beta_1 = \Sigma_{yx}\Sigma_x^{-1}$

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

Just put hats on.

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

where

$$\hat{\Sigma}_{yx} = \frac{1}{n} \sum_{i=1}^n (\mathbf{Y}_i - \bar{\mathbf{Y}})(\mathbf{X}_i - \bar{\mathbf{X}})^\top$$

$$\hat{\Sigma}_x = \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})^\top$$

# 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  $\theta$ ?
- Find the value of  $\theta$  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  $\theta$  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  $\theta$ ,
  - Put a hat on the answer.

## Example: Coffee taste test

A fast food chain is considering a change in the blend of coffee beans they use to make their coffee. To determine whether their customers prefer the new blend, the company plans to select a random sample of  $n = 100$  coffee-drinking customers and ask them to taste coffee made with the new blend and with the old blend, in cups marked “ $A$ ” and “ $B$ .” Half the time the new blend will be in cup  $A$ , and half the time it will be in cup  $B$ . Management wants to know if there is a difference in preference for the two blends.

## Statistical model for the taste test example

Letting  $\theta$  denote the probability that a consumer will choose the new blend, treat the data  $Y_1, \dots, Y_n$  as a random sample from a Bernoulli distribution. That is, independently for  $i = 1, \dots, n$ ,

$$f(y_i; \theta) = \theta^{y_i} (1 - \theta)^{1 - y_i}$$

for  $y_i = 0$  or  $y_i = 1$ , and zero otherwise.

# Find the MLE of $\theta$

Show your work

Maximize the log likelihood.

$$\begin{aligned}\frac{\partial}{\partial \theta} \ln L(\theta) &= \frac{\partial}{\partial \theta} \ln \left( \prod_{i=1}^n f(y_i; \theta) \right) \\ &= \frac{\partial}{\partial \theta} \ln \left( \prod_{i=1}^n \theta^{y_i} (1 - \theta)^{1 - y_i} \right) \\ &= \frac{\partial}{\partial \theta} \ln \left( \theta^{\sum_{i=1}^n y_i} (1 - \theta)^{n - \sum_{i=1}^n y_i} \right) \\ &= \frac{\partial}{\partial \theta} \left( \left( \sum_{i=1}^n y_i \right) \ln \theta + \left( n - \sum_{i=1}^n y_i \right) \ln(1 - \theta) \right) \\ &= \frac{\sum_{i=1}^n y_i}{\theta} - \frac{n - \sum_{i=1}^n y_i}{1 - \theta}\end{aligned}$$



Setting the derivative to zero,

$$\begin{aligned}\frac{\sum_{i=1}^n y_i}{\theta} &= \frac{n - \sum_{i=1}^n y_i}{1 - \theta} \Rightarrow (1 - \theta) \sum_{i=1}^n y_i = \theta(n - \sum_{i=1}^n y_i) \\ &\Rightarrow \sum_{i=1}^n y_i - \theta \sum_{i=1}^n y_i = n\theta - \theta \sum_{i=1}^n y_i \\ &\Rightarrow \sum_{i=1}^n y_i = n\theta \\ &\Rightarrow \theta = \frac{\sum_{i=1}^n y_i}{n}\end{aligned}$$

So it looks like the MLE is the sample proportion. Carrying out the second derivative test to be sure,

## Second derivative test

$$\begin{aligned}\frac{\partial^2 \ln \ell}{\partial \theta^2} &= \frac{\partial}{\partial \theta} \left( \frac{\sum_{i=1}^n y_i}{\theta} - \frac{n - \sum_{i=1}^n y_i}{1 - \theta} \right) \\ &= \frac{-\sum_{i=1}^n y_i}{\theta^2} - - - \frac{n - \sum_{i=1}^n y_i}{(1 - \theta)^2} \\ &= -n \left( \frac{1 - \bar{y}}{(1 - \theta)^2} + \frac{\bar{y}}{\theta^2} \right) < 0\end{aligned}$$

Concave down, maximum, verifying  $\hat{\theta} = \bar{y}$ .

## Numerical estimate

Suppose 60 of the 100 consumers prefer the new blend. Give a point estimate the parameter  $\theta$ . Your answer is a number.

```
> ybar = 60/100; ybar  
[1] 0.6
```

# 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  $\alpha$ 

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

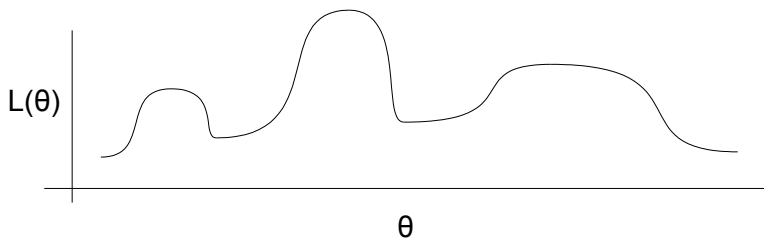
where

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

Nobody can do it.

# Maximize the likelihood numerically with software

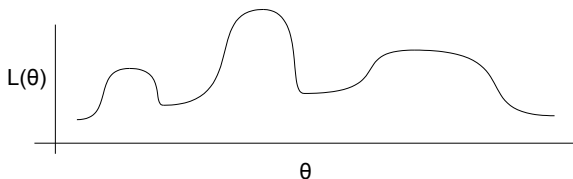
Usually this is in high dimension



- 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.
- Estimates of the variance are well known 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.

## 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/431s15>