

MAXIMUM LIKELIHOOD ESTIMATION  
of MEAN and VARIANCE  
as an ~~example~~ of MULTIVARIATE OPTIMIZATION

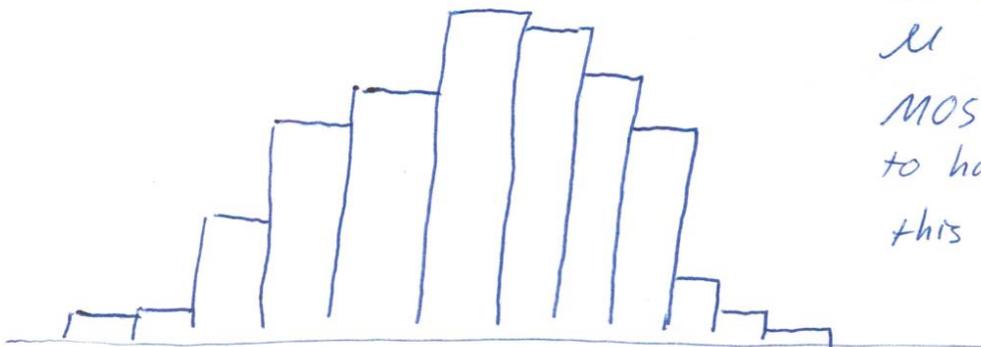
- Maximum Likelihood Estimation (MLE) is a statistical technique for estimating the parameters of a probability distribution that generated an observed set of data.
- In other words: "We see this data here. I'm willing to assume that it ~~was~~ comes from a NORMAL DISTRIBUTION. If that's the case, then what are the MEAN and VARIANCE of that distribution?"
- The ESTIMATION QUESTION is:

What values of  $\mu$  and  $\sigma$  are MOST LIKELY to have generated this data?

Optimization in two variables:  $\mu$  and  $\sigma$

p. 2

- The simplest example of MLE is to estimate the mean  $\mu$  and variance  $\sigma$  of a normal distribution
- This example is easily extensible to higher level analyses. For example, ~~the~~ ~~general~~ multivariate regression models, probability models and models with distributions other than ~~normal~~ normal (e.g., lognormal, logistic)
- And because this example only has two variables it serves as a good example for a course in quantitative methods.
- If your data is distributed normally then a histogram of the data may look like this:



"What values of  $\mu$  and  $\sigma$  are MOST LIKELY to have generated this data?"

(p.3)

→ If we assume that the true, underlying distribution is NORMAL, then we want to maximize the following likelihood function with respect to  $\mu$  and  $\sigma$ :

$$\text{MAX}_{\mu, \sigma} \mathcal{L} = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma}} \cdot e^{-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma}}$$

which is equivalent to maximizing the log of that same function:

$$\begin{aligned} \text{MAX}_{\mu, \sigma} \ln \mathcal{L} &= \sum_{i=1}^N \left\{ \ln \left[ (2\pi\sigma)^{-\frac{1}{2}} \right] - \frac{(x_i - \mu)^2}{2\sigma} \right\} \\ &= -\frac{N}{2} \ln(2\pi\sigma) - \frac{1}{2\sigma} \cdot \sum_{i=1}^N (x_i - \mu)^2 \end{aligned}$$

→ The FIRST-ORDER CONDITIONS are:

$$\frac{\partial \ln \mathcal{L}}{\partial \mu} = \frac{1}{\sigma} \sum_{i=1}^N (x_i - \mu) = 0 \quad \text{which implies: } \mu = \frac{1}{N} \sum x_i$$

$$\frac{\partial \ln \mathcal{L}}{\partial \sigma} = \frac{-N}{2\sigma} + \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2 = 0 \quad \text{which implies: } \sigma = \frac{1}{N} \sum (x_i - \mu)^2$$

ex. 4

→ The SECOND-ORDER CONDITIONS are

$$\frac{\partial^2 \ln L}{\partial \mu^2} = -\frac{N}{\sigma} < 0$$

"MEAN'S"  
OWN SECOND PARTIAL  
~~IS POSITIVE~~  
IS NEGATIVE



$$\frac{\partial^2 \ln L}{\partial \sigma^2} = \frac{N}{2\sigma^2} - \frac{1}{\sigma^3} \sum (x_i - \mu)^2$$

"VARIANCE'S"  
OWN SECOND PARTIAL  
IS NEGATIVE



$$= \frac{N}{2\sigma^2} - \frac{N}{\sigma^2} = \frac{N}{\sigma^2} \left( \frac{1}{2} - 1 \right) = -\frac{N}{2\sigma^2} < 0$$

↑ BY the  
1st Order Condition

AND THE PRODUCT OF THE OWN SECOND PARTIALS IS GREATER THAN THE PRODUCT OF THE CROSS PARTIALS (i.e. the square of the cross partial by Young's Theorem)

$$\frac{\partial^2 \ln L}{\partial \mu \partial \sigma} = -\frac{1}{\sigma^2} \sum (x_i - \mu)$$

$$= -\frac{1}{\sigma^2} [(\sum x_i) - N\mu] = 0$$

← BY THE FIRST ORDER CONDITION

$$H = \begin{bmatrix} \frac{\partial^2 \ln L}{\partial \mu^2} & \frac{\partial^2 \ln L}{\partial \mu \partial \sigma} \\ \frac{\partial^2 \ln L}{\partial \mu \partial \sigma} & \frac{\partial^2 \ln L}{\partial \sigma^2} \end{bmatrix} = \begin{bmatrix} -\frac{N}{\sigma} & 0 \\ 0 & -\frac{N}{2\sigma^2} \end{bmatrix}$$

★ HESSIAN MATRIX ★

7.5

→ The DETERMINANT of THE HESSIAN  $|H|$  is the product of the own second partials minus the square of the cross partial

$$|H| = \left( \frac{-N}{\gamma} \right) \left( \frac{-N}{2\gamma^2} \right) - 0 = 0$$

$$= \frac{N^2}{2\gamma^3} > 0$$

SATISFACTION OF THIS SECOND ORDER CONDITION ENSURES THAT WE ARE AT A MAXIMUM (as opposed to a SADDLE <sup>POINT</sup> ~~POINT~~)

→ The INFORMATION MATRIX is the negative of the INVERSE of the Hessian Matrix

$$I = -1 \cdot H^{-1} = \begin{bmatrix} \frac{\gamma}{N} & 0 \\ 0 & \frac{2\gamma^2}{N} \end{bmatrix}$$

The "information" contained ~~in~~ in this matrix are our measure of the precision with which we estimated the mean & variance.

Std. error of our estimate of the mean

$$\sqrt{\frac{\gamma}{N}}$$

variance of our estimate of the variance

$$\frac{2\gamma^2}{N}$$

p. 6

→ The true mean,  $\mu$  of  $x$  is unknown to us, what ~~was~~ MLE provides is an ESTIMATE of  $\mu$ . So it would be more accurate to write:

$$\hat{\mu} = \frac{1}{N} \sum x_i \quad \text{MEAN of } x$$

→ Similarly, we only obtain an estimate of the variance,  $\sigma^2$  of  $x$ , so it would be more accurate:

$$\hat{\sigma}^2 = \frac{1}{N} \sum (x_i - \hat{\mu})^2 \quad \text{VARIANCE of } x$$

→ Our estimate of the standard error of  $x$  is

$$\hat{\sigma} = \sqrt{\hat{\sigma}^2} = \sqrt{\frac{1}{N} \sum (x_i - \hat{\mu})^2}$$

→ BUT REMEMBER: We measure the mean of  $x$  and the variance of  $x$  with error.

→ The values of the own second partials (or more precisely the negative of the inverse of those values) tell us how precisely we have measured the mean of  $x$  and the variance of  $x$ .

p. 7

→ From the own second partials we get the STANDARD ERROR of the ESTIMATE of the MEAN (as opposed to the STANDARD ERROR of  $X$ )

$$\begin{array}{l} \text{std error} \\ \text{of estimate} \\ \text{of mean} \end{array} \sqrt{\frac{\hat{\sigma}^2}{N}} \quad \begin{array}{l} \leftarrow \text{increasing fn of} \\ \text{estimated variance} \\ \leftarrow \text{decreasing fn of} \\ \text{number of observations} \end{array}$$

→ We also obtain the variance of our estimate of the variance:

$$\frac{2\hat{\sigma}^2}{N} \quad \begin{array}{l} \leftarrow \text{INCREASING FN of} \\ \text{ESTIMATED VARIANCE} \\ \leftarrow \text{DECREASING FN of} \\ \text{NUMBER of OBSERVATIONS} \end{array}$$

→ SO the own SECOND PARTIALS ~~are~~ provide measures of ~~precision~~ HOW PRECISELY we have estimated the mean and variance.

→ The remaining question is:

"Why ~~does~~ do the measures of precision come from the own second partials?"

→ The answer can be seen by comparing the following functions

$$f(x) = -x^2 \quad \text{and} \quad g(x) = -2x^2$$

7.8

WIDER and LESS PRECISE

$$f(x) = -x^2$$

$$f'(x) = -2x$$

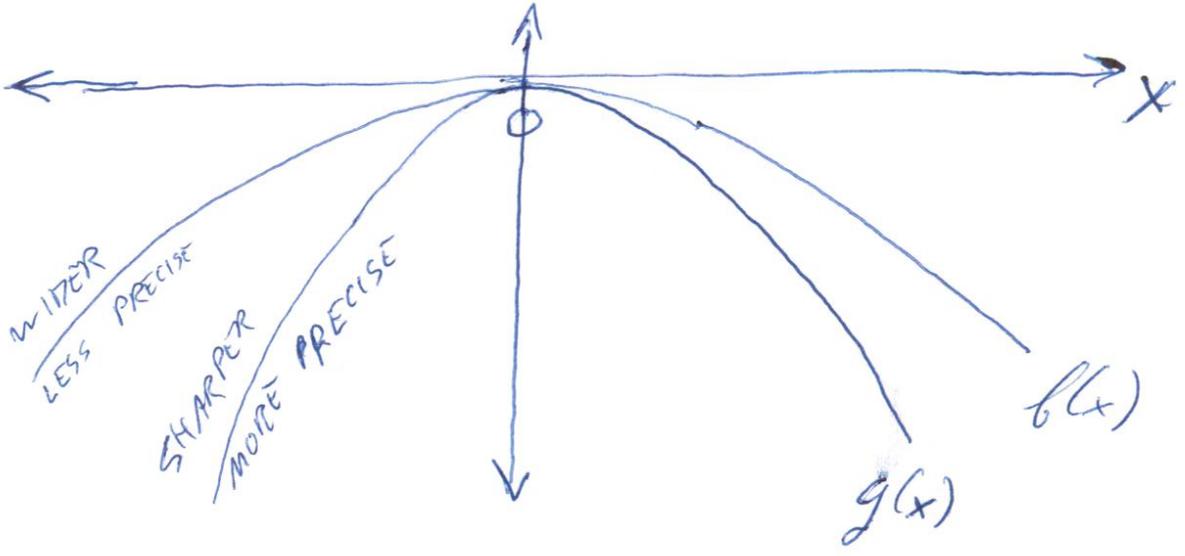
$$f''(x) = -2$$

$$g(x) = -2x^2$$

$$g'(x) = -4x$$

$$g''(x) = -4$$

SHARPER + therefore MORE PRECISE



→ Both  $f(x)$  and  $g(x)$  are maximized at  $x=0$ , but  $g(x)$  comes to a sharper peak, thus representing a more precise estimate

→ Comparing the negative of the inverses allows ~~that~~ us to compare ~~the~~ what would be the ~~the~~ measure of precision

Pretend that this is the ~~measure~~ ~~of the estimate~~ variance of the estimate

$$\frac{-1}{f''(x=0)} > \frac{-1}{g''(x=0)}$$

$$\frac{1}{2} > \frac{1}{4}$$

The  $g(x)$  function would yield lower variance of the estimate. It is a more precise measure.