



# Maximum Likelihood Estimation

PSYC 943: Fundamentals  
of Multivariate Modeling  
Lecture 4: September 18, 2013

# Today's Class

- The basics of maximum likelihood estimation
  - The engine that drives most modern statistical methods
- Additional information from maximum likelihood estimator (MLEs)
  - Likelihood ratio tests
  - Wald tests
  - Information criteria
- MLEs for GLMs
  - An introduction to SAS PROC MIXED

# Today's Example Data #1

- Imagine an employer is looking to hire employees for a job where IQ is important
  - We will only use 5 observations so as to show the math behind the estimation calculations
- The employer collects two variables:
  - IQ scores
  - Job performance
- Descriptive Statistics:

Variable	Mean	SD
IQ	114.4	2.30
Performance	12.8	2.28

Covariance Matrix		
IQ	5.3	5.1
Performance	5.1	5.2

Observation	IQ	Performance
1	112	10
2	113	12
3	115	14
4	118	16
5	114	12

# How Estimation Works (More or Less)

- Most estimation routines do one of three things:
  1. **Minimize Something**: Typically found with names that have “least” in the title. Forms of least squares include “Generalized”, “Ordinary”, “Weighted”, “Diagonally Weighted”, “WLSMV”, and “Iteratively Reweighted.” Typically the estimator of last resort...

Last Week's Class
  2. **Maximize Something**: Typically found with names that have “maximum” in the title. Forms include “Maximum likelihood”, “ML”, “Residual Maximum Likelihood” (REML), “Robust ML”. Typically the gold standard of estimators (and next week we'll see why).

Today's Class
  3. **Use Simulation to Sample from Something**: more recent advances in simulation use resampling techniques. Names include “Bayesian Markov Chain Monte Carlo”, “Gibbs Sampling”, “Metropolis Hastings”, “Metropolis Algorithm”, and “Monte Carlo”. Used for complex models where ML is not available or for methods where prior values are needed.

Future Class



# **AN INTRODUCTION TO MAXIMUM LIKELIHOOD ESTIMATION**

# Properties of Maximum Likelihood Estimators

- Provided several assumptions (“regularity conditions”) are met, maximum likelihood estimators have good statistical properties:
  1. Asymptotic Consistency: as the sample size increases, the estimator converges in probability to its true value
  2. Asymptotic Normality: as the sample size increases, the distribution of the estimator is normal (with variance given by “information” matrix)
  3. Efficiency: No other estimator will have a smaller standard error
- Because they have such nice and well understood properties, MLEs are commonly used in statistical estimation

## Maximum Likelihood: Estimates Based on Statistical Distributions

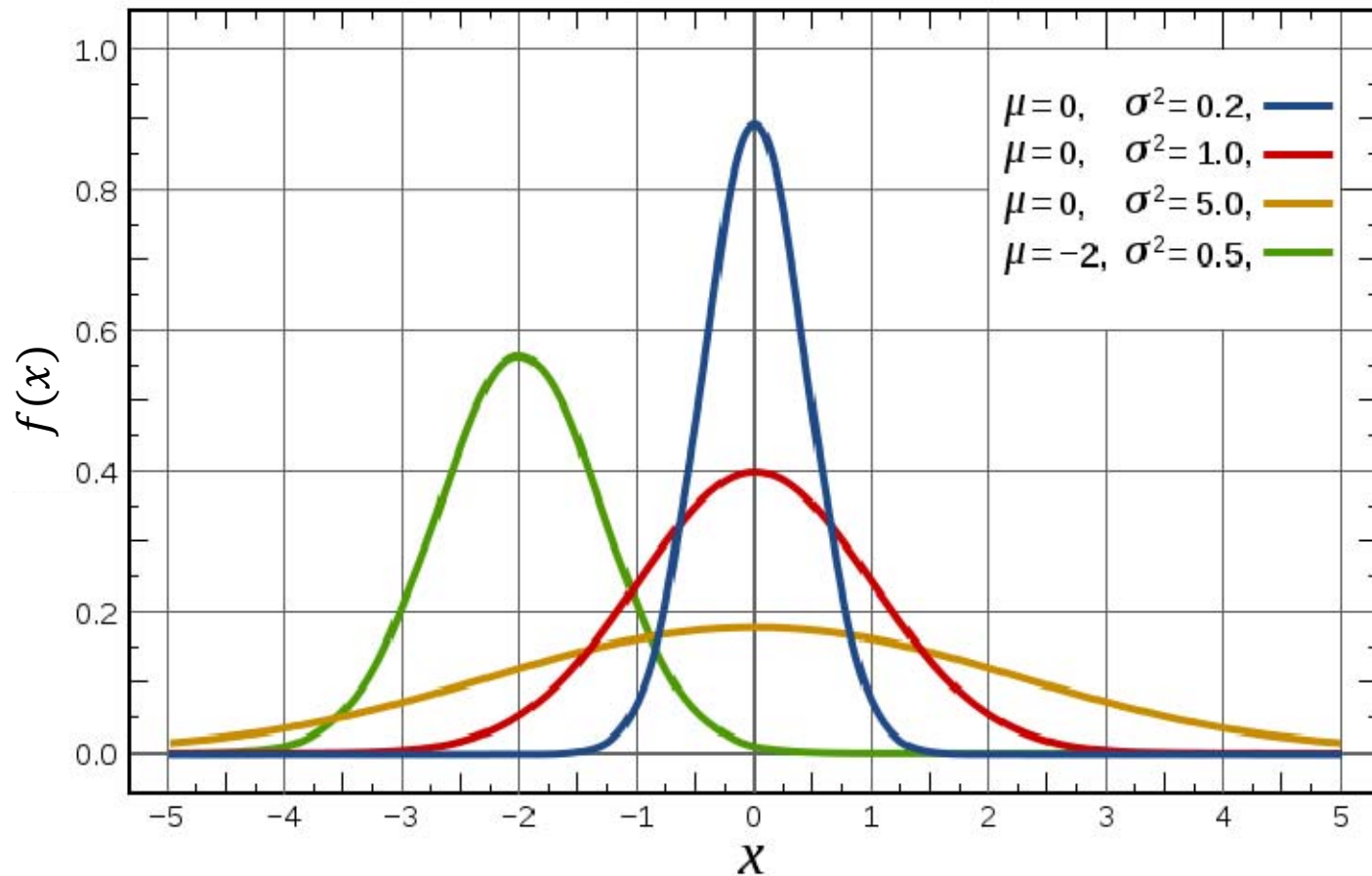
- Maximum likelihood estimates come from statistical distributions – assumed distributions of data
  - We will begin today with the univariate normal distribution but quickly move to other distributions (see this Friday's class)

- For a single random variable  $x$ , the univariate normal distribution is

$$f(x) = \frac{1}{\sqrt{2\pi\sigma_x^2}} \exp\left(-\frac{(x - \mu_x)^2}{2\sigma_x^2}\right)$$

- Provides the height of the curve for a value of  $x$ ,  $\mu_x$ , and  $\sigma_x^2$
- Last week we pretended we knew  $\mu_x$  and  $\sigma_x^2$ 
  - Today we will only know  $x$  (and maybe  $\sigma_x^2$ )

# Univariate Normal Distribution



For any value of  $x$ ,  $\mu_x$ , and  $\sigma_x^2$ ,  $f(x)$  gives the height of the curve (relative frequency)

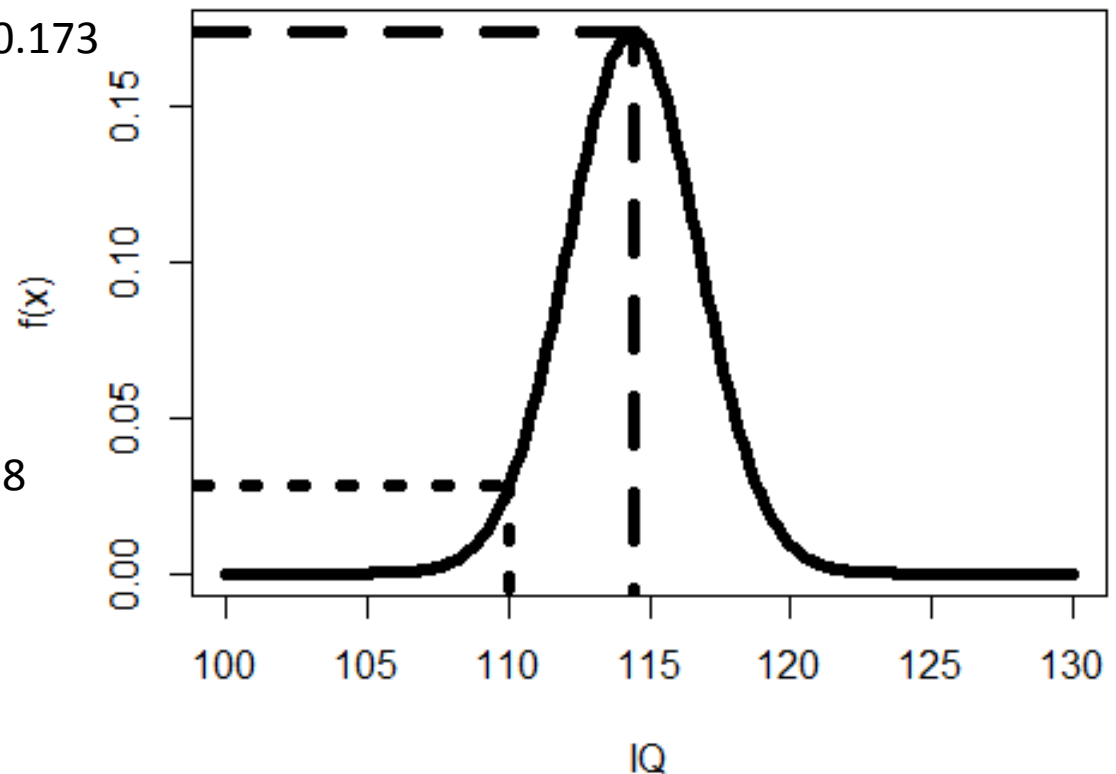


# Example Distribution Values

- Let's examine the distribution values for the IQ variable
  - We assume that we **know**  $\mu_x = 114.4$  and  $\sigma_x^2 = 5.29$  ( $\sigma_x = 2.30$ )
    - In reality we do not know what these values happen to be

For  $x = 114.4$ ,  $f(114.4) = 0.173$

For  $x = 110$ ,  $f(110) = 0.028$



# Constructing a Likelihood Function

- Maximum likelihood estimation begins by building a **likelihood function**
  - A likelihood function provides a value of a likelihood (think height of a curve) for a set of statistical parameters
- Likelihood functions start with probability density functions (PDFs)
  - Density functions are provided for each observation individually (marginal)
- The likelihood function for the entire sample is the function that gets used in the estimation process
  - The sample likelihood can be thought of as a joint distribution of all the observations, simultaneously
  - In univariate statistics, observations are considered independent, so the joint likelihood for the sample is constructed through a product
- To demonstrate, let's consider the likelihood function for one observation

# A One-Observation Likelihood Function

- Let's assume the following:
  - We have observed the first value of IQ ( $x = 112$ )
  - That IQ comes from a normal distribution
  - That the variance of  $x$  is known to be 5.29 ( $\sigma_x^2 = 5.29$ )
    - ◆ This is to simplify the likelihood function so that we only don't know one value
    - ◆ More on this later...empirical under-identification
- For this one observation, the likelihood function takes its assumed distribution and uses its PDF:

$$f(x, \mu_x, \sigma_x^2) = \frac{1}{\sqrt{2\pi\sigma_x^2}} \exp\left(-\frac{(x - \mu_x)^2}{2\sigma_x^2}\right)$$

- The PDF above now is expressed in terms of the three unknowns that go into it:  $x, \mu_x, \sigma_x^2$

# A One-Observation Likelihood Function

- Because we know two of these terms ( $x = 112$ ;  $\sigma_x^2 = 5.29$ ), we can create the likelihood function for the mean:

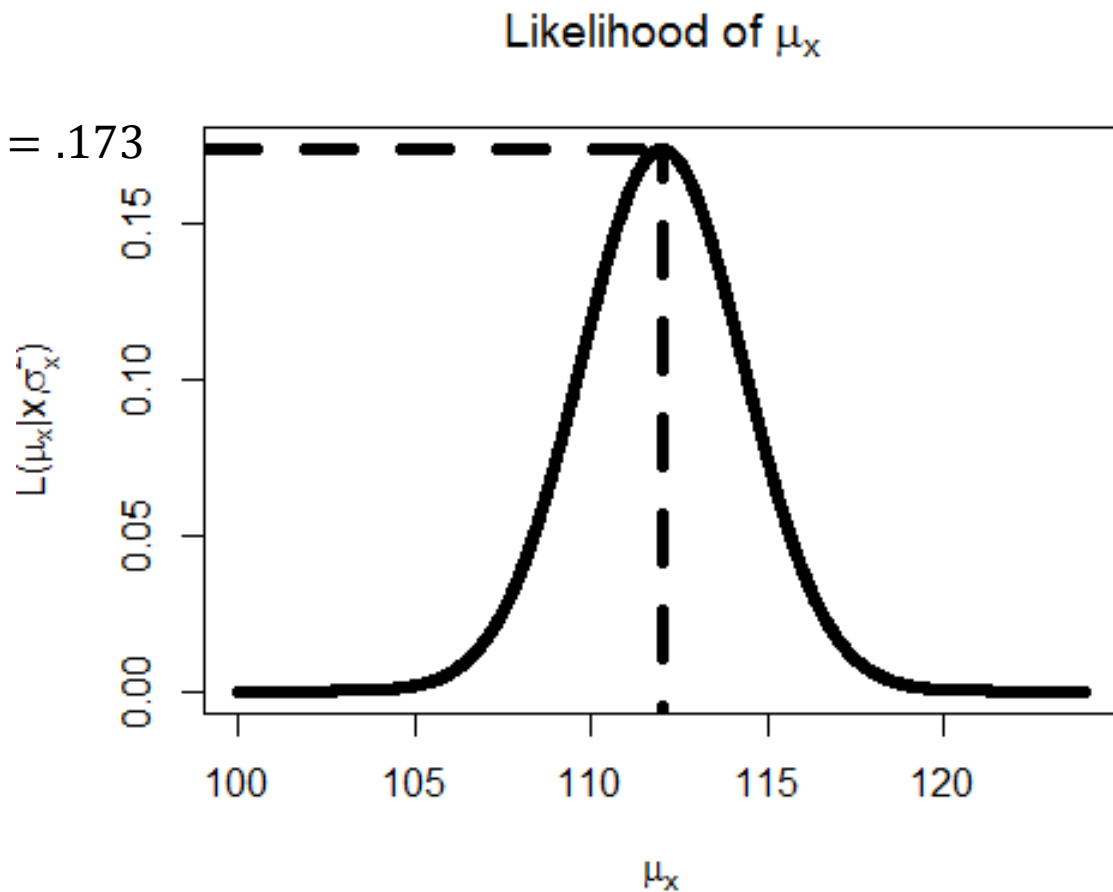
$$L(\mu_x | x = 112, \sigma_x^2 = 5.29) = \frac{1}{\sqrt{2\pi * 5.29}} \exp\left(-\frac{(112 - \mu_x)^2}{2 * 5.29}\right)$$

- For every value of  $\mu_x$  *could be*, the likelihood function now returns a number that is called **the likelihood**
  - The actual value of the likelihood is not relevant (yet)
- The value of  $\mu_x$  with the highest likelihood is called the **maximum likelihood estimate (MLE)**
  - For this one observation, what do you think the MLE would be?
  - This is asking: what is the most likely mean that produced these data?

# The MLE is...

- The value of  $\mu_x$  that maximizes  $L(\mu_x|x, \sigma_x^2)$  is  $\hat{\mu}_x = 112$ 
  - The value of the likelihood function at that point is  $L(112|x, \sigma_x^2) = .173$

For  $\hat{\mu}_x = 112$ ,  $L(112|x, \sigma_x^2) = .173$



# From One Observation...To The Sample

- The likelihood function shown previously was for one observation, but we will be working with a sample
  - Assuming the sample observations are independent and identically distributed, we can form the joint distribution of the sample
  - For normal distributions, this means the observations have the same mean and variance

Multiplication comes from independence assumption:

Here,  $L(\mu_x, \sigma_x^2 | x_p)$  is the univariate normal PDF for  $x_p$ ,  $\mu_x$ , and  $\sigma_x^2$

$$\begin{aligned} L(\mu_x, \sigma_x^2 | x_1, \dots, x_N) &= L(\mu_x, \sigma_x^2 | x_1) \times L(\mu_x, \sigma_x^2 | x_2) \times \dots \times L(\mu_x, \sigma_x^2 | x_N) \\ &= \prod_{p=1}^N f(x_p) = \prod_{p=1}^N \frac{1}{\sqrt{2\pi\sigma_x^2}} \exp\left(-\frac{(x_p - \mu_x)^2}{2\sigma_x^2}\right) = \\ &\quad (2\pi\sigma_x^2)^{-\frac{N}{2}} \exp\left(-\sum_{p=1}^N \frac{(x_p - \mu_x)^2}{2\sigma_x^2}\right) \end{aligned}$$

# The Sample Likelihood Function

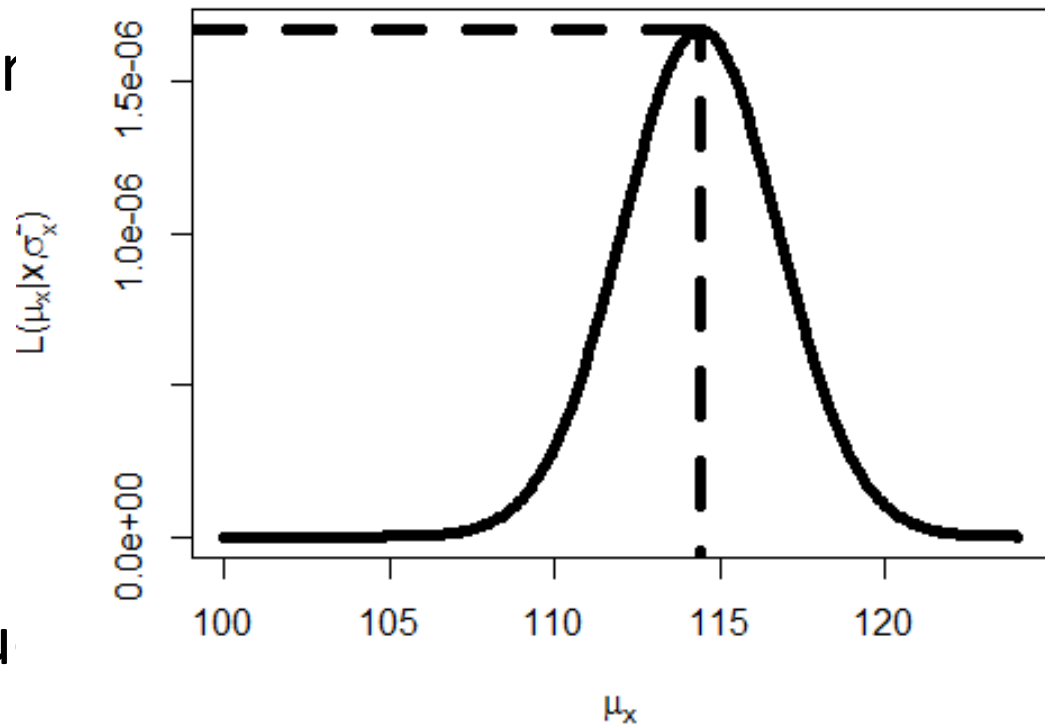
- From the previous slide:

$$L(x_1, \dots, x_N | \mu_x, \sigma_x^2) = L = (2\pi\sigma_x^2)^{-\frac{N}{2}} \exp \left( - \sum_{p=1}^N \frac{(x_p - \mu_x)^2}{2\sigma_x^2} \right)$$

- For this function, there is one mean ( $\mu_x$ ), one variance ( $\sigma_x^2$ ), and all of the data ( $x_1, \dots, x_N$ )
- If we observe the data but **do not know** the mean and/or variance, then we call this the sample likelihood function
- Rather than provide the height of the curve of any value of  $x$ , it provides the **likelihood** for any possible values of  $\mu_x$  **and**  $\sigma_x^2$ 
  - **Goal of Maximum Likelihood is to find values of  $\mu_x$  and  $\sigma_x^2$  that maximize this function**

# Likelihood Function for All Five Observations

- Imagine we know that  $\sigma_x^2 = 5.29$  but we do not know  $\mu_x$
- The likelihood function will give us the likelihood of a range of values of  $\mu_x$
- The value of  $\mu_x$  where  $L$  is the maximum is the MLE for  $\mu_x$ :
- $\hat{\mu}_x = 114.4$
- $L = 1.67e - 06$
- Note: likelihood value abbreviated as  $L$





# The Log-Likelihood Function

- The likelihood function is more commonly re-expressed as the log-likelihood:  $\log L = \ln(L)$

➤ The natural log of  $L$

$$\begin{aligned}\log L &= \log L(\mu_x, \sigma_x^2 | x_1, \dots, x_N) \\ &= \log(L(\mu_x, \sigma_x^2 | x_1) \times L(\mu_x, \sigma_x^2 | x_2) \times \dots \times L(\mu_x, \sigma_x^2 | x_N))\end{aligned}$$

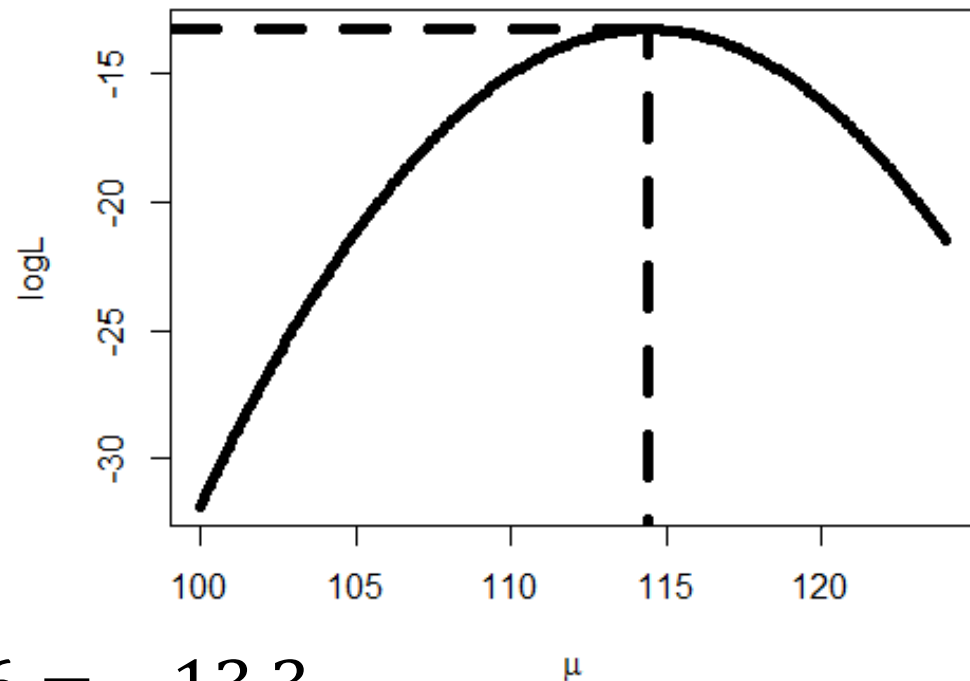
HW Hint:  $\log L$  can be found by taking the natural log of each observation's likelihood, then summing across observations

$$\begin{aligned}&= \sum_{p=1}^N \log L(\mu_x, \sigma_x^2 | x_p) \\ &= \log \left[ (2\pi\sigma_x^2)^{-\frac{N}{2}} \exp \left( - \sum_{p=1}^N \frac{(x_p - \mu_x)^2}{2\sigma_x^2} \right) \right] = \\ &\quad -\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma_x^2) - \sum_{p=1}^N \frac{(x_p - \mu_x)^2}{2\sigma_x^2}\end{aligned}$$

- The log-likelihood and the likelihood have a maximum at the same location of  $\mu_x$  and  $\sigma_x^2$

# Log-Likelihood Function In Use

- Imagine we know that  $\sigma_x^2 = 5.29$  but we do not know  $\mu_x$
- The log-likelihood function will give us the likelihood of a range of possible values of  $\mu_x$
- The value of  $\mu_x$  where  $\log L$  is the maximum is the MLE for  $\mu_x$ :
- $\hat{\mu}_x = 114.4$
- $\log L = \log 1.67e - 06 = -13.3$



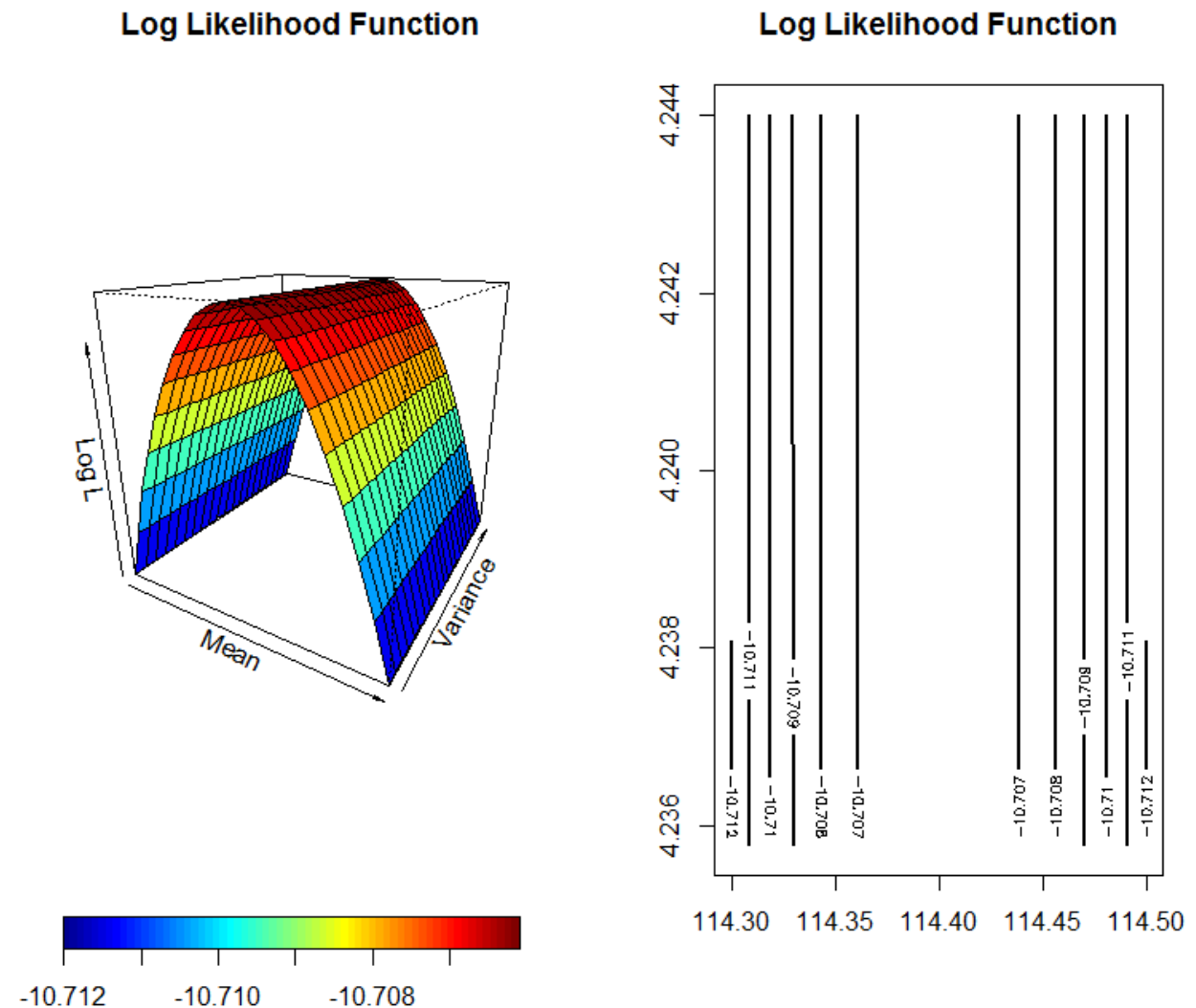
## But...What About the Variance?

- Up to this point, we have assumed the sample variance was known
  - Not likely to happen in practice
- We can jointly estimate the mean and the variance using the same log likelihood (or likelihood) function
  - The variance is now a parameter in the model
  - The likelihood function now will be with respect to two dimensions
    - ◆ Each unknown parameter is a dimension

$$\begin{aligned}\log L &= \log L(\mu_x, \sigma_x^2 | x_1, \dots, x_N) \\ &= -\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma_x^2) - \sum_{p=1}^N \frac{(x_p - \mu_x)^2}{2\sigma_x^2}\end{aligned}$$

# The Log Likelihood Function for Two Parameters

- The point where  $\log L$  is the maximum is the MLE for  $\mu_x$  and  $\sigma_x^2$
- $\log L = -10.7$
- $\hat{\mu} = 114.4$
- $\sigma_x^2 = 4.24$
- Wait... $\sigma_x^2 = 4$ .
  - It was 5.29 on slide 3
  - Why? Think

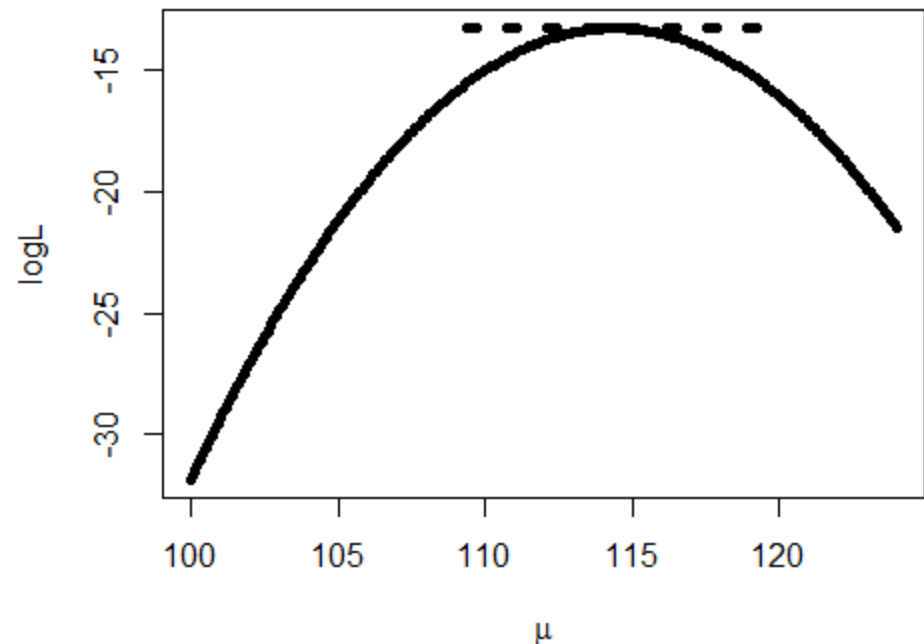


# Maximizing the Log Likelihood Function

- The process of finding the values of  $\mu_x$  and  $\sigma_x^2$  that maximize the likelihood function is complicated
  - What was shown was a grid search: trial-and-error process
- For relatively simple functions, we can use calculus to find the maximum of a function mathematically
  - Problem: not all functions can give closed-form solutions (i.e., one solvable equation) for location of the maximum
  - Solution: use efficient methods of searching for parameter (i.e., Newton-Raphson)

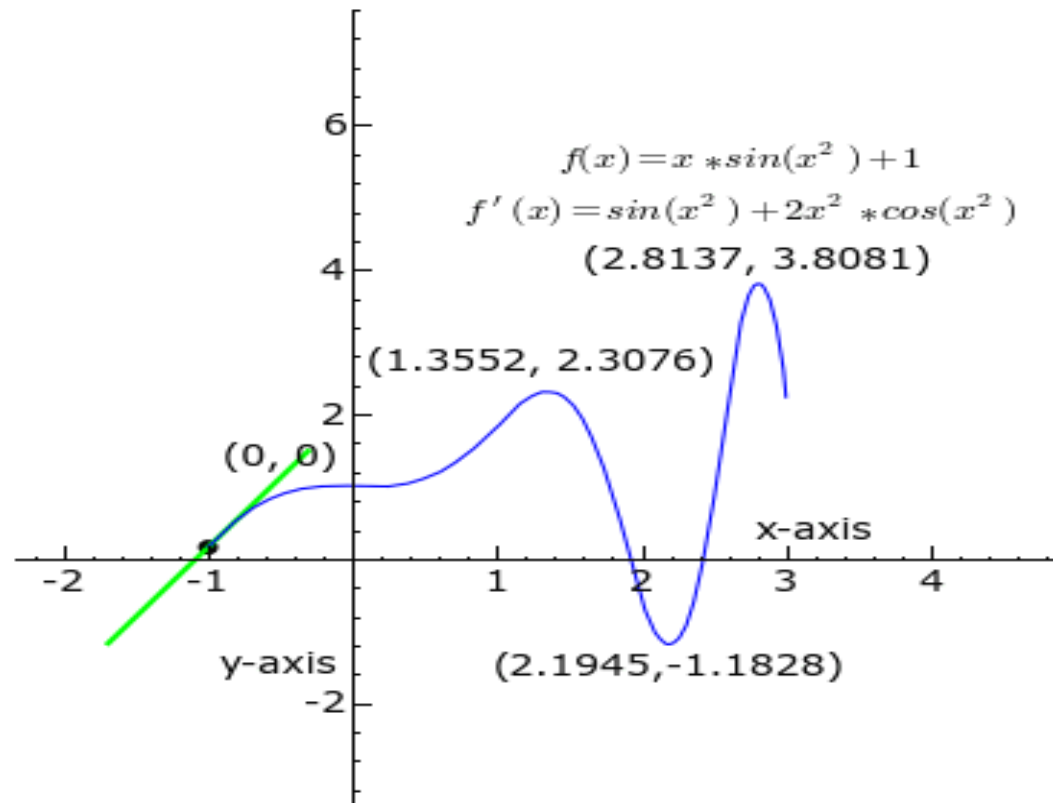
# Using Calculus: The First Derivative

- The calculus method to find the maximum of a function makes use of the first derivative
  - Slope of line that is tangent to a point on the curve
- When the first derivative is zero (slope is flat), the maximum of the function is found
  - Could also be at a minimum – but our functions will be inverted Us (convex)



# First Derivative = Tangent Line

From:  
Wikipedia



# The First Derivative for the Sample Mean

- Using calculus, we can find the first derivative for the mean from our normal distribution example (the slope of the tangent line for any value of  $\mu_x$ ):

$$\frac{\partial \log L}{\partial \mu_x} = \frac{1}{\sigma_x^2} \left( -N\mu_x + \sum_{p=1}^N x_p \right)$$

- To find where the maximum is, we set this equal to zero and solve for  $\mu_x$  (giving us an ML estimate  $\hat{\mu}_x$ ):

$$\frac{1}{\sigma_x^2} \left( -N\mu_x + \sum_{p=1}^N x_p \right) = 0 \rightarrow \hat{\mu}_x = \frac{1}{N} \sum_{p=1}^N x_p$$



# The First Derivative for the Sample Variance

- Using calculus, we can find the first derivative for the variance (slope of the tangent line for any value of  $\sigma_x^2$ ):

$$\frac{\partial \log L}{\partial \sigma_x^2} = -\frac{N}{2\sigma_x^2} + \sum_{p=1}^N \frac{(x_p - \mu_x)^2}{2\sigma_x^4}$$

- To find where the maximum is, we set this equal to zero and solve for  $\sigma_x^2$  (giving us an ML estimate  $\hat{\sigma}_x^2$ ):

$$-\frac{N}{2\sigma_x^2} + \sum_{p=1}^N \frac{(x_p - \mu_x)^2}{2\sigma_x^4} = 0 \rightarrow \hat{\sigma}_x^2 = \frac{1}{N} \sum_{p=1}^N (x_p - \mu_x)^2$$

- Where the  $\frac{1}{N}$  version of the variance/standard deviation comes from

## Standard Errors: Using the Second Derivative

- Although the estimated values of the sample mean and variance are needed, we also need the standard errors
- For MLEs, the standard errors come from the **information matrix**, which is found from the square root of -1 times the inverse matrix of second derivatives (only one value for one parameter)
  - Second derivative gives curvature of log-likelihood function
- Variance of the sample mean:

$$\frac{\partial^2 \log L}{\partial \mu_x^2} = \frac{-N}{\sigma_x^2} \rightarrow \text{Var}(\hat{\mu}_x) = \frac{\sigma_x^2}{N}$$



# **ML ESTIMATION OF GLMS: SAS PROC MIXED**

## Maximum Likelihood Estimation for GLMs in SAS: PROC MIXED

- Maximum likelihood estimation of GLMs can be performed in SAS using PROC MIXED
- PROC MIXED will grow in value to you as time goes on: most multivariate analyses assuming conditional normality can be run with PROC MIXED:
  - Multilevel models
  - Repeated measures
  - Some factor analysis models
- The **MIXED** part of PROC MIXED refers to the type of model it can estimate: **General Linear Mixed Models**
  - Mixed models *extend* the GLM to be able to model dependency between observations (either within a person or within a group, or both)

# Likelihood Functions in PROC MIXED

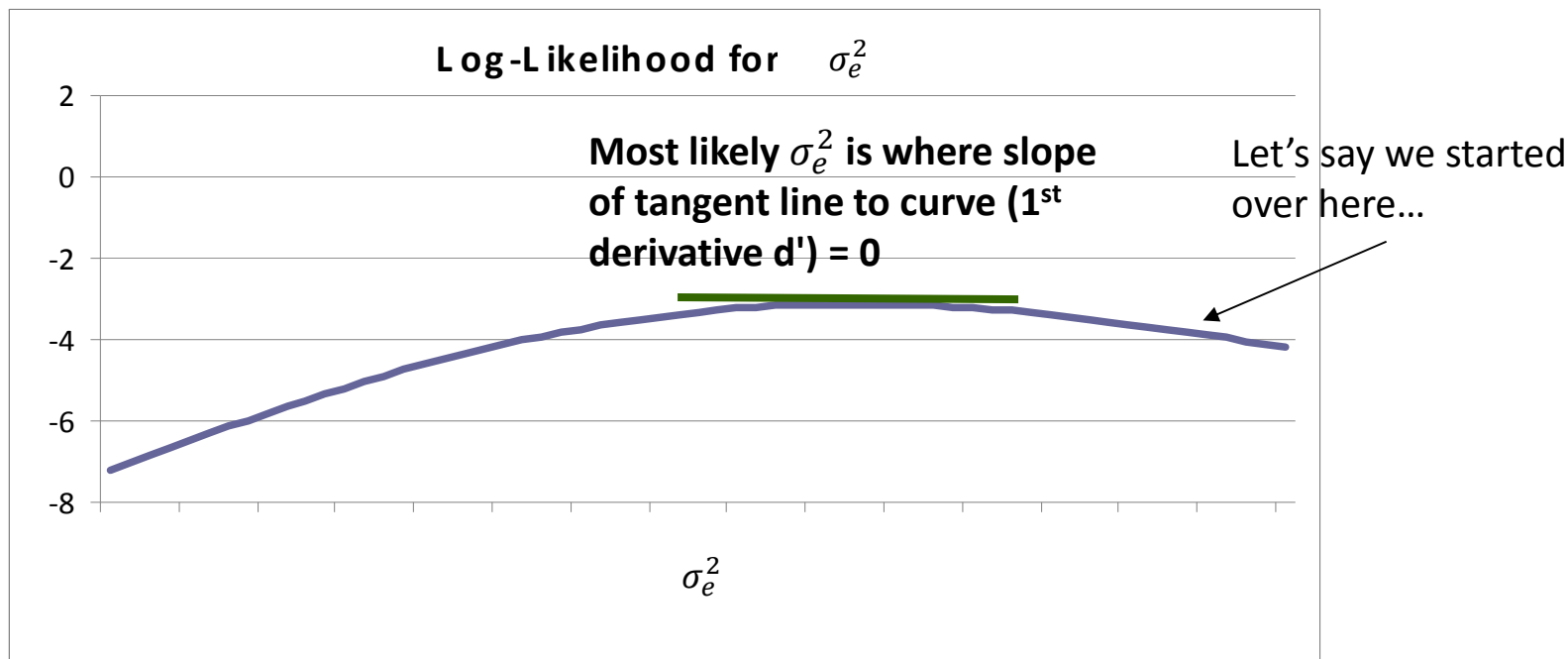
- PROC MIXED uses a common (but very general) log-likelihood function based on the GLM: the conditional distribution of Y given **X**

$$f(Y_p | X_p, Z_p) \sim N(\beta_0 + \beta_1 X_p + \beta_2 Z_p + \beta_3 X_p Z_p, \sigma_e^2)$$

- Y is normally distributed conditional on the values of the predictors
- The log likelihood for Y is then
$$\begin{aligned}\log L &= \log L(\sigma_e^2 | x_1, \dots, x_N) \\ &= -\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma_e^2) - \sum_{p=1}^N \frac{(Y_p - \hat{Y}_p)^2}{2\sigma_e^2}\end{aligned}$$
- Furthermore, there is a **closed form** (a set of equations) for the fixed effects (and thus  $\hat{Y}_p$ ) for any possible value of  $\sigma_e^2$ 
  - So...PROC MIXED seeks to find  $\sigma_e^2$  at the maximum of the log likelihood function – and after that finds everything else from equations
  - Begins with a naïve guess...then uses Newton-Raphson to find maximum

# $\sigma_e^2$ Estimation via Newton Raphson

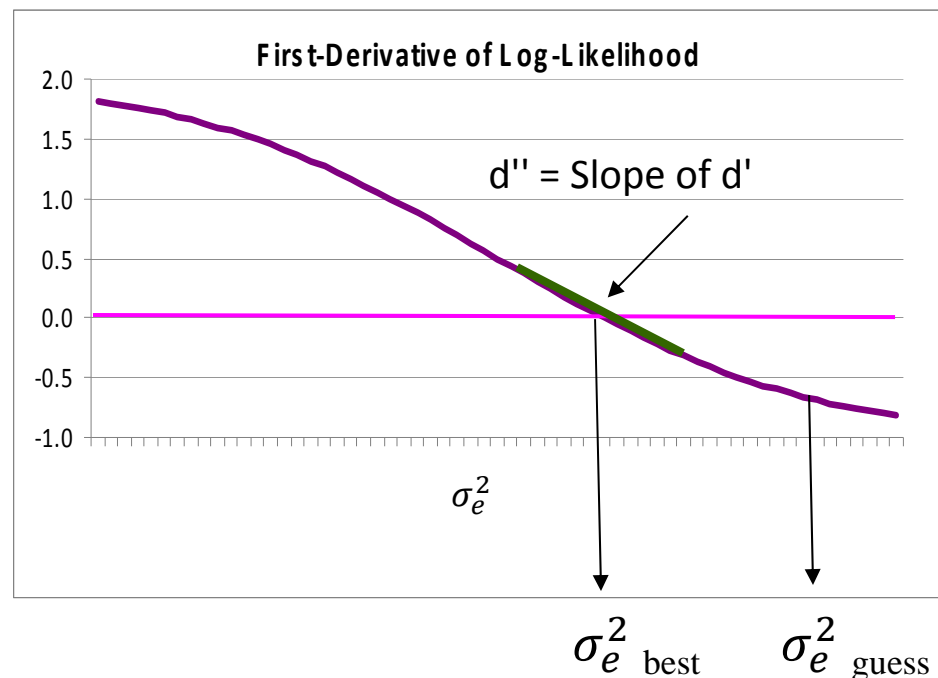
- We could calculate the likelihood over wide range of  $\sigma_e^2$  for each person and plot those log likelihood values to see where the peak is...
  - But we have lives to lead, so we can solve it mathematically instead by finding where the slope of the likelihood function (the 1<sup>st</sup> derivative,  $d'$ ) = 0 (its peak)
- Step 1: Start with a guess of  $\sigma_e^2$ , **calculate 1<sup>st</sup> derivative  $d'$**  of the log likelihood with respect to  $\sigma_e^2$  at that point
  - Are we there ( $d' = 0$ ) yet? Positive  $d'$  = too low, negative  $d'$  = too high



# $\sigma_e^2$ Estimation via Newton Raphson

- Step 2: **Calculate the 2<sup>nd</sup> derivative** (slope of slope,  $d''$ ) at that point
  - Tells us **how far off we are**, and is used to figure out how much to adjust by
  - $d''$  will always be negative as approach top, but  $d'$  can be positive or negative
- Calculate new guess of  $\sigma_e^2$  :  $\sigma_{e \text{ new}}^2 = \sigma_{e \text{ old}}^2 - (d'/d'')$ 
  - If  $(d'/d'') < 0 \rightarrow \sigma_e^2$  increases
  - If  $(d'/d'') > 0 \rightarrow \sigma_e^2$  decreases
  - If  $(d'/d'') = 0$  then you are done

- **2<sup>nd</sup> derivative  $d''$  also tells you how good of a peak you have**
  - Need to know where your best  $\sigma_e^2$  is (at  $d'=0$ ), as well as how precise it is (from  $d''$ )
  - If the function is flat,  $d''$  will be smallish
  - **Want large  $d''$  because  $1/\text{SQRT}(d'') = \sigma_e^2$ 's SE**



# Trying It Out: Using PROC MIXED with Our Example Data

- For now, we will know PROC MIXED to be largely like PROC GLM
  - Even the ESTIMATE statements work the same (another HW hint!)
- The first model will be the empty model where IQ is the DV
  - Linking PROC MIXED to our previous set of slides
  - After that, we will replicate the analysis from last Friday's class: predicting Performance from IQ
  - What we are estimating is  $\sigma_x^2 = \sigma_e^2$  (the variance of IQ, used in the likelihood function) and  $\beta_0^{IQ} = \mu_x$  (the mean IQ, found from equations)

- The PROC MIXED syntax for the empty model predicting IQ is:

```
*EMPTY MODEL PREDICTING IQ;  
PROC MIXED DATA=WORK.iqperf METHOD=ML COVTEST NOPROFILE ITDETAILS IC;  
MODEL iq = / SOLUTION;  
RUN;
```

- METHOD=ML: use Maximum Likelihood
- COVTEST: provide Wald test for  $\sigma_e^2$
- NOPROFILE: make sure  $\sigma_e^2$  is used in likelihood function
- ITDETAILS: list the iteration details
- IC: list the information criteria



# The Basics of PROC MIXED Output

- Dimensions: see Subjects and Max Obs Per Subject

```
Dimensions
Covariance Parameters      1
Columns in X               1
Columns in Z               0
Subjects                   5
Max Obs Per Subject        1
```

- Iteration History: provides starting and ending values of  $\sigma_e^2$

```
Iteration History
CovP1      Iteration      Evaluations      -2 Log Like      Criterion
1.0000          0              1          21.41220168
4.2400          1              1          21.41220168      0.00000000
```

- The best message: If you see this, you know things converged

```
Convergence criteria met.
```

- If you do not see this, you do not have the MLE (so all the good things about the MLE don't apply to your results)

## Further Unpacking Output

- The estimated  $\sigma_e^2$  is shown under the Covariance Parameter Estimates section (if you use the COVTEST option)

Covariance Parameter Estimates				
Cov Parm	Estimate	Standard Error	Z Value	Pr > Z
Residual	4.2400	2.6816	1.58	0.0569

- Note: PROC MIXED found the same estimate of  $\sigma_e^2$  as we did
  - Also: the SE of  $\sigma_e^2$  is the SD of a variance
- The Information Criteria section shows statistics that can be used for model comparisons

Information Criteria						
Neg2LogLike	Parms	AIC	AICC	HQIC	BIC	CAIC
21.4	2	25.4	31.4	23.3	24.6	26.6

- Note: Neg2LogLike is -2 times the log likelihood – our previous example estimating the mean and variance found the log likelihood to be -10.7
    - ♦ So  $-2 \times -10.7 = 21.4$

# Finally...the Fixed Effects

- The fixed effects are where the estimated regression slopes are listed – here  $\beta_0^{IQ} = \mu_x$

Solution for Fixed Effects					
Effect	Estimate	Standard Error	DF	t Value	Pr >  t
Intercept	114.40	0.9209	4	124.23	<.0001

- This also is the value we estimated in our example from before
- Not listed: traditional ANOVA table with Sums of Squares, Mean Squares, and F statistics
  - The Mean Square Error is no longer the estimate of  $\sigma_e^2$ : this comes directly from the model estimation algorithm itself
  - The traditional  $R^2$  change test also changes under ML estimation (see next section for what it becomes)



# **USEFUL PROPERTIES OF MAXIMUM LIKELIHOOD ESTIMATES**

# Useful Properties of MLEs

- Next, we demonstrate three useful properties of MLEs (not just for GLMs)
  - Likelihood ratio (aka Deviance) tests
  - Wald tests
  - Information criteria
- To do so, we will consider our example where we wish to predict job performance from IQ (but will now center IQ at its mean of 114.4)
- We will estimate two models, both used to demonstrate how ML estimation differs slightly from LS estimation for GLMs
  - Empty model predicting just performance:  $Y_p = \beta_0 + e_p$
  - Model where mean centered IQ predicts performance:
$$Y_p = \beta_0 + \beta_1(IQ - 114.4) + e_p$$

# PROC MIXED Syntax

- Syntax for the empty model predicting performance:

```
*EMPTY MODEL PREDICTING PERFORMANCE;  
PROC MIXED DATA=WORK.iqperf METHOD=ML COVTEST NOPROFILE ITDETAILS IC;  
MODEL perf = / SOLUTION;  
RUN;
```

- Syntax for the conditional model where mean centered IQ predicts performance:

```
*MEAN CENTERED IQ PREDICTING PERFORMANCE;  
PROC MIXED DATA=WORK.iqperf METHOD=ML COVTEST NOPROFILE ITDETAILS IC;  
MODEL perf = iq114 / SOLUTION;  
RUN;
```

- Questions in comparing between the two models:
  - How do we test the hypothesis that IQ predicts performance?
    - ♦ Likelihood ratio tests (can be multiple parameter/degree-of-freedom)
    - ♦ Wald tests (usually for one parameter)
  - If IQ does significantly predict performance, what percentage of variance in performance does it account for?
    - ♦ Relative change in  $\sigma_e^2$  from empty model to conditional model

# Likelihood Ratio (Deviance) Tests

- The likelihood value from MLEs can help to statistically test competing models assuming the models are nested
- Likelihood ratio tests take the ratio of the likelihood for two models and use it as a test statistic
- Using log-likelihoods, the ratio becomes a difference
  - The test is sometimes called a **deviance test**
$$D = \Delta - 2\log L = -2 \times (\log L_{H0} - \log L_{HA})$$
  - $D$  is tested against a Chi-Square distribution with degrees of freedom equal to the difference in number of parameters

# Deviance Test Example

- Imagine we wanted to test the null hypothesis that IQ did not predict performance:

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

- The difference between the empty model and the conditional model is one parameter
  - Null model: one intercept  $\beta_0$  and one residual variance  $\sigma_e^2$  estimated = 2 parameters
  - Alternative model: one intercept  $\beta_0$ , one slope  $\beta_1$ , and one residual variance  $\sigma_e^2$  estimated = 3 parameters
- Difference in parameters:  $3 - 2 = 1$  (will be degrees of freedom)



# LRT/Deviance Test Procedure

- Step #1: estimate null model (get  $-2 \times \log$  likelihood)

Information Criteria						
Neg2LogLike	Parms	AIC	AICC	HQIC	BIC	CAIC
21.3	2	25.3	31.3	23.2	24.5	26.5

- Step #2: estimate alternative model (get  $-2 \times \log$  likelihood)

Information Criteria						
Neg2LogLike	Parms	AIC	AICC	HQIC	BIC	CAIC
6.9	3	12.9	36.9	9.8	11.8	14.8

- Step #3: compute test statistic

$$D = -2 \times (\log L_{H0} - \log L_{HA}) = (21.3 - 6.9) = 14.4$$

- Step #4: calculate p-value from Chi-Square Distribution with 1 DF
  - I used =chidist(14.4,1) from Excel
  - p-value = 0.000148
- Inference: the regression slope for IQ was significantly different from zero -- we prefer our alternative model to the null model
- Interpretation: IQ significantly predicts performance

# Wald Tests (Usually 1 DF Tests in Software)

- For each parameter  $\theta$ , we can form the Wald statistic:

$$\omega = \frac{\hat{\theta}_{MLE} - \theta_0}{SE(\hat{\theta}_{MLE})}$$

- (typically  $\theta_0 = 0$ )
- As N gets large (goes to infinity), the Wald statistic converges to a standard normal distribution  $\omega \sim N(0,1)$ 
  - Gives us a hypothesis test of  $H_0: \theta = 0$
- If we divide each parameter by its standard error, we can compute the two-tailed p-value from the standard normal distribution (Z)
  - Exception: bounded parameters can have issues (variances)
- We can further add that variances are estimated, switching this standard normal distribution to a t distribution (SAS does this for us)
  - Note: some don't like calling this a "true" Wald test

# Wald Test Example

- We could have used a Wald test to compare between the empty and conditional model, or:

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

- SAS provides this for us in the Solution for Fixed Effects:

Solution for Fixed Effects					
Effect	Estimate	Standard Error	DF	t Value	Pr >  t
Intercept	12.8000	0.2163	3	59.17	<.0001
iq114	0.9623	0.1051	3	9.16	0.0028

➤ Note: these estimates are identical to the LS estimates from last week

- Here, the slope estimate has a t-test statistic value of 9.16 ( $p = .0028$ ), meaning we would reject our null hypothesis
- Typically, Wald tests are used for one additional parameter
  - Here, one slope

## Model Comparison with $R^2$

- To compute an  $R^2$ , we use the ML estimates of  $\sigma_e^2$ :
  - Empty model:  $\sigma_e^2 = 4.160$  (2.631)
  - Conditional model:  $\sigma_e^2 = 0.234$  (0.148)
- The  $R^2$  for variance in performance accounted for by IQ is:

$$R^2 = \frac{4.160 - 0.234}{4.160} = .944$$

- Hall of fame worthy

# Information Criteria

- Information criteria are statistics that help determine the relative fit of a model for non-nested models
  - Comparison is fit-versus-parsimony
- PROC MIXED reports a set of criteria (from conditional model)

Information Criteria						
Neg2LogLike	Parms	AIC	AICC	HQIC	BIC	CAIC
6.9	3	12.9	36.9	9.8	11.8	14.8

- Each uses  $-2 \times \log\text{-likelihood}$  as a base
    - ◆ Choice of statistic is **very** arbitrary and depends on field
- Best model is one with *smallest* value
- Note: don't use information criteria for nested models
  - LRT/Deviance tests are more powerful

# How ML and LS Estimation of GLMs Differ

- You may have recognized that the ML and the LS estimates of the fixed effects were identical
  - And for these models, they will be
- Where they differ is in their estimate of the residual variance  $\sigma_e^2$ :
  - From Least Squares (MSE):  $\sigma_e^2 = 0.390$  (no SE)
  - From ML (model parameter):  $\sigma_e^2 = 0.234$  (0.148)
- The ML version uses a **biased estimate** of  $\sigma_e^2$  (it is too small)
- Because  $\sigma_e^2$  plays a role in all SEs, the Wald tests differed from LS and ML
- Troubled by this? Don't be: a fix will come in a few weeks...



# **WRAPPING UP**

# Wrapping Up

- Today was our first pass at maximum likelihood estimation
- The topics discussed today apply to all statistical models, not just GLMs
- Maximum likelihood estimation of GLMs helps when the basic assumptions are obviously violated
  - Independence of observations
  - Homogeneous  $\sigma_e^2$
  - Conditional normality of Y (normality of error terms)