







f) Residuals and dispersion

The i^{th} **residual** = $x_i - \overline{x}$

The average residual length = $\sum |x_i - \overline{x}|/n$ is an intuitive measure of **dispersion:** the extent to which observations vary from each other.

It is rarely used because it is difficult to work with mathematically.





i) **Expected** value The **expected value** of a statistic is its average value from a very large number of experiments. The expected value of both x_i and \overline{x} is μ . The expected value of s^2 is σ^2 We write $E(s^2) = \sigma^2$ $E(\bar{x}) = \mu$ Unbiased estimate of a parameter j) A statistic is **unbiased estimate** of a parameter if its expected value equals the parameter. \overline{x} is an unbiased estimate of μ since $E(\overline{x}) = \mu$ The denominator of $s^2 = \sum (x_i - \overline{x})^2 / (n-1)$ is *n*-1 rather than *n* in order to make s^2 an unbiased estimate of σ^2 .

















3.	Simple Line	ar Regression
a)	The Model	
	We assume the	nat
	$y_i = \alpha$	$+\beta x_i + \varepsilon_i$
	Where x_i is a	variable observed on the $i^{ m th}$ patient
	ϵ_i	is assumed to be <mark>normally</mark> and <mark>independently</mark> distributed with mean 0 and standard deviation σ
	${\mathcal Y}_i$	is the response from the $i^{\rm th}$ patient
	α and	β are model parameters.
	The expected	value of y_i is $E(y_i) = \alpha + \beta x_i$.













5. The Stata Statistical Software Package

Stata is an excellent tool for the analysis of medical data. It is far easier to use than other software of similar sophistication. However, we will be using Stata for some complex analyses and you may be puzzled by some of its responses. If so please ask. I would very much like to minimize the time you spend struggling with Stata and maximize the time you spend learning statistics. I will be available to answer questions at most times on days, evenings and weekends.

If you have not used Stata since Biometry I you are probably very rusty. Here are a few reminders and aids that may help.

Punctuation a) Proper punctuation is mandatory. If Stata gives a confusing error message, the first thing to check is your punctuation. Stata commands are modified by **qualifiers** and **options**. Qualifiers precede options and there must be a comma before the first option. For example table age if treat==1 ,by(sex) Might produce a table showing the number of men and women of different ages receiving treatment 1. if treat==1 is a qualifier and by(sex) is an option. Without the comma, Stata will not recognize by (sex) as a valid option to the table command. Some command prefixes must be followed by a colon.

b) Capitalization

Stata variables and commands are case sensitive. That is, Stata considers age and Age to be two distinct variables. In general, I recommend that you always use lower case variables. Sometimes Stata will create variables for you that contain upper case letters. You must use the correct capitalization when referring to these variables.

c) Command summary

At the end of the text book is a summary of most of the commands that are needed in this course. These may be helpful in answering your class exercises.

d) GUI interface

You can avoid learning Stata syntax by using their pull down menus. These menus generate rather complex syntax but feel free to use them if it makes the exercises easier.

This interface is extensively documented in my text. See Section 1.3.8 on page 15.

d) Data files and log files

You may download the Stata data files, log files, do files and these lecture notes that you will need from this course from the web at:

http://biostat.mc.vanderbilt.edu/BiostatIILectureNotes

then click on the desired links for data files, Stata log files, do files or lecture notes. Pages for the class schedule, student names and exercises are password protected. The username and password for these pages is the same as for the other MPH courses except that the second letter of the username must be lower case while the other letters must be upper case. (This is due to a camel-case requirement for usernames on the Biostatistics wiki that I can't get around.)

The class exercises are very similar to the examples discussed in class. You can save yourself time by cutting and pasting commands from these log files into your Stata Command window, or by modifying these do files.



. * RosnerTable	:11.1.log		{1}
* Examine the * See Green 8	e Stata data set from Tabl & Touchstone 1963	e 11.1 of Rosner, p. 554	{2}
. <mark>use</mark> "C:\MyDoc	s\MPH\LectureNotes\rostab	o11.dta", clear	{3}
* Data > Desc	ribe data > Describe data	in memory	
. describe			
ontains data f	rom \\PMPC158\mph\analyse	s\linear_reg\stata\rostab11.d	ta
obs:	31	10 New 1000 10-01	
size:	3 496 (99.9% of memory fr	ee)	
1 id	float %0 Og		
1. 10	float %9.0g	Estriol (mg/24 hr)	
estriol		(5, , ,	

1	at: <u>http://biostat.mc.vanderbilt.edu/BiostatisticsTwoClassPage</u>
	and clicking on Example Logs and Data from Lecture Notes.
	Most of the class exercises may be completed by performing Stata sessions that are similar to those discussed in class.
{2}	• I have adopted the following color coding conventions for Stata log files throughout these notes.
•	Red is used for Stata comment statements. Also, red numbers in braces in the right margin refer to comments at the end of the program and are not part of the programming code. Red text on Stata output is my annotation rather than text printed by Stata.
•	Stata command words, qualifiers and options are written in blue
•	Variables and data set names are written in black , as are algebraic or logical formulas

describe - Describe data in memory
Variables: (leave empty for all variables)
Examples: yr* all variables starting with 'lyr'' syz-abc all variables between syz and abc Options: Display only variable names Display only general information Display additional details Do not abbreviate variable names Display additional details Display variable number along with name
C DK Cancel Submit

29.	25	39	
30. 31.	25 27	32 34	
	viiose esti.	ior values ar	e at least 20.

list - List values of variables		
Wain Epytion (Diplone) Summay (Advanced) Variables: [exticle weight] Column widths © © Default © © Compress width of columns in both table and display formats © Loss display format of each variable © Override minimum abbreviation of variable names ⑧ Haracters Truncate string variables 10 Haracters	Iist - List values of variables Main by///n Options Summary Advanced Repeat command by groups Variables that define groups Pesticit observations If: (expression) estinic) >= 25 Use a range of observations From: 11/1 to: 31/1	
Do not list observation numbers Display all k	O D Ba OK Cancel	Submit

Create	Overall By	 {5} Draw a scatter plot of birth weight against estriol levels.
Enable Move Up Move Down	Plot 1 Plot 1 Plot at/n Choose a plot category and type Baic plots Range plots Fit plots Advanced plots Plot type: (scatterplot) Y variable: X variable: biveight Marker properties Marker v	Batic plots: (select type) Bootter Connected Atea Bar Spike Spike Soft on x variable weights





regress - Linear regression Model by/t/n Weights SE/Robust Reporting
Dependent vanable: Independent vanables: bweight estrici
Suppress constant tem Has user-supplied constant Total SS with constant (advanced)



{4}	The Residual of It can be shown = SSM + SSE	or Error (Sum of Squares $\sum_{i=1}^{n} \sum_{j=1}^{n} (y_{i}^{j} - \overline{y})^{2} = \sum_{j=1}^{n} \sum_{j=1}^{n} (y_{j}^{j} - \overline{y})^{2} = \sum_{j=1$	(ESS) = 423.4 = $\sum (y_i - \hat{y}_i)^2$ $\sum (\hat{y}_i - \overline{y})^2 + \sum (y_i - \hat{y}_i)^2$
{5}	R-squared is the equals MSS/TS variation in b squared =1, s ² =	he square S and her weight t 0 and the	of the correlation nee measures the p hat is explained by data points fall o	proportion of the total y the model. When R- n the straight line
Source Model Residual Total	SS 250,574476 <mark>423,425524</mark> n-2 674.00	df 1 2 = 29 s ² = 30	MS 250.574476 = 14.6008801 {4} 22.4666667	Number of obs = 31 F(1, 29) = 17.16 Prob > F = 0.0003 R-squared = 0.3718 Adj R-squared = 0.3501 Root MSE = 3.8211
bweight	Coef.	Std. Er	r.tP> t	[95% Conf. Interval]
estriol _cons	b= .6081905 se(b a= 21.52343	o)= .14681 2.62041	17 4.14 P=0.000 7 8.21 0.000	.3079268 .9084541 16.16407 26.88278









predict - Prediction after estimat	ion 📃 🗶
Main if/in	
New variable name:	New variable type:
	float
Produce:	
Linear prediction (xb) C S	tandard error of the prediction
C Standardized residuals	tandard error of the recidual
C Studentized residuals C C	OVRATIO
C Cook's distance C D	FITS
C Leverage	/elsch distance
C Pr(yl (y <)	
C E(yl < y <)	
C E(y*), y* = max(, min(y,	
C DFBETA for variable:	
,	
00	OK Cancel Submit



|--|

lots If/n Yaxis ; Plot definitions: Plot 1 Plot 2 line y_hat estriol)	 X axis Titles Legend Overall By Flots I/In Y ax X axis Titles Flots I/In Y ax X axis Titles Title: Major tick/label properties Axis ine properties Reference lines Hide axis Place axis on opposite side of graph 	Is Image: Second Se
		Accept Cancel Submit

lots if/in Yaxis	X axis Titles Legend Overall By
Not definitions:	🗉 twoway - Twoway graphs
Plot 1 Plot 2	Plots if /in Y ax X axis titles Legend Overall By Title: Properties Major tick/label properties Axis fine properties Axis cale properties Axis cale properties
line y_hat estriol)	Hide axis Place axis on opposite side of graph Rule Labels Ticks Grid Axis rule Suggest # between major ticks Suggest # between major ticks 7 Minimum value 27 Minimum value 1 Deta Custom None The axis rule determines the number of ticks and their relative positions.

Etwoway - Twoway graphs Plots ii/n Y axis X axis Titles Legend Plot definitions: Create Flot 1 Create Edd Disable Enable Move Up Move Down Ine y_hat estriol)	rall By Title: Bith Weight (g/100) Major tick/label properties Axis scale properties Reference lines Hide axis	Propetties
	Place axis on opposite side of graph	Submit



b <u>+</u>	$t_{n-2,.025}$ se(b)	$= 0.608 \pm 0.608 \pm 0.608 \pm 0.608 \pm 0.608 \pm 0.608 \pm 0.308, 0$	$t_{29,.025} \times 0.1$ 2.045 × 0.1 0.300 0.908)	147 147		
Source	SS	df	MS		Number of obs	= 31
Model esidual	250.574476 423.425524	1 n-2 = 29 s ²	250.5744 ² = 14.6008	176 801	Prob > F R-squared =	= 0.0003 = 0.3718
Total	674.00	30	22.4666	667	Root MSE = s	= 3.8211
bweight	Coef.	Std. I	Err. t	P> t	[95% Conf. I	nterval]
estriol _cons	b= <mark>.6081905</mark> s a= 21.52343	se(b)= <mark>.1468</mark> 2.6204	<mark>3117</mark> 4.14 417 8.21	P=0.000 0.000	.3079268 16.16407	<mark>.9084541</mark> 26.88278



Let	$\hat{y}(x) = a + bx$			
The	e variance of $\hat{y}(x)$ var $(\hat{y}(x)) = [s^2]$	x) is $(n] + (x - \overline{x})$	$(b)^2 \operatorname{var}(b)$	{1.8}
Th	e 95% confidence $\hat{y} \pm t_{n-2,0.025} \sqrt{\text{var}}$	e interval f $\frac{1}{f(\hat{y}(x))}$	for $\hat{y}(x)$ is	{1.9}
		df	MS	Number of obs = 3
Source	55	u.		
Model Residual	55 250.574476 423.425524 n-:	1 2 = 29 s ²	250.574476 = 14.6008801	F(1, 29) = 17.10 Prob > F = 0.0003 R-squared = 0.3718 Adj P cquared = 0.3503
Source Model Residual Total	55 250.574476 423.425524 n-1 674.00	1 2 = 29 s ² 30	250.574476 = 14.6008801 22.4666667	F(1, 29) = 17.11 Prob > F = 0.000 R-squared = 0.3718 Adj R-squared = 0.350 Root MSE = s = 3.821
Model Residual Total bweight	55 250.574476 423.425524 n-1 674.00	1 2 = 29 s ² 30 Std. Err	250.574476 = 14.6008801 22.46666667	F(1, 29) = 17.11 Prob > F = 0.000 R-squared = 0.3718 Adj R-squared = 0.350 Root MSE = s = 3.821 [95% Conf. Interval]























2. Distribution of the Sum of Independent Variabl	es.
Suppose that x has mean μ_x and variance $\sigma_x^{\ 2}$	
y has mean μ_y and variance σ_y^{2} and	
<i>x</i> and <i>y</i> are independent.	
Then $x + y$ has mean $\mu_x + \mu_y$ and variance $\sigma_x^2 + \sigma_y^2$	
8. The 95% CI for the Forecasted Response of a No Patient	ew
3. The 95% CI for the Forecasted Response of a No Patient This response is $y = a + \beta x + \varepsilon_i \cong \hat{y}(x) + \varepsilon_i$	9 w
3. The 95% CI for the Forecasted Response of a Normalized Patient This response is $y = a + \beta x + \varepsilon_i \cong \hat{y}(x) + \varepsilon_i$ The variance of $y = var(\hat{y}(x)) + s^2$	ew
3. The 95% CI for the Forecasted Response of a Normalized Response of a Normalized Response of a Normalized Response is $y = a + \beta x + \varepsilon_i \cong \hat{y}(x) + \varepsilon_i$ This response is $y = a + \beta x + \varepsilon_i \cong \hat{y}(x) + \varepsilon_i$ The variance of $y = var(\hat{y}(x)) + s^2$ Therefore, a 95% confidence interval for y is	₹1,10}

 $\operatorname{var}(\hat{y}(x)) = [s^2/n] + (x - \overline{x})^2 \operatorname{var}(b)$ If x = 22 then $\operatorname{var}(\hat{y}(22)) = \frac{14.6009}{31} + (22 - 17.226)^2 \times \frac{0.1468}{0.1468^2} = \frac{0.9621}{0.9621}$ $\hat{y} = 21.523 + 0.6082 \times 22 = 34.903$ 95% C.I. for $\hat{y}(x) = \frac{34.903 \pm t_{29,.025} \times \sqrt{0.9621}}{1000}$ = 34.903 <u>+</u> 2.045 x 0.981 = (32.9, 36.9)Residual | 423.425524 n-2 = 29 s² = 14.6008801
 Total
 674.00
 30
 22.46666667
 Root MSE
 =
 3.8211
 -----bweight | Coef. Std. Err. t P>|t| [95% Conf. Interval] +----+ estriol |b= .6081905 se(b)= <mark>.1468117</mark> 4.14 P=0.000 .3079268 .9084541 _cons |a= 21.52343 2.620417 8.21 0.000 16.16407 26.88278

```
95% C.I. for y at x = \hat{y} \pm t_{n-2,0,025} \sqrt{\operatorname{var}(\hat{y}(x)) + s^2}
    If x = 22 then
    var(\hat{y}(22)) = 0.9621
    \hat{y} = 34.903
    95% C.I. for y at x = 34.903 \pm 2.045 \times \sqrt{0.9621 + 14.6009}
                    = (26.8, 43.0)
            SS df MS
                                              Number of obs =
 Source |
                                                                31
                                             F( 1, 29) = 17.16
Prob > F = 0.0003
R-squared = 0.3718
 ----+
  Model | 250.574476 1 250.574476
Residual | 423.425524 n-2 = 29 s<sup>2</sup> = 14.6008801
······
                                              Adj R-squared = 0.3501

        Total
        674.00
        30
        22.4666667
        Root MSE
        =
        3.8211

.....
bweight | Coef. Std. Err. t P>|t| [95% Conf. Interval]
       *
+-----
estriol |b= .6081905 se(b)= .1468117 4.14 P=0.000 .3079268 .9084541
_cons |a= 21.52343 2.620417 8.21 0.000 16.16407 26.88278
            _____
```



Plot definitions: Plot 1 Plot 2 Create East Disable Enable Move Up Move Down	Flot #/n Choose a plot category and type Range plots: [select type] Range plots Range plots (select type) Range plots Range plots (select type) Fit plots Range plots (select type) Choose a plot category and type Range splots (select type) Range plots Range splots (select type) Range plots Range splots (select type) Range splots Range splots (select type) Plot type: (range plot with lines) Range splots (select type)
(scatter bweight estriol)	Y1 variable: X variable: o_uf Y2 variable: o_lf Line properties
Line properties	Accept Cancel Sub

Plot 2 Plot 3	Edit	The second second		
	Enable	Plots if/in Yaxis Xa	xis Titles Legend Overall By	X
	Move Down	C Default		
(rline ci_uf ci_lf estriol)		Show legend Hide legend		
0 0 6	ОК	Ca Override default keys -		
		Specify order of keys and o	optionally change labels:	?
		Organization / Appearan	ce Placement	
		2 0 b	OK Cancel	Submit 1



15. Lowess Regression

Linear regression is a useful tool for describing a relationship that is linear, or approximately linear. It has the disadvantage that the linear relationship is assumed *a priori*. It is often useful to fit a line through a scatterplot that does not make any model assumptions. One such technique is **lowess regression**, which stands for locally weighted scatterplot smoothing. The idea is that for each observation (x_i, y_i) is fitted to a **separate linear regression** line based on adjacent observations. These points are **weighted** so that the farther away the *x* value is from x_i , the less effect it has on determining the estimate of \hat{y}_i . The proportion of the total data set considered for each \hat{y}_i is called the **bandwidth**. In Stata the default bandwidth is 0.8, which works well for small data sets. For **larger data sets** a bandwidth of **0.3 or 0.4** usually works better.

On large data sets lowess is computationally intensive.



Plots #/m Yaxis Xax Plot definitions: Plot 1	y graphs	×	1
(scatter bweight estriol)	Plot it/in Choose a plot category and type Basic plots Range plots Fit plots Fit plots Chimediate plots Advanced plots Plot type: (linear prediction plot) Y variable: V variable	(relect type) rediction to prediction to pre	
	00	Accept Cancel Submit	

Plots <i>il/in</i> Y exis X a Plot definitions: Plot 1 Plot 2	Ay graphs _ X is Titles Legend Overall By Create Edit Disable Emble
(lît bweight estriol)	Plot #/n Choose a plot category and type
	Cancel Submit





To avoid this problem, we usually calculate the studentized residual $t_i = e_i / s_{(i)} \sqrt{1 - h_i}$ {1.12} where $s_{(i)}$ denotes the root MSE estimate of σ with the i^{th} case deleted and $h_i = \operatorname{var}(\hat{y}(x)) / s^2_{(i)}$ is the variance of $\hat{y}(x)$ measured in units of s_i^2 . h_i is called the leverage of the i^{th} observation. t_i is sometimes referred to as the jackknife residual Plotting these studentized residuals against x_i assesses the homoscedasticity assumption.







18. Variance Stabilizing Transformations

a) Square root transform

Useful when the residual variance is proportional to the expected value.

b) <mark>Log transform</mark>

Useful when the residual standard deviation is proportional to the expected value.

N.B. Transformations that stabilize variance may cause non-linearity. In this case it may be necessary to use a non-linear regression technique.

19. Normalizing the Data Distribution

For skewed data we can often improve the quality of our model fit by transforming the data to give the residuals a more normal distribution.

For example, log transforms can normalize data that is right skewed.



























The pooled estimate of σ^2 is $s^2 = \left[\sum_i (y_{i1} - \hat{y}_1(x_{i1}))^2 + \sum_2 (y_{i2} - \hat{y}_2(x_{i2}))^2\right] / (n_1 + n_2 - 4)$ [1.20} Since $s_i^2 = \sum_i (y_{i1} - \hat{y}_1(x_{i1}))^2 / (n_1 - 2)$ and $s_2^2 = \sum_2 (y_{i2} - \hat{y}_2(x_{i2}))^2 / (n_2 - 2)$ var $(b_1 - b_2) = s^2 \left\{ \frac{1}{\sum_i (x_{i1} - \bar{x}_1)^2} + \frac{1}{\sum_2 (x_{i2} - \bar{x}_2)^2} \right\}$ [1.21} But var $(b_1) = s_i^2 / \sum_i (x_{i1} - \bar{x}_1)^2$ and hence $\sum_i (x_{i1} - \bar{x}_1)^2 = s_i^2 / \text{var}(b_1)$ Therefore var $(b_1 - b_2) = s^2(\text{var}(b_1)/s_1^2 + \text{var}(b_2)/s_2^2)$ $t = (b_1 - b_2) / \sqrt{\text{var}(b_1 - b_2)}$ [1.22] has a t distribution with $n_1 + n_2 - 4$ degrees of freedom. A 95% CI for $\beta_1 \cdot \beta_2$ is $(b_1 - b_2) \pm t_{n_1 + n_2 - 4, 0.025} \sqrt{\text{var}(b_1 - b_2)}$

To compute the preceding statistic we run two separate linear regressions on group 1 and 2.

For group 1, s_1 is the root MSE and $\sqrt{\operatorname{var}(b_1)}$ is the standard error of the slope estimate.

 s_2 and var(b_2) are similarly defined for group 2.

Substituting these values into the formulas from the previous slide gives the required test.



	Main it/in Options Languages
	Variables: (leave empty for all)
+	sex 🔹
	zxarupes. yr all variables starting with "yr" xvz-abc all variables between xvz and abc
	D C Cancel Submit

. regress sbp bmi if sex == 1 **{2}** Source | SS df MS Number of obs = 2047 F(1, 2045) = 121.09 Prob > F = 0.0000 R-squared = 0.0559 ------+ 44504.0296 1 44504.0296 Model | R-squared Residual | 751572.011 2045 367.516876 Adj R-squared = 0.0554 ----+ Total | 796076.041 2046 389.088974 Root MSE = 19.171 sbp | Coef. Std. Err. t P>|t| [95% Conf. Interval] bmi | <mark>1.375953</mark> .1250382 11.004 0.000 1.130738 1.621168 _cons | 96.43061 3.272571 29.466 0.000 90.01269 102.8485 _____ . predict yhatmen, xb (9 missing values generated) Comment **{2**} This command regresses sbp against bmi in men only. Note that Stata distinguishes between a = 1, which assigns the value 1 to a, and a==1 which is a logical expression that is true if the aequals 1 and is false otherwise.

Model by///m Weights SE/Robust Repo Dependent variable: Independent varia	nting] ables:
Treatment of constant Suppress constant term Has user-supplied constant Total SS with constant (advanced)	regress - Linear regression Model by/Win Weights SE/Robust Repeat command by groups Variables that define groups: Pesticit observations If: (expression) sex == 1 Obs. in range: 1 to: 4689 ±
0 0 b	Cincel Submit





Create Plot 1 Plot ii/n Choose a plot category and type Stain plots Basic plots Basic plots Stain plots Fit plots Chomeclad Advanced plots Bai Phot type: (scateplot) Y variable: Istp Istp Intermediate plots State Solution Solution Solution Solution Y variable: Variable: Solution Solution <th>Main Advanced Main Advanced Marker properties Color: Default</th>	Main Advanced Main Advanced Marker properties Color: Default
Ca Ca	Label color: Default Label size Label position: Default Label angle: Label angle: Label gap: V









Plot 1	
Plot ii/in	
Choose a plot category and type Basic plots Basic plots: (select type)	
C Range plots Scatter	Marker properties
C Immediate plots	Main Advanced
C Advanced plots Spike	Marker properties
- Dist luna: (so attained)	Symbol: Hollow circle
Y variable: X variable:	Color: Gray 10
sbp 💌 🛶 📷 💌 🔽 Sort on x variable	Size:
Marker properties Marker weights	
	Variable:
	Label color Default
	Labeleter
Accept	Cany Label actives
	Laber pusitum. Detault
	Laber angre:
	Label gap: 👻

Flots ii/n Y axis X axis Titles Legend Overall Plot definitions: Flot Create Creat Creat Creat </th <th>By By</th> <th></th> <th></th>	By By		
Plot 2 Plot (#/n Choose a plot category and type O Basic plots Range plots Fit plots C Investigate plots	Basic plot:: (select type)	×	
(scatter st C Advanced plots Plot type: (line plot) Y variable: V variable: Add a second y axis on right Line properties	Plot 2 Plot 2 Plot ii/in Restrict observations It (expression) sex == 1 Use a range of observations From 1 + to: 4633+	-	Create
Color: Default Color: Default Color: Default Pattern: Default Connecting method: Default Missing values: Default	-	Ŀ3	
Cancel Submit	00	Accept Cancel	Submit



Plot definitions:	axis rides Legend	Uveralit by				r
Plot 1 Plot 2 Plot 3	Edit Disable Enable Move Up	Plots if/in Y Plots Traw subgra Variables: Sex	Twoway graphs axis X axis Titles Lege phs for unique values of variat	nd Overall ^{By}		
(line yhatwom bmi,)	Move Down	Add a graph Add graphs Subgraph c Su	with totals or missing values Irganization <mark>twoway - Twoway</mark> ts ii/in Yaxis Xaxis	/ graphs Titles Legend Overall	By	×
			cheme: Default Graph size Width: (inches) 1/8	Name of graph:	F Replace	
			Scale text, markers, and li Scale multipler:	Aspect ratio:	pect ratio of plot re Placement Default	gion
		0	01	OK	Cancel	Submit

Plots if/in Plot definitions: Plot definitions: Plot 1 Plot 2 Plot 3	Twoway graphs Y axis X axis Titles Legend Directe Directe			
(line yhatwom	Itwoway - Twoway graphs Plots ii/in Y axis X axis Titles Legend Overal By Legend behavior Image: Constraint of the second of the seco	Legend organization a	and appearance properties	
2 D 🖻	Override default keys Specify order of keys and optionally change labels: T "Observed" 2 "Expected, Men" 3 "Expected, Women"	Organization Labers region liters Organization Rows/Columns: R Stack symbols and text. D Key sequence: De	Nws I ▼ 1 → Rows stault ▼ stault ▼	
ŝ		Symbol order: Du Symbol alignment: Du Row gap: Column gap:	slauk v	
		Generate keys for all symbols: (a	Ilow duplicates) Accept Cancel S	ubmit









 $\sum_{2} (x_{i2} - \overline{x}_2)^2 = s_2^2 / \operatorname{var}(b_2) = 534.75 / 0.09884^2 = 54738$ Women Source | SS df MS Number of obs = 2643 $\begin{array}{rcl} F(1, 2641) = & 428.48 \\ Prob > F & = & 0.0000 \\ R-squared & = & 0.1396 \end{array}$ -----+ Model | 229129.452 1 229129.452 Residual | 1412279.85 2641 534.751932 Adj R-squared = 0.1393 Total | 1641409.31 2642 621.275286 Root MSE = 23.125 _____ sbp | Coef. Std. Err. t P>|t| [95% Conf. Interval] bmi2.045966.098840320.7000.0001.8521542.239779_cons81.304352.54890931.8980.00076.3062986.30241 ---- $\operatorname{var}(b_1 - b_2) = s^2 \left\{ \frac{1}{\sum_1 (x_{i1} - \overline{x}_1)^2} + \frac{1}{\sum_2 (x_{i2} - \overline{x}_2)^2} \right\}$ $= 461.77 \times (1/23506 + 1/54738) = 0.02808$



24. Analyzing Subsets in Stata

```
The previous example illustrated how to restrict analyses to a
  subgroup such as men or women. This can be extended to more
  complex selections. Suppose that sex = 1 for males, 2 for females and
  that age = 1 for people < 10 years old,
           = 2 for people 10 to 19 years old, and
           = 3 for people \geq 20 years old. Then
sex==2 & age != 2
                         selects females who are not 10 to 19 years old. If
                         there are no missing values this is equivalent to
sex==2 & (age==1 | age == 3)
sex==1 | age==3
                         selects all men plus all women \geq 20.
 Logical expressions may be used to define new variables (generate
 command), to drop records from the data set (keep or drop command)
 or to restrict the data used by analysis commands such as regress.
  Logical expressions evaluate to 1 if true, 0 if false. Stata considers any
  non-zero value to be true.
```

```
Logical expressions may be used to keep or drop observations from the
data. For example
    . keep sex == 2 & age == 1
will keep young females and drop all other observations from memory
    . drop sex == 2 | age == 1
will drop women and young people keeping males ≥ 10 years old.
    . regress sbp bmi if age == 1 | age == 3
will regress sbp against bmi for people who are less than 10 or at
least 20 years of age.
```

24. What we have covered.

- Distinction between a parameter and a statistic
- \diamond The normal distribution
- Inference from a known sample about an unknown target population
- Simple linear regression: Assessing simple relationships between two continuous variables
- Interpreting the output from a linear regression program. Analyzing data with Stata
- Plotting linear regression lines with confidence bands
- Making inferences from simple linear regression models
- Lowess regression and residual plots. How do you know you have the right model?
- Transforming data to improve model fit
- Comparing slopes from two independent linear regressions

	Cited References
Greene JW, function.	Jr., Touchstone JC. Urinary estriol as an index of placental A study of 279 cases. <i>Am J Obstet Gynecol</i> 1963;85:1-9.
Gross, C. P. the Natio <i>Med</i> 340(, G. F. Anderson, et al. (1999). "The relation between funding by nal Institutes of Health and the burden of disease." <i>N Engl J</i> 24): 1881-7.
Levy D, Nat Commun National Hackensa	tional Heart Lung and Blood Institute., Center for Bio-Medical ication. 50 Years of Discovery : Medical Milestones from the Heart, Lung, and Blood Institute's Framingham Heart Study. ack, N.J.: Center for Bio-Medical Communication Inc.; 1999.
Rosner B. F	Fundamentals of Biostatistics. 6th ed. Belmont CA: Danbury 2006
	For additional references on these notes see.
Dupont WD Introducti U.K.: Can	. Statistical Modeling for Biomedical Researchers: A Simple ion to the Analysis of Complex Data. 2nd ed. Cambridge, nbridge University Press; 2009.