

3.2

Least Squares Regression Line (LSRL)

Regression Line

A regression line is a line that describes how a response variable y changes as an explanatory variable x changes. We often use a regression line to predict the value of y for a given value of x .

Regression Lines as Mathematical Models

Suppose that y is a response variable (plotted on the vertical axis) and x is an explanatory variable (plotted on the horizontal axis). A regression line relating y to x has an equation of the form

$$y = a + bx$$

In this equation, b is the **slope**, the amount by which y changes when x increases by one unit. The number a is the **y intercept**, the value of y when $x = 0$.

* **Example** to investigate the steps to develop an LSRL equation (ANSWER THE FOLLOWING QUESTIONS)

1. Enter L1 - Non-exercise activity
2. Enter L2 - Fat Gained
3. Plot the scatter plot. What is the association (direction, form, and strength)?
4. Find the mean and standard deviation for both variables in context.
5. Find the linear regression equation. What does it mean?
6. Plot the LSRL on the scatterplot. What are residuals?
7. Plot the residuals. What does this mean?
8. How do you assess the model? What does r^2 mean?
9. Use the LSRL equation to make predictions. When is it inappropriate to predict with LSRL?

NEA (calories)	Fat Gained (kilograms)
-94	4.2
-57	3.0
-29	3.7
135	2.7
143	3.2
151	3.6
245	2.4
355	1.3
392	3.8
473	1.7
486	1.6
535	2.2
571	1.0
580	0.4
620	2.3
690	1.1

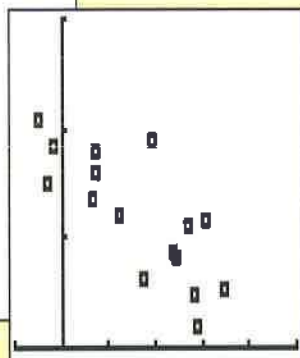
Review the Data

1ST STEP -
ALWAYS GRAPH!!!

THE SCATTERPLOT - The relationship between non-exercise activity and fat shows a negative association, with a linear form, and appears to have a moderately strong relationship.

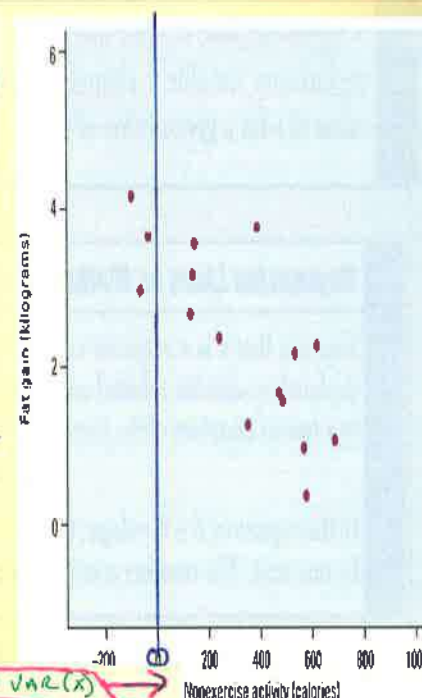
Plot1 Plot2 Plot3
Off
Type: ☒ WINDOW *
Xlist: L1
Ylist: L2
Mark: ☒

Xmin=-200
Xmax=1000
Xscl=200
Ymin=0
Ymax=6
Yscl=2
Xres=1



* ALWAYS SET YOUR OWN WINDOW TO SKETCH THE SCATTERPLOT

RESPONSE VAR (Y)



DESCRIPTIVE STATISTICS -

EXPLANATORY VAR (X)

- The mean for non-exercise activity is about 325 calories with a standard deviation of about 258 calories with a spread (based on range) of 794 calories.
- The mean for fat gain is 2.39 kilograms with a standard deviation of 1.14 kilograms and a spread (based on range) of 3.8 calories.

2-Var Stats
 $\bar{x}=324.75$
 $\Sigma x=5196$
 $\Sigma x^2=2683206$
 $Sx=257.6567484$
 $\sigma x=249.4750739$
 $n=16$

2-Var Stats
 $\bar{y}=2.3875$
 $\Sigma y=38.2$
 $\Sigma y^2=110.66$
 $Sy=1.138932248$
 $\sigma y=1.102766408$
 $\Sigma xy=8978.4$

2-Var Stats
 $\bar{y}=2.3875$
 $\Sigma y=38.2$
 $\Sigma y^2=110.66$
 $\Sigma xy=8978.4$
 $\min X=-94$
 $\max X=690$
 $\min Y=.4$
 $\max Y=4.2$

Equation of LSRL

Equation of the Least-Squares Regression Line

We have data on an explanatory variable x and a response variable y for n individuals. From the data, calculate the means \bar{x} and \bar{y} and the standard deviations s_x and s_y of the two variables and their correlation r . The least-squares regression line is

with slope

that passes through the point (\bar{x}, \bar{y}) .

$$\hat{y} = a + bx$$

$$\hat{Y} = b_0 + b_1 X$$

$$b = r \frac{s_y}{s_x}$$

$$b_1 = r \frac{s_y}{s_x}$$

DO NOT NEED TO MEMORIZE

- These are on your AP GREEN SHEET
- Check it out!

LinReg
 $y = a + bx$
 $a = 3.505122916$
 $b = -.003441487$
 $r^2 = .6061492049$
 $r = -.7785558457$

y INT (a)
 slope (b)
 COEFF OF DETERMINATION
 CORR COEF

The slope here $B = -.00344$ tells us that fat gained goes down by .00344 kg for each added calorie of NEA according to this linear model.

Our regression equation is the predicted RATE OF CHANGE in the response y as the explanatory variable x changes.

↑ FAT GAIN

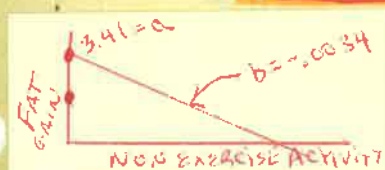
NONEXERCISE ACTIVITY (NEA)

The Y intercept $a = 3.505\text{kg}$ is the fat gain estimated by this model if NEA does not change when a person overeats.

LSRL EQUATION: $\hat{Y} = 3.51 - .0034X$

*** * (better to use words) * * ***

$\hat{Y}_{\text{HAT}} = \hat{Y}$
 \hat{Y} - predicted y



$$(\text{Fat Gain})_{\text{hat}} = 3.51 - .0034(\text{NEA})$$

PREDICTED VARIABLES MUST WEAR A HAT!

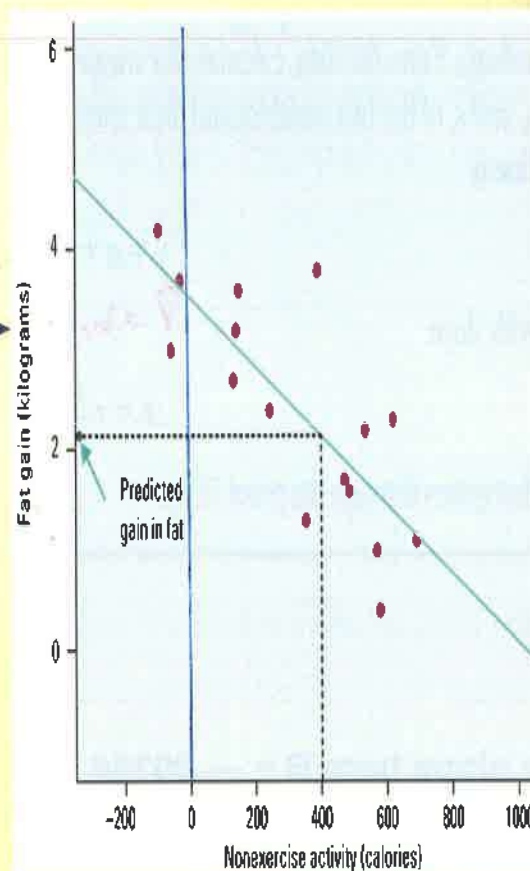
Graph the LSRL on our Scatterplot

LSRL EQUATION:

$$\hat{Y} = 3.51 - .0034X$$

$$(\text{Fat Gain})_{\text{hat}} = 3.51 - .0034(\text{NEA})$$

2nd 2 Plot2 Plot3
 Y1= 3.51 - .0034X
 Y2=
 Y3=
 Y4=
 Y5=
 Y6=
 Y7=



TIP:

But in fact, we can get the LSRL equation, from our calculator by saving the equation when we calculate the LSRL. Here are the steps if you want to try on your own.

LinReg(a+bx)

VARs Y=V1:2

Function

2:P1:Y1

3:P2:Y2

4:0:Y3

5:Y4

6:Y5

7:Y6

7↓Y7

LinReg(a+bx) Y1

2nd 2 Plot2 Plot3
 Y1= 3.505122915
 Y2=
 Y3=
 Y4=
 Y5=
 Y6=
 Y7=

The LSRL “Line”

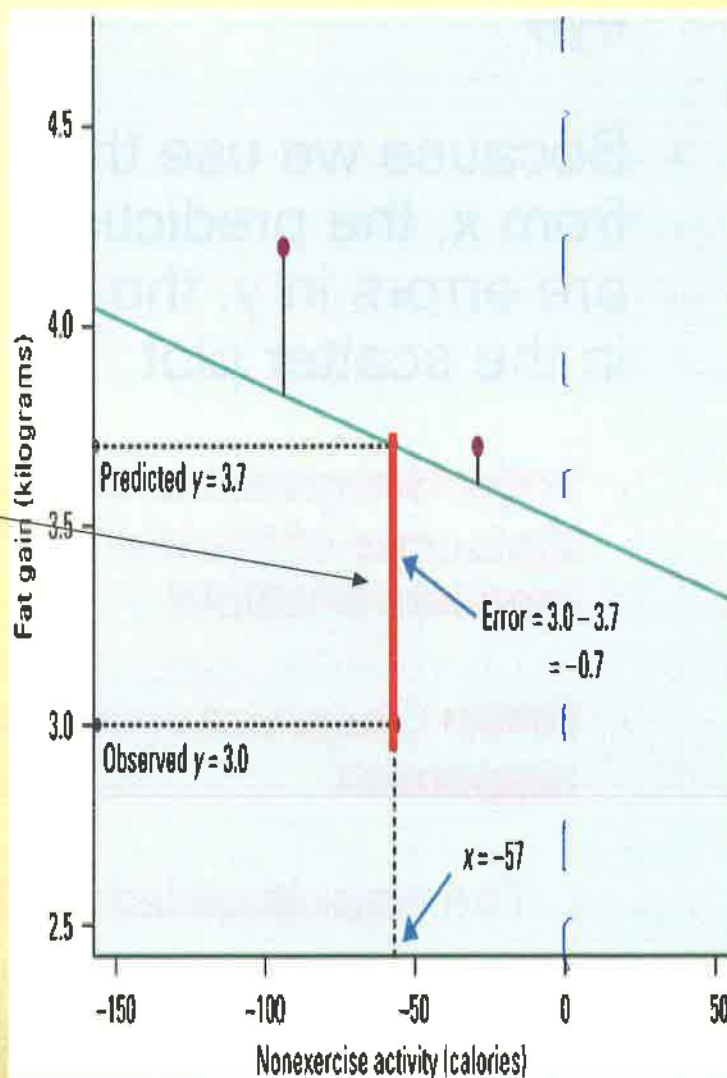
- In most cases, no line will pass exactly through all the points in a scatter plot and different people will draw different regression lines by eye.
- Because we use the line to predict y from x , the prediction errors we make are errors in y , the vertical direction in the scatter plot
 - A good regression line makes the vertical distances of the points from the line as small as possible
 - Error: Observed response - predicted response
 - The error is called RESIDUALS!

Goal of LSRL

Least-Squares Regression Line

The least-squares regression line of y on x is the line that makes the sum of the squared vertical distances of the data points from the line as small as possible.

- Goal of LSRL is to minimize error.
- The error is called residuals.
- Want to minimize the sum of the residuals squared.



Residuals

Residuals

A **residual** is the difference between an observed value of the response variable and the value predicted by the regression line. That is,

$$\begin{aligned}\text{residual} &= \text{observed } y - \text{predicted } y \\ &= y - \hat{y}\end{aligned}$$

- The error of our predictions, or vertical distance from predicted \hat{Y} to observed Y , are called residuals because they are “left-over” variation in the response.

EXAMPLE: One subject's NEA rose by 135 calories. That subject gained 2.7 KG of fat. The predicted gain for 135 calories is

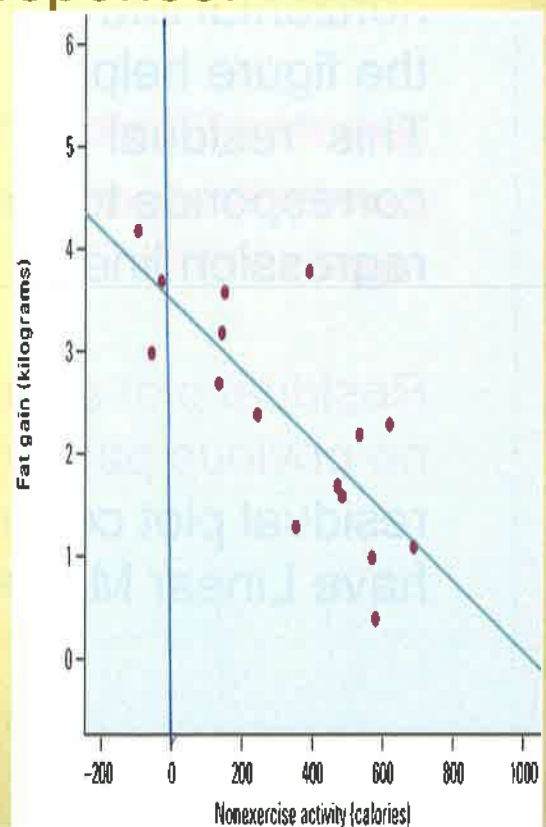
Predicted: (\hat{y})

$$\hat{Y} = 3.505 - .00344(135) = 3.04 \text{ kg}$$

Observed: 2.7 KG of fat

The **residual** for this subject is

$$\begin{aligned}y - \hat{y} \\ y - \hat{y} = 2.7 - 3.04 = -.34 \text{ kg}\end{aligned}$$

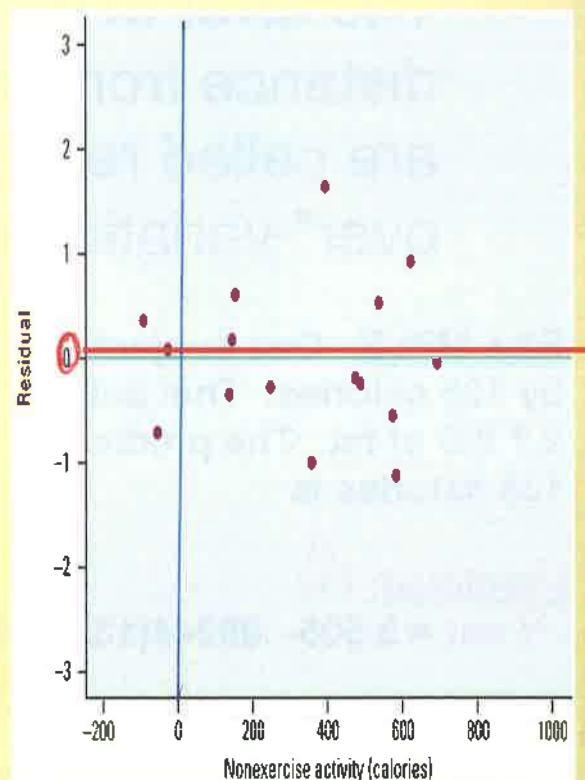


Residual Plot

Residual Plots

A residual plot is a scatterplot of the regression residuals against the explanatory variable (or equivalently, against the predicted y-values). Residual plots help us assess how well a regression line fits the data.

- The **sum of the least-squares residuals is always zero.**
- The mean of the residuals is always zero, the horizontal line at zero in the figure helps orient us. This "residual = 0" line corresponds to the regression line
- **Residual plot should show no obvious pattern.** Our residual plot confirms we have Linear Model.



Residuals Plots on Calc

- If you want to get all your residuals listed in L3 highlight L3 (the name of the list, on the top) and go to 2nd- stat- RESID then hit enter and enter and the list that pops out is your resid for each individual in the corresponding L1 and L2. (if you were to create a normal scatter plot using this list as your y list, so x list: L1 and Y list L3 you would get the exact same thing as if you did a residual plot defining x list as L1 and Y list as RESID as we had been doing).
- This is a helpful list to have to check your work when asked to calculate an individuals residual.

2nd STAT PLOT

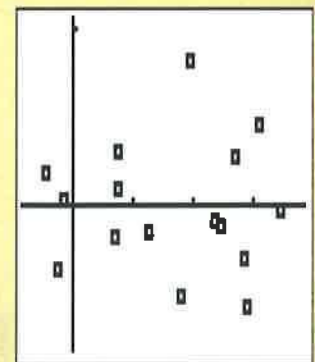
```
STAT PLOT
1:Plot1...On
  L1 RESID
2:Plot2...Off
  Plot2 Plot3
  Off
Type: [ ] [ ] [ ]
Xlist:L1
Ylist:RESID
Mark: [ ] +
```

2nd LIST

```
NAME OPS MATH
1:L1
2:L2
3:L3
4:L4
5:L5
6:L6
7:RESID
```

ZOOM 9

Residual Plot

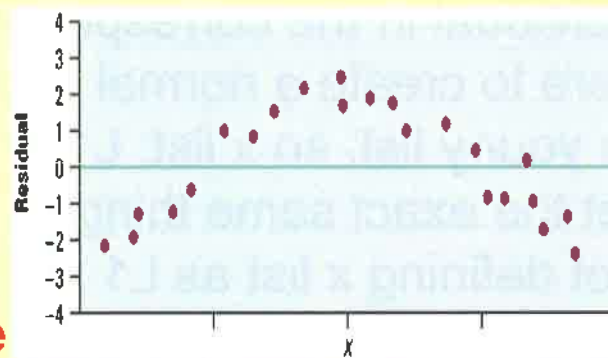


Assessing Models

Examining Residuals Plots and Residual Standard Deviation

- **Residual plot should show no obvious pattern.**

Residual plot should show no obvious pattern. A curved pattern shows that the relationship is not linear and a straight line may not be the best model.



- **Residuals should be** regression line in a model that fits the data well should come close" to most of the points.
- A commonly used measure of this is the standard deviation of the residuals, given by:

$$s = \sqrt{\frac{\sum \text{residuals}^2}{n-2}}$$

For the NEA and fat gain data, $s_e = \sqrt{\frac{7.66}{14}} = .740$

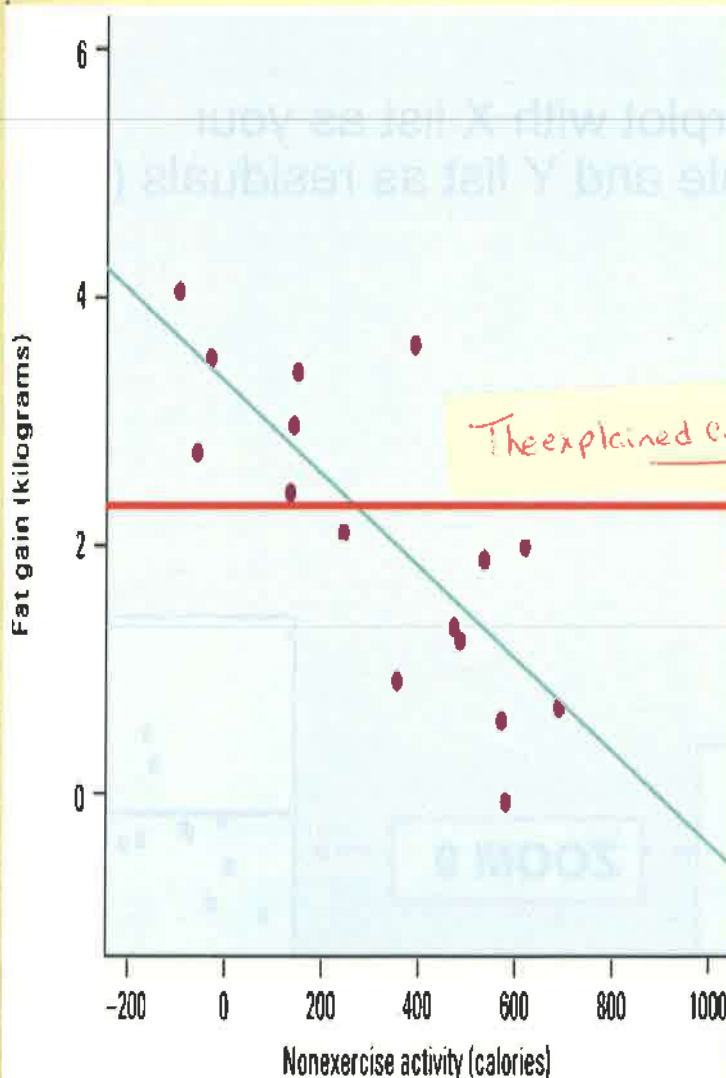
Assessing Models

R squared- Coefficient of determination

The Coefficient of Determination: r^2 in Regression

The coefficient of determination r^2 is the fraction of the variation in the values of y that is explained by the least-squares regression line of y on x . We can calculate r^2 using the following formula:

```
LinReg
y=a+bx
a=3.505122916
b=-.003441487
r2=.6061492049
r=-.7785558457
```



1. If all the points fall directly on the least-squares line, r squared = 1. Then all the variation in y is explained by the linear relationship with x .
2. So, if r squared = .606, that means that 61% of the variation in y among individual subjects is due to the influence of the other variable. The other 39% is "not explained".
3. r squared is a measure of how successful the regression was in explaining the response.

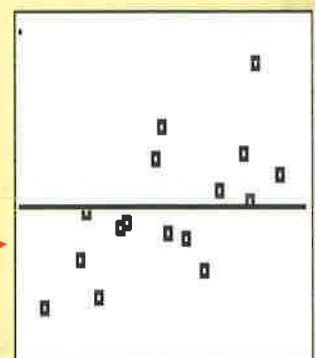
Plot Residual vs X or Y

- Produce Scatterplot and Regression line from data (lets use BAC if still in there)
- Turn all plots off
- Create new scatterplot with X list as your explanatory variable and Y list as residuals (2nd stat, resid)
- Zoom Stat

```
5:HI PLOT2
1:Plot1...On
  L2 RESID
2:Plot2...Off
  L5 1
3:Plot3...Off
  L1 L2
4↓PlotsOff
```

```
Plot2 Plot3
Off Off
Type: [ ] [ ] [ ]
      [ ] [ ] [ ]
Xlist:L2
Ylist:RESID
Mark: [ ] + .
```

ZOOM 9



NO pattern

NOT PART OF THIS EXAMPLE

Prediction

- We can use a regression line to predict the response y for a specific value of the explanatory variable x (but only for the range of x values used in our LSRL).

Extrapolation

Extrapolation is the use of a regression line for prediction outside the range of values of the explanatory variable x used to obtain the line. Such predictions are often not accurate.

