

LINEAR REGRESSION AND CORRELATION



Suppose we have the following paired data, also known as *bivariate data*.

$(1,3)$, $(2,2)$, $(3,4)$

We can enter this data into lists in our calculator and create a scatterplot.

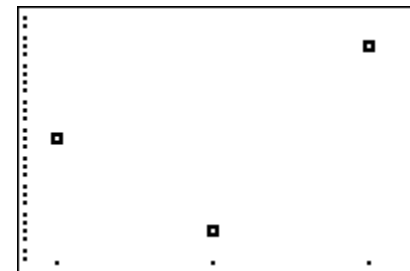
$(1,3)$, $(2,2)$, $(3,4)$

L1	L2	L3	2
1	3	-----	
2	2	-----	
3	4	-----	

L2(4) =

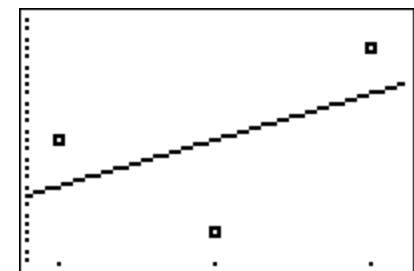
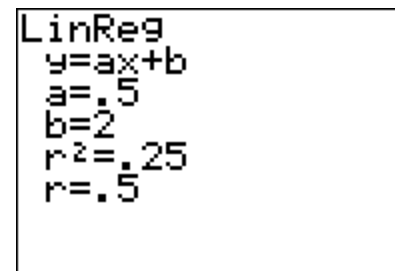
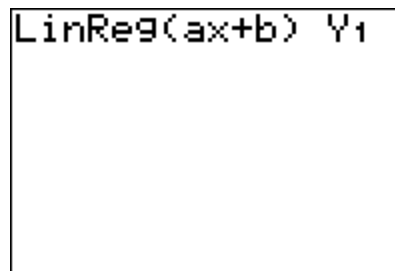
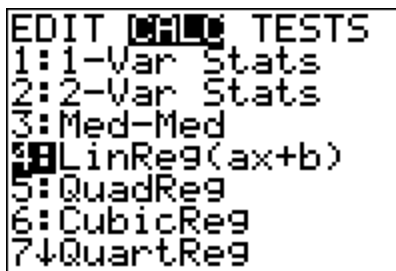
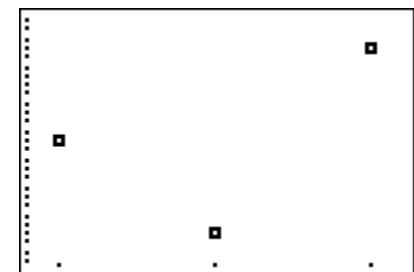
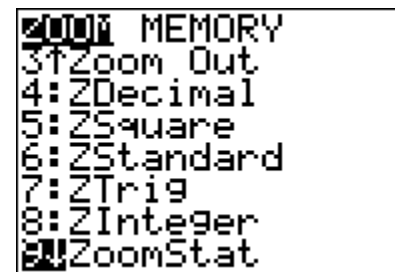
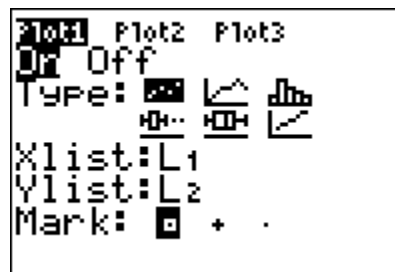
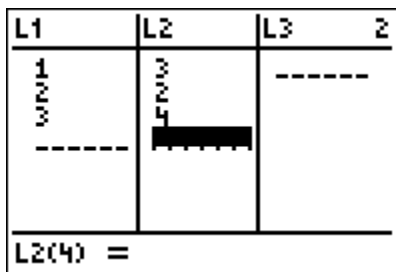
```
Plot2 Plot3
Off
Type: [ ] [ ] [ ]
Xlist:L1
Ylist:L2
Mark: [ ] + .
```

```
MEMORY
3:Zoom Out
4:ZDecimal
5:ZSquare
6:ZStandard
7:ZTrig
8:ZInteger
9:ZoomStat
```



We can also use our calculator to find the best fitting straight line for this data.

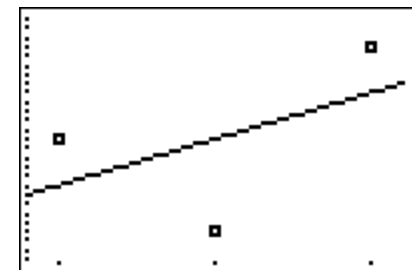
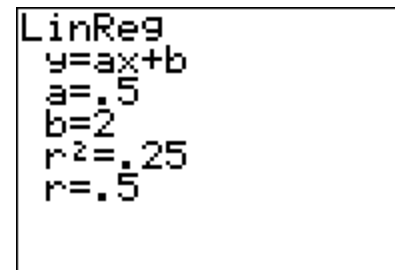
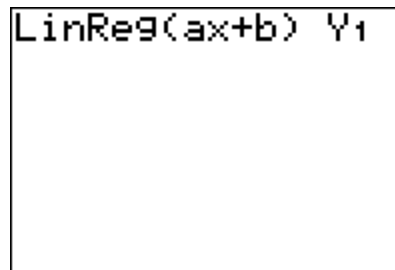
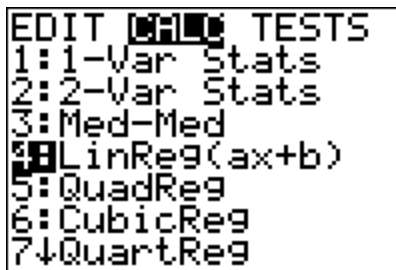
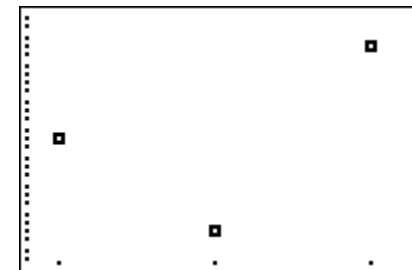
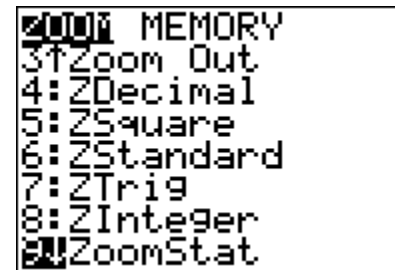
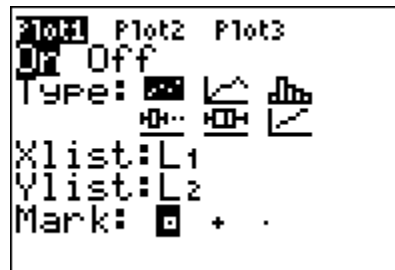
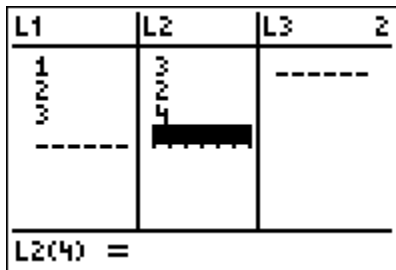
$(1,3)$, $(2,2)$, $(3,4)$



$VARS \rightarrow Y-VARS$
 $\rightarrow FUNCTION \rightarrow Y_1$

And now we have much to discuss.

$(1,3)$, $(2,2)$, $(3,4)$



$VARS \rightarrow Y-VARS$
 $\rightarrow FUNCTION \rightarrow Y_1$

If you don't see the values for r and r^2 on your screen, then do this:

```
LinReg  
y=ax+b  
a=.5  
b=2
```

If you don't see the values for r and r^2 on your screen, then do this:

```
LinReg
y=ax+b
a=.5
b=2
```

```
CATALOG
det(
DiagnosticOff
DiagnosticOn
dim(
Disp
DispGraph
DispTable
```

```
DiagnosticOn
```

```
DiagnosticOn
Done
```

```
EDIT TESTS
1:1-Var Stats
2:2-Var Stats
3:Med-Med
4:LinReg(ax+b)
5:QuadReg
6:CubicReg
7:QuartReg
```

```
LinReg(ax+b)
```

```
LinReg
y=ax+b
a=.5
b=2
r2=.25
r=.5
```

The number r is called the *coefficient of linear correlation*, or the *pearson product moment correlation coefficient* (among other things).

```
LinReg  
y=ax+b  
a=.5  
b=2  
r2=.25  
r=.5
```


This number measures the strength of the linear or straight line relationship between the two variables.

```
LinReg  
y=ax+b  
a=.5  
b=2  
r2=.25  
r=.5
```

We always have that $-1 \leq r \leq 1$.

```
LinReg  
y=ax+b  
a=.5  
b=2  
r2=.25  
r=.5
```

The closer $|r|$ is to 1, the stronger the linear relationship,
and the closer $|r|$ is to 0, the weaker the linear relationship.

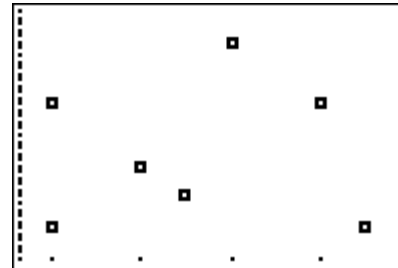
```
LinReg  
y=ax+b  
a=.5  
b=2  
r2=.25  
r=.5
```

Some statisticians use the following classification scheme:

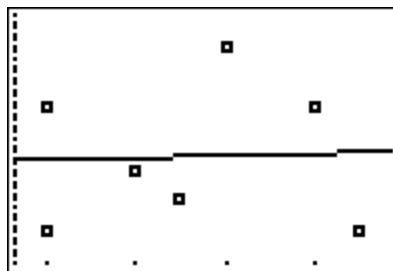
$$\begin{array}{ll} \text{Weak:} & 0 \leq |r| < .4 \\ \text{Moderate:} & .4 \leq |r| < .7 \\ \text{Strong:} & .7 \leq |r| \leq 1 \end{array}$$

Example 1:

L1	L2	L3	2
1	3	-----	
2	2		
3	4		
1	1		
2.5	1.5		
4	3		
4.5	1		
<hr/>			
L2(1)=3			

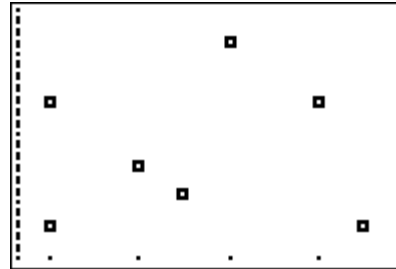


```
LinReg  
y=ax+b  
a=.0350318471  
b=2.124203822  
r2=.0017358122  
r=.0416630801
```

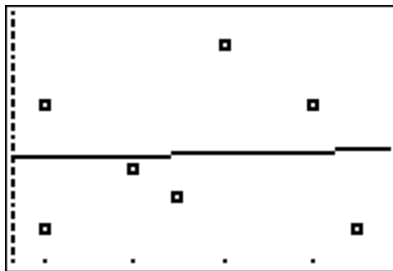


Example 1:

L1	L2	L3	2
1	3	-----	
2	2		
3	4		
1	1		
2.5	1.5		
4	3		
4.5	1		
<hr/>			
L2(1)=3			



```
LinReg  
y=ax+b  
a=.0350318471  
b=2.124203822  
r2=.0017358122  
r=.0416630801
```

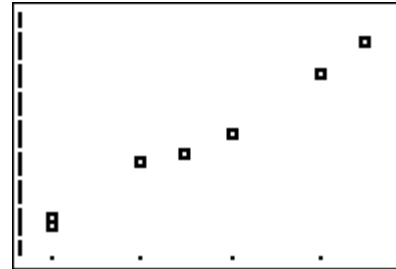


The linear correlation is very weak.

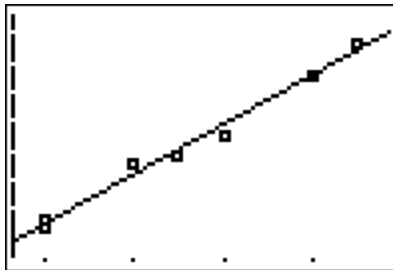
Example 2:

L1	L2	L3	2
1	2	-----	
2	3.2		
3	3.8		
4	4.1		
2.5	3.4		
4	5.1		
4.5	5.8		

L2(1)=2



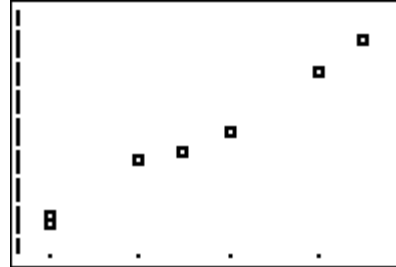
```
LinReg  
y=ax+b  
a=1.075159236  
b=.8210191083  
r2=.9874150108  
r=.9936875821
```



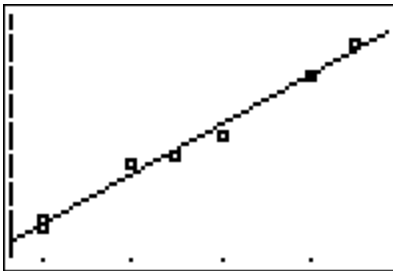
Example 2:

L1	L2	L3	2
1	2	-----	
2	3.2		
3	3.8		
4	4.8		
2.5	3.4		
4	5.1		
4.5	5.8		

L2(1)=2



```
LinReg  
y=ax+b  
a=1.075159236  
b=.8210191083  
r2=.9874150108  
r=.9936875821
```

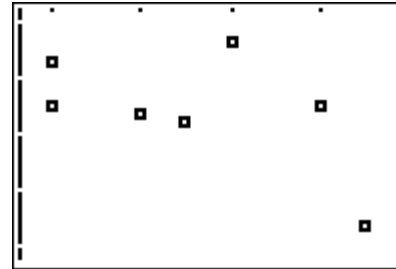


This is a strong, positive linear correlation.
As the input variable increases, the output variable also increases.

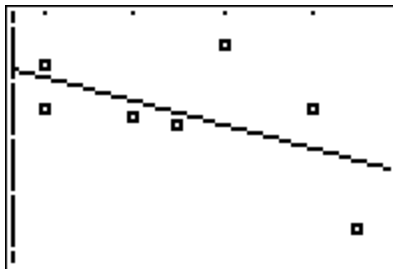
Example 3:

L1	L2	L3	2
1	-2	-----	
2	-3.2		
3	-1.5		
4	-3		
2.5	-3.4		
4	-3		
4.5	-5.8		

L2(1) = -2



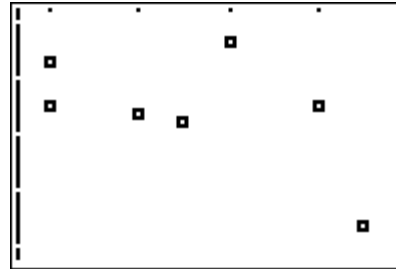
```
LinReg  
y=ax+b  
a=-.5515923567  
b=-1.710191083  
r2=.3053432504  
r=-.5525787278
```



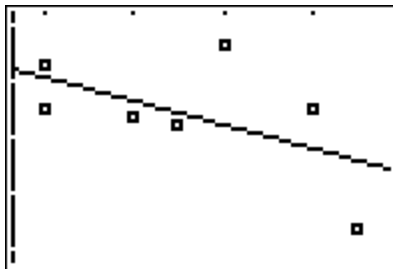
Example 3:

L1	L2	L3	2
1	-2	-----	
2	-3.2		
3	-1.5		
1	-3		
2.5	-3.4		
4	-3		
4.5	-5.8		

L2(1) = -2



```
LinReg  
y=ax+b  
a=-.5515923567  
b=-1.710191083  
r2=.3053432504  
r=-.5525787278
```

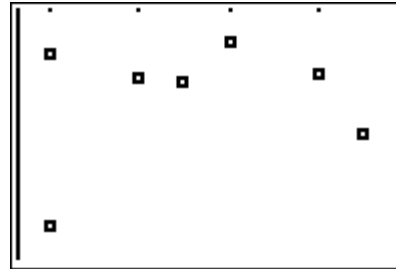


This is a moderate, negative linear correlation.
As the input variable increases, the output variable decreases.

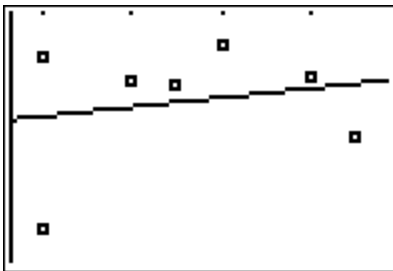
Example 4:

L1	L2	L3	2
1	-2	-----	
2	-3.2		
3	-1.5		
1	-10		
2.5	-3.4		
4	-3		
4.5	-5.8		

L2(1) = -2



```
LinReg  
y=ax+b  
a=.4292993631  
b=-5.232484076  
r2=.040229655  
r=.2005733158
```

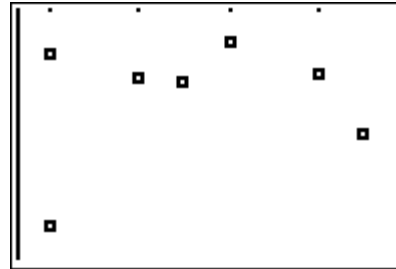


In this data set we changed just a single point, and it changes the correlation from moderate & negative to weak & positive.

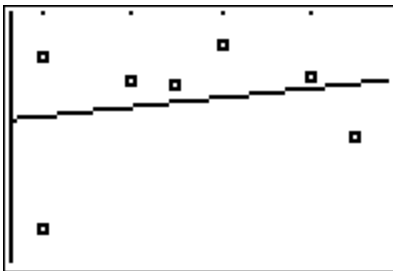
Example 4:

L1	L2	L3	2
1	-2	-----	
2	-3.2		
3	-1.5		
1	-10		
2.5	-3.4		
4	-3		
4.5	-5.8		

L2(1) = -2



```
LinReg  
y=ax+b  
a=.4292993631  
b=-5.232484076  
r2=.040229655  
r=.2005733158
```

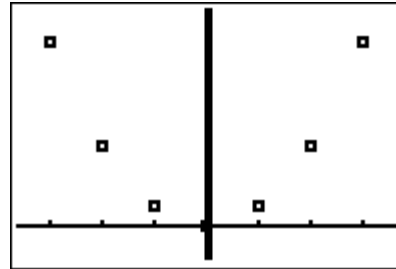


In this data set we changed just a single point, and it changes the correlation from moderate & negative to weak & positive.

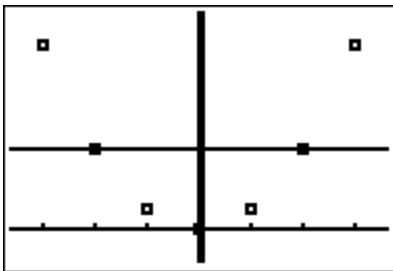
A single data point can strongly affect the degree of linear correlation.

Example 5:

L1	L2	L3	Z
1	9	-----	
2	1		
3	4		
4	1		
5	1		
6	1		
7	1		
8	1		
9	1		
10	1		
11	1		
12	1		
13	1		
14	1		
15	1		
16	1		
17	1		
18	1		
19	1		
20	1		
21	1		
22	1		
23	1		
24	1		
25	1		
26	1		
27	1		
28	1		
29	1		
30	1		
31	1		
32	1		
33	1		
34	1		
35	1		
36	1		
37	1		
38	1		
39	1		
40	1		
41	1		
42	1		
43	1		
44	1		
45	1		
46	1		
47	1		
48	1		
49	1		
50	1		
51	1		
52	1		
53	1		
54	1		
55	1		
56	1		
57	1		
58	1		
59	1		
60	1		
61	1		
62	1		
63	1		
64	1		
65	1		
66	1		
67	1		
68	1		
69	1		
70	1		
71	1		
72	1		
73	1		
74	1		
75	1		
76	1		
77	1		
78	1		
79	1		
80	1		
81	1		
82	1		
83	1		
84	1		
85	1		
86	1		
87	1		
88	1		
89	1		
90	1		
91	1		
92	1		
93	1		
94	1		
95	1		
96	1		
97	1		
98	1		
99	1		
100	1		
L2(1)=9			



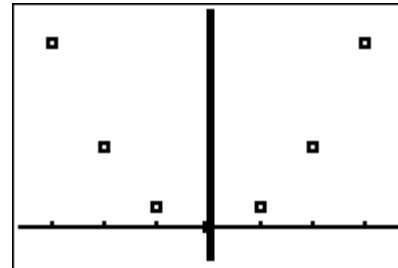
```
LinReg
y=ax+b
a=0
b=4
r=0
r=0
```



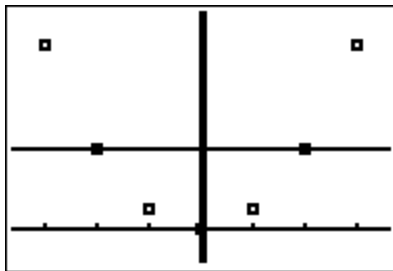
There is clearly a mathematical relationship between the variables, but it is not a linear one. Hence, $r = 0$.

Example 5:

L1	L2	L3	Z
1	9	-----	
2	5		
3	1		
4	4		
5	7		
6	1		
7	5		
8	9		
9	1		
10	4		
11	7		
12	1		
13	5		
14	9		
15	1		
16	4		
17	7		
18	1		
19	5		
20	9		
21	1		
22	4		
23	7		
24	1		
25	5		
26	9		
27	1		
28	4		
29	7		
30	1		
31	5		
32	9		
33	1		
34	4		
35	7		
36	1		
37	5		
38	9		
39	1		
40	4		
41	7		
42	1		
43	5		
44	9		
45	1		
46	4		
47	7		
48	1		
49	5		
50	9		
L2(1)=9			



```
LinReg  
y=ax+b  
a=0  
b=4  
r=0  
r=0
```



There is clearly a mathematical relationship between the variables, but it is not a linear one. Hence, $r = 0$.

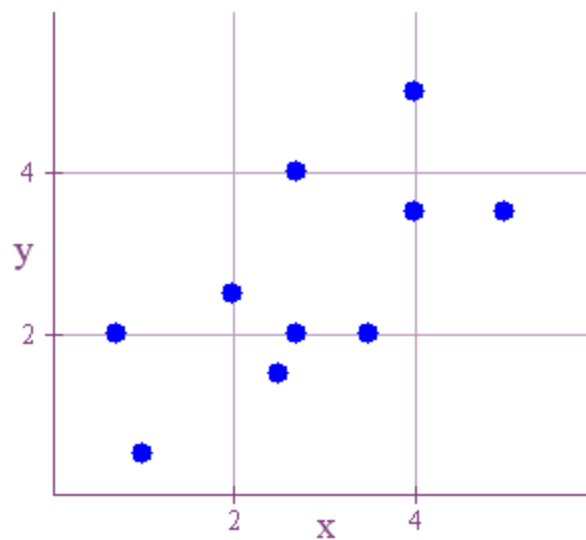
NOTE: If we are dealing with a population instead of a sample, then the linear correlation coefficient is denoted by $\rho = \rho$.

How do we compute the coefficient of linear correlation? We basically convert each x and y coordinate to a z score, multiply them together, add them up, and then divide by the number of ordered pairs in order to adjust the size of our sum. Of course, when we are dealing with sample data, we divide by $n-1$ instead of n .

$$r = \frac{\sum z_x z_y}{n - 1}$$

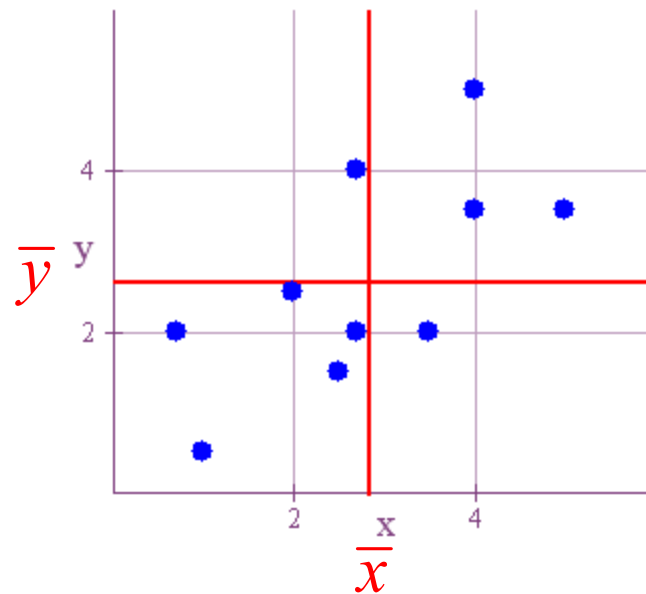
Here's how it works in practice. The scatterplot below clearly represents a positive linear correlation. Also, recall the formula for converting a raw score to a z score.

$$z = \frac{x - \mu}{\sigma}$$



If we draw some lines to represent the mean of the x coordinates and the mean of the y coordinates, then most of our coordinates are either both above their means or both below their respective means.

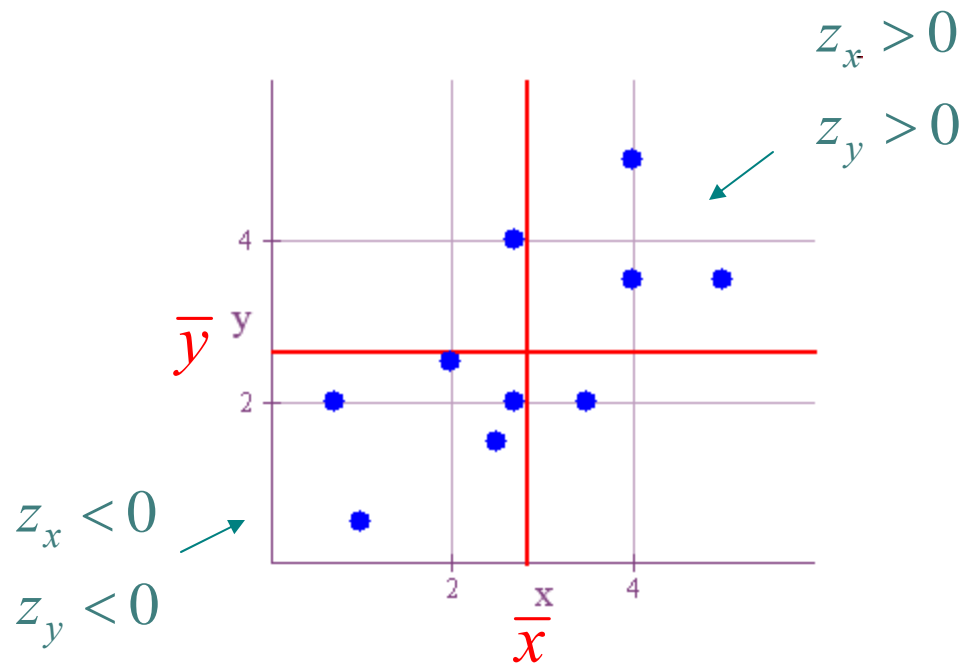
$$z = \frac{x - \mu}{\sigma}$$



This means that most of the products $z_x z_y$ will be positive, and, hence, r will be positive.

$$r = \frac{\sum z_x z_y}{n-1} > 0$$

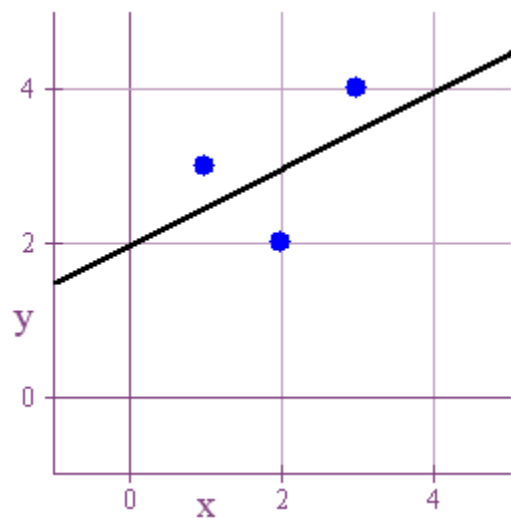
$$z = \frac{x - \mu}{\sigma}$$



Let's now go back to our first example.

```
LinReg  
y=ax+b  
a=.5  
b=2  
r2=.25  
r=.5
```

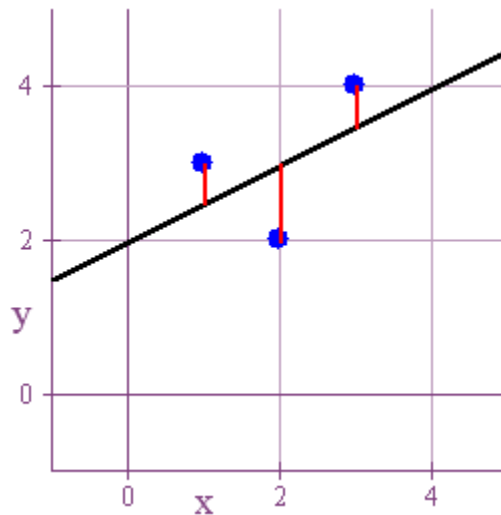
$(1,3)$, $(2,2)$, $(3,4)$



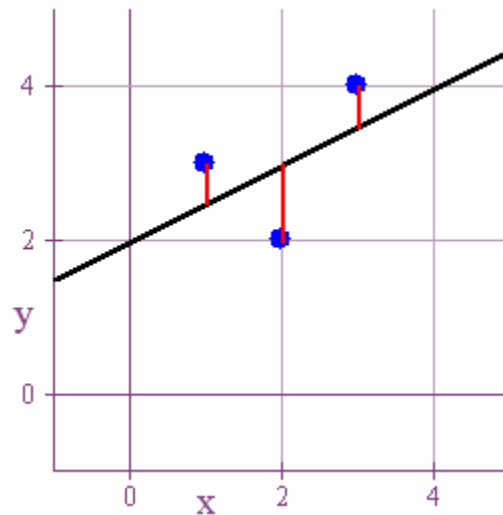
The line which best fits out data is called the *least squares regression line*.

```
LinReg  
y=ax+b  
a=.5  
b=2  
r2=.25  
r=.5
```

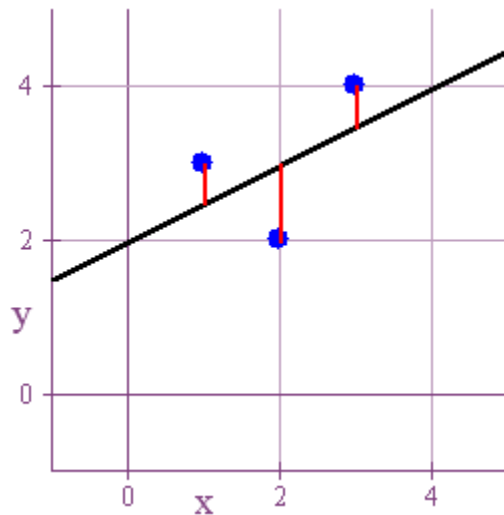
$(1,3)$, $(2,2)$, $(3,4)$



What this means is that if you look at the distance from the line to each data point, square that distance, and add everything up, then the sum of the squared distances will be minimized if we are using the *least squares regression line*.



Another way to look at it is this. Think of our data points as nails on a board, and think of the red lines as rubber bands connecting the nails to a broom stick. Then eventually the broom stick will settle down to some equilibrium position. That final position is the *least squares regression line*.

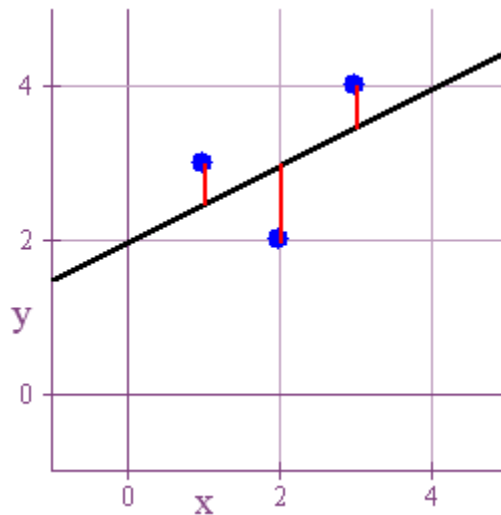


We often write the regression equation as \hat{y} to distinguish it from other y -values.

```
LinReg  
y=ax+b  
a=.5  
b=2  
r2=.25  
r=.5
```

$(1,3), (2,2), (3,4)$

$$\hat{y} = .5x + 2$$

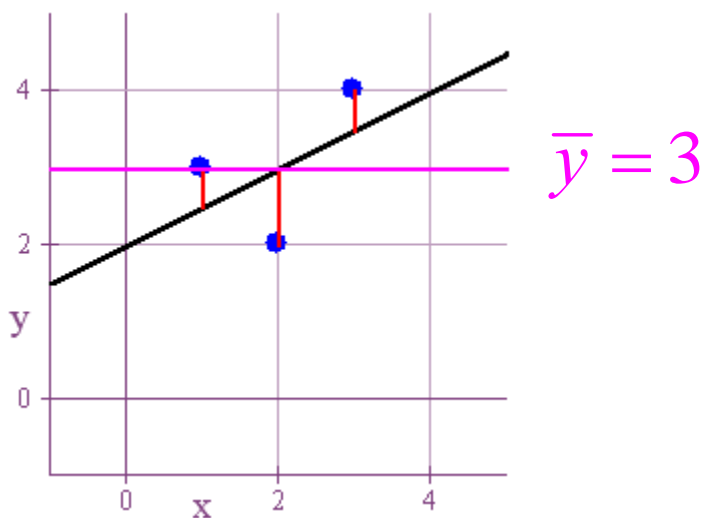


Now let's draw in a horizontal line at the position of the mean of the *y-values*.

```
LinReg  
y=ax+b  
a=.5  
b=2  
r2=.25  
r=.5
```

$(1,3), (2,2), (3,4)$

$$\hat{y} = .5x + 2$$

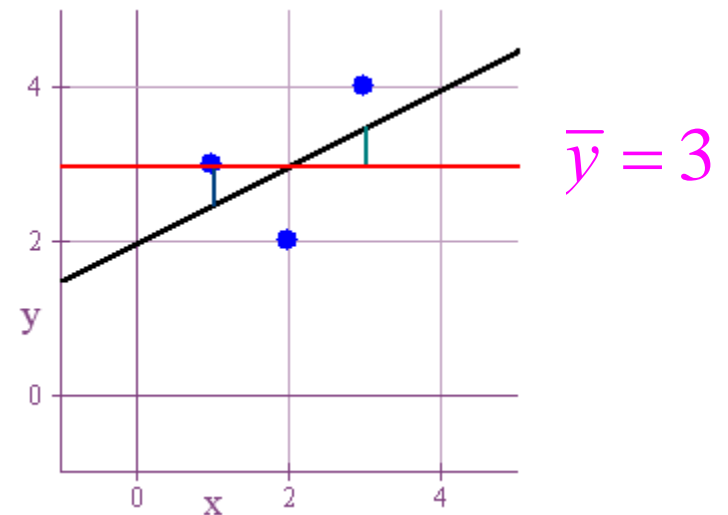


For each x value in our data points, the explained variance is the square of the distance between the line y -bar and the regression line.

$$(1,3), (2,2), (3,4)$$

$$\hat{y} = .5x + 2$$

$$\begin{aligned} \text{Explained Variance} &= \sum (\hat{y} - \bar{y})^2 \\ &= (2.5 - 3)^2 + (3 - 3)^2 + (3.5 - 3)^2 = 0.5 \end{aligned}$$



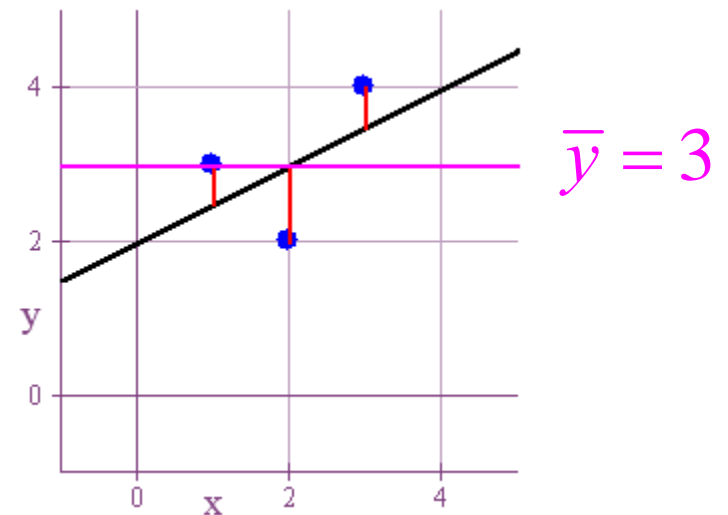
The unexplained variance is the square of the distance between the regression line and the y values of the data points.

$$(1,3), (2,2), (3,4)$$

$$\hat{y} = .5x + 2$$

$$\begin{aligned} \text{Explained Variance} &= \sum (\hat{y} - \bar{y})^2 \\ &= (2.5 - 3)^2 + (3 - 3)^2 + (3.5 - 3)^2 = 0.5 \end{aligned}$$

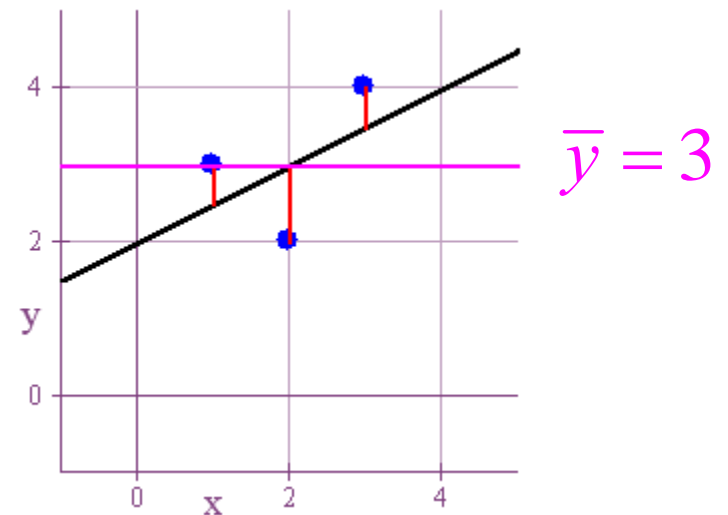
$$\begin{aligned} \text{Unexplained Variance} &= \sum (y - \hat{y})^2 \\ &= (3 - 2.5)^2 + (2 - 3)^2 + (4 - 3.5)^2 = 1.5 \end{aligned}$$



The *coefficient of determination* is the ratio of *explained variance* to *total variance*.

$(1,3), (2,2), (3,4)$

$$\hat{y} = .5x + 2$$

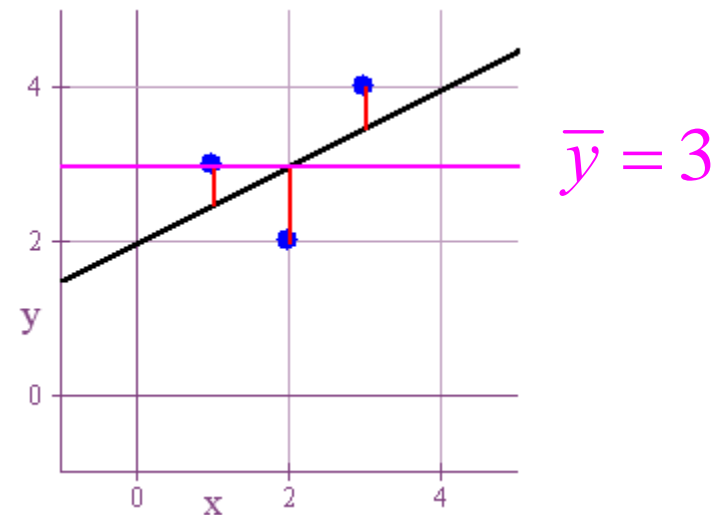


$$\text{Coefficient of Determination} = \frac{\text{Explained Variance}}{\text{Total Variance}} = \frac{0.5}{0.5 + 1.5} = \frac{0.5}{2} = 0.25$$

The *coefficient of determination* is equal to r^2 .

$(1,3), (2,2), (3,4)$

$$\hat{y} = .5x + 2$$



$$\text{Coefficient of Determination} = \frac{\text{Explained Variance}}{\text{Total Variance}} = \frac{0.5}{0.5 + 1.5} = \frac{0.5}{2} = 0.25$$

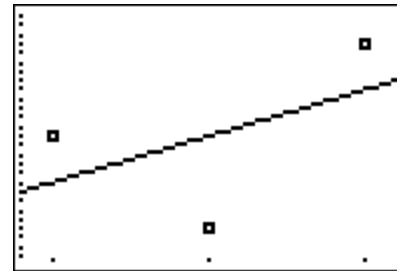
$$r^2 = 0.25$$

One of the more common uses of the linear regression equation is to make predictions.

$(1,3)$, $(2,2)$, $(3,4)$

```
LinReg(ax+b) Y1
```

```
LinReg  
y=ax+b  
a=.5  
b=2.5  
r2=.25  
r=.5
```



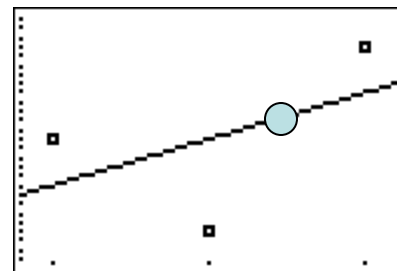
VARs → *Y - VARs*
→ *FUNCTION* → Y_1

If we use the equation to “fill in the gaps” within the range of our given x values, we call this *interpolation*.

$(1,3)$, $(2,2)$, $(3,4)$

```
LinReg(ax+b) Y1
```

```
LinReg  
y=ax+b  
a=.5  
b=2  
r2=.25  
r=.5
```



*VAR*S → *Y* - *VAR*S
→ *FUNCTION* → *Y*₁

```
Plot1 Plot2 Plot3  
Y1 = .5X+2  
Y2 =  
Y3 =  
Y4 =  
Y5 =  
Y6 =  
Y7 =
```

```
Y1(2.5) 3.25
```

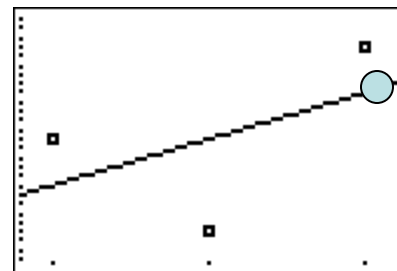
$(2.5, 3.25)$

If we use the equation to predict values outside the range of our given x values, we call this *extrapolation*.

$(1,3)$, $(2,2)$, $(3,4)$

```
LinReg(ax+b) Y1
```

```
LinReg
y=ax+b
a=.5
b=2
r2=.25
r=.5
```



VARs → *Y-VARS*
→ *FUNCTION* → Y_1

```
Plot1 Plot2 Plot3
\Y1 □ .5X+2
\Y2 =
\Y3 =
\Y4 =
\Y5 =
\Y6 =
\Y7 =
```

```
Y1(3.1)
3.55
```

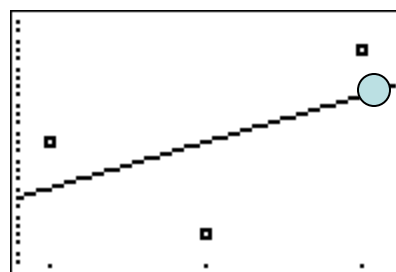
$(3.1, 3.55)$

As a general rule, we should be careful about extrapolating too far into either the future or the past.

$(1,3)$, $(2,2)$, $(3,4)$

```
LinReg(ax+b) Y1
```

```
LinReg
y=ax+b
a=.5
b=2
r2=.25
r=.5
```



VARs → *Y-VARS*
 → *FUNCTION* → Y_1

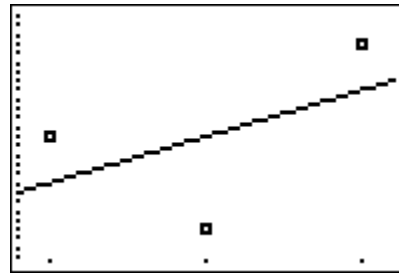
```
Plot1 Plot2 Plot3
Y1 .5X+2
Y2 =
Y3 =
Y4 =
Y5 =
Y6 =
Y7 =
```

```
Y1(3.1)
3.55
```

$(3.1, 3.55)$

Finally, there's one more question we should address. Below we found a linear correlation coefficient of 0.5. How do we know that this is significantly different from zero?

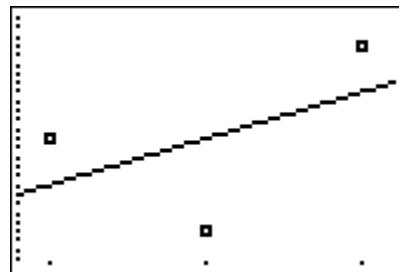
```
LinReg  
y=ax+b  
a=.5  
b=2  
r2=.25  
r=.5
```



$(1,3)$, $(2,2)$, $(3,4)$

Fortunately, we have a test for that! Let's do it at the .05 level of significance. Our null hypothesis is that $\rho=0$.

```
LinReg
y=ax+b
a=.5
b=2
r2=.25
r=.5
```



(1,3), (2,2), (3,4)

```
EDIT CALC TESTS
B†2-PropZInt...
C:X2-Test...
D:X2GOF-Test...
E:2-SampFTest...
FLinRegTTest...
G:LinRegTInt...
H:ANOVA(
```

```
LinRegTTest
Xlist:L1
Ylist:L2
Freq:1
B & P: EQ <0 >0
RegEQ:
Calculate
```

```
LinRegTTest
y=a+bx
B≠0 and ρ≠0
t=.5773502692
P=.666666667
df=1
↓a=2
```

```
LinRegTTest
y=a+bx
B≠0 and ρ≠0
↑b=.5
s=1.224744871
r2=.25
r=.5
```

Also, here is a template to follow when doing a hypothesis test to see if a linear correlation is significantly different from zero.

$$H_0 : \rho = 0$$

$$H_1 : \rho \neq 0$$

Level of Significance: $\alpha = .05$

P-Value: $P = 0.6667$

Decision: Accept H_0 (Fail to reject H_0)

Unfortunately, in this case we cannot reject the null hypothesis, and that also means we shouldn't use the regression equation for predictions. Instead, we use the mean of the y values as a point estimate for predictions.

$$H_0 : \rho = 0$$

$$H_1 : \rho \neq 0$$

$$(1,3), (2,2), (3,4)$$

Level of Significance: $\alpha = .05$

P-Value: $P = 0.6667$

Decision: Accept H_0 (Fail to reject H_0)

$$\bar{y} = \frac{3 + 2 + 4}{3} = \frac{9}{3} = 3$$