

Line of Best Fit aka (Linear Regression)

When data is displayed with a scatter plot, it is often useful to attempt to represent that data with the equation of a straight line for purposes of predicting values that may not be displayed on the plot.

Such a straight line is called the "line of best fit."

It may also be called a "trend" line.

A line of best fit is a straight line that **best** represents the data on a scatter plot. This line may pass through some of the points, none of the points, or all of the points.

approximation
estimation

Predicting:
- If you are looking for values that fall within the plotted values when using the line of best fit, you are **interpolating**.
- If you are looking for values that fall outside the plotted values when using the line of best fit, you are **extrapolating**. **Be careful** when extrapolating. The further away from the plotted values you go, the less reliable is your prediction.

Another term for line of best fit is

Choose:

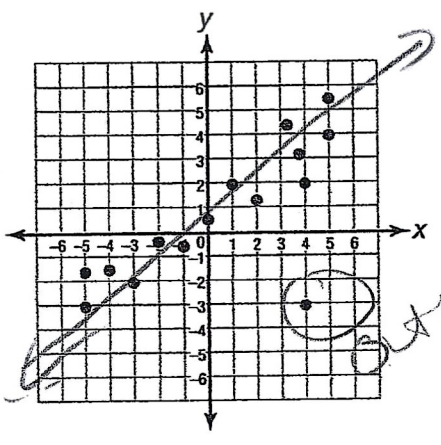
- a) scatter plot
- b) trend line
- c) tangent line
- d) slope



LOBF

Sketch a trend line for each scatter plot.

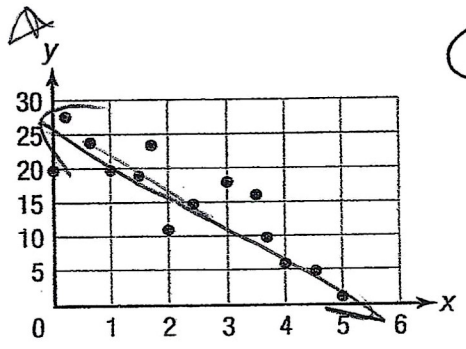
①



positive association / correlation
*line must be in the middle of the dots

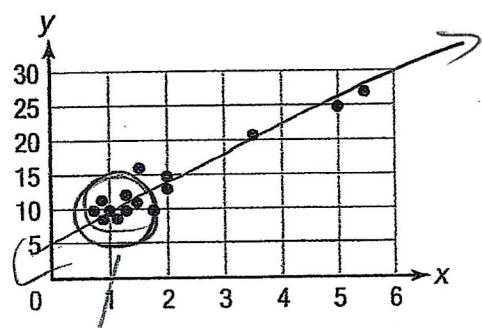
Ask Yourself
How can I draw a line that includes as many points as possible?

②



negative correlation

③

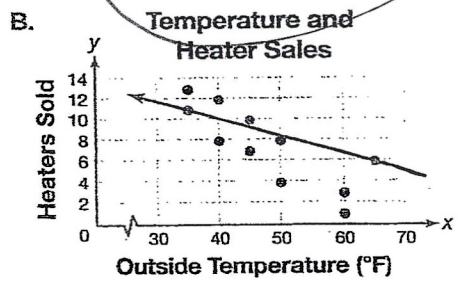
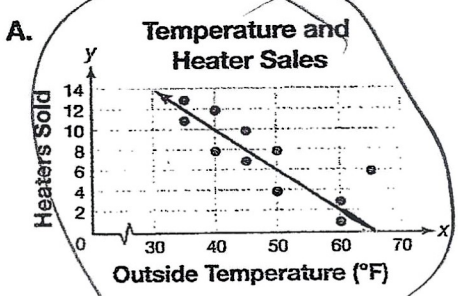


cluster

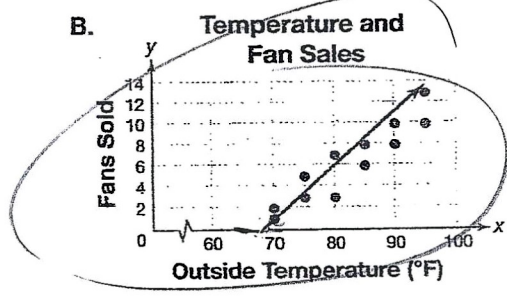
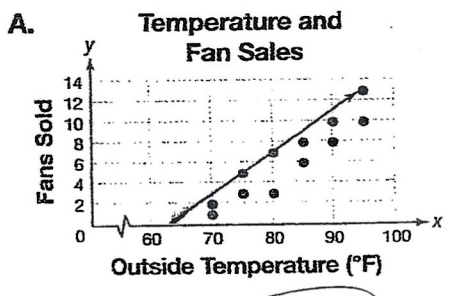
positive correlation

Consider each pair of identical scatter plots. Circle the letter of the plot that shows the better trend line. Explain your choice.

①



②



The line is in the middle of the dots

②

Correlation Coefficients

We know that the graphing calculator can find a "best fit" regression equation that can be used to predict new values. But, how reliable will these prediction be? Is there a way to determine how well our regression equation fits our data?

Yes! There is a way of measuring the "goodness of fit" of the ^{including} best fit line (least squares line), called the correlation coefficient. It is a number between -1 and 1, inclusive, which indicates the measure of linear association between the two variables, and also shows whether the correlation is positive or negative.

↳ how close to a straight line the dots are

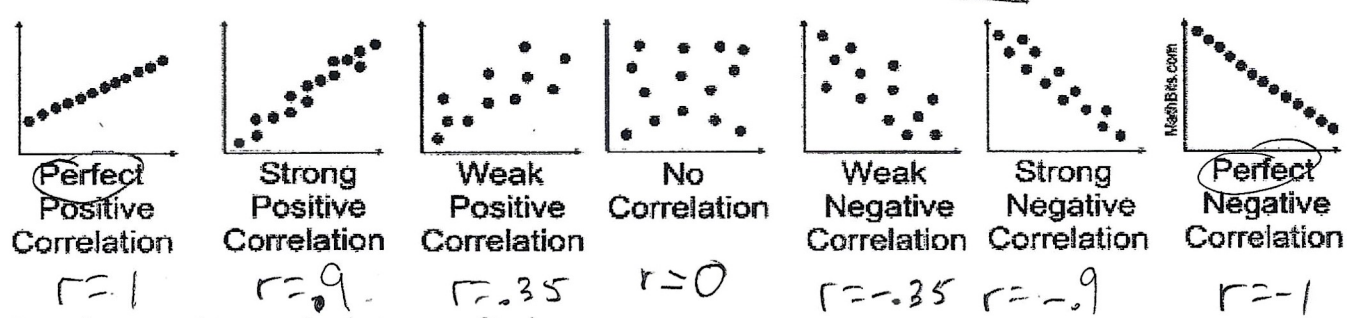
Definition:

A correlation coefficient, designated by r , is a number in the range $-1 \leq r \leq 1$, that indicates how well a regression equation truly represents data being examined.

If r is close to 1 (or -1), the model is considered a "good fit". (Close to a straight line)

- If r is close to 0, the model is "not a good fit". (Not a line)
- If $r = \pm 1$, the model is a "perfect fit" with all data points lying on the line. (perfect line)
- If $r = 0$, there is no linear relationship between the two variables. (No correlation)

A correlation greater than 0.8 is generally described as *strong*, whereas a correlation less than 0.5 is generally described as *weak*.



Using the graphing calculator to find r

Be sure the TI-84+'s "Diagnostics" are turned on. If not, you will not see the r -value.

Calc: (2nd) (0) (↓) Diagnostic on (enter) (enter) OR

When you choose a regression equation on the calculator, the correlation coefficient will be displayed on the screen with the regression equation information (assuming the Diagnostics are turned on).

The linear regression screen shown at the right shows an " r " value of 0.995970141, which implies a strong correlation.

The linear regression equation, in this case, will be a reliable model for future forecasts or predictions.



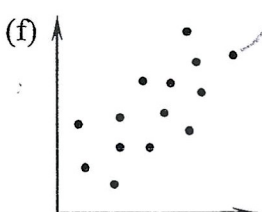
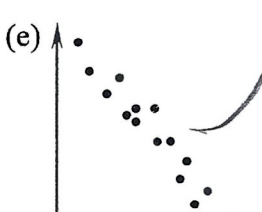
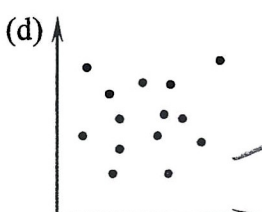
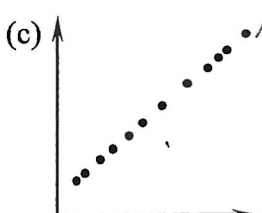
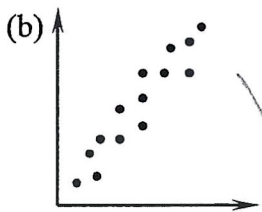
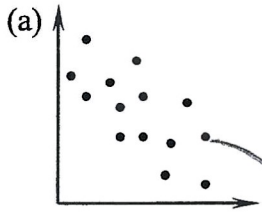
press (mode) + scroll down until you see STAT Diagnostics + (→) (2nd) + press (enter) then exit (6)

Name: _____

Date: _____

QUANTIFYING PREDICTABILITY COMMON CORE ALGEBRA I HOMEWORK

1. Below there are six scatter plots, six correlation coefficients, and six terms. Match the appropriate r -value with the scatter plot it most likely corresponds to. Then match the term you think is most appropriate to the r -value as well (not to the graph).



$r = 1.0$

$r = 0.35$

$r = -0.82$

$r = 0$

$r = -0.56$

$r = 0.93$

Weak Negative

Perfect Positive

Strong Positive

Weak Positive

Strong
Moderate Negative

No Correlation

More Examples

★ By hand: USC the smallest x-value & the largest x-value

7786 -

1) Given the data in the chart below:

x	4	5	6	7
y	8	10	12	14

$y = mx + b$

Determine a line of best fit. (equation)

$B = 0$ (4, 8) (7, 14)

$y = 2x$

$m = \frac{y_2 - y_1}{x_2 - x_1}$

$m = \frac{14 - 8}{7 - 4}$

$m = \frac{6}{3}$

$m = 2$

(4, 8) $m = 2$

$y = mx + b$

$8 = (2)(4) + b$

$8 = 8 + b$

$-8 - 8$
 $0 = b$

$r = 1$
A perfect pos. linear relationship

3)

James uses data that he collected in a science experiment to calculate a line of best fit. He determines the equation of the line to be $y = 7x + 2.25$.

Use this equation to calculate the value of y when $x = 6$.

- A) 15.25
- B) 39.75

- C) 44.25
- D) 42

$y = 7x + 2.25$

$y = 7(6) + 2.25$

$y = 42 + 2.25$

$y = 44.25$

2) The chart below shows the number of minutes studied and the grade received on a test.

Minutes Studied (x)	Test Grade (y)
15	50
40	67
45	75
60	75
70	73
75	89

must put #5 in the calc

Determine a line of best fit for this data.

★ Round to the nearest thousandths

$y = .519x + 45.115$

$r = .903 \rightarrow$ strong pos. close to pos. straight line

(15, 50) + (75, 89)
 x_1, y_1 x_2, y_2

$m = \frac{y_2 - y_1}{x_2 - x_1}$

$m = \frac{89 - 50}{75 - 15}$

$m = \frac{39}{60}$

$m = .65$

(15, 50) $m = .65$

$y = mx + b$

$50 = .65(15) + b$

$50 = 9.75 + b$

$-9.75 - 9.75$

$40.25 = b$

$y = mx + b$
 $m = .65$
 $b = 40.25$

$y = .65x + 40.25$

CALC: 1) STAT 2) Edit put #5 in L1, L2

3) ZOOM 9: STAT5 to graph (make sure stat plot is on)

3) STAT → CALC 4: LinReg (ax + b)

↓ calculate enter to get the equation

Equations are similar but the calculator will give you

Line of best fit

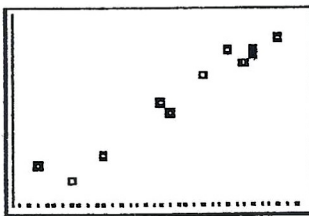
Can we predict the number of total calories based upon the total fat grams?

1. Enter the data in the calculator lists. Place the data in L₁ and L₂.
STAT, #1Edit, type values into the lists

L1	L2	L3	3
9	260		
13	320		
21	420		
30	520		
31	560		
31	550		
34	590		

L3(1)=

2. Prepare a scatter plot of the data. Set up for the scatterplot.
2nd StatPlot - choices shown at right.
Choose ZOOM #9 ZoomStat. Graph shown below.



```

2nd [STAT] Plot2 Plot3
Type: Off
Type: [ ] [ ] [ ]
Xlist: L1
Ylist: L2
Mark: [ ] [ ]
    
```

3. Have the calculator determine the line of best fit.
STAT → CALC #4 LinReg(ax+b)
Include the parameters L₁, L₂, Y₁.
(Y₁ comes from VARS → YVARS, #Function, Y₁)

```

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

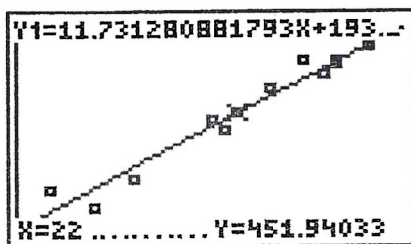
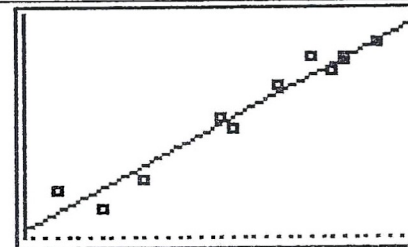
LinReg(ax+b) L1,
L2, Y1

```

LinReg
y=ax+b
a=11.73128088
b=193.8521475
r^2=.9498583012
r=.9746067418
    
```

You now have the values of a and b needed to write the equation of the line of best fit. See values at the right.
 $y = 11.73128088x + 193.8521475$

4. Graph the line of best fit. Simply hit GRAPH. To get a predicted value within the window, hit TRACE, up arrow, and type the desired value.



Question: Predict the total calories based upon 22 grams of fat.

ANS: 451.940 calories

70pt result or 91pt

Y= → VARS/STATS → EQ 1: Reg EQ ZOOM 9: ZoomStat

CALC Steps →

make sure you make a Y1

34 The data given in the table below show some of the results of a study comparing the height of a certain breed of dog, based upon its mass.

must put the #'s in the calc in the order they are in in the table

x	L1	Mass (kg)	4.5	5	4	3.5	5.5	5	5	4	4	6	3.5	5.5
y	L2	Height (cm)	41	40	35	38	43	44	37	39	42	44	31	30

LDBF

Write the linear regression equation for these data, where x is the mass and y is the height. Round all values to the nearest tenth.

STAT \square put #'s into
L1 + L2 then go
to *STAT* *CALC* \square
to get the equation

$$y = 1.9x + 29.8$$

State the value of the correlation coefficient to the nearest tenth, and explain what it indicates.

Model \square \downarrow do
Stat diagnostics
ON to get the
'r'

$$r = 0.3$$

The correlation coefficient suggests a weak positive linear relationship between the mass and height of a certain breed of dog b/c the correlation coefficient is not close to +1.

31 At Mountain Lakes High School, the mathematics and physics scores of nine students were compared as shown in the table below.

must put the #'s in the calc in the order they are in the table

L1 X	Mathematics	55	93	89	60	90	45	64	76	89
L2 Y	Physics	66	89	94	52	84	56	66	73	92

State the correlation coefficient, to the nearest hundredth, for the line of best fit for these data.

$$y = .81x + 15.19$$

$$r = 0.92$$

*mode 2 to stat diagnostics
2nd to get the "r"*

*STAT 1 put #'s into
L1 + L2 + then go
to STAT CALC 4
to get the equation*

Explain what the correlation coefficient means with regard to the context of this situation.

The correlation coefficient suggests a strong positive linear relationship between the math scores and physics scores of nine students because the correlation coefficient is close to +1.

35 Stephen collected data from a travel website. The data included a hotel's distance from Times Square in Manhattan and the cost of a room for one weekend night in August. A table containing these data appears below.

Must put the #'s in the calc in the order they are in the tabs

Distance From Times Square (city blocks) (x)	0	0	1	1	3	4	7	11	14	19
Cost of a Room (dollars) (y)	293	263	244	224	185	170	219	153	136	111

Write the linear regression equation for this data set. Round all values to the nearest hundredth.

$$y = -7.76x + 246.34$$

STAT 1 put #'s in h
L1 + L2 + then go
to STAT CALC 4
to get the equation

State the correlation coefficient for this data set, to the nearest hundredth.

$$r = -0.84$$

mode $\sqrt{\square}$ to stat
diagnostics on to
get the r's

Explain what the sign of the correlation coefficient suggests in the context of the problem.

The correlation coefficient suggests a strong negative linear relationship between the distance from Times Square and the cost of a room. B/c the correlation coefficient is close to -1.

The negative sign suggests a negative correlation and means as the distance from Times Square increases, the cost of a room decreases.

you can
right →
this
BUT you
MUST have →
B/c it's asking
about the sign

↑ MUST have B/c it asks about the sign

- 36 The percentage of students scoring 85 or better on a mathematics final exam and an English final exam during a recent school year for seven schools is shown in the table below.

Must put the #'s in the calc in the order they are in the table

Percentage of Students Scoring 85 or Better	
Mathematics, x_L	English, y_L
27	46
12	28
13	45
10	34
30	56
45	67
20	42

Write the linear regression equation for these data, rounding all values to the nearest hundredth.

$$y = 0.96x + 23.95$$

STAT \square put #'s into x_L & y_L +
 then go to STAT \square \square \square
 to get the equation.

State the correlation coefficient of the linear regression equation, to the nearest hundredth. Explain the meaning of this value in the context of these data.

$$r = 0.92$$

mode \square to
 stat diagnostic
 \square to get the "r"

The correlation coefficient suggests a strong positive linear relationship between the percentage of students scoring 85 or better on a math final exam and an English final exam. The correlation coefficient is close to +1.