Describing Relationships Between Variables (Ch. 4)

Will Landau

Iowa State University

April 11, 2013

Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

s the model Jseful?

Outline

Introduction

Fitting a regression line

Is the model useful?

Is the model valid?

Describing Relationships Between Variables (Ch. 4)

Will Landau

Introduction

Fitting a regression line

s the model 1seful?



Pressing pressures and specimen densities for a ceramic compound

A mixture of Al_2O_3 , polyvinyl alcohol, and water was prepared, dried overnight, crushed, and sieved to obtain 100 mesh size grains. These were pressed into cylinders at pressures from 2,000 psi to 10,000 psi, and cylinder densities were calculated.

(pressure in psi)	y (density in g/cc)
2000.00	2.49
2000.00	2.48
2000.00	2.47
4000.00	2.56
4000.00	2.57
4000.00	2.58
6000.00	2.65
6000.00	2.66
6000.00	2.65
8000.00	2.72
8000.00	2.77
8000.00	2.81
10000.00	2.86
10000.00	2.88
10000.00	2.86

Describing Relationships Between Variables (Ch. 4)

Will Landau

Introduction

Fitting a regression line

s the model useful?

Is the model valid?

>

Scatterplot: ceramics data



Describing Relationships Between Variables (Ch. 4)

Will Landau

Introduction

Fitting a regression line

s the model iseful?



Describing Relationships Between Variables (Ch. 4)

Will Landau

Introduction

Fitting a regression line

s the model Iseful?

Is the model valid?

► The line, y ≈ 2.375 + 4.867 × 10⁻⁵x, is the regression line fit to the data.

Why fit a regression line?

- 1. To predict future values of y based on x.
 - ▶ I.e., a new ceramic under pressure x = 5000 psi should have a density of $2.375 + 4.867 \times 10^{-5} \cdot 5000 = 2.618$ g/cc.
- 2. To characterize the relationship between x and y in terms of strength, direction, and shape.
 - In the ceramics data, density has a strong, positive, linear association with x.
 - ➤ On average, the density increases by 4.867 × 10⁻⁵ g/cc for every increase in pressure of 1 psi.

Describing Relationships Between Variables (Ch. 4)

Will Landau

Introduction

Fitting a regression line

s the model useful?

Outline

Introduction

Fitting a regression line

Is the model useful?

Is the model valid?

Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

s the model 1seful?



Fitting a linear regression line

For a response variable y and a predictor variable x, we declare:

$$y \approx b_0 + b_1 x$$

- ► and then calculate the intercept b₀ and slope b₁ using least squares.
 - ▶ We apply the principle of least squares: that is, the best-fit line is given by minimizing the loss function in terms of b₀ and b₁:

$$S(b_0,b_1)=\sum_{i=1}^n(y_i-\widehat{y}_i)^2$$

• Here,
$$\hat{y}_i = b_0 + b_1 x_i$$

Will Landau

ntroduction

Fitting a regression line

ls the model useful?

Minimize $\sum_{i=1}^{n} (y_i - \hat{y}_i)^2$ to get the line as close as possible to the points.



Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

s the model 1seful?

How to apply least squares to get the regression line

From the principle of least squares, one can derive the normal equations:

$$nb_0 + b_1 \sum_{i=1}^n x_i = \sum_{i=1}^n y_i$$
$$b_0 \sum_{i=1}^n x_i + b_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i$$

• and then solve for b_0 and b_1 :

$$b_1 = rac{\sum (x_i - \overline{x})(y_i - \overline{y})}{\sum (x_i - \overline{x})^2}$$
 $b_0 = \overline{y} - b_1 \overline{x}$

Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

s the model useful?

Example: plastics hardness data

Eight batches of plastic are made. From each batch one test item is molded. At a given time (in hours), it hardness is measured in units (assume freshly-melted plastic has a hardness of 0 units). The following are the 8 measurements and times.

time	hardness									
32.00	230.00		о г							0
72.00	323.00	ts)	35						~	-
64.00	298.00	(uni	<u> </u>					° c	,	
48.00	255.00	Jess	0				~			
16.00	199.00	lardr	- 25(0	0	0			
40.00	248.00	±	8 -0							
80.00	359.00			 20	30	40	50	 60	70	 80
56.00	305.00			20	00	-10		00	10	00
					-	Time	(hou	urs)		

Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

s the model useful?

Fitting the line

- ▼ x = 51
- $ightharpoonup \overline{y} = 277.125$

х	У	$x_i - \overline{x}$	$y_i - \overline{y}$	$(x_i - \overline{x})(y_i - \overline{y})$	$(x_i - \overline{x})^2$
32.00	230.00	-19.00	-47.12	895.38	361.00
72.00	323.00	21.00	45.88	963.38	441.00
64.00	298.00	13.00	20.88	271.38	169.00
48.00	255.00	-3.00	-22.12	66.38	9.00
16.00	199.00	-35.00	-78.12	2734.38	1225.00
40.00	248.00	-11.00	-29.12	320.38	121.00
80.00	359.00	29.00	81.88	2374.38	841.00
56.00	305.00	5.00	27.88	139.38	25.00

▶
$$\sum (x_i - \overline{x})(y_i - \overline{y}) = 895.38 + 963.38 + \cdots 139.38 = 7765$$

▶ $\sum (x_i - \overline{x})^2 = 361 + 441 + \cdots 25 = 3192$
▶ $b_1 = \frac{7765}{3192} = 2.43$
▶ $b_0 = \overline{y} - b_1 \overline{x} = 277.125 - 2.43 \cdot 51 = 153.19$

Describing Relationships Between Variables (Ch. 4)

Will Landau

Introduction

Fitting a regression line

s the model Jseful?

Plot the line to check the fit.



Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

s the model useful?

Interpret the model terms

- ▶ b₁ = 2.43 means that on average, the plastic hardens 2.43 more units for every additional hour it is allowed to harden.
- ▶ b₀ = 153.19 means that at the very beginning of the hardening process (time = 0 hours), the plastics had a hardness of 153.19 on average, IF the model is still correct around time 0.
 - But we know that the plastics were completely molten at the very beginning, with a hardness of 0.
 - Don't extrapolate: i.e., predict y values beyond the range of the x data.

Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

s the model useful?

Checking a fitted line

- 1. Is the model useful?
 - How closely do the points cluster around the line?
 - ▶ How strong is the linear relationship between *x* and *y*?
 - How much variation in y can be explained by the fitted line?
 - ► How well can the fitted line predict future values of *y*?
 - Is the model precise?
- 2. Is the model valid?
 - Should we really be using a straight line to explain y using x, or would some other equation (like a parabola) be better?
 - Does y deviate from the fitted line in some systematic way?
 - Is the model valid?

Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

s the model useful?

Outline

Introduction

Fitting a regression line

Is the model useful?

Is the model valid?

Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

ls the model useful?

Linear correlation: a measure of the usefulness of a fitted line

Linear correlation:

$$r = \frac{\sum (x_i - \overline{x})(y_i - \overline{y})}{\sqrt{\sum (x_i - \overline{x})^2 \sum (y_i - \overline{y})^2}}$$

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

where s_x is the standard deviation of the x_i 's and x_y is the standard deviation of the y_i 's.

Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

Is the model useful?

Facts about linear correlation

• $-1 \le r \le 1$

- r < 0 means a negative slope, r > 0 means a positive slope
- ► High |r| means x and y have a strong linear relationship (high correlation), and low |r| implies a weak linear relationship (low correlation).



Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

Is the model useful?

Correlation in the ceramics data



•
$$s_x = 2927.7002, \ s_y = 0.1438 \ b_1 = 4.867 \cdot 10^{-5}$$

• $r = b_1 \frac{s_x}{s_y} = 4.867 \times 10^{-5} \ \frac{2927.7002}{0.1438} = 0.9911$

© Will Landau

April 11

Describing

Relationships Between Variables (Ch. 4) Will Landau

Is the model useful?

Correlation in the plastics data

▶ <u>x</u> = 51

 $\overline{y} = 277.125$

х	У	$x_i - \overline{x}$	$y_i - \overline{y}$	$(x_i - \overline{x})^2$	$(y_i - \overline{y})^2$	$\Delta x \Delta y$
32.00	230.00	-19.00	-47.12	361.00	2220.77	895.38
72.00	323.00	21.00	45.88	441.00	2104.52	963.38
64.00	298.00	13.00	20.88	169.00	435.77	271.38
48.00	255.00	-3.00	-22.12	9.00	489.52	66.38
16.00	199.00	-35.00	-78.12	1225.00	6103.52	2734.38
40.00	248.00	-11.00	-29.12	121.00	848.27	320.38
80.00	359.00	29.00	81.88	841.00	6703.52	2374.38
56.00	305.00	5.00	27.88	25.00	777.02	139.38

Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

Is the model useful?

CAUTION: the data may be highly correlated even if the *linear* correlation, r, is low.



Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression ine

Is the model useful?

Coefficient of determination

Coefficient of determination: another measure of the usefulness of a fitted line, defined by:

$$R^2 = \frac{\sum (y_i - \overline{y})^2 - \sum (y_i - \widehat{y}_i)^2}{\sum (y_i - \overline{y})^2}$$

where
$$y_i = b_0 + b_1 x_i$$
.

Fortunately,

 $R^2 = r^2$

- Interpretation: R² is the fraction of variation in the response variable (y) explained by the fitted line.
- Ceramics data: $R^2 = r^2 = 0.9911^2 = 0.9823$, so 98.23% of the variation in density is explained by a linear equation in terms of pressure. Hence, the line is useful for predicting density from pressure.
- ▶ Plastics data: $R^2 = r^2 = 0.9796^2 = 0.9596$, so 95.96% of the variation in hardness is explained by a linear equation in terms of time. Hence, so the line is useful for predicting hardness from time.

Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

Is the model useful?

 R^2 measures usefulness (or precision), not validity.

• x and y can have a true linear relationship despite a low R^2



х

Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

Is the model useful?

Is the model valid?



April 11, 2013 23 / 31

Outline

Introduction

Fitting a regression line

Is the model useful?

Is the model valid?

Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

s the model 1seful?

CAUTION: Sometimes, the true relationship between x and y is not linear, despite a high R^2



х

Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

s the model useful?



Residuals: a way to check the validity of a fitted line

Residuals: numbers *e_i* of the form:

$$e_i = y_i - \widehat{y}_i$$

= $y_i - (b_0 + b_1 x_i)$

Instead of:

$$y_i \approx b_0 + b_1 x_i$$

or:

$$\widehat{y}_i = b_0 + b_1 x_i$$

you can now write:

$$y_i = b_0 + b_1 x_i + e_i$$

© Will Landau

Will Landau

ntroduction

Fitting a regression line

s the model .seful?

Is the model valid?

April 11, 2013 26 / 31

What do residuals mean? (Scatterplot: heights and weights of 10 elderly men)



Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

s the model Iseful?

Is the model valid?

 Residuals are the vertical distances between the points and the fitted line.

Residuals: heights and weights of elderly men data

x_i (height in cm)	<i>y</i> i (weight in kg)	ŷi	$e_i = y_i - \widehat{y}_i$
172.70	65.00	74.19	-9.19
165.00	57.00	65.32	-8.32
172.50	77.00	73.96	3.04
182.20	89.00	85.13	3.87
177.60	93.00	79.83	13.17
181.00	73.00	83.75	-10.75
182.50	83.00	85.48	-2.48
182.50	86.00	85.48	0.52
162.80	70.00	62.79	7.21
177.80	83.00	80.06	2.94

Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

ls the model useful?

Plots of residuals



Describing Relationships Between Variables (Ch. 4) Will Landau

ntroduction

Fitting a regression line

s the model iseful?

Is the model valid?

The model fits well since there is no discernible pattern in the residuals when plotted.

© Will Landau

Iowa State University

April 11, 2013 29 / 31

Residual plots and validity

- Left: data that don't fit a line
- Right: the plot of residuals on x
 - The residuals show a nonlinear pattern in the residual plot.
 - Hence, the fitted line is not a valid model.



Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

s the model useful?

More residual plots and patterns

All patterns are bad in plots of residual vs. fitted values, x, time, etc.



When we get to inference, we want to make sure the residuals have a bell-shaped distribution:



This normal QQ plot shows that the residuals are roughly bell-shaped, which is good. Describing Relationships Between Variables (Ch. 4)

Will Landau

ntroduction

Fitting a regression line

s the model useful?