

Regression

Cal State Northridge

Ψ320

Andrew Ainsworth PhD

What is regression?

- How do we predict one variable from another?
- How does one variable change as the other changes?
- Cause and effect

Psy 320 - Cal State Northridge

2

Linear Regression

- A technique we use to predict the most likely score on one variable from those on another variable
- Uses the *nature of the relationship* (i.e. correlation) between two (or more; next chapter) variables to *enhance* your prediction

Psy 320 - Cal State Northridge

3

Linear Regression: Parts

- Y - the variables you are predicting
 - i.e. dependent variable
- X - the variables you are using to predict
 - i.e. independent variable
- \hat{Y} - your predictions (also known as Y')

Psy 320 - Cal State Northridge

4

Why Do We Care?

- We may want to make a prediction.
- More likely, we want to understand the relationship.
 - How fast does CHD mortality rise with a one unit increase in smoking?
 - Note: we speak about predicting, but often don't actually predict.

Psy 320 - Cal State Northridge

5

An Example

- Cigarettes and CHD Mortality from Chapter 9
- Data repeated on next slide
- We want to predict level of CHD mortality in a country averaging 10 cigarettes per day.

Psy 320 - Cal State Northridge

6

The Data

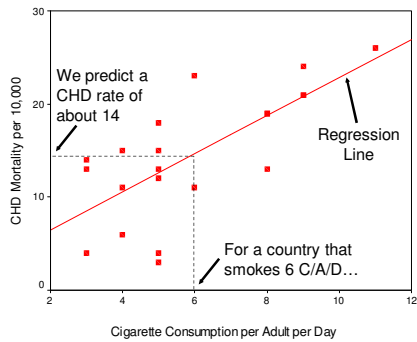
Based on the data we have what would we predict the rate of CHD be in a country that smoked 10 cigarettes on average?

First, we need to establish a prediction of CHD from smoking...

Country	Cigarettes	CHD
1	11	26
2	9	21
3	9	24
4	9	21
5	8	19
6	8	13
7	8	19
8	6	11
9	6	23
10	5	15
11	5	13
12	5	4
13	5	18
14	5	12
15	5	3
16	4	11
17	4	15
18	4	6
19	3	13
20	3	4
21	3	14

Psy 320 - Cal State Northridge

7



Psy 320 - Cal State Northridge

8

Regression Line

- Formula

$$\hat{Y} = bX + a$$

- \hat{Y} = the predicted value of Y (e.g. CHD mortality)
- X = the predictor variable (e.g. average cig./adult/country)

Psy 320 - Cal State Northridge

9

Regression Coefficients

- “Coefficients” are a and b
- b = slope
 - Change in predicted Y for one unit change in X
- a = intercept
 - value of \hat{Y} when $X = 0$

Psy 320 - Cal State Northridge

10

Calculation

- Slope

$$b = \frac{\text{cov}_{XY}}{s_X^2} \text{ or } b = r \left[\frac{s_Y}{s_X} \right]$$

$$\text{or } b = \frac{N \sum XY - \sum X \sum Y}{[N \sum X^2 - (\sum X)^2]}$$

- Intercept

$$a = \bar{Y} - b\bar{X}$$

11

For Our Data

- $\text{Cov}_{XY} = 11.12$
- $s_X^2 = 2.33^2 = 5.447$
- $b = 11.12/5.447 = 2.042$
- $a = 14.524 - 2.042*5.952 = 2.32$
- See SPSS printout on next slide

Answers are not exact due to rounding error and desire to match SPSS.

Psy 320 - Cal State Northridge

12

SPSS Printout

Coefficients^a

Model	Unstandardized Coefficients		Standardized Coefficients	t	Sig.
	B	Std. Error	Beta		
1	(Constant)	2.367	2.941	.805	.431
	Cigarette Consumption per Adult per Day	2.042	.461	.713	4.426

a. Dependent Variable: CHD Mortality per 10,000

Psy 320 - Cal State Northridge

13

Note:

- The values we obtained are shown on printout.
- The intercept is the value in the *B* column labeled “constant”
- The slope is the value in the *B* column labeled by name of predictor variable.

Psy 320 - Cal State Northridge

14

Making a Prediction

- Second, once we know the relationship we can predict
- $$\hat{Y} = bX + a = 2.042X + 2.367$$
- $$\hat{Y} = 2.042 * 10 + 2.367 = 22.787$$
- We predict 22.77 people/10,000 in a country with an average of 10 C/A/D will die of CHD

Psy 320 - Cal State Northridge

15

Accuracy of Prediction

- Finnish smokers smoke 6 C/A/D
- We predict:

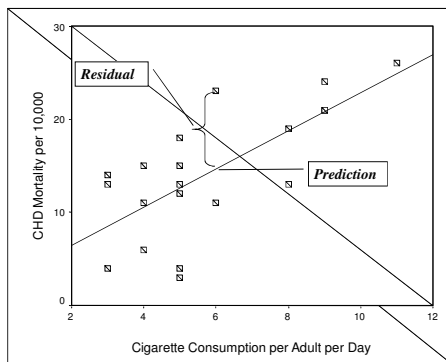
$$\hat{Y} = bX + a = 2.042X + 2.367$$

$$\hat{Y} = 2.042 * 6 + 2.367 = 14.619$$
- They actually have 23 deaths/10,000
- Our error ("residual") =

$$23 - 14.619 = 8.38$$
 – a large error

Psy 320 - Cal State Northridge

16



Psy 320 - Cal State Northridge

17

Residuals

- When we predict \hat{Y} for a given X , we will sometimes be in error.
- $Y - \hat{Y}$ for any X is a an **error of estimate**
- Also known as: a **residual**
- We want to $\Sigma(Y - \hat{Y})$ as small as possible.
- BUT, there are infinitely many lines that can do this.
- Just draw ANY line that goes through the mean of the X and Y values.
- Minimize Errors of Estimate... How?

Psy 320 - Cal State Northridge

18

Minimizing Residuals

- Again, the problem lies with this definition of the mean:

$$\sum (X - \bar{X}) = 0$$

- So, how do we get rid of the 0's?
- Square them.

Psy 320 - Cal State Northridge

19

Regression Line: A Mathematical Definition

- The regression line is the line which when drawn through your data set produces the smallest value of:

$$\sum (Y - \hat{Y})^2$$

- Called the Sum of Squared Residual or $SS_{residual}$
- Regression line is also called a "least squares line."

Psy 320 - Cal State Northridge

20

Summarizing Errors of Prediction

- Residual variance
 - The variability of predicted values

$$s_{Y-\hat{Y}}^2 = \frac{\sum (Y_i - \hat{Y}_i)^2}{N - 2} = \frac{SS_{residual}}{N - 2}$$

Psy 320 - Cal State Northridge

21

Partitioning Variability

- Sums of square deviations

– Total

$$SS_{total} = \sum (Y - \bar{Y})^2$$

– Regression

$$SS_{regression} = \sum (\hat{Y} - \bar{Y})^2$$

– Residual we already covered

$$SS_{residual} = \sum (Y - \hat{Y})^2$$

- $SS_{total} = SS_{regression} + SS_{residual}$

Psy 320 - Cal State Northridge

25

Partitioning Variability

- Degrees of freedom

– Total

- $df_{total} = N - 1$

– Regression

- $df_{regression} = \text{number of predictors}$

– Residual

- $df_{residual} = df_{total} - df_{regression}$

- $df_{total} = df_{regression} + df_{residual}$

Psy 320 - Cal State Northridge

26

Partitioning Variability

- Variance (or Mean Square)

– Total Variance

- $s^2_{total} = SS_{total} / df_{total}$

– Regression Variance

- $s^2_{regression} = SS_{regression} / df_{regression}$

– Residual Variance

- $s^2_{residual} = SS_{residual} / df_{residual}$

Psy 320 - Cal State Northridge

27

r^2 for our example

- $r = .713$
- $r^2 = .713^2 = .508$
- or $r^2 = \frac{SS_{\text{regression}}}{SS_Y} = \frac{454.307}{895.247} = .507$
- Approximately 50% in variability of incidence of CHD mortality is associated with variability in smoking.

Psy 320 - Cal State Northridge

31

Coefficient of Alienation

- It is defined as $1 - r^2$ or

$$1 - r^2 = \frac{SS_{\text{residual}}}{SS_Y}$$

- Example

$$1 - .508 = .492$$

$$1 - r^2 = \frac{SS_{\text{residual}}}{SS_Y} = \frac{440.757}{895.247} = .492$$

Psy 320 - Cal State Northridge

32

r^2 , SS and $s_{Y-\hat{Y}}$

- $r^2 * SS_{\text{total}} = SS_{\text{regression}}$
- $(1 - r^2) * SS_{\text{total}} = SS_{\text{residual}}$
- We can also use r^2 to calculate the standard error of estimate as:

$$s_{Y-\hat{Y}} = s_y \sqrt{(1-r^2) \left(\frac{N-1}{N-2} \right)} = 6.690 * \sqrt{(.492) \left(\frac{20}{19} \right)} = 4.816$$

Psy 320 - Cal State Northridge

33

Hypothesis Testing

- Test for overall model
- Null hypotheses
 - $b = 0$
 - $a = 0$
 - population correlation (ρ) = 0
- We saw how to test the last one in Chapter 9.

Psy 320 - Cal State Northridge

34

Testing Overall Model

- We can test for the overall prediction of the model by forming the ratio:
$$\frac{s_{regression}^2}{s_{residual}^2} = F \text{ statistic}$$
- If the calculated F value is larger than a tabled value (Table D.3 $\alpha = .05$ or Table D.4 $\alpha = .01$) we have a significant prediction

Psy 320 - Cal State Northridge

35

Testing Overall Model

- Example
$$\frac{s_{regression}^2}{s_{residual}^2} = \frac{454.307}{23.198} = 19.594$$
- Table D.3 – F critical is found using 2 things $df_{regression}$ (numerator) and $df_{residual}$ (denominator)
- Table D.3 our $F_{crit}(1, 19) = 4.38$
- $19.594 > 4.38$, significant overall
- Should all sound familiar...

Psy 320 - Cal State Northridge

36

SPSS output

Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.713 ^a	.508	.482	4.81640

a. Predictors: (Constant), CIGARETT

ANOVA^b

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	454.482	1	454.482	19.592	.000 ^a
	Residual	440.757	19	23.198		
	Total	895.238	20			

a. Predictors: (Constant), CIGARETT

b. Dependent Variable: CHD

Psy 320 - Cal State Northridge

37

Testing Slope and Intercept

- The regression coefficients can be tested for significance
- Each coefficient divided by it's standard error equals a t value that can also be looked up in a table (Table D.6)
- Each coefficient is tested against 0

Psy 320 - Cal State Northridge

38

Testing Slope

- With only 1 predictor, the standard error for the slope is:

$$se_b = \frac{s_{Y-\hat{Y}}}{s_X \sqrt{N-1}}$$

- For our Example:

$$se_b = \frac{4.816}{2.334\sqrt{21-1}} = \frac{4.816}{10.438} = .461$$

Psy 320 - Cal State Northridge

39

Testing Slope and Intercept

- With only 1 predictor, the standard error for the intercept is:

$$se_a = s_{y-\hat{y}} \sqrt{\frac{1}{N} + \frac{\bar{X}^2}{\sum (X - \bar{X})^2}} = s_{y-\hat{y}} \sqrt{\frac{1}{N} + \frac{\bar{X}^2}{SS_X}}$$

- For our Example:

$$= 4.816 \sqrt{\frac{1}{21} + \frac{5.952^2}{2.334^2 * 20}} = 4.816 \sqrt{\frac{1}{21} + \frac{35.426}{108.951}} = 2.94$$

Psy 320 - Cal State Northridge

40

Testing Slope

- These are given in computer printout as a *t* test.

Coefficients ^a						
Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.
		B	Std. Error	Beta		
1	(Constant)	2.367	2.941		.805	.431
	Cigarette Consumption per Adult per Day	2.042	.461	.713	4.426	.000

a. Dependent Variable: CHD Mortality per 10,000

Psy 320 - Cal State Northridge

41

Testing

- The *t* values in the second from right column are tests on slope and intercept.
- The associated *p* values are next to them.
- The slope is significantly different from zero, but not the intercept.
- Why do we care?

Psy 320 - Cal State Northridge

42

Testing

- What does it mean if slope is not significant?
 - How does that relate to test on r ?
- What if the intercept is not significant?
- Does significant slope mean we predict quite well?

Psy 320 - Cal State Northridge

43
