

SOS3003
**Applied data analysis for social
science**
Collected lecture notes spring 2010

Erling Berge
Department of sociology and political science
NTNU

Spring 2010

© Erling Berge 2010

1

Required reading SOS3003

- Hamilton, Lawrence C. 1992. *Regression with graphics*. Belmont: Duxbury. Ch 1-8
- Hamilton, Lawrence C. 2008. *A Low-Tech Guide to Causal Modelling*.
<http://pubpages.unh.edu/~lch/causal2.pdf>
- Allison, Paul D. 2002. *Missing Data*. Sage University Paper: QASS 136. London: Sage.

Spring 2010

© Erling Berge 2010

2

Background to the sciences

- In the history of civilization there are 2 unrivalled accelerators:
 - The invention of writing about 5-6000 years ago
 - The invention of the scientific method for separating facts from fantasy about 5-600 years ago
- There is no topic more important to learn than the basics of the scientific method
- That does not mean that it is not – at times – rather boring

Spring 2010

© Erling Berge 2010

3

Basics of causal beliefs

- First: doubt what you believe is a causal link until you can give good valid reasons justifying your belief
- Second: there are usually many types of good valid reasons for believing in a particular causal link, for example, scientific consensus
 - If the overwhelming majority of certified scientists says that human activities **contribute** to global warming, then we are justified believing that by changing our activities we could contribute less to global warming
- Third: random conjunctures (“correlation”) are not good valid reasons for believing in a causal link

Spring 2010

© Erling Berge 2010

4

Causal mechanism

- Elster 2007 *Explaining Social Behaviour*:
- ”mechanisms are frequently occurring and easily recognizable causal patterns that are triggered under generally unknown conditions or with indeterminate consequences” (page 36)
- Also sometimes limited to “causal chains”

Spring 2010

© Erling Berge 2010

5

Causal correlations

- This class will focus on how to distinguish between random conjunctures and that which might be a valid causal correlation
- That which might be a valid causal correlation will need a *causal mechanism* explaining how the cause can produce the effect before we have a valid reason to believe in the causal link

Spring 2010

© Erling Berge 2010

6

Primacy of theory

- To say it more bluntly: If you do not have a believable theory (and this may well start as a fantasy) then regression techniques will tell you nothing even if you find a seemingly non-random correlation
- But without a valid and believable empirical analysis any believable fantasy will remain just that: a fantasy (assuming you cannot find other valid verifications)

Spring 2010

© Erling Berge 2010

7

Types of causal mechanisms I

- Structural causation
 - A Social structure has causal impacts that are not well understood. In a framework of methodological individualism one may say that it limits and orders the options that actors can choose from. Hence, variables such as age, sex, and place of living can be used as proxies for poorly understood causal factors.
 - Budget constraints (time and income constraints) have the same character. They limit and orders the options that actors can choose from. However, they enter the model more through the way the dependent variable is constructed, and the kind of link function (linear or logistic) used to mediate between observations and dependent variable.

Spring 2010

© Erling Berge 2010

8

Types of causal mechanisms II

- Individual causation
 - Preferences (norms, values, attitudes) may be difficult to observe in detail but are assumed to be present
 - Resources (income/ capital, education/ human capital, access to networks/ social capital) are usually measured extensively even if unevenly. Resources represent budget constraints
 - Perception of opportunities will often depend on position in social structure
 - Beliefs about resources and opportunities are important. They may be based on both fact and fiction

Spring 2010

© Erling Berge 2010

9

Preliminaries

- Prerequisite: SOS1002 or equivalent
- Goal: to read critically research articles using quantitative methods in your field of interest
- Required reading ... see above
- Term paper: this is part of the examination and evaluation procedure

Spring 2010

© Erling Berge 2010

10

Goals for the class

- The goal is that each of you shall be able to read critically research articles discussing quantitative data. This means
 - You are to know the pitfalls so you can evaluate the validity of an article
- You are to learn how to perform straightforward analyses of co-variation in "quantitative" and "qualitative" data (nominal scale data in regression analysis), and in particular:
 - Also here you have to demonstrate that you know the pitfalls

Spring 2010

© Erling Berge 2010

11

Lecture I

Basics of what you are assumed to know

- The following is basically repeating known stuff
- Variable distributions
 - Ringdal Ch 12 p251-270
 - Hamilton Ch 1 p1-23
- Bivariat regression
 - Ringdal Ch 17-18 p361-387
 - Hamilton Ch 2 p29-59

Spring 2010

© Erling Berge 2010

12

Some basic concepts

- Cause
- Model
- Population
- Sample
- Variable: level of measurement
- Variable: measure of centralization
- Variable: measure of dispersion

Spring 2010

© Erling Berge 2010

13

Data analysis

- Descriptive use of data
 - Developing classifications
- Analytical use of data
 - Describe phenomena that cannot be observed directly (inference)
 - Causal links between directly eller indirectly observable phenomena (theory or model development)

Spring 2010

© Erling Berge 2010

14

Causal analysis: from co-variation to causal connection

- From colloquial speech to theory
 - Fantasy and intuition, established science tradition
- From theory to model
 - Operationalisation
- From observation to generalisation
 - Causal analysis

Spring 2010

© Erling Berge 2010

15

THREE BASIC DIVISIONS

<u>Observed</u>		<u>Real interest</u>
THEORY/ MODEL	-	REALITY
SAMPLE	-	POPULATION
CO-VARIATION	-	CAUSE

On the one hand we have what we are able to observe and record, on the other hand, we have what we would like to discuss and know more about

Spring 2010

© Erling Berge 2010

16

Basic sources of error

- Errors in theory / model
 - Model specification: valid conclusions require a correct (true) model
- Errors in the sample
 - Selection bias
- Measurement problems
 - Missing cases and measurement errors
 - Validity og reliability
- Multiple comparisons
 - Conclusions are valid only for the sample

Spring 2010

© Erling Berge 2010

17

From population to sample

- POPULATION (all units)

Simple random sampling

- SAMPLE (selected units)

Spring 2010

© Erling Berge 2010

18

Unit and variable

- A unit, as a carrier of data, will be contextually defined
 - SUPER - UNIT: e.g. the local community
 - UNIT: e.g. household
 - SUB - UNIT: e.g. person
- Variable: empirical concept used to characterize units under investigation. Each unit is characterized by being given a variable value

Spring 2010

© Erling Berge 2010

19

Data matrix and level of measurement

- Matrix defined by Units * Variables
 - A table presenting the characteristics of all investigated units ordered so that all variable values are listed in the same sequence for all units
- Level of measurement for a variable
 - Nominal scale *classification
 - Ordinal scale *classification and rank
 - Interval scale *classification, rank and distance
 - Ratio scale *classification, rank, distance and absolute zero

Spring 2010

© Erling Berge 2010

20

Variable analysis

- Description
 - Central tendency and dispersion
 - Form of distribution
 - Frequency distributions and histograms
- Comparing distributions
 - Quantile plots
 - Box plots

Spring 2010

© Erling Berge 2010

21

VARIABLE: central tendency

- **Mean**
sum of all values of the variable for all units divided by the number of units $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$
- **MEDIAN**
The variable value in an ordered distribution that has half the units on each side $\sum_{i=1}^n (X_i - \bar{X})$
- **MODUS**
The typical value. The value in a distribution that has the highest frequency $\sum_{i=1}^n (X_i - \bar{X})^2 < \sum_{i=1}^n (X_i - C)^2$
if
 $C \neq \bar{X}$

Spring 2010

© Erling Berge 2010

22

VARIABLE: measures of dispersion I

- **MODAL PERCENTAGE**
- The percentage of units with value like the mode
- **RANGE OF VARIATION**
- The difference between highest and lowest value in an ordered distribution
- **QUARTILE DIFFERENCE**
- Range of variation of the 50% of units closest to the median ($Q_3 - Q_1$)
- **MAD - Median Absolute Deviation**
- Median of the absolute value of the difference between median and observed value:
 - $MAD(x_i) = \text{median } |x_i - \text{median}(x_i)|$

Spring 2010

© Erling Berge 2010

23

VARIABLE: measures of dispersion II

- **STANDARD DEVIATION**
 - Square root of mean squared deviation from the mean
 - $s_y = \sqrt{[(\sum_i (Y_i - \tilde{Y})^2)/(n - 1)]}$
 - **MEAN DEVIATION**
 - Mean of the absolute value of the deviation from the mean
 - **VARIANCE**
 - Standard deviation squared:
 - $s_y^2 = (\sum_i (Y_i - \tilde{Y})^2)/(n - 1)$
- (nb: here \tilde{Y} is the mean of Y)

Spring 2010

© Erling Berge 2010

24

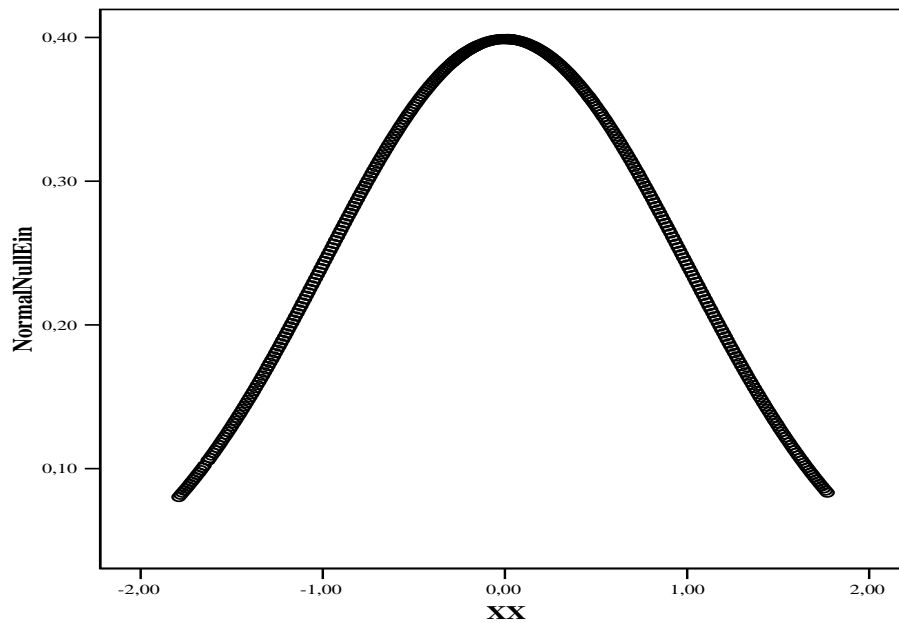
Variable: form of distribution I

- Symmetrical distributions
- Skewed distributions
 - "Heavy" and "Light" tails
- Normal distributions
 - Are not "normal"
 - Are unambiguously determined by mean and variance (μ og σ^2)

Spring 2010

© Erling Berge 2010

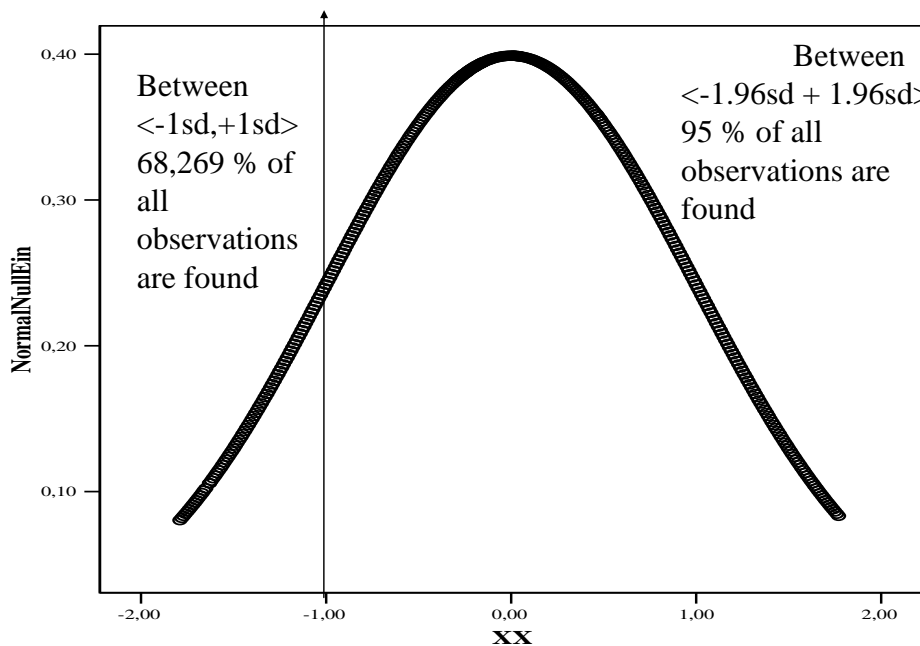
25



Spring 2010

© Erling Berge 2010

26



Spring 2010

© Erling Berge 2010

27

Skewed distributions

- Positively skewed has $\tilde{Y} > Md$
- Negatively skewed has $\tilde{Y} < Md$
- Symmetric distributions has $\tilde{Y} \approx Md$

- nb: here $\tilde{Y} = \text{mean of } Y$

Spring 2010

© Erling Berge 2010

28

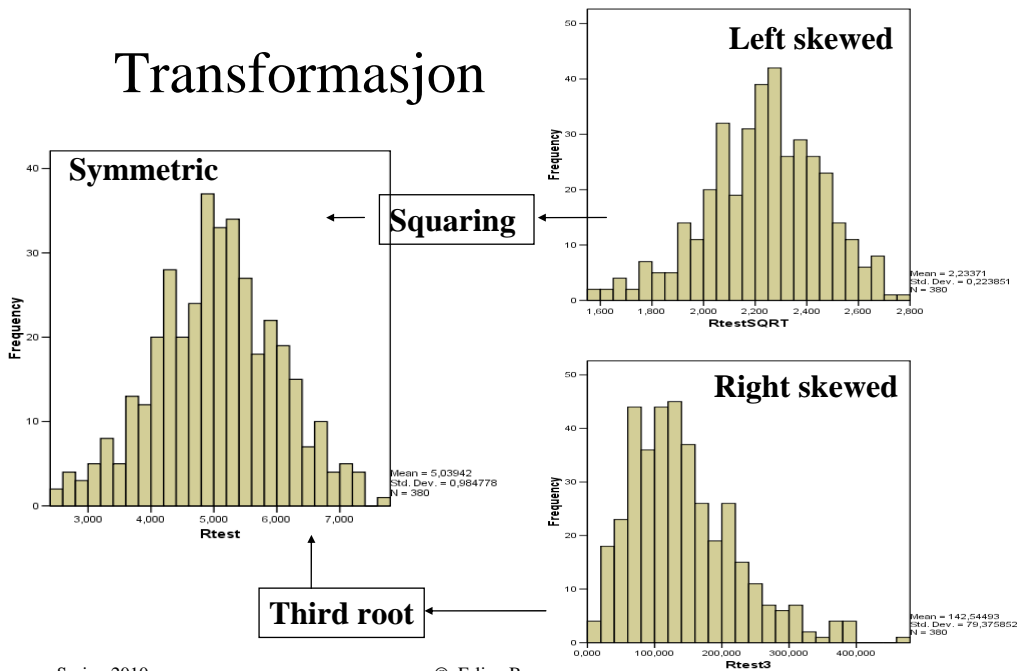
Symmetric distributions

- Median and IQR are resistant against the impact of extreme values
- Mean and standard deviation are not
- In the normal distribution (ND) $s_y \approx IQR/1.35$
- If we in a symmetric distribution find
 - $s_y > IQR/1.35$ then the tails are heavier than in the ND
 - $s_y < IQR/1.35$ then the tails are lighter than in the ND
 - $s_y \approx IQR/1.35$ then the tails are about similar to the ND

Spring 2010

© Erling Berge 2010

29



Spring 2010

© Erling Berge

Variable: analyzing distributions I

- Box plot
 - The box is constructed based on the quartile values Q_1 og Q_3 . Observations within $\langle Q_1, Q_3 \rangle$ are in the box-
 - Adjacent large values are defined as those outside the box but inside $Q_3 + 1.5 \cdot \text{IQR}$ or $Q_1 - 1.5 \cdot \text{IQR}$
 - Outliers (seriously extreme values) are those outside of $Q_3 + 1.5 \cdot \text{IQR}$ or $Q_1 - 1.5 \cdot \text{IQR}$

Spring 2010

© Erling Berge 2010

31

Variables: analyzing distributions II

- Quantiles is a generalisation of quartiles and percentiles
- Quantile values are variable values that correspond to particular fractions of the total sample or observed data, e.g.
 - Median is 0.5 quantile (or 50% percentile)
 - Lower quartile is 0.25 quantile
 - 10% percentile is 0.1 quantile ...

Spring 2010

© Erling Berge 2010

32

Variables: analyzing distributions III

- Quantile plots
 - Quantile values against value of variable
 - The Lorentz curve is a special case of this (it gives us the Gini-index)
- Quantile-Normal plot
 - Plot of quantile values on one variable against quantile values of a Normal distribution with the same mean and standard deviation

Spring 2010

© Erling Berge 2010

33

Example: Randaberg 1985

- Questionnaire: (the number of decares land you own / 10 da = 1 ha)

Q: ANTALL DEKAR GRUNN DU
eier: _____

(Number of decares you own: ____)

Spring 2010

© Erling Berge 2010

34

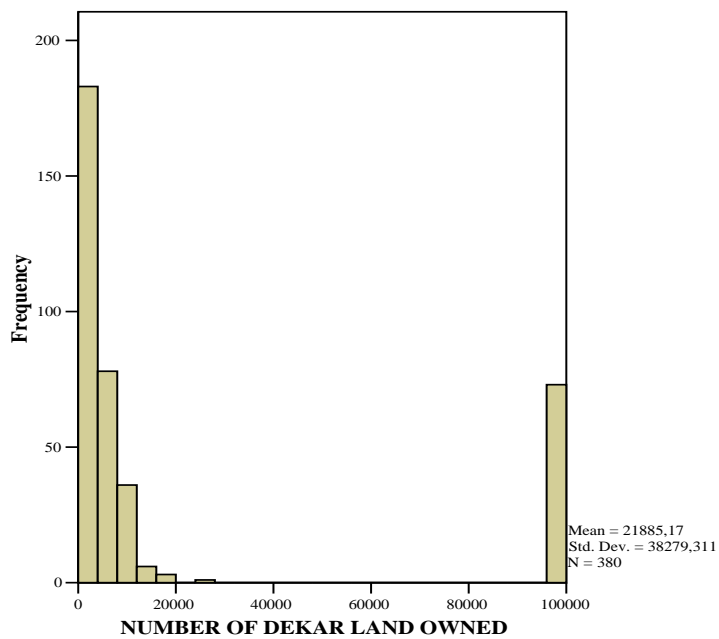
NUMBER OF DEKARE LAND OWNED

	NUMBER OF DEKARE LAND OWNED	Valid N (listwise)
N	380	380
Minimum	0	
Maximum	99900	
Mean	21885.17	
Std. Deviation	38279.311	

Spring 2010

© Erling Berge 2010

35



Spring 2010

© Erling Berge 2010

36

XAreaOwned (NUMBER OF DEKARE LAND OWNED)

	XAreaOwned	Valid N (listwise)
N	307	307
Minimum	.00	
Maximum	25000.00	
Mean	3334.4104	
Std. Deviation	4201.54943	

Spring 2010

© Erling Berge 2010

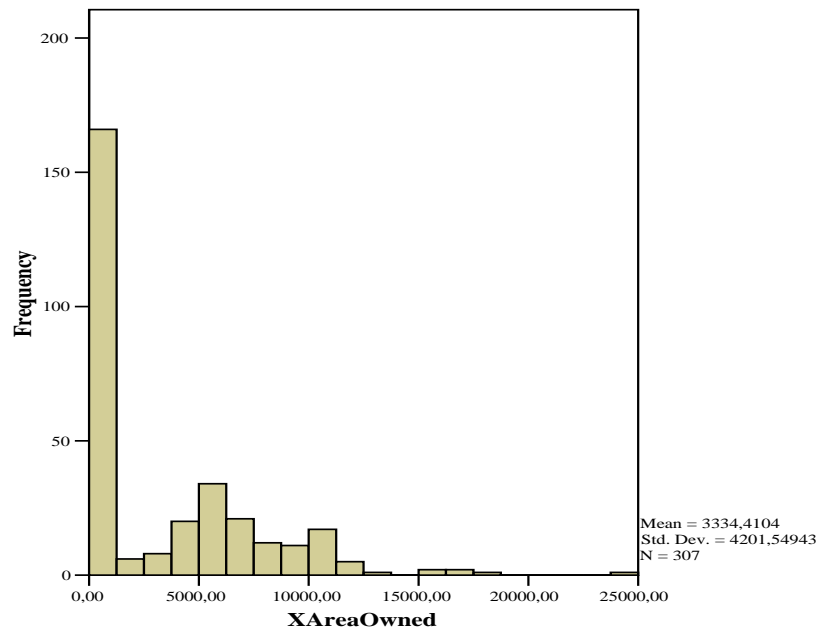
37

		XAreaOwned	Valid N (listwise)
N	Statistic	307	307
Range	Statistic	25000.00	
Minimum	Statistic	.00	
Maximum	Statistic	25000.00	
Sum	Statistic	1023664.00	
Mean	Statistic	3334.4104	
	Std. Error	239.79509	
Std. Deviation	Statistic	4201.54943	
Variance	Statistic	17653017.596	
Skewness	Statistic	1.352	
	Std. Error	.139	
Kurtosis	Statistic	2.194	
	Std. Error	.277	

Spring 2010

© Erling Berge 2010

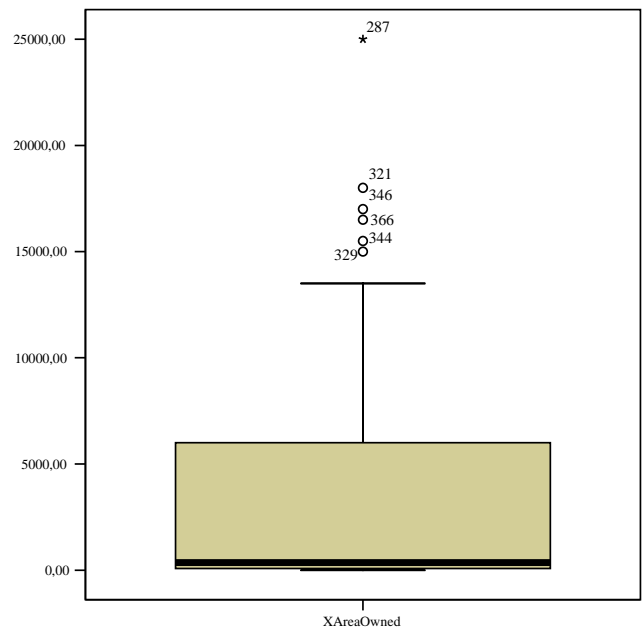
38



Spring 2010

© Erling Berge 2010

39



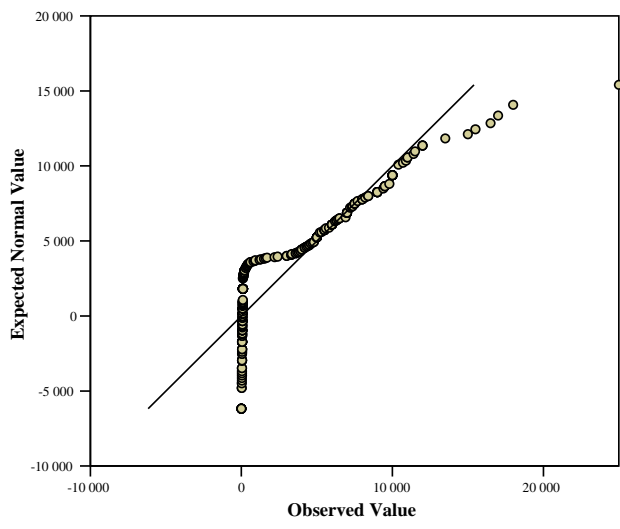
Spring 2010

© Erling Berge 2010

40

NB
Figures
from SPSS
are mirrors
of figures
in
Hamilton

Normal Q-Q Plot of XAreaOwned

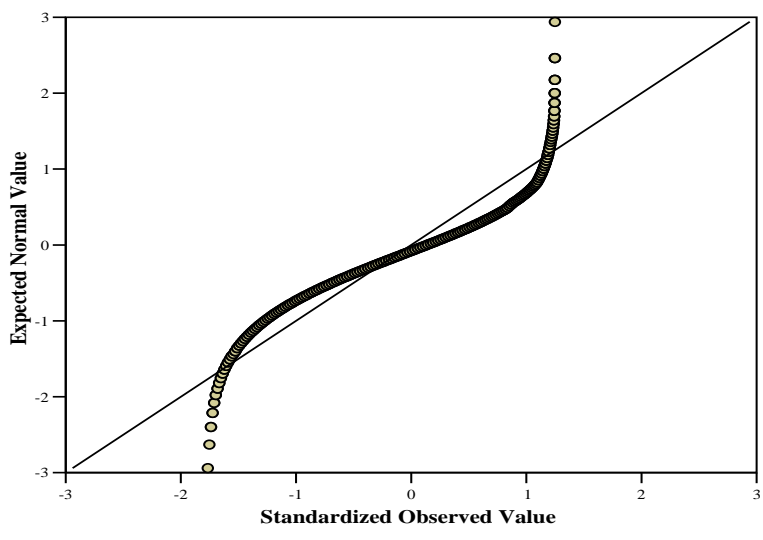


Spring 2010

© Erling Berge 2010

41

Normal Q-Q Plot of NormalNullEin



Spring 2010

© Erling Berge 2010

42

Questionnaire:

- **Hvor viktig er det at myndighetene kontrollerer og regulerer bruken av arealer gjennom for eksempel kontroll av**

- av tomtetildelinger (kommunal formidl.)

1 2 3 4 5 6 7 8

- avkjørsler fra hus til vei

1 2 3 4 5 6 7 8

- kjøp og salg av landbrukseiendommer

1 2 3 4 5 6 7 8

Spring 2010

© Erling Berge 2010

43

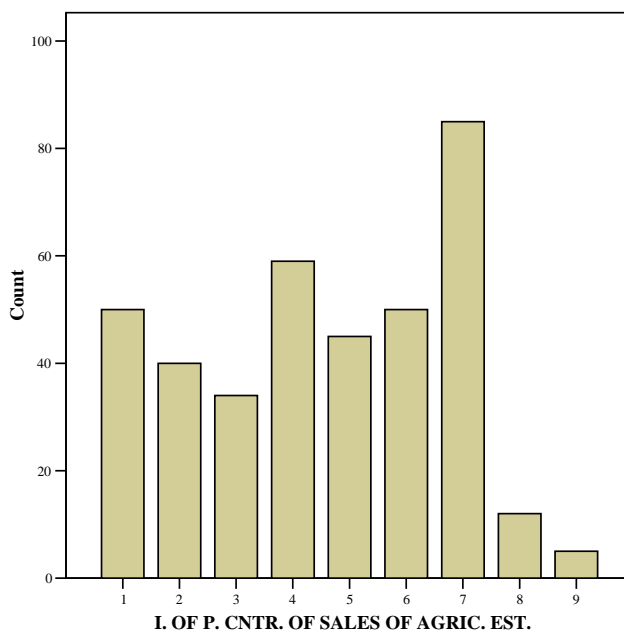
Importance of public control of sales of agric. estates

	Frequency	Percent	Valid Percent	Cumulative Percent
Valid 1	50	13.2	13.2	13.2
2	40	10.5	10.5	23.7
3	34	8.9	8.9	32.6
4	59	15.5	15.5	48.2
5	45	11.8	11.8	60.0
6	50	13.2	13.2	73.2
7	85	22.4	22.4	95.5
8	12	3.2	3.2	98.7
9	5	1.3	1.3	100.0
Total	380	100.0	100.0	

Spring 2010

© Erling Berge 2010

44



Spring 2010

© Erling Berge 2010

45

Questionnaire: coding

Ved utfylling: sett ring rundt et tall som synes å gi passelig uttrykk for viktigheten når 1 betyr svært lite viktig og 7 særdeles viktig, eller sett et kryss inne i parantesene () som står bak svaret du velger

På noen spørsmål kan du krysse av flere svar

	lykkes dårlig/ lite viktig						lykkes godt/ svært viktig	vet ikke
Kodeverdi	1	2	3	4	5	6	7	8

Dei som ikkje kryssar av noko svar vert koda 9 (ie. missing)

Spring 2010

© Erling Berge 2010

46

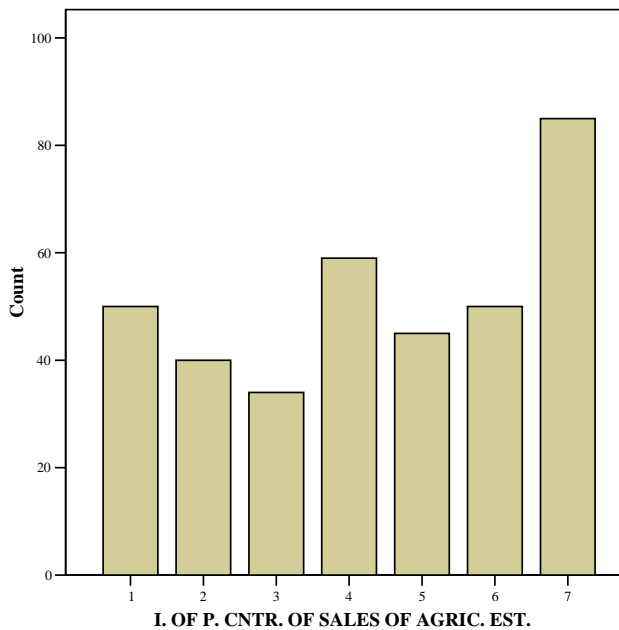
I. OF P. CNTR. OF SALES OF AGRIC. EST.

		Frequency	Percent	Valid Percent	Cumulative Percent
Valid	1	50	13.2	13.8	13.8
	2	40	10.5	11.0	24.8
	3	34	8.9	9.4	34.2
	4	59	15.5	16.3	50.4
	5	45	11.8	12.4	62.8
	6	50	13.2	13.8	76.6
	7	85	22.4	23.4	100.0
	Total	363	95.5	100.0	
Missing	8	12	3.2		
	9	5	1.3		
	Total	17	4.5		
Total		380	100.0		

Spring 2010

© Erling Berge 2010

47



Spring 2010

© Erling Berge 2010

48

		I. OF P. CNTR. OF SALES OF AGRIC. EST.	Y regressed on ControlSalesAgricEstate Valid N (listwise)
N	Statistic	380	363
Range	Statistic	8	6.00
Minimum	Statistic	1	1.00
Maximum	Statistic	9	7.00
Sum	Statistic	1729	1588.00
Mean	Statistic	4.55	4.3747
	Std. Error	.114	.11045
Std. Deviation	Statistic	2.213	2.10435
Variance	Statistic	4.897	4.428
Skewness	Statistic	-.171	-.234
	Std. Error	.125	.128
Kurtosis	Statistic	-1.148	-1.267
	Std. Error	.250	.255

Spring 2010

© Erling Berge 2010

49

Distributions with or without missing?

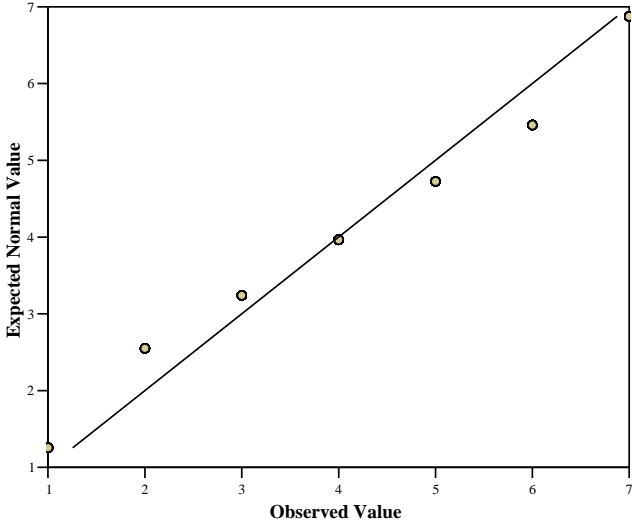
- What difference do the 17 missing observations make in the
 - Quantile-Normal plot?
 - Box plot?

Spring 2010

© Erling Berge 2010

50

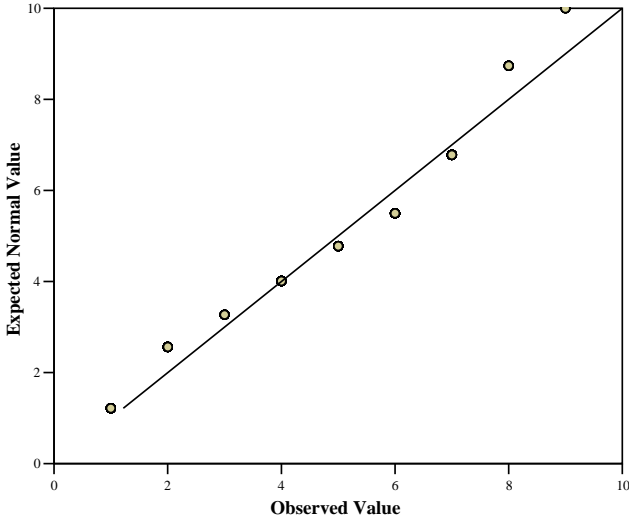
Normal Q-Q Plot of I. OF P. CNTR. OF SALES OF AGRIC. EST.



Spring 2010

51

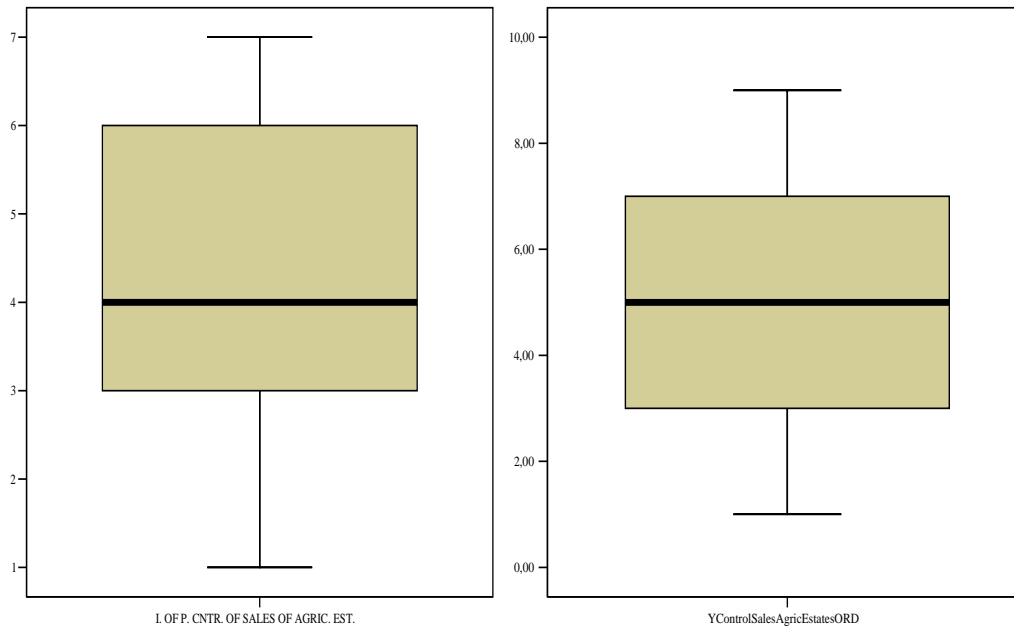
Normal Q-Q Plot of I. OF P. CNTR. OF SALES OF AGRIC. EST.



Spring 2010

© Erling Berge 2010

52



Spring 2010

© Erling Berge 2010

53

Data collection and data quality I

- Questions – techniques for asking questions will not be discussed
- Sample
 - From sampling to final data matrix: selection of cases, refusing to participate, and missing answers on questions
- Variables: Data on cases collected as variable values for each case
- Statistics: Data on samples collected as statistics (Norwegian: “observatorer” where values are estimated for each sample)
- Statistics is also the science of assessing the quality of each statistic

Spring 2010

© Erling Berge 2010

54

Data collection and data quality II

- What is important for the quality of the data?
 - Validity of questions asked and reliability of the procedures used.
 - Selection bias: A possible causal link between missing observations and the topic studied
- What can be done if data are faulty?
 - Not much!

Spring 2010

© Erling Berge 2010

55

Writing up a model

- Defining the elements of the model
 - Variables, error term, population, and sample
- Defining the relations among the elements of the model
 - Sampling procedure, time sequence of the events and observations, the functions that links the elements into an equation
- Specification of the assumptions stipulated to be true in order to use a particular method of estimation
 - Relationship to substance theory (specification requirement)
 - Distributional characteristics of the error term

Spring 2010

© Erling Berge 2010

56

Elements of a model

- Population (who or what are we interested in?)
- Sample (simple random sample or exact specification of how each case came into the sample)
- Variables (characteristics of cases relevant to the questions we are investigating)
- Error terms (the sum of impacts from all other causes than those explicitly included)

Spring 2010

© Erling Berge 2010

57

Relations among elements of a model

- Sampling: biased sample?
- Time sequence of events and observations (important to aid causal theory)
- Co-variation (genuine vs spurious co-variation)
 - Conclusions about causal impacts require genuine co-variation
- Equations and functions

Spring 2010

© Erling Berge 2010

58

Bivariat Regression: Modelling a population

- $Y_i = \beta_0 + \beta_1 x_{1i} + \varepsilon_i$
- $i=1, \dots, n$ $n = \#$ cases in the population
- Y and X must be defined unambiguously, and Y must be interval scale (or ratio scale) in ordinary regression (OLS regression)

Spring 2010

© Erling Berge 2010

59

Bivariat Regression: Modelling a sample

- $Y_i = b_0 + b_1 x_{1i} + e_i$
- $i=1, \dots, n$ $n = \#$ cases in the sample
- e_i is usually called the residual (not the error term as in the population model)
- Y and X must be defined unambiguously, and Y must be interval scale (or ratio scale) in ordinary regression (OLS regression)

Spring 2010

© Erling Berge 2010

60

An example of a bad regression

- The example following contains a series of errors. If you present such a regression in your term paper you will fail
- Your task is to identify the errors as quickly as possible and then never do the same
- Clue: look again at the distributions of the variables above

Spring 2010

© Erling Berge 2010

61

Importance of public control of sales of agric. Estates Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.047(a)	.002	.000	2.213

a Predictors: (Constant), NUMBER OF DEKAR LAND OWNED

Spring 2010

© Erling Berge 2010

62

Importance of public control of sales of agric. Estates ANOVA(b)

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	4.145	1	4.145	.846	.358(a)
	Residual	1851.905	378	4.899		
	Total	1856.050	379			

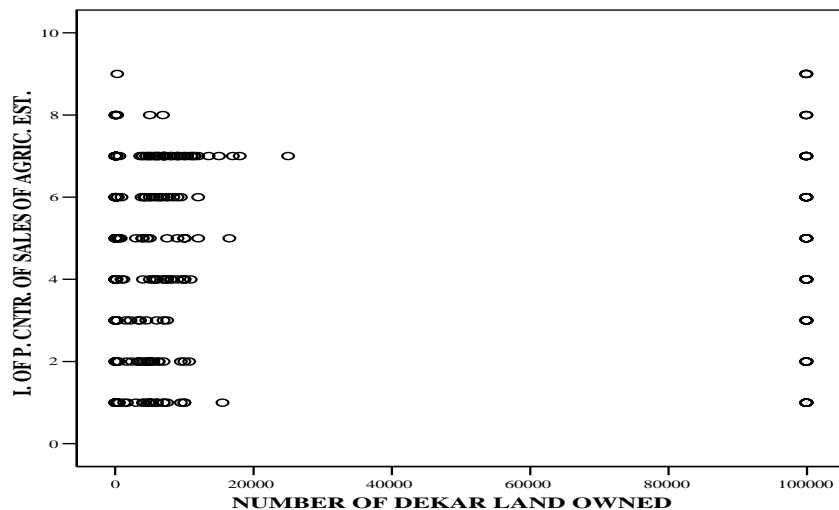
a Predictors: (Constant), NUMBER OF DEKAR LAND OWNED
 b Dependent Variable: I. OF P. CNTR. OF SALES OF AGRIC. EST.

Importance of public control of sales of agric. Estates Coefficients (a)

Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.
		B	Std. Error	Beta		
1	(Constant)	4.610	.131		35.233	.000
	NUMBER OF DEKAR LAND OWNED	.000	.000	-.047	-.920	.358

a Dependent Variable: I. OF P. CNTR. OF SALES OF AGRIC. EST.

Scatterplot

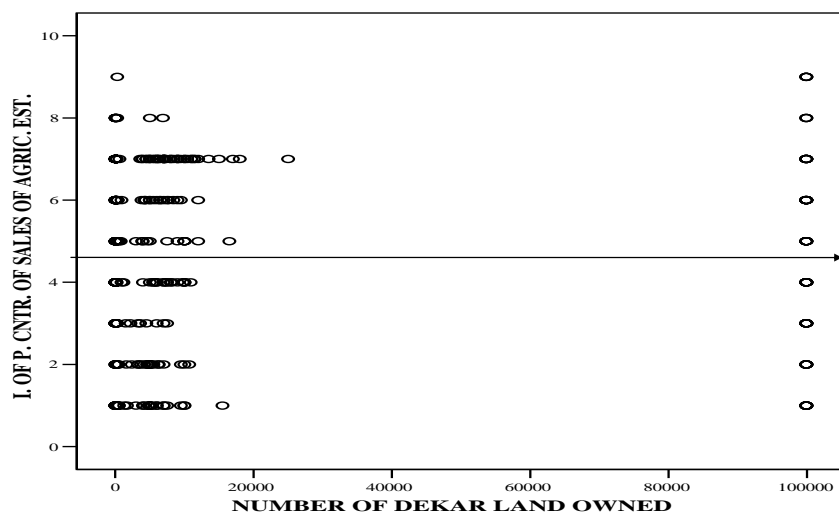


Spring 2010

© Erling Berge 2010

65

Scatterplot with regression line



Spring 2010

© Erling Berge 2010

66

Assumptions needed for the use of OLS to estimate a regression model

OLS: ordinary least squares (minste kvadrat metoden)

Requirements for OLS estimation of a regression model can shortly be summed up as

- We assume that the linear model is correct (true) with independent, and identical normally distributed error terms ("normal i.i.d. errors")

Spring 2010

© Erling Berge 2010

67

Estimation method: OLS

- Model $Y_i = b_0 + b_1 x_{1i} + e_i$

The observed error (the residual) is

- $e_i = (Y_i - b_0 - b_1 x_{1i})$

Squared and summed residual

- $\sum_i (e_i)^2 = \sum_i (Y_i - b_0 - b_1 x_{1i})^2$

Find b_0 and b_1 that minimizes the squared sum

Spring 2010

© Erling Berge 2010

68

Relationship sample - population (1)

- A new mathematical operator: $E[\varpi]$ meaning the expected value of $[\varpi]$ where ϖ stands for some expression containing at least one variable or unknown parameter, e.g.
- $E[Y_i] = E[b_0 + b_1 x_{1i} + e_i]$

$$= \beta_0 + \beta_1 x_{1i}$$
- Note in particular that in our model
 - $E[b_0] = \beta_0$
 - $E[b_1] = \beta_1$
 - $E[e_i] = \varepsilon_i$

Spring 2010

© Erling Berge 2010

69

Relationship sample – population (2)

- Relationship sample - population is determined by the characteristics that the error term has been given in the sampling and observation procedure
- In a simple random sample with complete observation
 $E[\varepsilon_i] = 0$ for all i , and
 $\text{var}[\varepsilon_i] = \sigma^2$ for all i

NB: $\text{var}(\varpi)$ is a new mathematical operator meaning "the procedure that will find the variance of some algebraic expression " ϖ "

Spring 2010

© Erling Berge 2010

70

Complete observation

- Make it possible to make a completely specified model. This means that all variables that causally affects the phenomenon we study (Y) have been observed, and are included in the model equation
- This is practically impossible. Therefore the error term will include also unobserved factors affecting (Y)

Spring 2010

© Erling Berge 2010

71

Testing hypotheses I

	In reality H_0 is true	In reality H_0 is untrue
We conclude that H_0 is true	Our method gives the correct answer with probability $1 - \alpha$	<u>Error of type II</u> (probability $1 - \beta$)
We conclude that H_0 is untrue	<u>Error of type I</u> The test level α is the probability of errors of type I	Our method gives the correct answer with probability β (= power of the test)

Spring 2010

© Erling Berge 2010

72

Testing hypotheses II

- A test is always constructed based on the assumption that H_0 is true
- The construction leads to a
 - **Test statistic**
- The test statistic is constructed so that it has a known probability distribution, usually called a
 - **Sampling distribution**

Spring 2010

© Erling Berge 2010

73

Testing hypotheses III

- It is easier to construct tests based on the assumption that it is true that a particular test statistic is zero, [H_0 stating that a parameter is 0], than any particular other value
- In regression this means that we assume a particular parameter $\beta = 0$ in order to evaluate how large the probability is for this to be true given the sample we have observed

Spring 2010

© Erling Berge 2010

74

The p-value of a test

- The p-value of a test gives the estimated probability for observing the values we have in our sample or values that are even more in accord with a conclusion that **H_0 is untrue**; assuming that our sample is a simple random sample from the population where H_0 in reality is true
- Very low p-values suggest that we cannot believe that H_0 is true. We conclude that $\beta \neq 0$

Spring 2010

© Erling Berge 2010

75

T-test and F-test

- Sums of squares
 - $TSS = ESS + RSS$
 - $RSS = \sum_i (e_i)^2 = \sum_i (Y_i - \hat{Y}_i)^2$ distance observed- estimated value
 - $ESS = \sum_i (\hat{Y}_i - \tilde{Y})^2$ distance estimated value - mean
 - $TSS = \sum_i (Y_i - \tilde{Y})^2$ distance observed value – mean
- Test statistic
 - $t = (\mathbf{b} - \boldsymbol{\beta}) / SE_b$ SE = standard error
 - $F = [ESS/(K-1)]/[RSS/(n-K)]$ K = number of model parameters

Spring 2010

© Erling Berge 2010

76

Confidence interval for β

- Picking a t_α - value from the table of the t-distribution with $n-K$ degrees of freedom makes the interval

$$\langle b - t_\alpha(\text{SE}_b), b + t_\alpha(\text{SE}_b) \rangle$$
into a two-tailed test giving a probability of α for committing error of type I
- This means that $b - t_\alpha(\text{SE}_b) \leq \beta \leq b + t_\alpha(\text{SE}_b)$ with probability $1 - \alpha$

Spring 2010

© Erling Berge 2010

77

Coefficient of determination

Coefficient of determination:

- $R^2 = \text{ESS}/\text{TSS} = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$
 - Tells us how large a fraction of the variation around the mean we can "explain by" (attribute to) the variables included in the regression ($\hat{Y}_i =$ predicted y)
- In bi-variate regression the coefficient of determination equals the coefficient of correlation:

$$r_{yu}^2 = s_{yu} / s_y s_u$$
- Co-variance: $s_{yu} = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})(U_i - \bar{U})$

Spring 2010

© Erling Berge 2010

78

Detecting problems in a regression

- Take a second look at the example presented above where
 - Y = IMPORTANCE OF PUBLIC CONTROL OF SALES OF AGRICULURAL ESTATES
 - X = NUMBER OF DEKAR LAND OWNED
 - $Y_i = b_0 + b_1 x_{1i} + e_i$

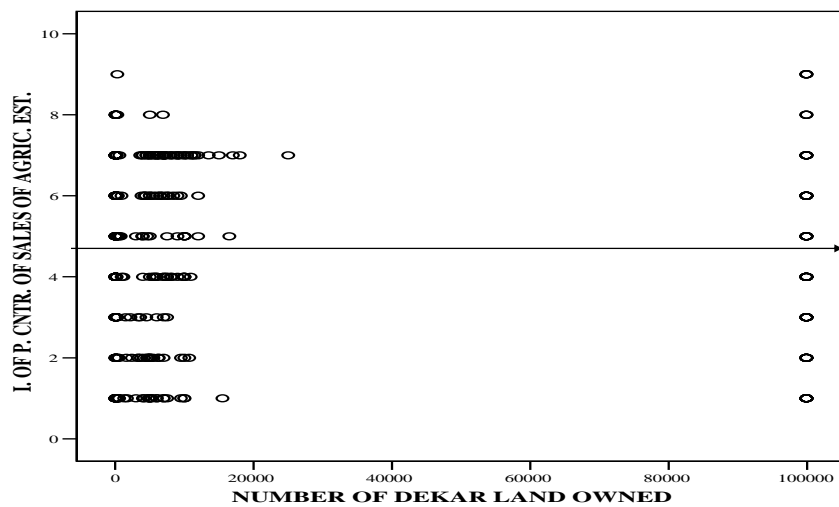
What was the problem in this example?

Spring 2010

© Erling Berge 2010

79

What is wrong in this scatter plot with regression line?



Spring 2010

© Erling Berge 2010

80

In general: what can possibly cause problems?

- Omitted variables (specification error)
- Non-linear relationships (specification error)
- Non-constant error term (heteroskedastisitet)
- Correlation among error terms (autocorrelation)
- Non-normal error terms

Spring 2010

© Erling Berge 2010

81

Problems also from

- High correlations among included variables (multicollinearity)
- High correlation between an included and an excluded variable (spurious correlation in the model)
- Cases with high influence
- Measurement errors

Spring 2010

© Erling Berge 2010

82

Non-normal errors:

- Regression **DO NOT need assumptions about the distribution of variables**
- But to test hypotheses about the parameters we need to assume that the **error terms are normally distributed** with the same mean and variance
- **If the model is correct** (true) and n (number of cases) is large the central limit theorem demonstrates that the error terms approach the normal distribution
- **But usually a model will be erroneously or incompletely specified.** Hence we need to inspect and test residuals (observed error term) to see if they actually are normally distributed

Spring 2010

© Erling Berge 2010

83

Residual analysis

- This is the most important starting point for diagnosing a regression analysis

Useful tools:

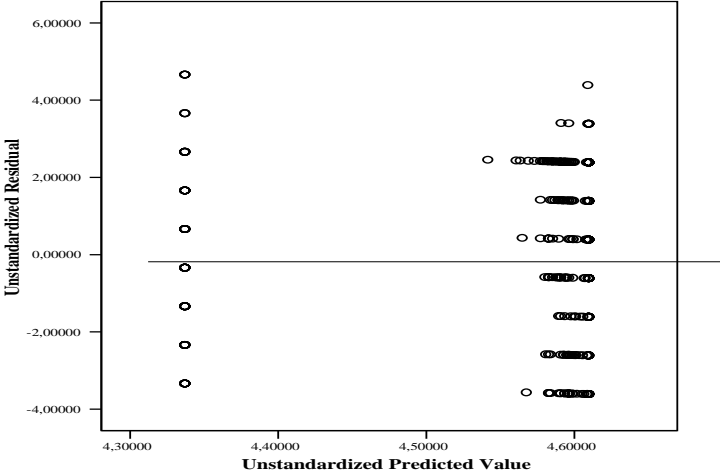
- Scatter plot
- Plot of residual against predicted value
- Histogram
- Box plot
- Symmetry plot
- Quantil-normal plot

Spring 2010

© Erling Berge 2010

84

What went wrong? (1) residual-predicted value plot



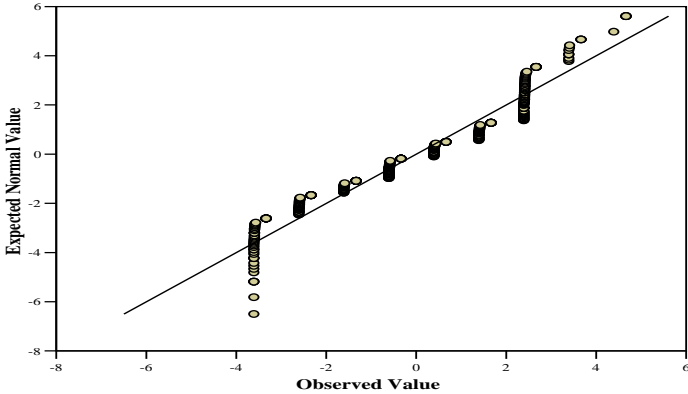
Spring 2010

© Erling Berge 2010

85

What went wrong? (1) normal-quantile plot

Normal Q-Q Plot of Unstandardized Residual



Spring 2010

© Erling Berge 2010

86

Power transformations

May solve problems related to

- Curvilinearity in the model
- Outliers
- Influential cases
- Non-constant variance of the error term (heteroscedasticity)
- Non-normal error term

NB: Power transformations are used to solve a problem. If you do not have a problem do not solve it.

Spring 2010

© Erling Berge 2010

87

Power transformations (see H:17-22)

Y^* : read

“transformed Y ”

(transforming Y to Y^*)

- $Y^* = Y^q \quad q > 0$
- $Y^* = \ln[Y] \quad q = 0$
- $Y^* = - [Y^q] \quad q < 0$

Inverse

transformation

(transforming Y^* to Y)

- $Y = [Y^*]^{1/q} \quad q > 0$
- $Y = \exp[Y^*] \quad q = 0$
- $Y = [- Y^*]^{1/q} \quad q < 0$

Spring 2010

© Erling Berge 2010

88

Power transformations: consequences

- $X^* = X^q$
 - $q > 1$ increases the weight of the right hand tail relative to the left hand tail
 - $q = 1$ produces identity
 - $q < 1$ reduces the weight of the right hand tail relative to the left hand tail
- If $Y^* = \ln(Y)$ the regression coefficient of an interval scale variable X can be interpreted as % change in Y per unit change in X

E.g. if $\ln(Y) = b_0 + b_1 x + e$

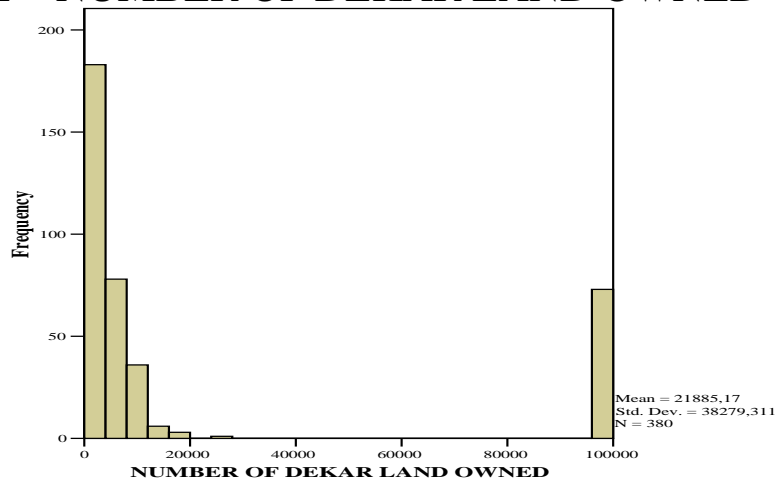
b_1 can be interpreted as % change in Y pr unit change in X

Spring 2010

© Erling Berge 2010

89

Point of departure $X = \text{NUMBER OF DEKAR LAND OWNED}$

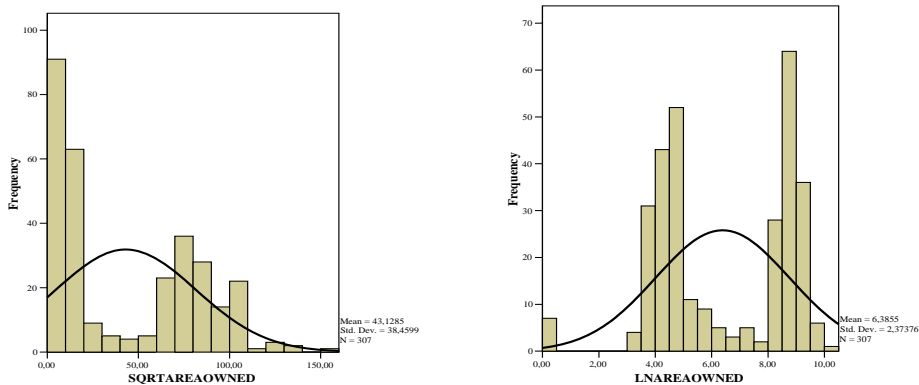


Spring 2010

© Erling Berge 2010

90

Power transformed X = NUMBER OF DEKAR LAND OWNED



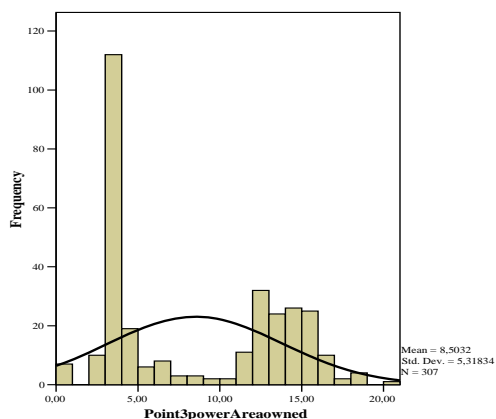
SQRT=square root of areaowned – LN= natural logarithm of (areaowned+1)

Spring 2010

© Erling Berge 2010

91

Power transformed X = NUMBER OF DEKAR LAND OWNED



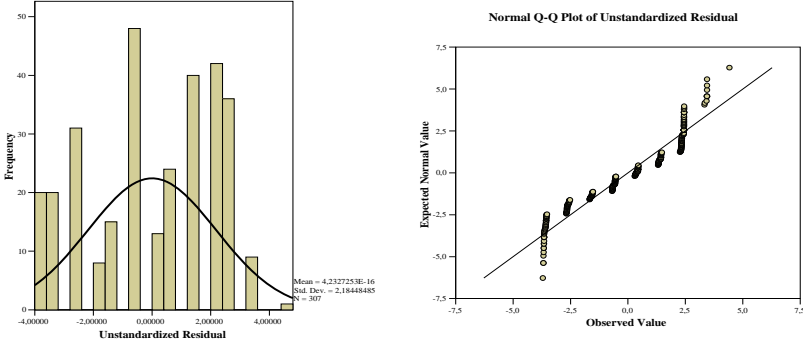
Point3power = 0,3 power of areaowned

Spring 2010

© Erling Berge 2010

92

Does power transformation help?



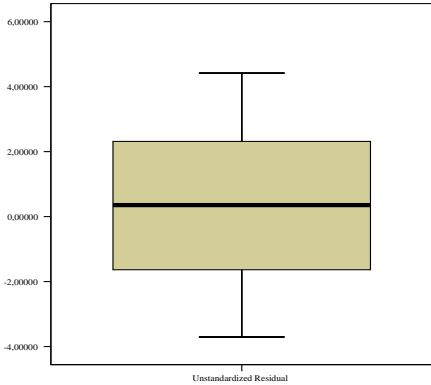
0.3 power-transformation gives lighter tails and no outliers

Spring 2010

© Erling Berge 2010

93

Box plot of the residual shows approximate symmetry and no outliers



Spring 2010

© Erling Berge 2010

94

Curvilinear regression

- The example above used the variable "Point3powerAreaowned", or 0.3 power of number of dekar land owned:

- $\text{Point3powerAreaowned} = (\text{NUMBER OF DEKAR LAND OWNED})^{0.3}$

The model estimated is thus

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

$$y_i = b_0 + b_1 (\text{Point3powerAreaowned}_i) + e_i$$

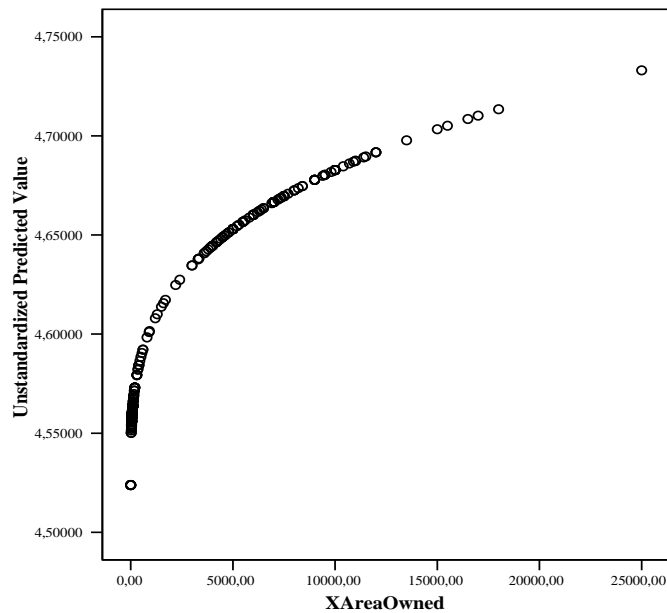
$$\hat{y}_i = 4.524 + 0.010 * (\text{NUMBER OF DEKAR LAND OWNED}_i)^{0.3}$$

Spring 2010

© Erling Berge 2010

95

Use of power transformed variables means that the regression is curvilinear



Spring 2010

© Erling Berge 2010

96

Summary

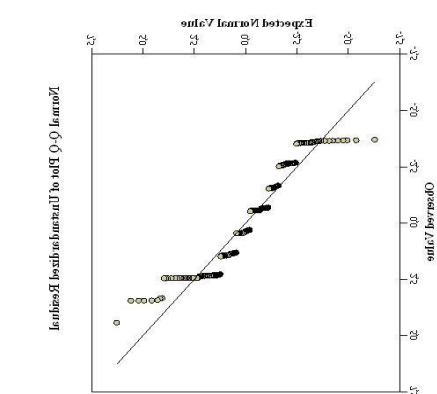
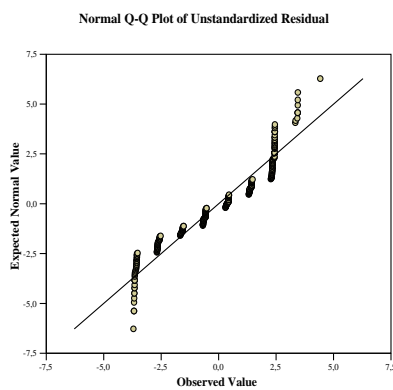
- In bivariate regression the OLS method finds the "best" LINE or CURVE in a two dimensional scatter plot
- Scatter-plot and analysis of residuals are tools for diagnosing problems in the regression
- Transformations are a general tool helping to mitigate several types of problems, such as
 - Curvilinearity
 - Heteroscedasticity
 - Non-normal distributions of residuals
 - Case with too high influence
- Regression with transformed variables are always curvilinear. Results can most easily be interpreted by means of graphs

Spring 2010

© Erling Berge 2010

97

SPSS printout vs the book (see p16)



Spring 2010

© Erling Berge 2010

98

Reading printout from SPSS (1)

Descriptive Statistics					Mean	Std. Deviation ¹	N ²		
I. OF P. CNTR. OF SALES OF AGRIC. EST.					4.61	2.185	307		
Point3powerAreaowned					8.5032	5.31834	307		

Model	R	R Square ³	Adjusted R Square ⁴	Std. Error of the Estimate ⁵	Change Statistics				
					R Square Change	F Change	df1	df2	Sig. F Change
1	.024(a)	.001	-.003	2.188	.001	.182	1	305	.670

a Predictors: (Constant), Point3powerAreaowned

b Dependent Variable: I. OF P. CNTR. OF SALES OF AGRIC. EST.

Footnotes to the table above (1)

1. Standard deviation of the mean
2. Number of cases used in the analysis
3. Coefficient of determination
4. The adjusted coefficient of determination (see Hamilton page 41)
5. Standard deviation of the residual

$$s_e = \text{SQRT} (\text{RSS}/(n-K)),$$

where SQRT (*) = square root of (*)

Reading printout from SPSS (2)

Model		Sum of Squares ³	df	Mean Square	F ¹	Sig. ²
1	Regression	.870	1	.870	.182	.670(a)
	Residual	1460.224	305	4.788		
	Total	1461.094	306			

•Sums of squares: $TSS = ESS + RSS$

• $RSS = \sum_i (e_i)^2 = \sum_i (Y_i - \hat{Y}_i)^2$: sum of squared (distance observed – estimated value)

•Mean Square = RSS / df For RSS it is known that $df=n-K$

K equals number of parameters estimated in the model (b_0 og b_1)

Here we have $n=307$ and $K=2$, hence $Df = 305$

Spring 2010

© Erling Berge 2010

101

Footnotes to the table above (2)

1. F-statistic for the null hypothesis $\beta_1 = 0$ (see Hamilton p45)
2. p-value of the F-statistic: the probability of finding a F-value this large or larger assuming that the null hypothesis is correct
3. Sums of squares
 1. $TSS = ESS + RSS$
 2. $RSS = \sum_i (e_i)^2 = \sum_i (Y_i - \hat{Y}_i)^2$ distance observed value – estimated value
 3. $ESS = \sum_i (\hat{Y}_i - \bar{Y})^2$ distance estimated value – mean
 4. $TSS = \sum_i (Y_i - \bar{Y})^2$ distance observed value – mean

Spring 2010

© Erling Berge 2010

102

Reading printout from SPSS (3)

M o d e l		Unstandardized Coefficients		Standardized Coefficients	t ⁴	Sig. ⁵	95% Confidence Interval for B	
		B ¹	Std. Error ²	Beta ³			Lower Bound	Upper Bound
1	(Constant)	4.524	.236		19.187	.000	4.060	4.988
	Point3-powerA re-owned	.010	.024	.024	.426	.670	-.036	.056

Spring 2010

© Erling Berge 2010

103

Footnotes to the table above (3)

1. Estimates of the regression coefficients b_0 og b_1
2. Standard error of the estimates of b_0 og b_1
3. Standardized regression coefficients: $b_1^{st} = b_1 * (s_x / s_y)$ see Hamilton pp38-40
4. t-statistic for the null hypothesis $\beta_1 = 0$ (see Hamilton p44)
5. p-value of the t-statistic: the probability of finding a t-value this large or larger assuming that the null hypothesis is correct

Spring 2010

© Erling Berge 2010

104

Multiple regression

Hamilton Ch 3 p65-101

Spring 2010

© Erling Berge 2010

105

Recall from first lecture:

Bivariate regression: Modelling a sample

- $Y_i = b_0 + b_1 x_{1i} + e_i$
 - $i=1, \dots, n$ $n = \#$ cases in the sample
- e_i is usually called the residual (**not** the error term as in the population model)
- Y and X must be defined unambiguously, and Y must be interval scale (or ratio scale) in ordinary regression (OLS regression)

Spring 2010

© Erling Berge 2010

106

Recall from first lecture:
Bivariate regression: Modelling a population

- $Y_i = \beta_0 + \beta_1 x_{1i} + \varepsilon_i$
 - $i=1, \dots, n$ $n = \#$ cases in the population
 - ε_i is the error term for case no i
- Y and X must be defined unambiguously, and Y must be interval scale (or ratio scale) in ordinary regression (OLS regression)

Spring 2010

© Erling Berge 2010

107

Summary on bivariate regression

- In bivariate regression the OLS method finds the "best" LINE or CURVE in a two dimensional scatter plot
- Best is defined as the "a" and "b" that minimizes the sum of squared deviations between the line/ curve and observed variable values
- **Scatter-plot and analysis of residuals** are tools for diagnosing problems in the regression
- Transformation (by powers) is a general tool helping to mitigate several types of problems, such as
 - Curvilinearity
 - Heteroscedasticity
 - Non-normal distributions of residuals
 - Cases with too high influence
- Regression with (power) transformed variables are always curvilinear. Results can most easily be interpreted by means of graphs

Spring 2010

© Erling Berge 2010

108

Multiple regression: model (1)

- The goal of multiple regression is to find the net impact of one variable controlled for the impact of all other variables
- Let K = number of parameters in the model (this means that $K-1$ is the number of variables)
- Then the population model can be written
- $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \dots + \beta_{K-1} x_{i,K-1} + \varepsilon_i$

Spring 2010

© Erling Berge 2010

109

Multiple regression: model (2)

- This can also be written

$$y_i = E[y_i] + \varepsilon_i ,$$

this means that

- $E[y_i]$ is read as “the expected value of y_i ”
- $E[y_i] = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \dots + \beta_{K-1} x_{i,K-1}$

Spring 2010

© Erling Berge 2010

110

Multiple regression: model (3)

- We will find the OLS estimates of the model parameters as the b-values in

$$\hat{y}_i = b_0 + b_1 x_{i1} + b_2 x_{i2} + b_3 x_{i3} + \dots + b_{K-1} x_{i,K-1}$$

(\hat{y}_i is read as "estimated" or "predicted" value of y_i)

that minimizes the squared sum of the residuals

$$RSS = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n e_i^2$$

Spring 2010

© Erling Berge 2010

111

Estimation methods

- The OLS method means that parameters are found by minimizing RSS (residual sum of squares)
- But this is not the only method for finding suitable b-values. Two alternatives are:
 - WLS: Weighted least squares
 - ML: maximum likelihood

Spring 2010

© Erling Berge 2010

112

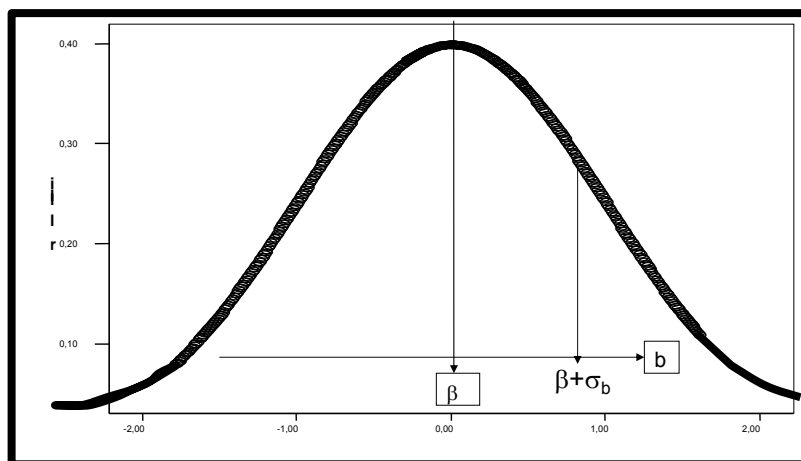
More on testing hypotheses

- We can draw many samples from a population
- In every new sample we can estimate new values (a new b_k -value) of the same population regression parameter (β_k)
- If we make a histogram of the many estimates of e.g. b_k we will see that b_k has a distribution. This distribution is called the sampling distribution of b_k
- Different types of parameters have different types of sampling distributions
- Regression parameters (OLS regression b_k) have t-distributions (Student's t-distribution)

Spring 2010

© Erling Berge 2010

113



Sampling distribution of the regression parameter b :

$$E[b] = \beta$$

Spring 2010

© Erling Berge 2010

114

On partial effects (1)

- Example with 2 variables
- If we estimate a model with 2 x-variables

$$y_i = b_0 + b_1 x_{i1} + b_2 x_{i2} + e_i$$

it will in principle involve 3 different correlations:

- Between y and x_1
- Between y and x_2
- Between x_1 and x_2

Spring 2010

© Erling Berge 2010

115

On partial effects (2)

- This might have been represented by 3 different bivariate regressions where the third variable was kept constant

$$(1) y = a_{y|x_1} + b_{y|x_1} x_1 + e_{y|x_1} \quad x_2 \text{ constant}$$

$$(2) y = a_{y|x_2} + b_{y|x_2} x_2 + e_{y|x_2} \quad x_1 \text{ constant}$$

$$(3) x_1 = a_{x_1|x_2} + b_{x_1|x_2} x_2 + e_{x_1|x_2} \quad y \text{ constant}$$

the index "y|x1" is read "from the regression of y on x1"

- Equations (2) and (3) can be rewritten as:

Spring 2010

© Erling Berge 2010

116

On partial effects (3)

$$(2) e_{y|X_2} = y - (a_{y|X_2} + b_{y|X_2}x_2)$$

$$(3) e_{x_1|X_2} = x_1 - (a_{x_1|X_2} + b_{x_1|X_2}x_2)$$

We may interpret this as a removal of the effect of x_2 from y and from x_1

We also see that $e_{y|X_2}$ and $e_{x_1|X_2}$ become the new y and x_1 variables where the effect of x_2 has been removed

Spring 2010

© Erling Berge 2010

117

On partial effects (4)

- If we, based on this, make a new regression

$$\hat{e}_{y|X_2} = a + b e_{x_1|X_2}$$

we find that

$$a = 0$$

$$b = b_1 \text{ from the regression}$$

$$y_i = b_0 + b_1 x_{i1} + b_2 x_{i2} + e_i$$

- b_1 is in other words the effect of x_1 on y **after we have removed the effect of x_2**

Spring 2010

© Erling Berge 2010

118

Experiments and partial effects

- Experiments investigate the causal connection between two variables controlled for all other causal impacts
- Multiple regression is a kind of half-way replication of experiments – the next best solution – and is a close relative of quasi-experimental research designs

Spring 2010

© Erling Berge 2010

119

Partial effects

A leverage plot for y and x_k is a plot where

- y-axis is the residual from the regression of y on all x -variables except x_k , and
- x-axis is the residual from regression of x_k on all the other x -variables

The regression line in such a plot will always go through $y=0$ and will have a slope coefficient equal to b_k

Spring 2010

© Erling Berge 2010

120

An example with 2 independent variables

Table 2.2 Dependent: Summer 1981 Water Use	B	Std. Error	t	Sig.
(Constant)	1201.124	123.325	9.740	.000
Income in Thousands	47.549	4.652	10.221	.000

Table 3.1 Dependent: Summer 1981 Water Use	B	Std. Error	t	Sig.
(Constant)	203.822	94.361	2.160	.031
Income in Thousands	20.545	3.383	6.072	.000
Summer 1980 Water Use	.593	.025	23.679	.000

From the table 2.2 (p46) and 3.1 (p68) in Hamilton. In the tables in the book the constant is on the last line. SPSS put it on the first line.

Question: What does it mean that the coefficient of income declines when we add a new variable?

Spring 2010

© Erling Berge 2010

121

On the addition of new variables

- It is not common that existing theory will give precise prescriptions for what variables to include in a model. Usually there is an element of trial and error in developing a model
- When new variables are added to a model several things happen
 - The explanatory force increase: R^2 increase, but will the increase be significant?
 - The coefficient of the regression shows the effect on y . Is this effect significantly different from 0?
 - If the coefficient is significantly different from 0, is it also so big that it is of substantial interest?
 - Spurious coefficients can decline. Do the new variable change the interpretation of the effect of the other variables?

Spring 2010

© Erling Berge 2010

122

Parsimony

- Parsimony is what might be called an aesthetic criterion of a good model. We want to explain as much as possible of the variation in y by means of as few variables as possible
- The adjusted coefficient of determination, Adjusted R^2 , is based on parsimony in the sense that it takes into consideration the complexity of the data relative to the complexity of the model by the difference between n and K
($n-K$ is the degrees of freedom in the residual,
 n = number of observations, K = number of estimated parameters)

Spring 2010

© Erling Berge 2010

123

Irrelevant variable

- Including irrelevant variables
 - A variable is irrelevant if the real effect (β) is 0; or more pragmatically, if it is so small that it has no substantive interest
 - **Inclusion of an irrelevant variable** makes the model unnecessarily complex and will have the consequence that coefficient estimates on all variables have larger variance (coefficients varies more from sample to sample)
- Including an irrelevant variable in OLS model estimation is probably **the least damaging error** we can do

Spring 2010

© Erling Berge 2010

124

Relevant variable

- A variable is relevant if
 - Its real effect (β) is significantly different from 0, and
 - Large enough to have substantive interest, and
 - It is **correlated with other included x-variables**
- If we exclude a relevant variable all results from our regression will be unreliable. The model is unrealistically simple
- Not including a relevant variable is **the most damaging error** we can do. But consider requirement 2 and 3. This makes it a lot easier to avoid this problem.

Sample specific results?

- Choice of variables is a trade-off among risks. Which risk is worse depends on the purpose of the study and the strength of relations
- With a test level of 0.05 one may easily find sample specific results. In about 5% of all samples a coefficient that show up as not significantly different from 0 will in "reality" be different from 0 ($\beta \neq 0$) and vice versa for those we find to be significantly different from 0 may in reality be 0
- The best defence against this is the theoretical argument for finding an effect different from 0

Hamilton (s74) example

y_i	Post shortage water use (1981)
x_{i1}	Household income, in thousands of dollars
x_{i2}	Pre-shortage water use, in cubic feet (1980)
x_{i3}	Education of household head, in years
x_{i4}	Retirement (coded 1 if household head is retired and 0 otherwise)
x_{i5}	Number of people living in household at time of water shortage (summer 1981)
x_{i6}	Change in number of people, summer 1981 minus summer 1980

Spring 2010

© Erling Berge 2010

127

Table 3.2 (Hamilton p74)

Dependent Variable: Summer 1981 Water Use	B	Std. Error	t	Sig.	Beta
(Constant)	242.220	206.864	1.171	.242	
Income in Thousands	20.967	3.464	6.053	.000	.184
Summer 1980 Water Use	.492	.026	18.671	.000	.584
Education in Years	-41.866	13.220	-3.167	.002	-.087
Head of house retired?	189.184	95.021	1.991	.047	.058
# of People Resident, 1981	248.197	28.725	8.641	.000	.277
Increase in # of People	96.454	80.519	1.198	.232	.031

How do we interpret the coefficient of "Increase in # of People" ?

What leads to less water use after the crisis?

Spring 2010

© Erling Berge 2010

128

Standardized coefficients

- Standardized variables (z-scores) have standard deviation as unit of measurement and a mean of 0

$$Z_{iX} = \frac{(X_i - \bar{X})}{s_X}$$

- Standardized regression coefficients (beta-weights, or path coefficients)
 $b_k^s = b_k(s_k/s_y)$ (varies between -1 and +1)
- Predicted standard score of y_i (\hat{z}_{iy}) = $0.18z_{i1} + 0.58z_{i2} - 0.09z_{i3} + 0.06z_{i4} + 0.28z_{i5} + 0.03z_{i6}$

Spring 2010

© Erling Berge 2010

129

t-test

- The difference between the observed coefficient (b_k) and the unobserved coefficient (β_k) standardized by the standard deviation of the observed coefficient (SE_{b_k}) will usually be very close to zero if the observed b_k is close to the population value. This means that if we in the formula
- $t = (b_k - \beta_k) / SE_{b_k}$ substitutes $\beta_k = 0$ (H_0) and find that "t" is small we will believe that the population value β_k in reality equals 0 (we cannot refute H_0)
- How big "t" has to be before we stop believing that $\beta_k = 0$ we can find from knowing the sampling distribution of b_k and SE_{b_k}

Spring 2010

© Erling Berge 2010

130

360 Appendix 4 Statistical Tables

Table A4.1 Critical values for student's t-distribution

Probability $P(|t| \leq c)$ for confidence intervals

Probability $P(|t| \geq c)$ for two-sided tests

Probability $P(t \geq c)$ for one-sided tests

df	Probability									Confidence intervals		
	.50	.40	.30	.25	.20	.15	.10	.05	.025		.01	.005
1	1.000	3.078	6.314	12.706	31.821	63.637	127.32	318.31	636.62			
2	.816	1.886	2.920	4.303	6.965	9.925	14.069	22.326	31.598			
3	.765	1.638	2.353	3.182	4.541	5.841	7.453	10.213	12.924			
4	.741	1.533	2.132	2.776	3.747	4.604	5.598	7.173	8.610			
5	.727	1.476	2.015	2.571	3.365	4.032	4.773	5.893	6.869			
6	.718	1.440	1.943	2.447	3.143	3.707	4.317	5.208	5.959			
7	.711	1.415	1.895	2.365	2.998	3.499	4.020	4.785	5.408			
8	.706	1.397	1.860	2.306	2.896	3.355	3.833	4.501	5.041			

"t" has a sampling distribution called the t-distribution. The t-distribution varies with the number of degrees of freedom (n-K) and is listed according to level of significance α .

Spring 2010

© Erling Berge 2010

131

Confidence interval for β (1)

- We have defined $t = (b_k - \beta_k) / SE_{b_k}$. This means that
- $t * (SE_{b_k}) = b_k - \beta_k$ or $\beta_k = b_k - t * (SE_{b_k})$ where t follows the t distribution with n-K degrees of freedom
- Choosing a t_α -value from the table of the t-distribution with n-K degrees of freedom then it is true that
- $Pr\{b_k - t * (SE_{b_k}) < \beta_k < b_k + t * (SE_{b_k})\} = 1 - \alpha$
- Then if $\beta_k = b_k$ is correct, a two-tailed test will have a probability of α to reject $H_0 : \beta_k = 0$ when H_0 in reality is correct (type I error)

Spring 2010

© Erling Berge 2010

132

Confidence interval for β (2)

- This means that there is a probability of α that β_k in reality is outside the interval

$$\langle \mathbf{b}_k - t_\alpha(\mathbf{SE}_{\mathbf{b}_k}), \mathbf{b}_k + t_\alpha(\mathbf{SE}_{\mathbf{b}_k}) \rangle$$
- This is equivalent to saying that

$$\mathbf{b}_k - t_\alpha(\mathbf{SE}_{\mathbf{b}_k}) \leq \beta_k \leq \mathbf{b}_k + t_\alpha(\mathbf{SE}_{\mathbf{b}_k})$$
 is correct with probability $1 - \alpha$ (our confidence of this result is $1 - \alpha$)
- $\Pr\{\mathbf{b}_k - t^*(\mathbf{SE}_{\mathbf{b}_k}) < \beta_k < \mathbf{b}_k + t^*(\mathbf{SE}_{\mathbf{b}_k})\} = 1 - \alpha$

Spring 2010

© Erling Berge 2010

133

F-test: big model against small (1)

Define:

$$F_{n-K}^H = \frac{\frac{RSS_{[K-H]} - RSS_{[K]}}{H}}{\frac{RSS_{[K]}}{n-K}}$$

$RSS_{[*]}$ = residual sum of squares with index $[*]$ where $*$ stands for number of parameters in the model

Spring 2010

© Erling Berge 2010

134

F-test: big model against small (2)

- Big model: $RSS_{[K]}$
- Small model: $RSS_{[K-H]}$
- H is the difference in number of parameters in the two models

F^H_{n-K} will have a sampling distribution called the **F-distribution** with H and n-K degrees of freedom

Spring 2010

© Erling Berge 2010

135

Example (Hamilton table 3.1 and 3.2)

Small model Table 3.1	Sum of Squares	df	Mean Square	F	Sig.
Regression (Model) (Explained)	671025350.237	2	335512675.119	391.763	.000(a)
Residual	422213359.440	493	856416.551		
Total	1093238709.677	495			

Large model Table 3.2	Sum of Squares	df	Mean Square	F	Sig.
Regression	740477522.059	K - 1 = 6	123412920.343	171.076	.000(a)
Residual	352761187.618	n - K = 489	721393.022		
Total	1093238709.677	n - 1 = 495			

Test if the big model (7 parameters) is better than the small (3 parameters)

Spring 2010

© Erling Berge 2010

136

Notes to the example

- K = number of parameters of the big model (6 variables plus constant) = 7
- $H = K - [\text{number of parameters in the small model (2 variables plus constant)}] = 7 - 3 = 4$
- $RSS_{[K-H]} = 422213359.440$
- $RSS_{[K]} = 352761187.618$
- $n = 496$
- $n - K = 496 - 7 = 489$
- $(RSS_{[K-H]} - RSS_{[K]})/H = (422213359.440 - 352761187.618)/4 = 17363042.9555$
- $RSS_{[K]}/(n-K) = 352761187.618/489 = 721393.0217$

Spring 2010

© Erling Berge 2010

137

Testing all parameters in one test

- If the big model has K parameters and we let the small model be as small as possible with only 1 parameter (the constant = the mean) our test will have $H=K-1$. Inserting this into our formula we have

$$F_{n-K}^{K-1} = \frac{\frac{RSS_{[1]} - RSS_{[K]}}{K-1}}{\frac{RSS_{[K]}}{n-K}}$$

This is the F-value we find in the ANOVA tables from SPSS
[note: $\{RSS[1] - RSS[K]\} = \text{ESS (explained sum of squares)}$]

Spring 2010

© Erling Berge 2010

138

Multicollinearity (1)

- Multicollinearity only involves the x-variables, not y, and is about linear relationships between two or more x-variables
- If there is a perfect correlation between 2 explanatory variables, e.g. x and w ($r_{xw} = 1$) the multiple regression model breaks down
- The same will happen if there is perfect correlation between two groups of x-variables

Spring 2010

© Erling Berge 2010

139

Multicollinearity (2)

- Perfect correlation is rarely a practical problem
- But high correlations between different x-variables or between groups of x-variables will make estimates of their effect unreliable.
- The effects of two highly correlated variables (like x and x^2) may be arbitrarily assigned to one, the other, or both
- Individual regression coefficients will have large standard deviations and t-tests will practically speaking have no interest whatsoever
- **F-tests of groups of variables will not be affected by this**

Spring 2010

© Erling Berge 2010

140

Search strategies

- There are methods for automatic searches for explanatory variables in a large set of data
- The best advice to give on this is to avoid using it
- One problem is that the p-values of the tests from such searches are wrong and too "kind". The the probability of making Type I errors increase with the number of tests
- This difficulty is called "the problem of multiple comparisons"
- Another problem is that such searches do not work well if the variables are highly correlated

Spring 2010

© Erling Berge 2010

141

Dummy variables: group differences

- Dichotomous variables taking the values of 0 or 1 are called dummy variables, or more generally binary variables
- In the example in table 3.2 (p74) x_{i4} (Head of house retired?) is a dummy variable
- First put into the equation $x_{i4} = 1$ then $x_{i4} = 0$
 $y_i = 242 + 21x_{i1} + 0.49x_{i2} - 42x_{i3} + 189x_{i4} + 248x_{i5} + 96x_{i6}$ og
- Explain what the two equations tell us

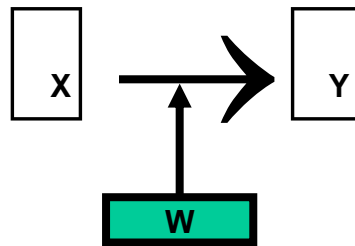
Spring 2010

© Erling Berge 2010

142

Interaction

- There is interaction between two variables if the effect of one variable changes or varies depending on the value of the other variable



Spring 2010

© Erling Berge 2010

143

Interaction effects in regression (1)

- If we do a non-linear transformation of y all estimated effects will implicitly be interaction effects
- Simple additive interaction effects can be included in a linear model by means of product terms where two x -variables are multiplied
- $\hat{y}_i = b_0 + b_1x_i + b_2w_i + b_3x_iw_i$
- Conditional effect plots will be able to illustrate what interaction means

Spring 2010

© Erling Berge 2010

144

Interaction effects in regression (2)

- An interaction effect involving x and w can be included in a regression model by means of an auxiliary variable equal to the product of the two variables, i.e.
- Auxiliary variable $H=x*w$
- $y_i = b_0 + b_1*x_i + b_2*w_i + b_3*H_i + e_i$
- $y_i = b_0 + b_1*x_i + b_2*w_i + b_3*x_i*w_i + e_i$

Spring 2010

© Erling Berge 2010

145

Example from Hamilton(p85-91)

Let

- y = natural logarithm of chloride concentration
- x = depth of well (1=deep, 0=shallow)
- w = natural logarithm of distance from road
- xw = interaction term between distance and depth (product $x*w$). Then
- $\hat{y}_i = b_0 + b_1x_i + b_2w_i + b_3x_iw_i$

First take a look at the simple models without interaction

Spring 2010

© Erling Berge 2010

146

Figures 3.3 and 3.4 (Hamilton p85-86)

Figure 3.3 is based on

Dependent Variable: lnChlorideConcentra	B	Std. Error	Beta	t	Sig.
(Constant)	3.775	.429		8.801	.000
x= BEDROCK OR SHALLOW WELL?	-.706	.477	-.205	-1.479	.145

Figure 3.4 is based on

Dependent Variable: lnChlorideConcentra	B	Std. Error	Beta	t	Sig.
(Constant)	4.210	.961		4.381	.000
w= lnDistanceFromRoad	-.091	.180	-.071	-.506	.615
x= BEDROCK OR SHALLOW WELL?	-.697	.481	-.202	-1.449	.154

Spring 2010

© Erling Berge 2010

147

Figure 3.3

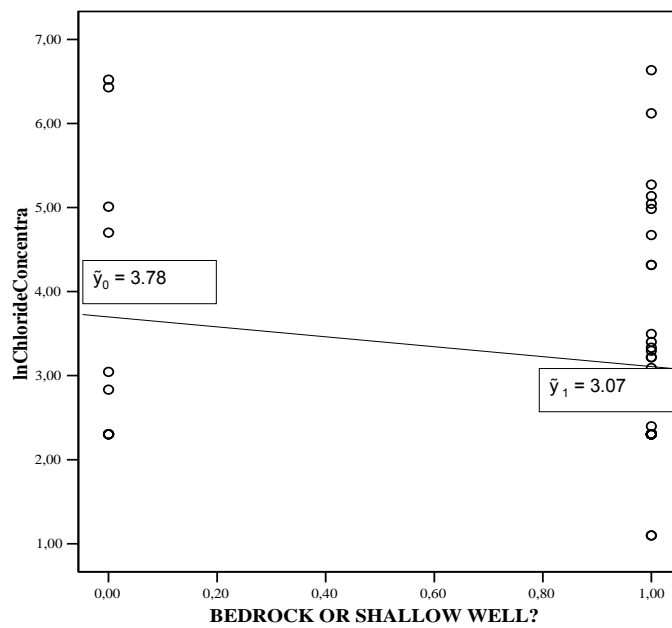
$$\hat{y}_i = 3.78 - .71x_i$$

Let

$$x_i = 1 \text{ (deep)}$$

and

$$x_i = 0 \text{ (shallow)}$$



Spring 2010

© Erling Berge 2010

148

Figure 3.4

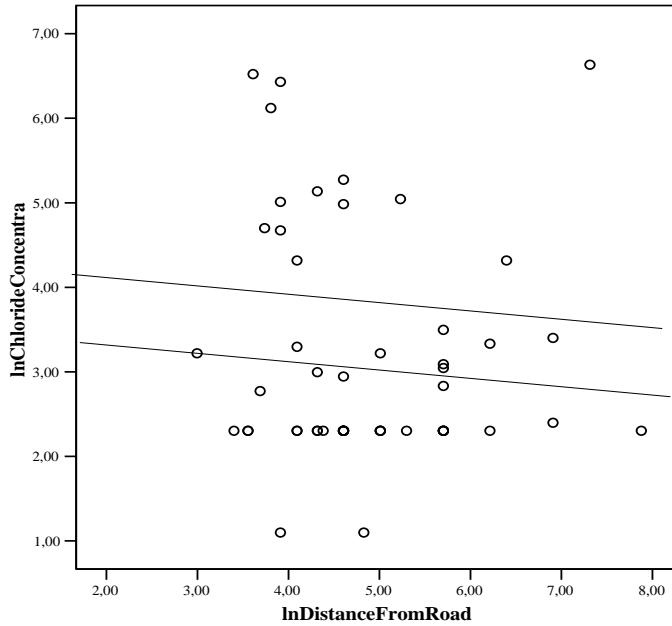
$$\hat{y}_i = 4.21 - .70x_i - .09w_i$$

Let

$x_i = 1$ (deep)

and

$x_i = 0$ (shallow)



Spring 2010

© Erling Berge 2010

147

Figures 3.5 and 3.6 (Hamilton p89-91)

Take note of significance changes

Figure 3.5 is based on

Dependent Variable: lnChlorideConcentra	B	Std. Error	Beta	t	Sig.
(Constant)	3.666	.905		4.050	.000
w= lnDistanceFromRoad	-.029	.202	-.022	-.144	.886
x*w= lnDroadDeep	-.081	.099	-.128	-.819	.417

Figure 3.6 is based on

Also see Table 3.4 in Hamilton p90 Dependent Variable: lnChlorideConcentra	B	Std. Error	Beta	t	Sig.
(Constant)	9.073	1.879		4.828	.000
w= lnDistanceFromRoad	-1.109	.384	-.862	-2.886	.006
x= BEDROCK OR SHALLOW WELL?	-6.717	2.095	-1.948	-3.207	.002
x*w= lnDroadDeep	1.256	.427	1.979	2.942	.005

Spring 2010

© Erling Berge 2010

150

Figure 3.5

$$\hat{y}_i = 3.67 - .03w_i - .08x_i$$

For

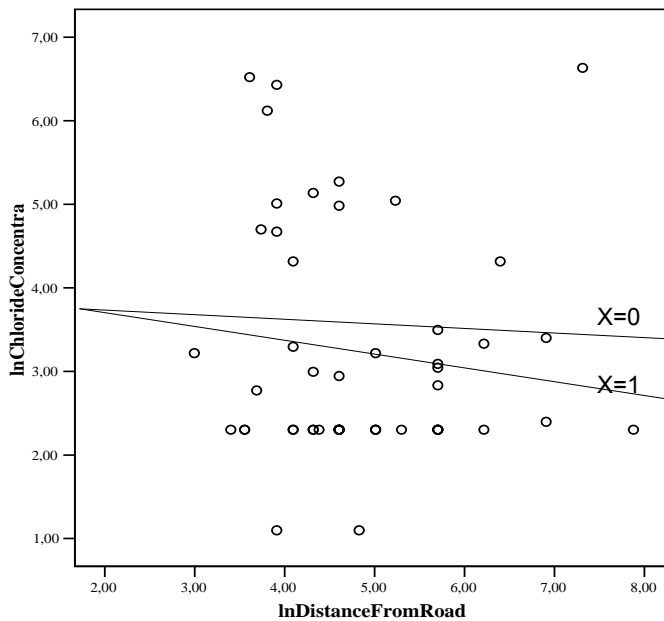
$$x_i = 1 \text{ (deep)}$$

$$\hat{y}_i = 3.67 - .11w_i$$

and for

$$x_i = 0 \text{ (shallow)}$$

$$\hat{y}_i = 3.67 - .03w_i$$



Spring 2010

© Erling Berge 2010

151

Figure 3.6

$$\hat{y}_i = 9.07 - 6.72x_i - 1.11w_i + 1.26x_iw_i$$

For

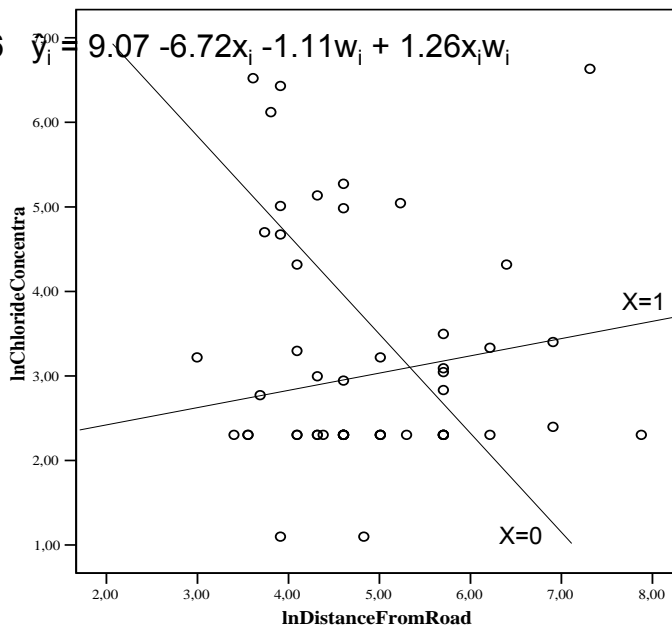
$$x_i = 1 \text{ (deep)}$$

$$\hat{y}_i = 2.35 + .15w_i$$

and for

$$x_i = 0 \text{ (shallow)}$$

$$\hat{y}_i = 9.07 - 1.11w_i$$



Spring 2010

© Erling Berge 2010

152

Multicollinearity

- Taking all three variables, x , w , and $x*w$ will introduce an element of multicollinearity. This means that we cannot trust tests of single coefficients
- But as shown in the previous example one can not drop any one of the variables without dropping a relevant variable
- F-test of e.g. w and $z*w$ simultaneously circumvents the test problem, and with some experimentation with different models one may see if excluding w or $x*w$ changes the relations substantially

Spring 2010

© Erling Berge 2010

153

Testing in the presence of multicollinearity

- To specify a model correctly we may need to add terms containing variables already in the equation. This applies to
 - Interaction terms
 - Curvilinear relations (use of squared variables in addition to the one present)
- Let us take a look at curvilinear relations:

Spring 2010

© Erling Berge 2010

154

Test for Curvilinear Relations

- Testing for curvilinearity in “age”
 - Set age squared = “age2”
- Remember:
 - Age is one substance variable that may be represented either by one technical variable or by two technical variables (somewhat like one variable being represented by different ways of coding)
- Substance variable Age is represented by
 - age
 - or
 - age + age2

Spring 2010

© Erling Berge 2010

155

- ## Testing for curvilinearity
- Model 0
 - (some variables)
 - Model 1
 - (some variables) + age
 - Model 2
 - (some variables) + age + age2
 - In model 1 the impact of Age is tested by the t-test and the corresponding p-value (there is no difference between the substance variable and its technical representation)

Spring 2010

© Erling Berge 2010

156

Testing for curvilinearity 2

- In model 1 the test may conclude that Age does not contribute to the model. If so we go to model 2
- In model 2 the testing of the impact of the substance variable Age (represented by age and age2) is done by an F-test of Model 2 against Model 0
- The F-test may conclude that Age does not contribute to the model. Then we drop both age and age2.
- The F-test may conclude that Age (represented by age and age2) contributes significantly to the model. Then we keep both age and age2

Spring 2010

© Erling Berge 2010

157

Testing for curvilinearity 3

- In model 1 the test may conclude that Age does contribute to the model. If so we may still go to Model 2
- If either the t-test of model 1, or the F-test of model 2, or both show that Age contributes significantly to the model, there are several possibilities
 - T-test significant, F-test not significant: drop age2, keep age
 - T-test significant, F-test significant, p-value of age is unchanged or higher (compared to model 1) while p-value of age2 is clearly insignificant: drop age2, keep age

Spring 2010

© Erling Berge 2010

158

Testing for curvilinearity 4

- (continued)
 - T-test significant, F-test significant, p-value of age improves (compared to model 1): keep age2 no matter what p-value for age2 is
 - T-test significant, F-test significant, p-value of age shows no significance (compared to model 1) while p-value of age2 shows clear significance: keep age2 no matter what p-value for age is
 - T-test significant, F-test significant, p-value of both age and age2 show no significance but are fairly close. Then the F-test decides. Keep age2.
- And remember: age2 never appears alone, always with age

Spring 2010

© Erling Berge 2010

159

Nominal scale variables

- Can be included in regression models by the use of new auxiliary variables: one for each category of the nominal scale variable. J categories implies $H(j)$, $j=1, \dots, J$ new auxiliary variables
- If the dependent variable is interval scale and the the only independent variable is nominal scale analysis of variance (ANOVA) is the most common approach to analysis
- By introducing auxiliary variables the same type of analysis can be done in a regression model

Spring 2010

© Erling Berge 2010

160

Analysis of variance - ANOVA

- Analysing an interval scale dependent variable with one or more nominal scale independent variables, often called factors
 - One way ANOVA uses one nominal scale variable
 - Two way ANOVA uses two nominal scale variable
 - And so on ...
- Tests of differences between groups are based on an evaluation of whether the variation within a group (defined by the "factors") is large compared to the variation between groups

Spring 2010

© Erling Berge 2010

161

Nominal scale variables in regression (1)

- If the nominal scale has J categories a maximum of $J-1$ auxiliary variables can enter the regression
If $H(j)$, $j=1, \dots, J-1$ are included $H(J)$ have to be excluded
- The excluded auxiliary variable is called the **reference category** and is the most important category in the interpretation of the results from the regression

Spring 2010

© Erling Berge 2010

162

Nominal scale variables in regression (2)

Dummy coding of a nominal scale variable

- The auxiliary variable $H(j)$ for a person i is coded 1 if the person belongs to category j on the nominal scale variable, it is coded 0 if the person do not belong to category j
- NB: The mean of a dummy coded variable is the proportion in the sample with value 1 (i.e. that belongs in the category)

Spring 2010

© Erling Berge 2010

163

Nominal scale variables in regression (3)

The reference category

(the excluded auxiliary variable)

- The chosen reference category ought to be large and clearly defined
- The estimated effect of an included auxiliary variable measures the effect of being in the included category relative to being in the reference category

Spring 2010

© Erling Berge 2010

164

Nominal scale variables in regression (4)

- This means that the regression parameter for an included dummy coded auxiliary variable tells us about additions or subtractions from the expected Y-value a person gets by being in this category rather than in the reference category
- When all auxiliary variables are zero the effect of being in the reference category is included in the constant

Spring 2010

© Erling Berge 2010

165

Nominal scale variables in regression (5)

Testing I

- Testing if a regression coefficient for an included auxiliary variable equals 0 answers the question whether the persons in this group have a mean Y value different from the mean value of the persons in the reference category

Spring 2010

© Erling Berge 2010

166

Nominal scale variables in regression (6)

Testing II

- Testing whether a Nominal scale variable contributes significantly to a regression model have to be done by testing if all auxiliary variables in sum contributes significantly to the regression
- For this we use the F-test as explained above. See formula 3.28 in Hamilton (p80)

Spring 2010

© Erling Berge 2010

167

Nominal scale variables in regression (7)

Interaction

- When dummy coded nominal scale variables are entered into an interaction all included auxiliary variables have to be multiplied with the variable suspected of interacting with it

Spring 2010

© Erling Berge 2010

168

On terminology (1)

- **Dummy coding** of nominal scale variables are called different names in different textbooks. For example it is
 1. Dummy coding in Hamilton, Hardy, and Weisberg
 2. Indicator coding in Menard (and also Weisberg)
 3. Reference coding or partial method in Hosmer&Lemeshow

Spring 2010

© Erling Berge 2010

169

On terminology (2)

- To reproduce results from the analysis of variance (ANOVA) by means of regression techniques Hamilton introduces a coding of the auxiliary variables he calls effect coding. Other authors call it differently:
 - It is called effect coding by Hardy
 - It is called deviance coding by Menard
 - It is called the marginal method or deviance method by Hosmer&Lemeshow
- To highlight particular group comparisons Hardy (Ch5) introduces a coding scheme called contrast coding

Spring 2010

© Erling Berge 2010

170

Ordinal scale variables

- Can be included as an interval scale if the unobserved theoretical dimension is continuous and distance measures seems reasonable
- Also it may be used directly as dependent variable if the program allows ordinal dependent variables
 - In that case parameters are estimated for every level above the lowest as cumulative effects relative to the lowest level

Spring 2010

© Erling Berge 2010

171

Nominal scale variables

TYPE OF GROUP	Frequency	Percent	Valid Percent	Cumulative Percent
POLITICIAN	48	12.6	12.6	12.6
FARMER	132	34.7	34.7	47.4
PEOPLE not Farmers or Pol	200	52.6	52.6	100.0
Total	380	100.0	100.0	

Spring 2010

© Erling Berge 2010

172

Example of dummy coding

Nominal scale			Auxiliary	variables	H (*)	
Type of group	Code	N	H(1)= Pol	H(2)= Farmer	H(3)= People	
Politicians	1	48	1	0	0	
Farmers	2	132	0	1	0	
Other People	3	200	0	0	1	Reference

A variable with 3 categories leads to 2 dummy coded variables in a regression with the third used as reference

Spring 2010

© Erling Berge 2010

173

Example of effect coding

Nominal scala			Auxiliary			
Type of group	Code	N	H(1)= Pol	H(2)= Farmer		
Politicians	1	48	1	0		
Farmers	2	132	0	1		
Other People	3	200	-1	-1		Reference category

In effect coding the reference category is coded -1. Effect coding makes it possible to duplicate all F-tests of ordinary ANOVA analyses.

Spring 2010

© Erling Berge 2010

174

Contrast coding

- Is used to present just those comparisons that are of the highest theoretical interest
- Contrast coding requires
 - That with J categories there have to be J-1 contrasts
 - The values of the codes on each auxiliary variable have to sum to 0
 - The values of the codes on any two auxiliary variables have to be orthogonal (their vector product has to be 0)

Spring 2010

© Erling Berge 2010

175

Use of dummy coded variables(1)

Dependent Variable: I. of political contr. of sales of agric. est.	B	Std. Error	Beta	t	Sig.
(Constant)	4.106	.152		26.991	.000
Pol	.914	.337	.147	2.711	.007
Farmer	.421	.240	.096	1.758	.080

- The constant shows the mean of the dependent variable for those who belong to the reference category
- The mean of the dependent variable for politicians are 0.91 opinion score points above the mean of the reference category
- The mean on the dependent variable for farmers are 0.42 opinion score points above the mean of the reference category

Spring 2010

© Erling Berge 2010

176

Use of dummy coded variables (2)

Dependent Variable: I. of political control of sales of agricultural estates	B	Std. Error	t	Sig.
(Constant)	4.264	.186	22.954	.000
Number of decares land Owned	.000	.000	2.176	.030
Pol	.566	.382	1.482	.139
Farmer	-.309	.338	-.913	.362

Compare this table with the previous. What has changed?

How do we interpret the coefficient on "Pol" and "Farmer"?

Spring 2010

© Erling Berge 2010

177

Recall:

Multiple regression: model

Let K = number of parameters in the model

(then $K-1$ = number of variables)

Population model

- $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \dots + \beta_{K-1} x_{i,K-1} + \varepsilon_i$
 $i = 1, \dots, N$; where N = number of case in the population

Sample model

- $y_i = b_0 + b_1 x_{i1} + b_2 x_{i2} + b_3 x_{i3} + \dots + b_{K-1} x_{i,K-1} + e_i$
 $i = 1, \dots, n$; where n = number of case in the sample

Spring 2010

© Erling Berge 2010

178

A note on the dependent variable in OLS regression:

- The requirement is that Y in OLS regression has to be interval scale. It has to be able to take any value between minus infinity and plus infinity.
- Deviations from this may cause problems
- It is not, I repeat NOT, most emphatically **NOT** required that it shall have any particular distribution such as a normal distribution
- In some other types of models this is different. Maximum likelihood factor analysis for example assumes a multivariate normal distribution
- Normal distributions are assumed in order to be able to do tests

Spring 2010

© Erling Berge 2010

179

Conclusions (1)

- Linear regression can easily be extended to use 2 or more explanatory variables
- If the assumptions of the regression is satisfied (that the error terms are normally distributed with independent and identically distributed errors – “normal i.i.d. errors”) the regression will be a versatile and strong tool for analytical studies of the connection between a dependent and one or more independent variables

Spring 2010

© Erling Berge 2010

180

Conclusions (2)

- The most common method of estimating coefficients for a regression model is called OLS (ordinary least squares)
- Coefficients computed based on a sample are seen as estimates of the population coefficient
- Using the t-test we can judge how good such coefficient estimates are
- Using the F-test we may evaluate several coefficient estimates in one test (dummy coded variables, interaction terms, curvilinear variables)

Spring 2010

© Erling Berge 2010

181

Conclusions (3)

- Dummy variables are useful in several ways
 - A single dummy coded x-variable will give a test of the difference in means for two groups (coded 0 and 1)
 - Nominal scale variables with more than 2 categories can be recoded by means of dummy coding and included in regression analysis
 - By using effect coding we can perform analysis of variance of the ANOVA type

Spring 2010

© Erling Berge 2010

182

Logistic regression

- Hamilton Ch 7 p217-234

Spring 2010

© Erling Berge 2010

183

LOGIT REGRESSION

- **Should be used if the dependent variable (Y) is a nominal scale**
- Here it is assumed that Y has the values 0 or 1
- The model of the conditional probability of Y, $E[Y | X]$, is based on the logistic function ($E[Y | X]$ is read “the expected value of Y given the value of X”)
- But
Why cannot $E[Y | X]$ be a linear function also in this case?

Spring 2010

© Erling Berge 2010

184

The linear probability model: LPM

- The linear probability model (LPM) of y_i when y_i can take only two values (0, 1) assumes that we can interpret $E[y_i | \mathbf{X}_i]$ as a probability
- $\mathbf{X}_i = \{x_{1i}, x_{2i}, x_{3i}, \dots, x_{(K-1)i}\}$
- $E[y_i | \mathbf{X}_i] = b_0 + \sum_j b_j x_{ji} = \Pr[y_i = 1]$
- This leads to severe problems:

Spring 2010

© Erling Berge 2010

185

Are the assumptions of a linear regression model satisfied for the LPM?

- One assumption of the LPM is that the residual, e_i satisfies the requirements of OLS
- The residual must be either
 - $e_i = 1 - (b_0 + \sum_j b_j x_{ji})$ or
 - $e_i = 0 - (b_0 + \sum_j b_j x_{ji})$
- This means that there is heteroscedasticity (the residual varies with the size of the values on the x-variables)
- There are estimation methods that can get around this problem (such as 2-stage weighted least squares method)
- One example of LPM:

Spring 2010

© Erling Berge 2010

186

OLS regression of a binary dependent variable on the independent variable "years lived in town"

ANOVA tabell	Sum of Squares	df	Mean Square	F	Sig.
Regression	3,111	1	3,111	13,648	,000(a)
Residual	34,418	151	,228		
Total	37,529	152			

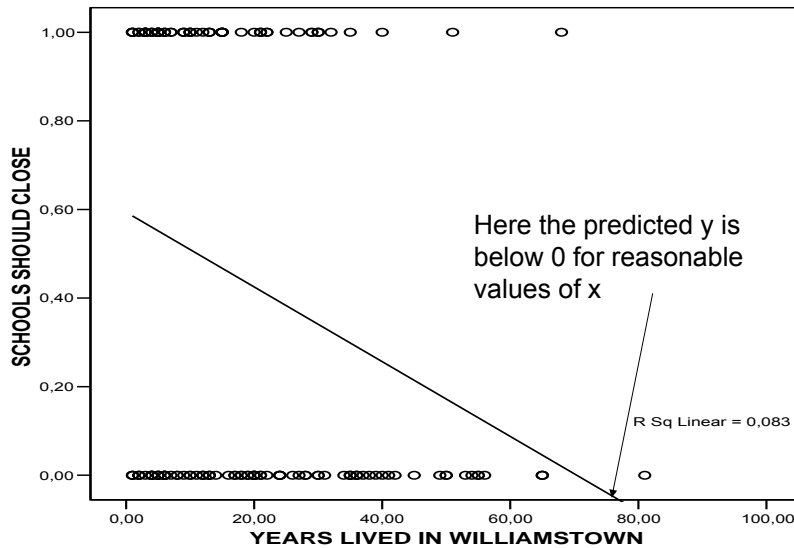
Dependent Variable: SCHOOLS SHOULD CLOSE	B	Std. Error	t	Sig.
(Constant)	,594	,059	10,147	,000
YEARS LIVED IN TOWN	-,008	,002	-3,694	,000

The regression looks OK in these tables

Spring 2010

© Erling Berge 2010

187



Scatter plot with line of regression. Figure 7.1 Hamilton

Spring 2010

© Erling Berge 2010

188

Conclusion: LPM model is wrong

- The example shows that for reasonable values of the x variable we can get values of the predicted y where $E[y_i | \mathbf{X}_i] > 1$ or $E[y_i | \mathbf{X}_i] < 0$,
- For this there is no remedy
- LPM is for substantial reasons a wrong model
- We need a model where we always will have $0 \leq E[y_i | \mathbf{X}_i] \leq 1$
- The logistic function can provide such a model

Spring 2010

© Erling Berge 2010

189

The logistic function

The general logistic function is written

$$y_i = \alpha / (1 + \gamma \cdot \exp[-\beta x_i]) + \varepsilon_i$$

$\alpha > 0$ provides an upper limit for y_i

this means that $0 < y_i < \alpha$

γ determines the horizontal point for rapid growth

If we determine that $\alpha = 1$ and $\gamma = 1$ one will always find that

$$0 < 1 / (1 + \exp[-\beta x_i]) < 1$$

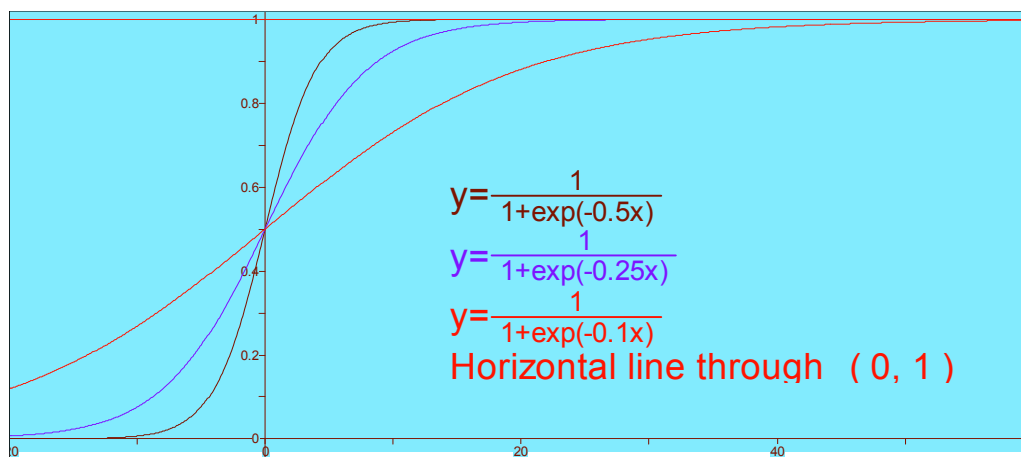
The logistic function will for all values of x_i lie between 0 and 1

Spring 2010

© Erling Berge 2010

190

Logistic curves for different β



β determines how rapidly the curve grows

Spring 2010

© Erling Berge 2010

191

MODEL (1)

Definitions:

- The probability that person no i shall have the value 1 on the variable y_i will be written $\Pr(y_i = 1)$.
- Then $\Pr(y_i \neq 1) = 1 - \Pr(y_i = 1)$
- The odds that person no i shall have the value 1 on the variable y_i , here called O_i , is the ratio between two probabilities

$$O_i(y_i = 1) = \frac{\Pr(y_i = 1)}{1 - \Pr(y_i = 1)} = \frac{p_i}{1 - p_i}$$

Spring 2010

© Erling Berge 2010

192

MODEL (2)

Definitions:

- The LOGIT , L_i , for person no i (corresponding to $\Pr(y_i=1)$) is the natural logarithm of the odds, O_i , that person no i has the value 1 on variable y_i , is written:

$$L_i = \ln(O_i) = \ln\{p_i/(1-p_i)\}$$
- The model assumes that L_i is a linear function of the explanatory variables x_j ,
- i.e.:
- $L_i = \beta_0 + \sum_j \beta_j x_{ji}$, where $j=1,\dots,K-1$, and $i=1,\dots,n$

Spring 2010

© Erling Berge 2010

193

MODEL (3)

- Let $X =$ (the collection of all x_j), then the probability of $Y_i = 1$ for person no i

$$\Pr(y_i = 1) = E[y_i | X_i] = \frac{1}{1 + \exp(-L_i)} = \frac{\exp(L_i)}{1 + \exp(L_i)}$$

$$\text{where } L_i = \beta_0 + \sum_{j=1}^{K-1} \beta_j X_{ji}$$

The graph of this relationship is useful for the interpretation what a change in x means

Spring 2010

© Erling Berge 2010

194

MODEL (4)

In the model $Y_i = E[y_i | X_i] + \varepsilon_i$ the error is either

- $\varepsilon_i = 1 - E[y_i | X_i]$ with probability $E[y_i | X_i]$
(since $\Pr(y_i = 1) = E[y_i | X_i]$),

or the error is

- $\varepsilon_i = -E[y_i | X_i]$ with probability $1 - E[y_i | X_i]$

- Meaning that the error has a distribution known as the binomial distribution with

$$p_i = E[y_i | X_i]$$

Estimation by the ML method

- The method used to estimate the parameters in the model is Maximum Likelihood
- The ML-method gives us the parameters that maximize the likelihood of finding just the observations we have got
- This Likelihood we call \mathcal{L}
- The criterion for choosing regression parameters is that the Likelihood becomes as large as possible

Maximum Likelihood (1)

- The Likelihood equals the product of the probability of each observation. For a dichotomous variable where $\Pr(Y_i = 1) = P_i$ this can be written

$$\mathcal{L} = \prod_{i=1}^n \left\{ P_i^{Y_i} (1 - P_i)^{(1-Y_i)} \right\}$$

Spring 2010

© Erling Berge 2010

197

Maximum Likelihood (2)

- It is easier to maximize the likelihood \mathcal{L} if one uses the natural logarithm of \mathcal{L} :

$$\ln(\mathcal{L}) = \sum_{i=1}^n \left\{ y_i \ln P_i + (1 - y_i) \ln(1 - P_i) \right\}$$

- The natural logarithm of \mathcal{L} is called the LogLikelihood, It will be written \mathcal{LL} .
- \mathcal{LL} has a central role in logistic regression.

Spring 2010

© Erling Berge 2010

198

Maximum Likelihood (3)

- The LogLikelihood \mathcal{L} will always be negative
- Maximizing \mathcal{L} is the same as minimizing the **positive LogLikelihood**; i.e. minimizing $-\mathcal{L}$
- Finding parameter values that minimizes $-\mathcal{L}$ can be done only by "trial and error", i.e. using an iterative procedure

Spring 2010

© Erling Berge 2010

199

Iterative estimation

From Hamilton Tabell 7.1	Iteration	-2 Log Likelihood	Coefficients	
			Constant	lived
Initial	0	209,212	-,276	
Step	1	195,684	,376	-,034
	2	195,269	,455	-,041
	3	195,267	,460	-,041
	4	195,267	,460	-,041

Note the column titled -2 LogLikelihood

Spring 2010

© Erling Berge 2010

200

Footnotes to the tables

- Step 0: Point of departure is a model with only a constant and no variables
- **Iterative estimation**
 - Estimation ends at iteration no 4 since the parameter estimates changed less than 0.001

For the next slide:

- The Wald statistic that SPSS provides equals the square of the “t” that Hamilton (and STATA) provides (Wald = t²)

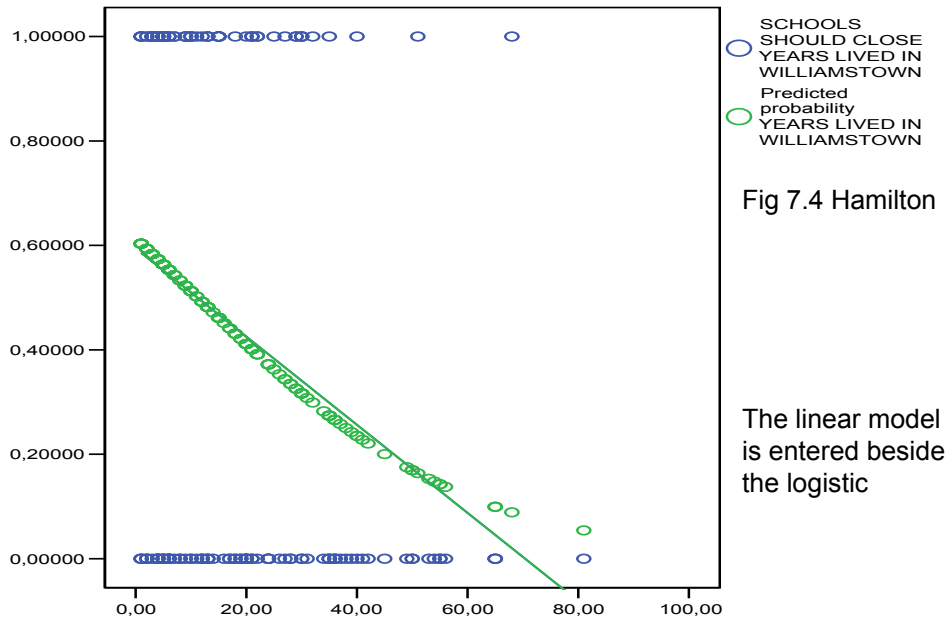
Logistic model instead of LPM

OLS regression (slide 6 above)

Dependent Variable: SCHOOLS SHOULD CLOSE	B	Std. Error	t	Sig.
(Constant)	,594	,059	10,147	,000
YEARS LIVED IN TOWN	-,008	,002	-3,694	,000

Logistic regression

Dependent: Schools should close	B	S.E.	Wald	df	Sig.	Exp(B)
Lived in town	-,041	,012	11,399	1	,001	,960
Constant	,460	,263	3,069	1	,080	1,584



Spring 2010

© Erling Berge 2010

203

TESTING

Two tests are useful

- (1) The Likelihood ratio test
 - This can be used analogous to the F-test (e.g. comparing two NESTED models)
- (2) Wald test
 - The square root of this can be used analogous to the t-test

Spring 2010

© Erling Berge 2010

204

Interpretation (1)

- The difference between the linear model and the logistic is large in the neighbourhood of 0 and 1
- LPM is easy to interpret: $Y_i = \beta_0$ when $x_{1i}=0$, and when x_{1i} increases with one unit Y_i increases with β_1 units
- The logistic model is more difficult to interpret. It is non-linear both in relation to the odds and the probability

Spring 2010

© Erling Berge 2010

205

ODDS and ODDS RATIOS

- The Logit, L_i , ($L_i = \beta_0 + \sum_j \beta_j x_{ji}$) is defined as the natural logarithm of the odds

This means that

- $\text{odds} = O_i(Y_i=1) = \exp(L_i) = e^{L_i}$

and

- **Odds ratio** = $O_i(Y_i=1 | L_i') / O_i(Y_i=1 | L_i)$
– where L_i' and L_i have different values on only one variable x_j .

Spring 2010

© Erling Berge 2010

206

Interpretation (2)

- When all x equals 0 then $L_i = \beta_0$ This means that the odds for $y_i = 1$ in this case is $\exp\{\beta_0\}$
- If all x -variables are kept fixed (they sum up to a constant) while x_1 increases with 1, the odds for $y_i = 1$ will be multiplied by $\exp\{\beta_1\}$
- This means that it will change with $100(\exp\{\beta_1\} - 1) \%$
- The probability $\Pr\{y_i = 1\}$ will change with a factor affect by all elements in the logit

Spring 2010

© Erling Berge 2010

207

Logistic regression: assumptions

- The model is correctly specified
 - The logit is linear in its parameters
 - All relevant variables are included
 - No irrelevant variables are included
- x -variables are measured without error
- Observations are independent
- No perfect multicollinearity
- No perfect discrimination
- Sufficiently large sample

Spring 2010

© Erling Berge 2010

208

Assumptions that cannot be tested

- Model specification
 - All relevant variables are included
- x-variables are measured without error
- Observations are independent

Two will be tested automatically.

- If the model can be estimated by SPSS there is
 - No perfect multicollinearity and
 - No perfect discrimination

Assumptions that can be tested

- Model specification
 - logit is linear in the parameters
 - no irrelevant variables are included
- Sufficiently large sample
 - What is “sufficiently large” depends on the number of different patterns in the sample and how cases are distributed across these
- Testing implies an assessment of whether statistical problems leads to departure from the assumptions

LOGISTIC REGRESSION

Statistical problems may be due to

- Too small a sample
- High degree of **multicollinearity**
 - Leading to large standard errors (imprecise estimates)
 - Multicollinearity is discovered and treated in the same way as in OLS regression
- High degree of **discrimination** (or separation)
 - Leading to large standard errors (imprecise estimates)
 - Will be discovered automatically by SPSS

Spring 2010

© Erling Berge 2010

211

Discrimination in Hamilton table 7.5

- Odds for weaker requirements is $44/202 = 0,218$ among women without small children
- Odds for weaker requirement is $0/79 = 0$ among women with small children
- Odds rate is $0/0,218 = 0$ hence $\exp\{b_{\text{woman}}\}=0$
- This means that $b_{\text{woman}} = \text{minus infinity}$

Y = Strength of water quality standards	Women without small children	Women with small children
Not weaker	202	79
Weaker OK	44	0

Spring 2010

© Erling Berge 2010

212

Discrimination/ separation

- Problems with discrimination appear when we for a given x-value get almost perfect prediction of the y-value (nearly all with a given x-value have the same y-value)
- In SPSS it may produce the following message:

Warnings

- | |
|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <ul style="list-style-type: none">• There is possibly a quasi-complete separation in the data. Either the maximum likelihood estimates do not exist or some parameter estimates are infinite. |
| <ul style="list-style-type: none">• The NOMREG procedure continues despite the above warning(s). Subsequent results shown are based on the last iteration. Validity of the model fit is uncertain. |

Spring 2010

© Erling Berge 2010

213

The LikeLihood Ratio test (1)

- The ratio between two Likelihoods equals the difference between two **LogLikelihoods**
- The difference between the **LogLikelihood** (\mathcal{LL}) of two **nested** models, estimated on **the same data**, can be used to test which of two models fits the data best, just like the F-statistic is used in OLS regression
- The test can also be used for single regression coefficients (single variables). In small samples it has better properties than the Wald statistic

Spring 2010

© Erling Berge 2010

214

The LikeLihood Ratio test (2)

The LikeLihood Ratio test statistic

- $\chi^2_H = -2[\mathcal{LL}(\text{model1}) - \mathcal{LL}(\text{model2})]$

will, if the null hypothesis of no difference between the two models is correct, be distributed approximately (for large n) as the chi-square distribution with number of degrees of freedom equal to the difference in number of parameters in the two models (H)

Spring 2010

© Erling Berge 2010

215

Example of a Likelihood Ratio test

- Model 1: just constant
- Model 2: constant plus one variable
- $\chi^2_H = -2[\mathcal{LL}(\text{model1}) - \mathcal{LL}(\text{model2})]$
 $= -2\mathcal{LL}(\text{model1}) + 2\mathcal{LL}(\text{model2})$
- Find the value of the ChiSquare and the number of degrees of freedom
- e.g.: LogLikelihood (mod1) = 209,212/(-2)
- LogLikelihood (mod2) = 195,267/(-2)

From Tab 7.1: -2 Log Likelihood
209,212
195,684
195,269
195,267
195,267

Spring 2010

© Erling Berge 2010

216

The Wald test (1)

- The Wald (or chisquare) test statistic provided by SPSS = $t^2 = (b_k / SE(b_k))^2$ (where t is the t used by Hamilton) can be used for testing single parameters similarly to the t-statistic of the OLS regression
- If the null hypothesis is correct, t will (for large n) in logistic regression be approximately normally distributed
- If the null hypothesis is correct, the Wald statistic will (for large n) in logistic regression be approximately chisquare distributed with df=1

Spring 2010

© Erling Berge 2010

217

Excerpt from Hamilton Table 7.2

Iterasjon	-2 Log likelihood					
0	209,212					
1	152,534					
2	149,466					
3	149,382					
4	149,382					
5	149,382					
Variables	B	S.E.	Wald	df	Sig.	Exp(B)
Lived	-,046	,015	9,698	1	,002	,955
Educ	-,166	,090	3,404	1	,065	,847
Contam	1,208	,465	6,739	1	,009	3,347
Hsc	2,173	,464	21,919	1	,000	8,784
Constant	1,731	1,302	1,768	1	,184	5,649

Spring 2010

© Erling Berge 2010

218

Confidence interval for parameter estimates

- Can be constructed based on the fact that the square root of the Wald statistic approximately follows a normal distribution with 1 degree of freedom
- $b_k - t_\alpha * SE(b_k) < \beta_k < b_k + t_\alpha * SE(b_k)$
where t_α is a value taken from the table of the **normal distribution** with level of significance equal to α

Spring 2010

© Erling Berge 2010

219

Can be constructed based on the t-distribution
(1)

- If a table of the normal distribution is missing one may use the **t-distribution** since the t-distribution is approximately normally distributed for large $n-K$ (e.g. for $n-K > 120$)

Spring 2010

© Erling Berge 2010

220

Excerpt from Hamilton Table 7.3 (from SPSS)

STATA SPSS		B	S.E.	t ² Wald	df	Prob>t Sig.	Exp(B)
Step 1	lived	-,047	,017	7,550	1	,006	,954
	educ	-,206	,093	4,887	1	,027	,814
	contam	1,282	,481	7,094	1	,008	3,604
	hsc	2,418	,510	22,508	1	,000	11,223
	female	-,052	,557	,009	1	,926	,950
	kids	-,671	,566	1,406	1	,236	,511
	nodad	-2,226	,999	4,964	1	,026	,108
	Constant	2,894	1,603	3,259	1	,071	18,060

Spring 2010

© Erling Berge 2010

221

More from Hamilton Table 7.3

Iteration		-2 Log likelihood	Coefficients							
			Const	lived	educ	contam	hsc	female	kids	nodad
Step0		209,212	-0,276							
Step1	1	147,028	1,565	-,027	-,130	,782	1,764	-,015	-,365	-1,074
	2	141,482	2,538	-,041	-,187	1,147	2,239	-,037	-,580	-1,844
	3	141,054	2,859	-,046	-,204	1,269	2,401	-,050	-,662	-2,184
	4	141,049	2,893	-,047	-,206	1,282	2,418	-,052	-,671	-2,225
	5	141,049	2,894	-,047	-,206	1,282	2,418	-,052	-,671	-2,226

Spring 2010

© Erling Berge 2010

222

Is the model in table 7.3 better than the model
in table 7.2 ?

- $\mathcal{LL}(\text{model in 7.3}) = 141,049/(-2)$
- $\mathcal{LL}(\text{model in 7.2}) = 149,382/(-2)$
- $\chi^2_{\text{H}} = -2[\mathcal{LL}(\text{model 7.2}) - \mathcal{LL}(\text{model 7.3})]$
- Find χ^2_{H} value
- Find H
- Look up the table of the chisquare distribution

Spring 2010

© Erling Berge 2010

223

The model of the probability of observing
 $y=1$ for person i

$$\Pr(y_i = 1) = E[y_i | x] = \frac{1}{1 + \exp(-L_i)} = \frac{\exp(L_i)}{1 + \exp(L_i)}$$

where the logit $L_i = \beta_0 + \sum_{j=1}^{K-1} \beta_j X_{ji}$ is a linear function
of the explanatory variables

It is not easy to interpret the meaning of the β
coefficients just based on this formula

Spring 2010

© Erling Berge 2010

224

The odds ratio

- The odds ratio, **O**, can be interpreted as the relative effect of having one variable value rather than another
- e.g. if $x_{ki} = t+1$ in L_i' and $x_{ki} = t$ in L_i
- $O = O_i(Y_i=1 | L_i') / O_i(Y_i=1 | L_i)$
 $= \exp[L_i'] / \exp[L_i]$
 $= \exp[\beta_k]$
- Why β_k ?

Spring 2010

© Erling Berge 2010

225

The odds ratio : example I

- The Odds for answering yes =
 $e^{b_0 + b_1 * Alder + b_2 * Kvinne + b_3 * E.utd + b_4 * Barn_i_HH}$
- The odds ratio for answering yes between women and men =

$$\frac{e^{b_0 + b_1 * Alder + b_2 * 1 + b_3 * E.utd + b_4 * Barn_i_HH}}{e^{b_0 + b_1 * Alder + b_2 * 0 + b_3 * E.utd + b_4 * Barn_i_HH}} = e^{b_2}$$

Remember the rules of power exponents

Spring 2010

© Erling Berge 2010

226

The odds ratio : example II

- The Odds for answering yes given one year of extra education

$$\frac{e^{b_0+b_1*Alder+b_2*Kvinne+b_3*(E.utd+1)+b_4*Barn_i_HH}}{e^{b_0+b_1*Alder+b_2*Kvinne+b_3*E.utd+b_4*Barn_i_HH}} = e^{b_3}$$

Remember the rules of power exponents

Spring 2010

© Erling Berge 2010

227

Example from Hamilton table 7.2

- What is the odds ratio for yes to closing the school from one year extra education?
- The odds ratio is the ratio of two odds where one odds is the odds for a person with one year extra education

$$\frac{e^{b_0+b_1*ÅrBuddIByen+b_2*(Utdanning+1)+b_3*UreiningEigEigedom+b_4*MangeHSCmøter}}{e^{b_0+b_1*ÅrBuddIByen+b_2*Utdanning+b_3*UreiningEigEigedom+b_4*MangeHSCmøter}} = \frac{e^{b_2*(Utdanning+1)}}{e^{b_2*Utdanning}} = e^{b_2}$$

Spring 2010

© Erling Berge 2010

228

Example from Hamilton table 7.2 cont.

- Odds ratio = $\text{Exp}\{b_2\} = \exp(-0,166) = 0,847$
- One extra year of education implies that the odds is reduced with a factor of 0.847
- One may also say that the odds has increased with a factor of
 $100(0,847-1)\% = -15,3\%$
- Meaning that it has declined with 15,3%

Spring 2010

© Erling Berge 2010

229

Concluding on logistic regression

- If the assumptions are satisfied logistic regression will provide normally distributed, unbiased and efficient (minimal variance) estimates of the parameters

Spring 2010

© Erling Berge 2010

230

Regression criticism

- Hamilton Ch 4 p109-123

Spring 2010

© Erling Berge 2010

231

Analyses of models are based on assumptions

- OLS is a simple technique of analysis with very good theoretical properties. But:
- The good properties are based on certain assumptions
- If the assumptions do not hold the good properties evaporates
- Investigating the degree to which the assumptions hold is the most important part of a regression analysis

Spring 2010

© Erling Berge 2010

232

OLS-REGRESSION: assumptions

- I SPECIFICATION REQUIREMENT
 - The model is correctly specified
- II GAUSS-MARKOV REQUIREMENTS
 - Ensures that the estimates are “BLUE”
- III NORMALLY DISTRIBUTED ERROR TERM
 - Ensures that the tests are valid

Spring 2010

© Erling Berge 2010

233

I SPECIFICATION REQUIREMENT

- The model is correctly specified if
 - The expected value of y , given the values of the independent variables, is a linear function of the parameters of the x -variables
 - All included x -variables have an impact on the expected y -value
 - No other variable has an impact on expected y -value *at the same time as they correlate with included x -variables*

Spring 2010

© Erling Berge 2010

234

II GAUSS-MARKOV REQUIREMENTS

(i)

- (1) x is known, without stochastic variation
(2) Errors have an expected value of 0 for all i

$$\bullet E(\varepsilon_i) = 0 \quad \text{for all } i$$

Given (1) and (2) ε_i will be independent of x_k for all k
and OLS provides **unbiased estimates** of β
(unbiased = forventningsrett)

Spring 2010

© Erling Berge 2010

235

II GAUSS-MARKOV REQUIREMENTS (ii)

- (3) Errors have a constant variance for all i

$$\bullet \text{Var}(\varepsilon_i) = \sigma^2 \quad \text{for all } i$$

This is called homoscedasticity

- (4) Errors are uncorrelated with each other

$$\bullet \text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \quad \text{for all } i \neq j$$

This is called no autocorrelation

Spring 2010

© Erling Berge 2010

236

II GAUSS-MARKOV REQUIREMENTS (iii)

Given (3) and (4) in addition to (1) and (2) provides:

- a. Estimates of standard errors of regression coefficients are unbiased and

- b. The **Gauss-Markov theorem**:

OLS estimates have **less variance** than any other linear unbiased estimate (including ML estimates)

OLS gives “BLUE”
(Best Linear Unbiased Estimate)

Spring 2010

© Erling Berge 2010

237

II GAUSS-MARKOV REQUIREMENTS (iv)

(1) - (4) are called the GAUSS-MARKOV requirements

- Given (2) - (4) with an additional requirement that errors are uncorrelated with x-variables:

$$\bullet \text{cov}(x_{ik}, \varepsilon_i) = 0 \quad \text{for all } i, k$$

The coefficients and standard errors are consistent (converging in probability to the true population value as sample size increases)

Spring 2010

© Erling Berge 2010

238

Footnote 1: Unbiased estimators

- Unbiased means that

$$E[b_k] = \beta_k$$
- In the long run we are bound to find the population value - β_k - if we draw sufficiently many samples, calculate b_k and average these

Spring 2010

© Erling Berge 2010

239

Footnote 2: Consistent estimators

- An estimator is consistent if we as sample size (n) grows towards infinity, find that b approaches β and s_b [or SE_b] approaches σ_β
- b_k is a consistent estimator of β_k if we for any small value of c have

$$\lim_{n \rightarrow \infty} [\Pr\{ |b_k - \beta_k| < c \}] = 1$$

Spring 2010

© Erling Berge 2010

240

Footnote 3: In BLUE "Best" means
minimal variance estimator

- Minimal variance or efficient estimator means that
 $\text{var}(b_k) < \text{var}(a_k)$ for all estimators a different from b
- Equivalent:
 $E[b_k - \beta_k]^2 < E[a_k - \beta_k]^2$ for all estimators a unlike b

Spring 2010

© Erling Berge 2010

241

Footnote 4:

Biased estimators

- Even if the requirements ensuring that our estimates are BLUE one may at times find biased estimators with less variance such as in
- Ridge Regression

Non-linear estimators

- There may be **non-linear estimators** that are unbiased and with less variance than BLUE estimators

Spring 2010

© Erling Berge 2010

242

III NORMALLY DISTRIBUTED ERROR TERM

- (5) If all errors are normally distributed with expectation 0 and standard deviation of σ^2 , that is if

$$\varepsilon_i \sim N(0, \sigma^2) \quad \text{for all } i$$

- Then we can test hypotheses about β and σ , and
- OLS estimates will have less variance than estimates from all other unbiased estimators
- OLS results are “BUE”

(Best Unbiased Estimate)

Spring 2010

© Erling Berge 2010

243

Problems in regression analysis that cannot be tested

- If all relevant variables are included
- If x-variables have measurement errors
- If the expected value of the error is 0

This means that we are unable to check if the correlation between the error term and x-variables actually is 0

OLS constructs residuals so that $\text{cov}(x_{ik}, e_i) = 0$

This is in reality saying the same as the first point that we are unable to test if all relevant variables are included

Spring 2010

© Erling Berge 2010

244

Problems in regression analysis that can be tested (1)

- Non-linear relationships
- Inclusion of an irrelevant variable
- Non-constant variance of the error term (heteroscedasticity)
- Autocorrelation for the error term
- Correlations among error terms
- Non-normal error terms
- Multicollinearity

Spring 2010

© Erling Berge 2010

245

Consequences of problems (Hamilton, p113)

Requirement	Problem	Unwanted properties of estimates			
		Biased estimate of b	Biased estimate of SE_b	Invalid t&F-tests	High var[b]
Specification	Non-linear relationship	X	X	X	-
-"-	Excluded relevant variable	X	X	X	-
-"-	Included irrelevant variable	0	0	0	X
Gauss-Markov	X with measurement error	X	X	X	-
-"-	Heteroscedasticity	0	X	X	X
-"-	Autocorrelation	0	X	X	X
-"-	X correlated with ϵ	X	X	X	-
Normal distribution	ϵ not normally distributed	0	0	X	X
... no requirement	Multicollinearity	0	0	0	X

Spring 2010

© Erling Berge 2010

246

Problems in regression analysis that can be discovered (2)

- Outliers (extreme y-values)
- Influence (cases with large influence: unusual combinations of y and x-values)
- Leverage (potential for influence)

Spring 2010

© Erling Berge 2010

247

Tools for discovering problems

- Studies of
 - One-variable distributions (frequency distributions and histogram)
 - Two-variable co-variation (correlation and scatter plot)
 - Residual (distribution and covariation with predicted values)

Spring 2010

© Erling Berge 2010

248

Correlation and scatter plot

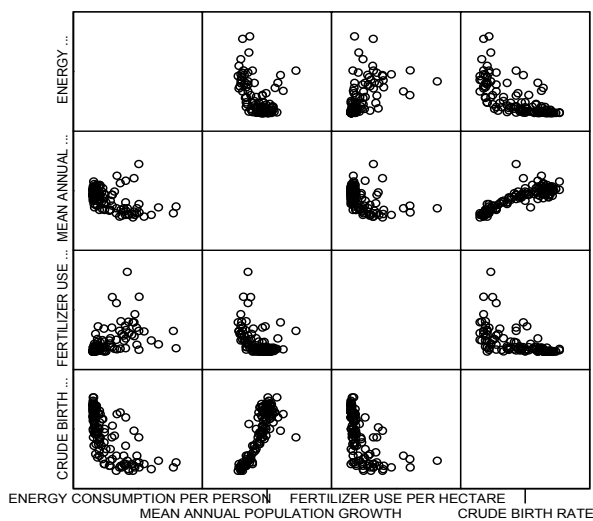
Data from 122 countries		ENERGY CONSUMPTION PER PERSON	MEAN ANNUAL POPULATION GROWTH	FERTILIZER USE PER HECTARE	CRUDE BIRTH RATE
ENERGY CONSUMPTION PER PERSON	Pearson Correlation	1	-,505	,533	-,689
	N	125	122	125	122
MEAN ANNUAL POPULATION GROWTH	Pearson Correlation	-,505	1	-,469	,829
	N	122	125	125	125
FERTILIZER USE PER HECTARE	Pearson Correlation	,533	-,469	1	-,589
	N	125	125	128	125
CRUDE BIRTH RATE	Pearson Correlation	-,689	,829	-,589	1
	N	122	125	125	125

Spring 2010

© Erling Berge 2010

249

Correlation and scatter plot



Spring 2010

© Erling Berge 2010

250

Heteroscedasticity

(non-constant variance of error term) can arise from:

- Measurement error (e.g. y more accurate the larger x is)
- Outliers
- If ε_i contains an important variable that varies with both x and y (specification error)
- Specification error is the same as the wrong model and may cause heteroscedasticity
- An important diagnostic tool is a plot of the residual against predicted value (\hat{Y})

Spring 2010

© Erling Berge 2010

251

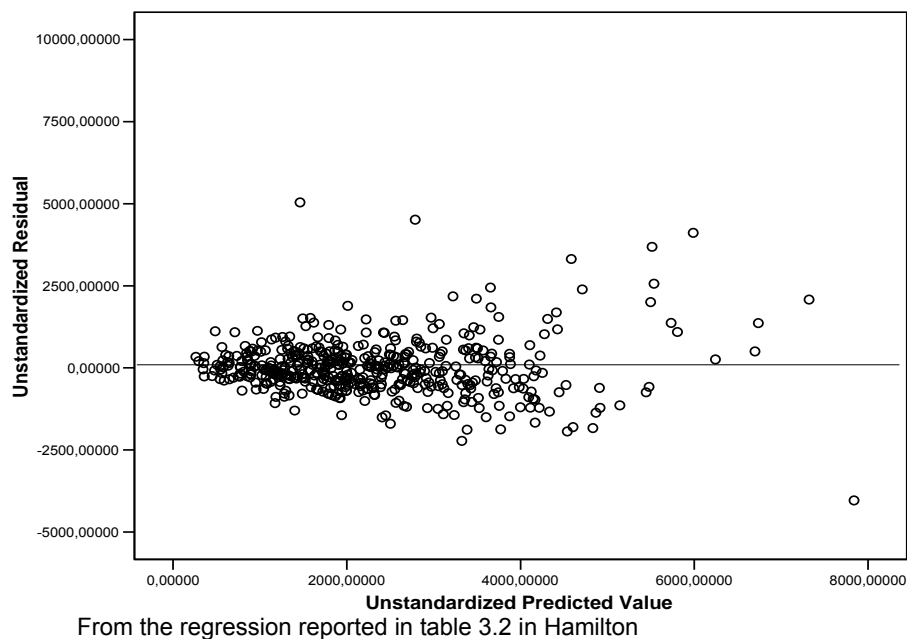
Example: Hamilton table 3.2

Dependent Variable: Summer 1981 Water Use	Unstandardized Coefficients		t	Sig.
	B	Std. Error		
(Constant)	242,220	206,864	1,171	,242
Income in Thousands	20,967	3,464	6,053	,000
Summer 1980 Water Use	,492	,026	18,671	,000
Education in Years	-41,866	13,220	-3,167	,002
head of house retired?	189,184	95,021	1,991	,047
# of People Resident 1981	248,197	28,725	8,641	,000
Increase in # of People	96,454	80,519	1,198	,232

Spring 2010

© Erling Berge 2010

252



Spring 2010

© Erling Berge 2010

253

Footnote for the previous figure

- There is heteroscedasticity if the variation of the residual (variation around a typical value) varies systematically with the value of one or more x-variables
- The figure shows that the variation of the residual increases with increasing predicted y : \hat{y}
- Predicted y (\hat{y}) is in this case an index showing high average x-values
- When the variation of the residual varies systematically with the values of the x-variables like this, we conclude with heteroscedasticity

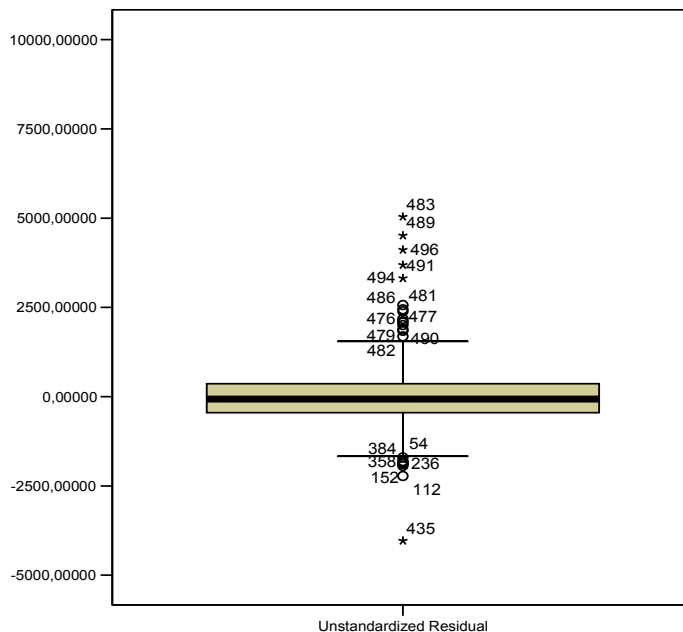
Spring 2010

© Erling Berge 2010

254

- Box-plot of the residual shows
- Heavy tails
 - Many outliers
 - Weakly positively skewed distribution

Will any of the outliers affect the regression?

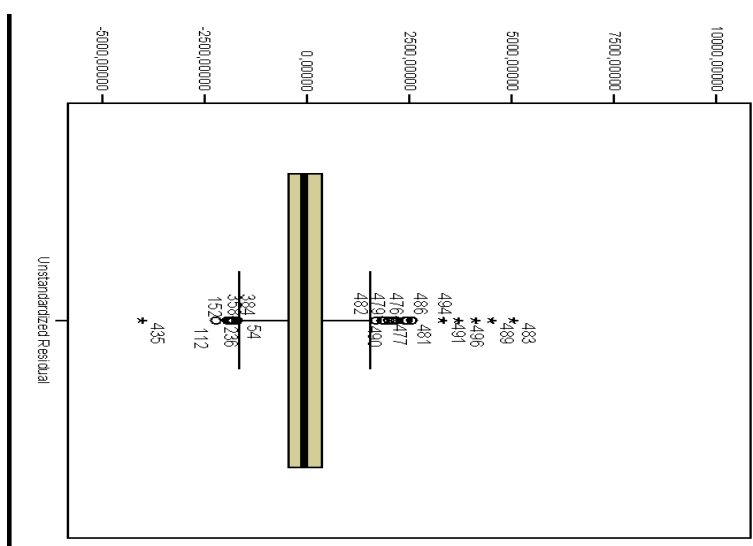


Spring 2010

© Erling Berge 2010

255

The distribution seen from another angle



Spring 2010

© Erling Berge 2010

256

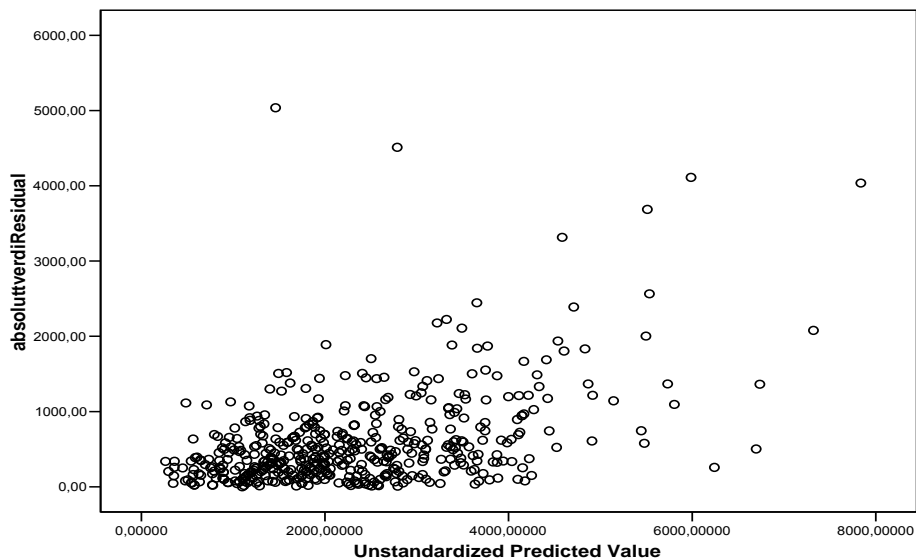
Band-regression

- Homoscedasticity means that the median (and the average) of the absolute value of the residual, i.e.: $\text{median}\{|e_i|\}$, should be about the same for all values of the predicted y_i
- If we find that the median of $|e_i|$ for given predicted values of y_i changes systematically with the value of predicted y_i (\hat{y}_i) it indicates heteroscedasticity
- Such analyses can easily be done in SPSS

Spring 2010

© Erling Berge 2010

257



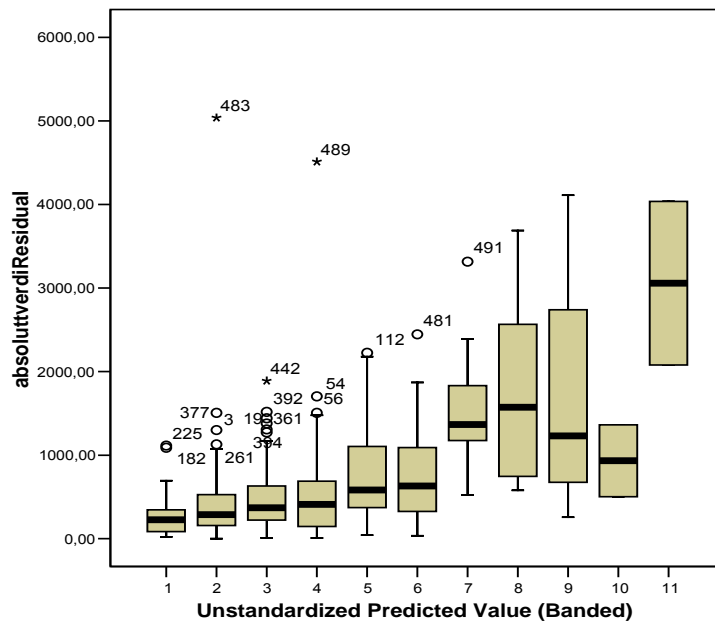
Absolute value of e_i (Based on regression in table 3.2 in Hamilton)

Spring 2010

© Erling Berge 2010

258

Approximate band regression (cpr figure 4.4 in Hamilton)



Spring 2010

© Erling Berge 2010

259

Band regression in SPSS

- Start by saving the residual and predicted y from the regression
- Compute a new variable by taking the absolute value of the residual (Use “compute” under the “transform” menu)
- Then partition the predicted y into bands by using the procedure ”Visual bander” under the ”Transform” menu
- Then use ”Box plot” under ”Graphs” where the absolute value of the residual is specified as variable and the band variable as category axis

Spring 2010

© Erling Berge 2010

260

Footnote to Eikemo and Clausen 2007

- Page 121 describes White's test of Heteroscedasticity
- The description **is wrong**
- They say to replace y with e^2 in the regression on all the x variables
- That is not sufficient.
- The x -variables have to be replaced by all unique cross products of x with x (including x^2)
- Unique elements of the Kronecker product of x with x (where x is the vector of x -variables)

Spring 2010

© Erling Berge 2010

261

Autocorrelation (1)

- Correlation among variable values on the same variable across different cases
(e.g. between ε_i and ε_{i-1})
- Autocorrelation leads to larger variance and biased estimates of the standard error - similar to heteroscedasticity
- In a simple random sample from a population autocorrelation is improbable

Spring 2010

© Erling Berge 2010

262

Autocorrelation (2)

- Autocorrelation is the result of a wrongly specified model. A variable is missing
- Typically it is found in time series and geographically ordered cases
- Tests (e.g. Durbin-Watson) is based on the sorting of the cases. Hence:
- A hypothesis about autocorrelation needs to specify the sorting order of the cases

Spring 2010

© Erling Berge 2010

263

Durbin-Watson test (1)

$$d = \frac{\sum_{i=2}^n (e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2}$$

Should not be used for **autoregressive models**, i.e. models where the y-variable also is an x-variable, see table 3.2

Spring 2010

© Erling Berge 2010

264

Durbin-Watson test (2)

- The sampling distribution of the d-statistic is known and tabled as d_L and d_U (table A4.4 in Hamilton), the number of degrees of freedom is based on n and $K-1$
- Test rule:
 - Reject if $d < d_L$
 - Do not reject if $d > d_U$
 - If $d_L < d < d_U$ the test is inconclusive
- $d=2$ means uncorrelated residuals
- Positive autocorrelation results in $d < 2$
- Negative autocorrelation results in $d > 2$

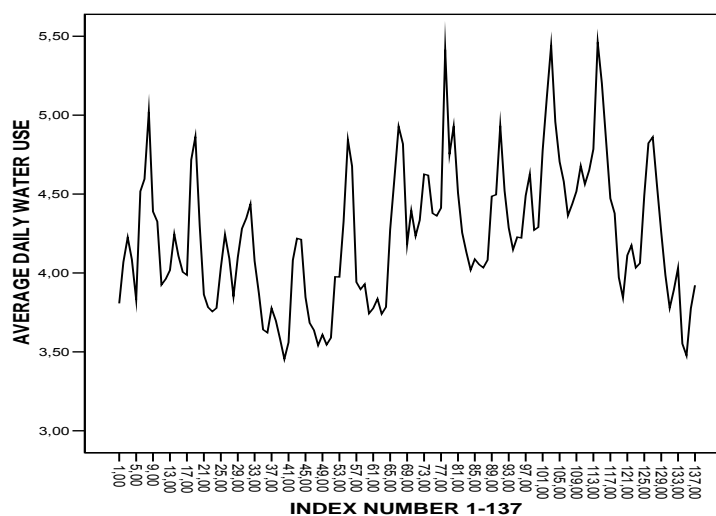
Spring 2010

© Erling Berge 2010

265

Daily water use, average pr month

Example:



Spring 2010

© Erling Berge 2010

266

Ordinary OLS-regression where the case is month

Dependent Variable: AVERAGE DAILY WATER USE	Unstandardized Coefficients		t	Sig.
	B	Std. Error		
(Constant)	3,828	,101	38,035	,000
AVERAGE MONTHLY TEMPERATURE	,013	,002	7,574	,000
PRECIPITATION IN INCHES	-,047	,021	-2,234	,027
CONSERVATION CAMPAIGN DUMMY	-,247	,113	-2,176	,031

Predictors: (Constant), CONSERVATION CAMPAIGN DUMMY, AVERAGE MONTHLY TEMPERATURE, PRECIPITATION IN INCHES

Spring 2010

© Erling Berge 2010

267

Test of autocorrelation

Dependent Variable: AVERAGE DAILY WATER USE	R	R Square	Adjusted R Square	Std. Error of the Estimate	Durbin-Watson
1	,572(a)	,327	,312	,36045	,535

Predictors: (Constant), CONSERVATION CAMPAIGN DUMMY, AVERAGE MONTHLY TEMPERATURE, PRECIPITATION IN INCHES

N = 137, K-1 = 3

Find limits for rejection / acceptance of the null hypothesis of no autocorrelation with level of significance 0,05

Tip: Look up table A4.4 in Hamilton, p355

Spring 2010

© Erling Berge 2010

268

Autocorrelation coefficient

m-th order autocorrelation coefficient

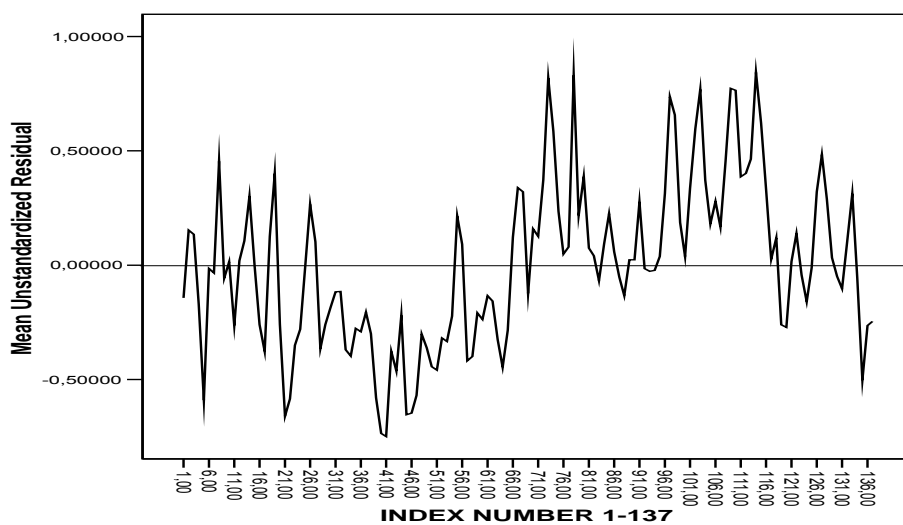
$$r_m = \frac{\sum_{t=1}^{T-m} (e_t - \bar{e})(e_{t+m} - \bar{e})}{\sum_{t=1}^T (e_t - \bar{e})^2}$$

Spring 2010

© Erling Berge 2010

269

Residual "Daily water use", month



Spring 2010

© Erling Berge 2010

270

Smoothing with 3 points

- Sliding average

$$e_t^* = \frac{e_{t-1} + e_t + e_{t+1}}{3}$$

- "Hanning"

$$e_t^* = \frac{e_{t-1}}{4} + \frac{e_t}{2} + \frac{e_{t+1}}{4}$$

- Sliding median

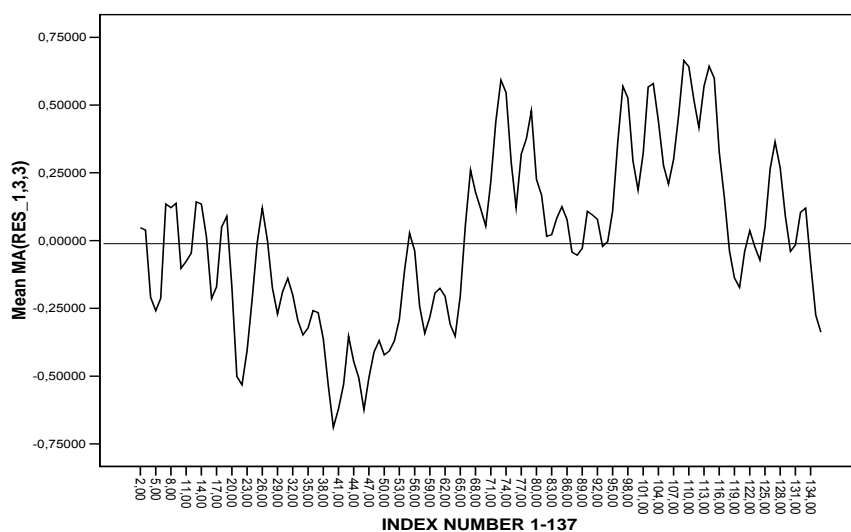
$$e_t^* = \text{median}\{e_{t-1}, e_t, e_{t+1}\}$$

Spring 2010

© Erling Berge 2010

271

Residual, smoothing once

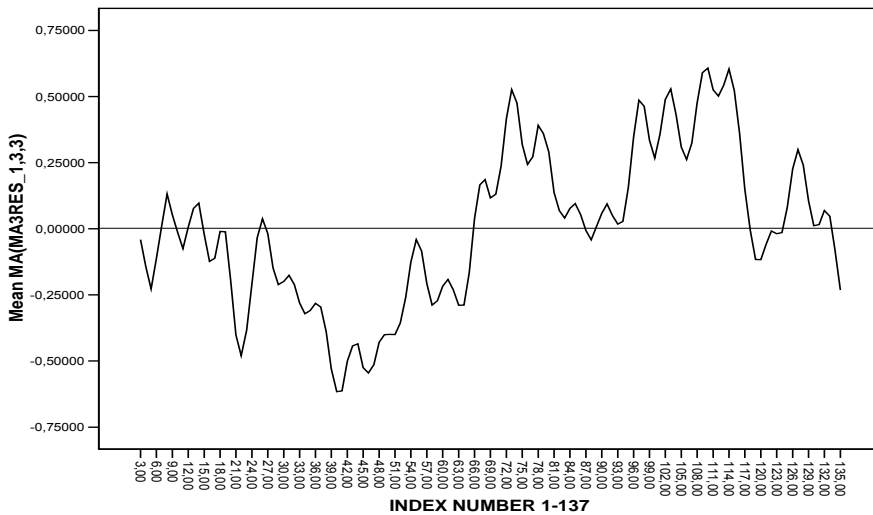


Spring 2010

© Erling Berge 2010

272

Residual, smoothing twice

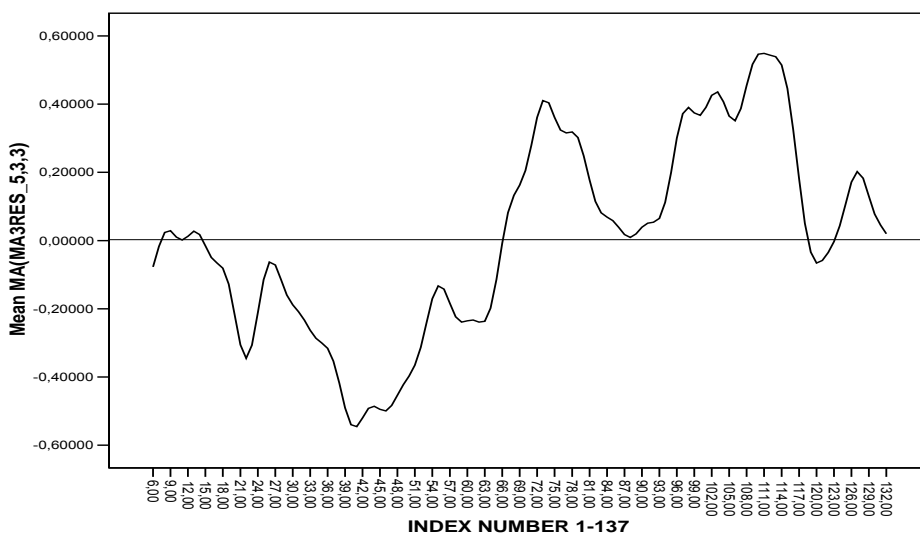


Spring 2010

© Erling Berge 2010

273

Residual, smoothing five times



Spring 2010

© Erling Berge 2010

274

Consequences of autocorrelation

- Tests of hypotheses and confidence intervals are unreliable. Regressions may nevertheless provide a good description of the sample. Parameters are unbiased
- Special programs can estimate standard errors consistently
- Include in the model variables affecting neighbouring cases
- Use techniques developed for time series analysis (e.g.: analyse the difference between two points in time, Δy)

Spring 2010

© Erling Berge 2010

275

Concluding on Autocorrelation

- Correlation among variable values on the same variable across different cases (e.g. between ε_i and ε_{i-1})
- Autocorrelation leads to larger variance and biased estimates of the standard error - similar to heteroscedasticity
- Autocorrelation is the result of a wrongly specified model
- Typically it is found in time series and geographically ordered cases. In a simple random sample from a population autocorrelation is improbable
- Tests (e.g. Durbin-Watson) is based on the sorting of the cases. Hence: hypotheses about autocorrelation need to specify the sorting order of the cases

Spring 2010

© Erling Berge 2010

276

Analyses of models are based on assumptions

- OLS is a simple technique of analysis with very good theoretical properties. But
- The good properties are based on certain assumptions
- If the assumptions do not hold the good properties evaporates
- Investigating the degree to which the assumptions hold is the most important part of the analysis

Spring 2010

© Erling Berge 2010

277

OLS-REGRESSION: assumptions

- I SPECIFICATION REQUIREMENT
 - The model is correctly specified
- II GAUSS-MARKOV REQUIREMENTS
 - (1) x is known, without stochastic variation
 - (2) Errors have an expected value of 0 for all i
 - (3) Errors have a constant variance for all i
 - (4) Errors are uncorrelated with each other

(Ensures that the estimates are “BLUE”)
- III NORMALLY DISTRIBUTED ERROR TERM
 - Ensures that the tests are valid

Spring 2010

© Erling Berge 2010

278

Problems in regression analysis that cannot be tested

- If all relevant variables are included
- If x-variables have measurement errors
- If the expected value of the error is 0
- (This means that we are unable to check if the correlation between the error term and x-variables actually is 0 and is actually the same as the first point that we are unable to test if the model is correctly specified)

Spring 2010

© Erling Berge 2010

279

The most important problems in regression analysis that can be tested

- Non-linear relationships
- Non-constant error of the error term (heteroscedasticity)
- Autocorrelation for the error term
- Non-normal error terms

Spring 2010

© Erling Berge 2010

280

More on Heteroscedasticity

- Is present if the variance of the error term varies with the size of x-values
- Predicted y is an indicator of the size of x-values (hence scatter plot of residual against predicted y)
- Heteroscedasticity (non-constant variance of error term) can arise from
 - Measurement error (e.g. y more accurate the larger x is)
 - Outliers
 - The wrong functional form
 - If ε_i contain an important variable that varies with one or more x and y. The error term ε_i is not independent of the x-es. Hence the Gauss-Markov requirements 1 and 2 cannot be correct.

Spring 2010

© Erling Berge 2010

281

Indicators of heteroscedasticity

- Inspection of the scatter plot of residual against predicted value of y
- Band regression of the scatter plot

An interesting option here is:

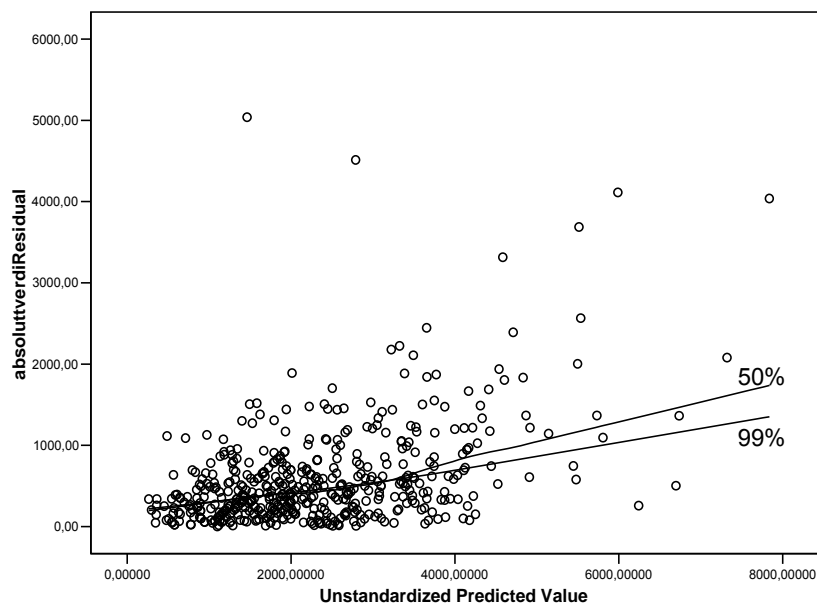
- Locally weighted / "sliding" regression on the central part of the sample

Spring 2010

© Erling Berge 2010

282

"Sliding"
adapted line
by means of
locally
weighted
OLS
regression
The
procedure is
called
LOESS (see
next slide)



Spring 2010

© Erling Berge 2010

283

A footnote: SPSS explains

Fit Lines

- In a fit line, the data points are fitted to a line that usually does not pass through all the data points. The fit line represents the trend of the data. Some fit lines are regression based. Others are based on iterative weighted least squares.
- Fit lines apply to scatter plots. You can create fit lines for all of the data values on a chart or for categories, depending on what you select when you create the fit line.

Loess

- Draws a fit line using iterative weighted least squares. At least 13 data points are needed. This method fits a specified percentage of the data points, with the default being 50%. In addition to changing the percentage, you can select a specific kernel function. The default kernel (probability function) works well for most data.

Spring 2010

© Erling Berge 2010

284

Non-normal residuals

- Imply that t- and F-tests cannot be used
- Since OLS estimates of parameters are easily affected by outliers, heavy tails in the distribution of the residual will indicate large variation in estimates from sample to sample
- We can test the assumption of normally distributed error term by inspecting the distribution of the residual, e.g. by inspecting
 - Histogram, box plot, or quantile-normal plot
 - There are also more formal tests (but not very useful) based on skewness and kurtosis

Spring 2010

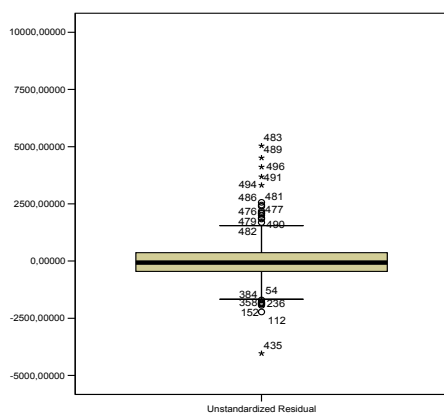
© Erling Berge 2010

285

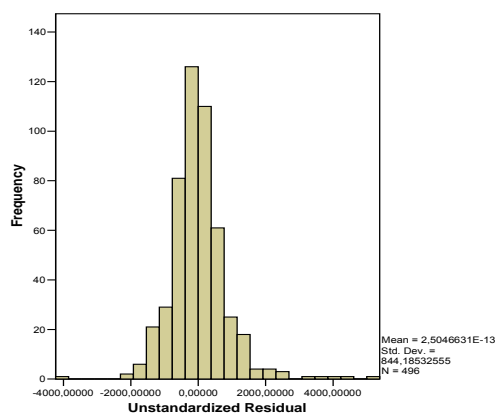
Diagram of the residual shows:

Heavy tails, many outliers, and weakly positively skewed distribution

BOX PLOT



HISTOGRAM

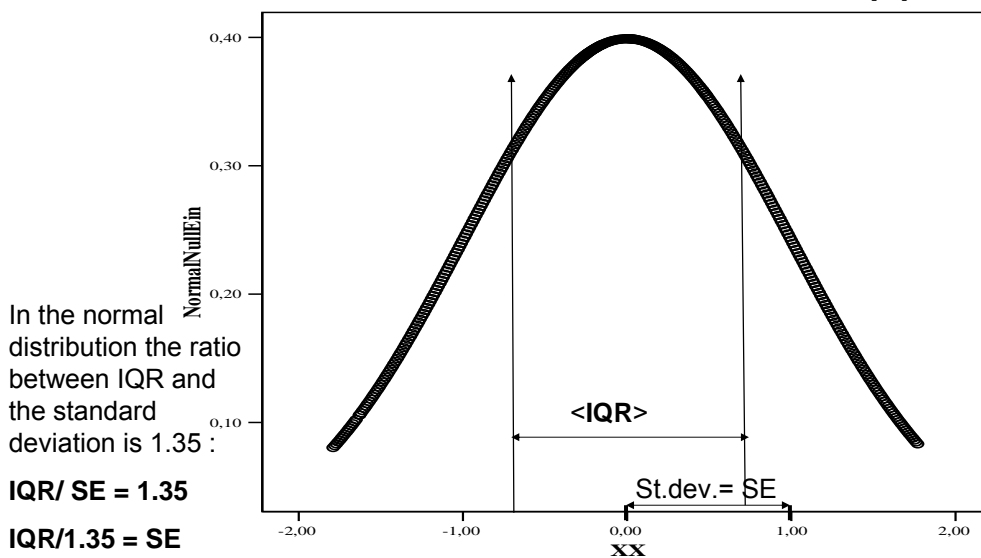


Spring 2010

© Erling Berge 2010

286

Skewed distribution of the residual (1)



Spring 2010

© Erling Berge 2010

287

Skewed distribution of the residual (2)

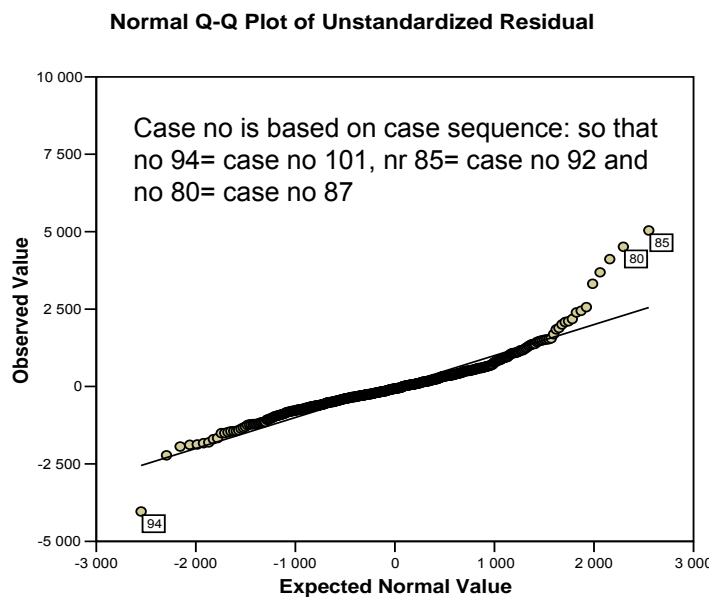
- Since the average of the residuals (e_i) always equals 0, the distribution will be skewed if the median is unequal to 0
- It is known that for the normal distribution the standard deviation (or the standard error) equals approximately $IQR/1.35$
- If the distribution of the residual is symmetric we can compare SE_e to $IQR/1.35$. If
 - $SE_e > IQR/1.35$ the tails are heavier than the normal distribution
 - $SE_e \approx IQR/1.35$ the tails are approximately equal to the normal distribution
 - $SE_e < IQR/1.35$ the tails are lighter than the normal distribution

Spring 2010

© Erling Berge 2010

288

Quantile-
Normal plot
of residual
from
regression
in table 3.2
in Hamilton



Spring 2010

© Erling Berge 2010

289

Options if non-normality is found

- Test out if the right function has been used
- Test out if some important variable has been excluded
 - If the model cannot be improved substantially, we may try transforming the dependent variable to symmetry
- Test out if lack of normality is caused by outliers or influential cases
 - If there are outliers, transforming of the variable where the case is outlier may help

Spring 2010

© Erling Berge 2010

290

Influence (1)

- A case (or observation) has influence if the regression result changes when the case is excluded
- Some cases have unusually large influence because of
 - Unusually large y-value (outliers)
 - Unusually large value on an x-variable
 - Unusual combinations of variable values

Spring 2010

© Erling Berge 2010

291

Influence (2)

- We can see if a case has influence by comparing regressions with and without a particular case. One may for example
- Inspect the difference between b_k and $b_{k(i)}$ where case no i has been excluded in the estimation of the last coefficient
- This difference measured relative to the standard error of $b_{k(i)}$ is called $DFBETAS_{ik}$

Spring 2010

© Erling Berge 2010

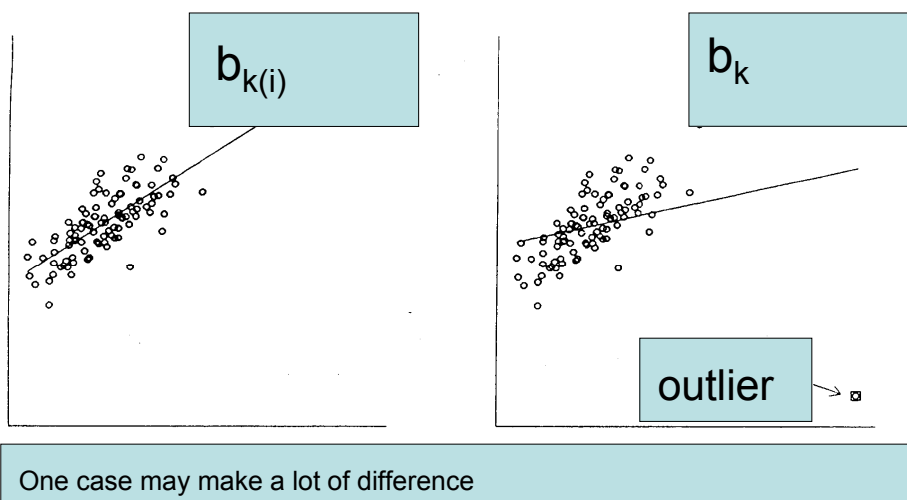
292

DFBETAS_{ik}

$$DFBETAS_{ik} = \frac{b_k - b_{k(i)}}{\frac{s_{e(i)}}{\sqrt{RSS_k}}}$$

$s_{e(i)}$ is the standard deviation of the residual when case no i has been excluded from the analysis RSS_k is Residual Sum of Squares from the regression of x_k on all other x -variables

DFBETAS_{ik} :



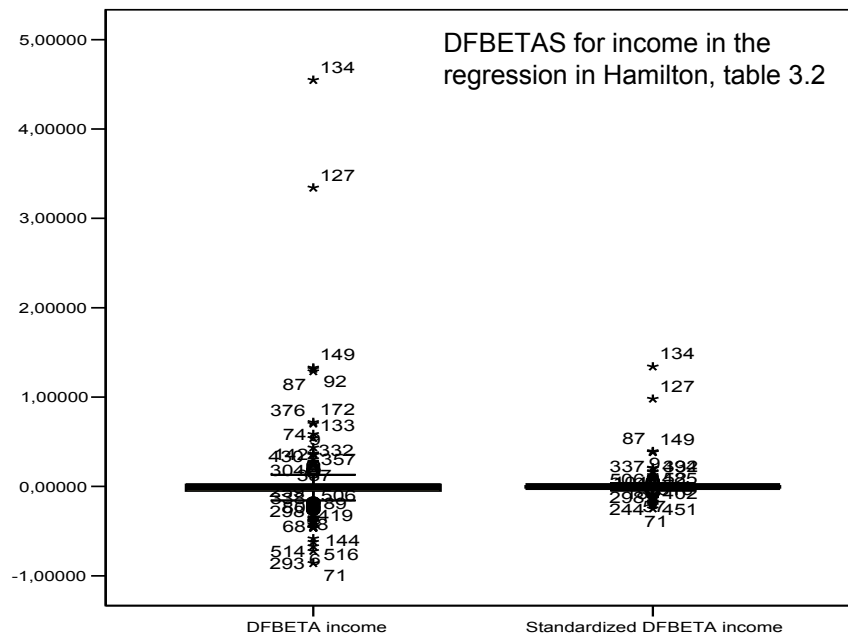
What is a large DFBETAS?

- $DFBETAS_{ik}$ is calculated for every independent variable for every case. We do not want to inspect all values for it
- Three criteria for finding large values we need to inspect are
 - External scaling. $|IDFBETAS_{ik}| > 2/\text{SQRT}(n)$
 - Internal scaling. Look for **severe outliers** in the box plot of $DFBETAS_{ik}$:
 $DFBETAS_{ik} < Q_1 - 3IQR$
 $Q_3 + 3IQR < DFBETAS_{ik}$
 - Gap in the distribution of $DFBETAS_{ik}$
- None of the $DFBETAS_{ik}$ needs to be problematic

Spring 2010

© Erling Berge 2010

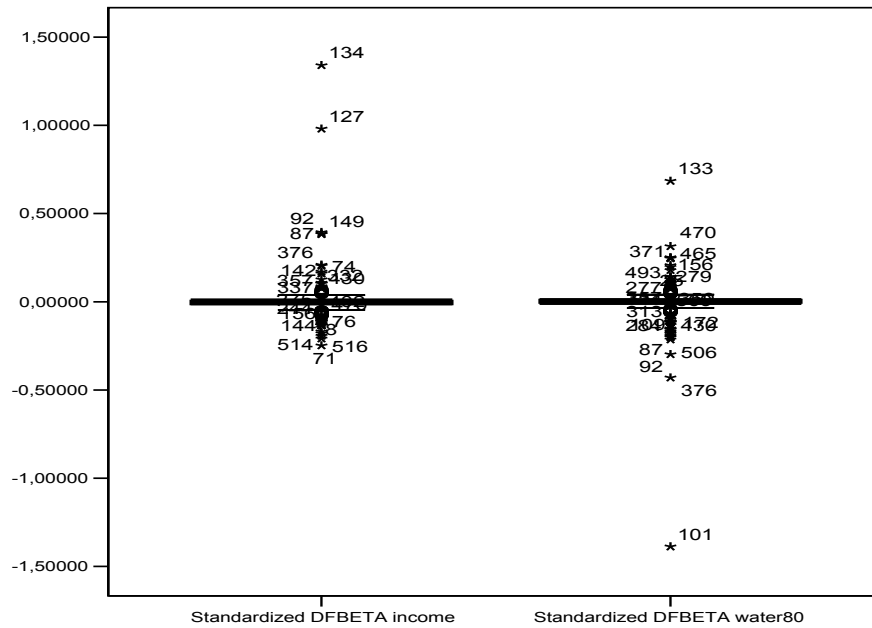
295



Spring 2010

© Erling Berge 2010

296



Spring 2010

© Erling Berge 2010

297

Sequence in the data set and case no is not the same.
Case no is fixed. Variable values.

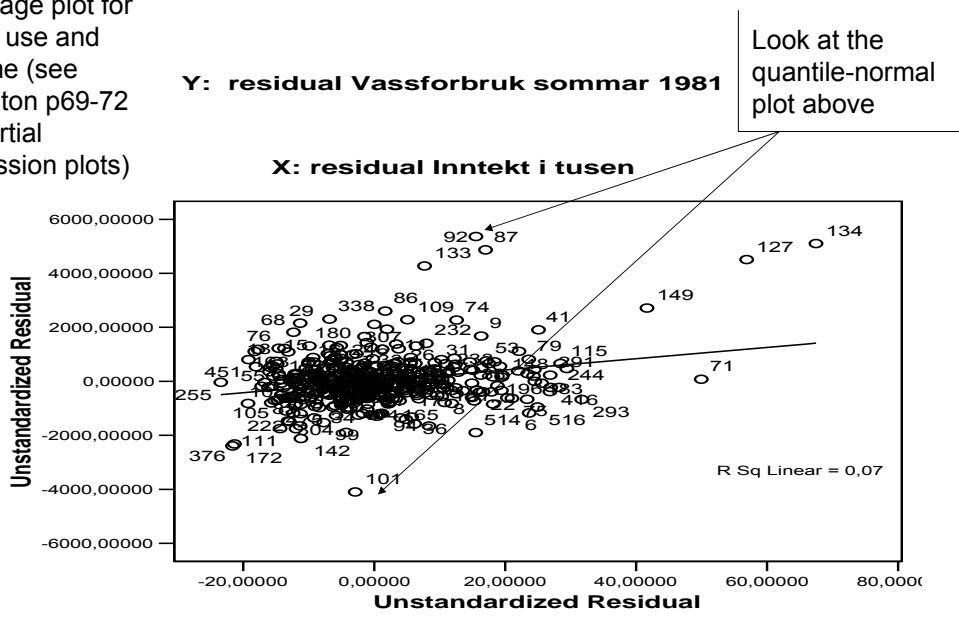
Sequence no	Case nr	water81	water80	water79	educat	retire	peop81	cpeop
91	98	1500	1300	1500	16	0	2	0
92	99	3500	6500	5100	14	0	6	0
93	100	1000	1000	2700	12	1	1	0
94	101	3800	12700	4800	20	0	5	0
95	102	4100	4500	2600	20	0	5	0
96	103	4200	5600	5400	16	0	5	-1
97	104	2400	2700	800	16	0	6	0
98	105	1600	2300	2200	14	0	4	0
99	107	2300	2300	3100	16	0	4	-2

Spring 2010

© Erling Berge 2010

298

Leverage plot for water use and income (see Hamilton p69-72 on partial regression plots)



Spring 2010

© Erling Berge 2010

299

Consequences of case with large influence

- If we discover cases with large influence we should not remove them from the analysis unless they contain serious errors
- Take a careful look at influential cases, maybe there are measurement errors
- When influential cases are outliers their influence can be reduced by transformation
- Use robust regression not so easily affected as OLS regression
- **If no errors are found report results both with and without one or two of the most influential cases**

Spring 2010

© Erling Berge 2010

300

Potential influence: leverage

- The potential for influence of a case from a particular combination of x-values is measured by the hat statistic h_i
- h_i varies from $1/n$ to 1. It has an average of K/n ($K = \#$ parameters)
- SPSS reports **the centred h_i**
 - i.e. $(h_i - K/n)$, we may call this for h_i^c
 - We must compute the normal $h_i = h_i^c + K/n$ to judge the size by the criteria supplied by Hamilton

Spring 2010

© Erling Berge 2010

301

What is a large value of leverage?

- As for DFBETAS different criteria can be suggested. They all depend on the sample size n
 - If $h_i > 2K/n$ (or $h_i^c > K/n$) we find the ca 5% largest h_i ; alternatively
 - If $\max(h_i) \leq 0.2$ there is no problem
 - If $0.2 \leq \max(h_i) \leq 0.5$ there is some risk for a problem
 - If $0.5 \leq \max(h_i)$ probably there is a problem

Spring 2010

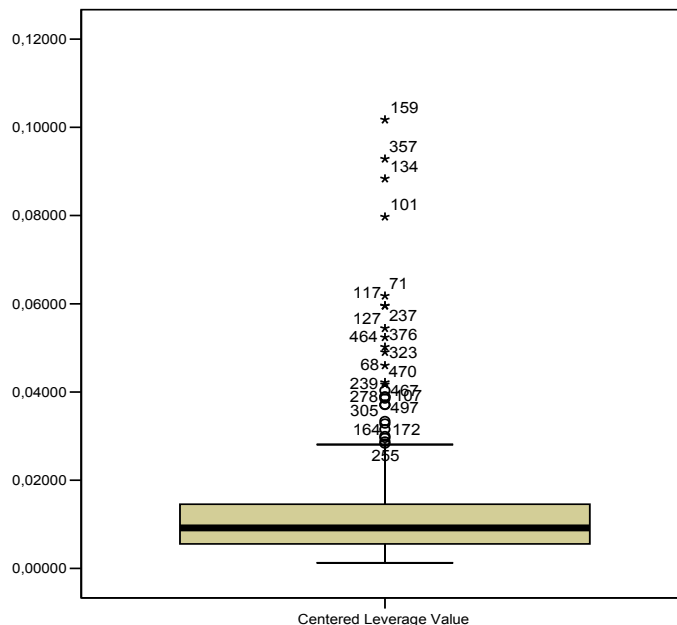
© Erling Berge 2010

302

Centred leverage (h^c_i) from the regression in table 3.2 in Hamilton

Max av h^c_i er 0.102

Or max of $h_i = 0.102 + K/n = 0.102 + 7/496 = 0.116 < 0.2$

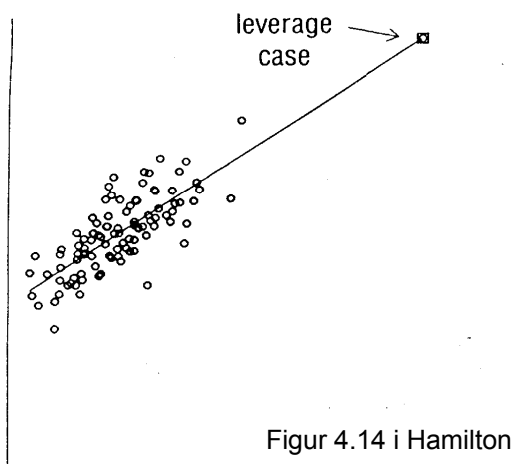


Spring 2010

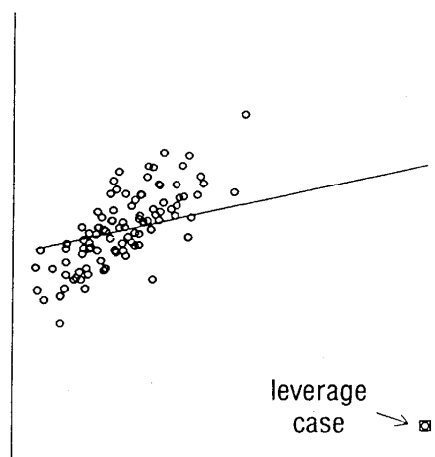
© Erling Berge 2010

303

The difference between influence and leverage



High Leverage, Low Influence



High Leverage, High Influence

Spring 2010

© Erling Berge 2010

304

The leverage statistic is found in many other case statistics

- Variance of the i -th residual

$$\text{var}[e_i] = s_e^2[1 - h_i]$$

- Standardized residual (*ZRESID in SPSS)

$$z_i = \frac{e_i}{s_e \sqrt{1 - h_i}}$$

- Studentized residual (*SRESID in SPSS)

$$t_i = \frac{e_i}{s_{e(i)} \sqrt{1 - h_i}}$$

- And remember that the standard deviation of the residual is

$$s_e = \sqrt{RSS / (n - K)}$$

Spring 2010

© Erling Berge 2010

305

Total influence: Cook's D_i

- Cook's distance D_i measure influence on the model as a whole, not on a specific coefficient as $DFBETAS_{ik}$

$$D_i = \frac{z_i^2 h_i}{K(1 - h_i)}$$

where z_i is the standardized residual

and h_i is the hat statistic (leverage)

Spring 2010

© Erling Berge 2010

306

What is a large D_i ?

- One might want to take a look at all
 - $D_i > 1$ or
 - $D_i > 4/n$ these are about the 5% largest D_i
- Even if a case has low D_i it may still be the case that it affects the size of single coefficients (it has a large $DFBETAS_{ik}$)

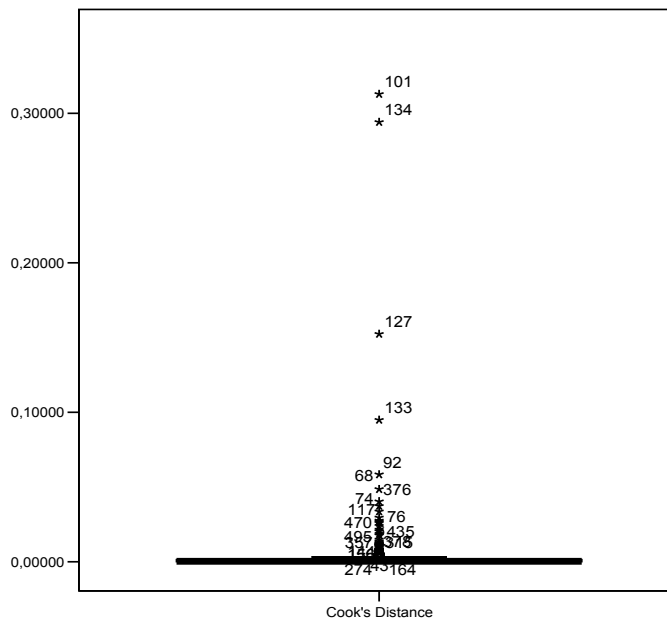
Spring 2010

© Erling Berge 2010

307

Cook's distance D_i
from the regression
in table 3.2 in
Hamilton

Also see table 4.4
(p133) in Hamilton



Spring 2010

© Erling Berge 2010

308

Summarizing

What can be done with outliers and cases with large influence? We can

- Investigate if data are erroneous. If data are wrong the case can be removed from the analysis
- Investigate if transformation to symmetry helps
- Report two equations: with and without cases with unreasonably large influence
- Get more data

Spring 2010

© Erling Berge 2010

309

Multicollinearity

- Means very high intercorrelations among x-variables
- Check if parameter estimates are correlated
- Check if tolerance (the part of the variation of x that is not shared with other variables) is less than say 0.1. If so there may be a problem
- $VIF = \text{variance inflation factor} = 1/\text{tolerance}$
- If multicollinearity is caused by squaring of variables or interaction terms it should not be seen as problematic

Spring 2010

© Erling Berge 2010

310

Tolerance

- The amount of variation in a variable x_k unique to that variable is called the tolerance of the variable
- Let R_k^2 be the coefficient of determination in the regression of x_k on all the rest of the x -variables. The other x -variables explain the proportion R_k^2 of the variation in x_k .
- Then $1 - R_k^2$ is the unique variation:
 - Tolerance = $1 - R_k^2$
- Perfect multicollinearity means that
 - $R_k^2 = 1$ and tolerance = 0
- Low values of tolerance make regression results less precise (larger standard errors)

Spring 2010

© Erling Berge 2010

311

Variance Inflation Factor (VIF)

- $1/\text{tolerance} = 1/(1 - R_k^2) = \text{VIF}$
- The standard error of the regression coefficient b_k can be written

$$SE_{b_k} = \frac{s_e}{\sqrt{RSS_k}} = \frac{s_e}{\sqrt{(1 - R_k^2)TSS_k}} = \sqrt{\text{VIF}} \frac{s_e}{\sqrt{TSS_k}}$$

- Other things being equal lower tolerance (larger VIF) for x_k will give higher standard error for b_k [SE increase with a factor equal to square root of VIF]

Spring 2010

© Erling Berge 2010

312

Indicators of multicollinearity

- The best indicator is tolerance or VIF (both are based on R^2_k)
- Other indicators are
 - Correlation among single variables (not reliable)
 - Inclusion/ exclusion of single variables give large changes in the effect of other variables
 - Unexpected signs on the effects of some variable
 - Standardized regression coefficients larger than 1 or less than -1
 - Correlation among parameter estimates

Spring 2010

© Erling Berge 2010

313

Tolerance and VIF from regression in table 3.2 in Hamilton

Dependent Variable: Summer 1981 Water Use	Unstandardized Coefficients		t	Sig.	Collinearity Statistics	
	B	Std. Error			Tolerance	VIF
(Constant)	242,220	206,864	1,171	,242		
Summer 1980 Water Use	,492	,026	18,671	,000	,675	1,482
Income in Thousands	20,967	3,464	6,053	,000	,712	1,404
Education in Years	-41,866	13,220	-3,167	,002	,873	1,145
head of house retired?	189,184	95,021	1,991	,047	,776	1,289
# of People Resident, 1981	248,197	28,725	8,641	,000	,643	1,555
Increase in # of People	96,454	80,519	1,198	,232	,957	1,045

Spring 2010

© Erling Berge 2010

314

What is low tolerance?

When $R^2_k > 0,9$
tolerance is $< 0,1$
and $VIF > 10$

Factor of
multiplication for the
standard error is the
square root of VIF
(ca 3.2 for $R^2_k = 0,9$)

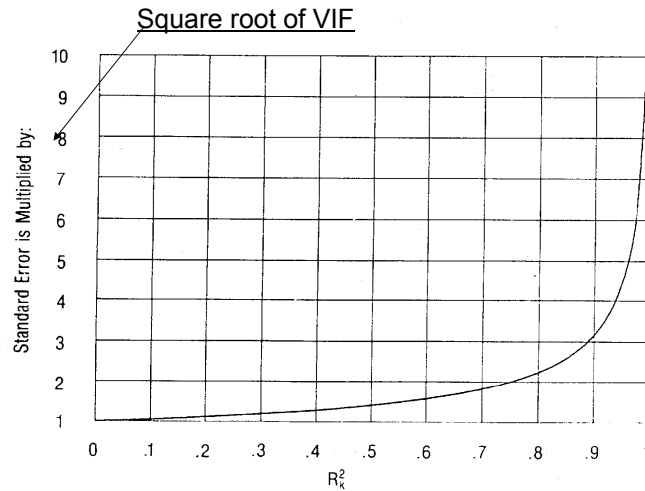


Figure 4.15 Effect of multicollinearity on standard errors (simplified).

Spring 2010

© Erling Berge 2010

315

When is multicollinearity a problem?

- It is not a problem if the reason is curvilinearity or interaction terms in the model. But in testing we need to take account of the fact that if VIF is high parameter estimates are imprecise (high standard errors). They are tested as a group by the F-test
- If the reason is that two variables measure the same concept one of them should be dropped, or they can be combined in an index
- It is a problem if we need estimates of the separate effects of two highly correlated variables (if a test of their joint effect is not sufficient)

Spring 2010

© Erling Berge 2010

316

Summarizing (1)

- When errors are independent and identically normally distributed OLS estimates are as good or better than other possible estimates
- But the assumptions are rarely satisfied completely, we have to test the degree to which they are satisfied
- Many problems can be corrected if we learn about them
- Check early on if curvilinearity, outliers or heteroscedasticity are problems (for example by use of scatter plots)

Spring 2010

© Erling Berge 2010

317

Summarizing (2)

- Do more exact investigations using residual/predicted Y plots and leverage plots
 - Curvilinearity (leverage plot, residual vs predicted Y plot)
 - Heteroscedasticity (leverage plot, [absolute value of] residual against predicted Y plot)
 - Non-normal residuals (quantile-normal plot, box-plot with analysis of median and IQR/1.35)
 - Influence (check DFBETAS and Cook's D)
 - When we do not find serious problems we can have more confidence in our conclusions

Spring 2010

© Erling Berge 2010

318

Fitting Curves Robust Regression

- Hamilton Ch 5 p145-173
- Hamilton Ch 6 p183-212

Spring 2010

© Erling Berge 2010

319

Ch 5 Fitting Curves

- A correctly specified model require that the function linking x-variables and y-variable is true to what really exist: is the relationship linear?
- Data can be inspected by means of band regression or smoothing
- The theory of causal impact can specify a non-linear relationship
- For phenomena that cannot be represented by a line we shall present some alternatives
 - Curvilinear regression
 - Non-linear regression

Spring 2010

© Erling Berge 2010

320

Band regression

- Can be used to explore how the relationship among the variables actually appears
- If we can see a non-linear underlying trend of the data we must through transformations or use of curves find a form for the function better representing the relationship

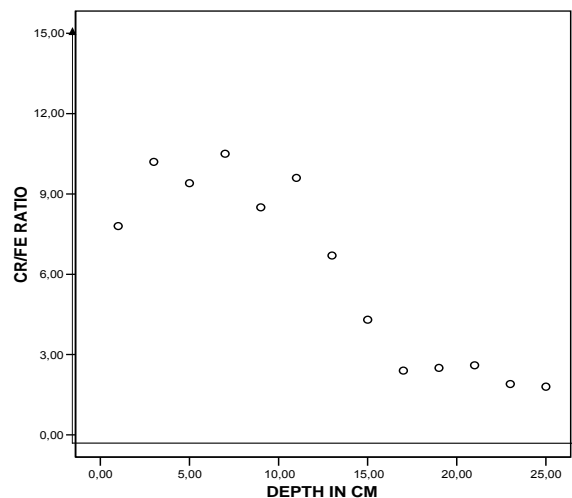
Spring 2010

© Erling Berge 2010

321

Pollution at different depths in sediments outside the coast of NH

- Pollution measured by the ratio chromium/iron at different depths of various sediment samples
- Is the relationship linear?

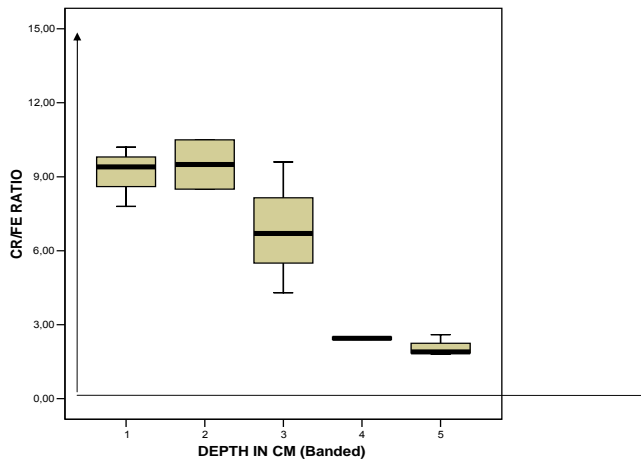


Spring 2010

© Erling Berge 2010

322

Medians of 5 bands: rate of chromium/iron in sediments outside the coast of NH



The relationship is obviously non-linear

Spring 2010

© Erling Berge 2010

323

Transformed variables

- Using transformed variables makes a regression curvilinear. The transformation makes the original curve relationship into a linear relationship
- This is the most important reason for a transformation
- At the same time transformations may rectify several other types of statistical problems (outliers, heteroscedasticity, non-normal errors)
- Procedure:
 - Choose an appropriate transformation and make new transformed variables
 - Do a standard regression analysis with the transformed variables
 - To interpret the results one usually will have to transform back to the original measurement scale

Spring 2010

© Erling Berge 2010

324

The linear model

$$y_i = \beta_0 + \sum_{j=1}^{K-1} \beta_j X_{ji} + \varepsilon_i$$

- In the linear model we can transform both x- and y-variables without any consequences for the properties of OLS estimates of the parameters
- OLS is a valid method as long as the model is linear in the parameters

Spring 2010

© Erling Berge 2010

325

Curvilinear Models

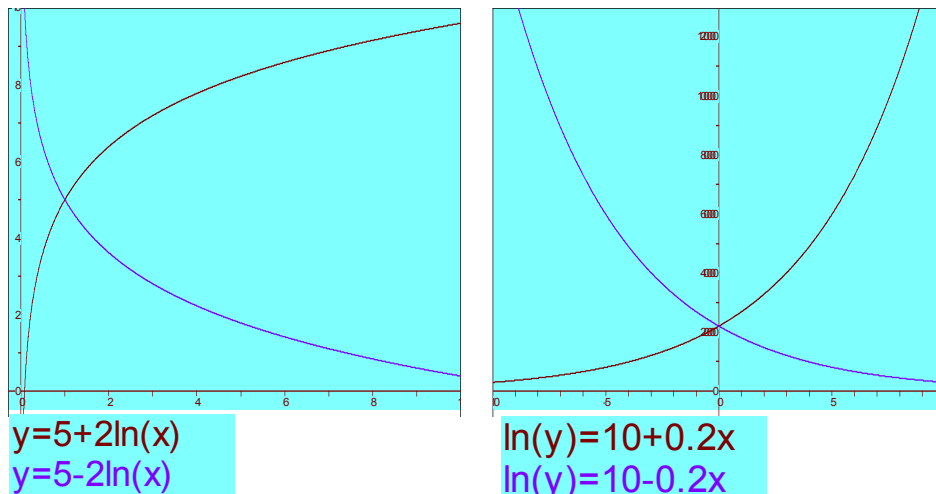
- Practically speaking this is regression with transformed variables
- We shall take a look at how different transformations provide different forms for the variable relations
 - Semi-logarithmic curves
 - Log-Log curves
 - Log-reciprocal curves
 - Polynomials (2 and 3 order)

Spring 2010

© Erling Berge 2010

326

Semilog curves Fig 5.2 in Hamilton



$y=5+2\ln(x)$
 $y=5-2\ln(x)$

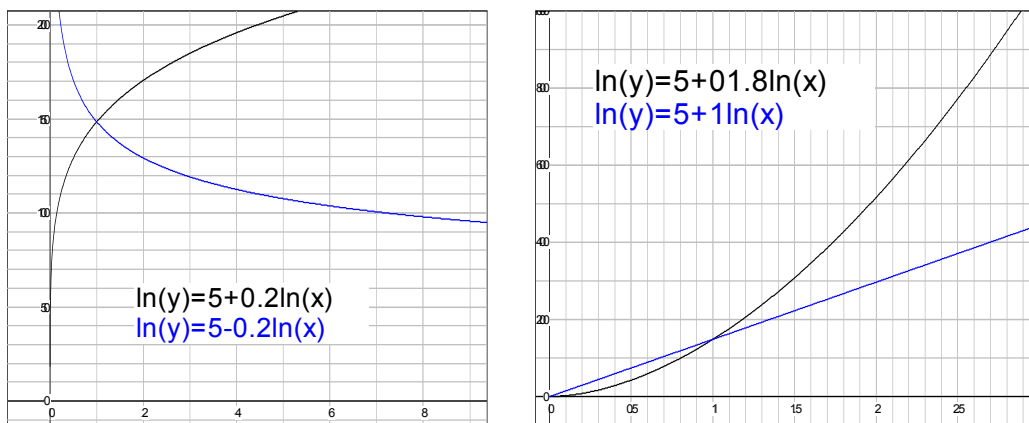
$\ln(y)=10+0.2x$
 $\ln(y)=10-0.2x$

Spring 2010

© Erling Berge 2010

327

Log-log curves Fig 5.3 in Hamilton



$\ln(y)=5+0.2\ln(x)$
 $\ln(y)=5-0.2\ln(x)$

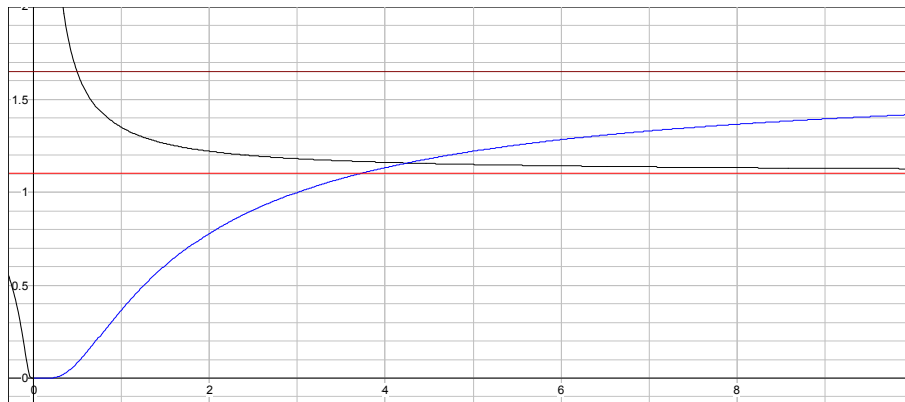
$\ln(y)=5+0.18\ln(x)$
 $\ln(y)=5+1\ln(x)$

Spring 2010

© Erling Berge 2010

328

Log-reciprocal curves Fig 5.4 in Hamilton



$\ln(y)=0.1+0.2/x$

$\ln(y)=0.5-1.5/x$

Horizontal line through (0, 1.105)

Horizontal line through (0, 1.649)

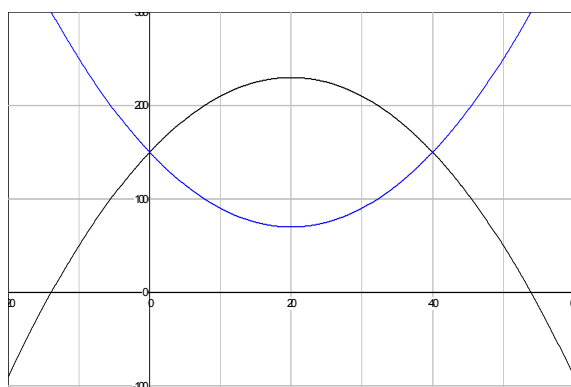
The horizontal lines give the value of y when x grows towards infinity: the asymptote for y

Spring 2010

© Erling Berge 2010

329

Second order polynomials Fig 5.5 in Hamilton



$y=150+8x-0.2x^2$

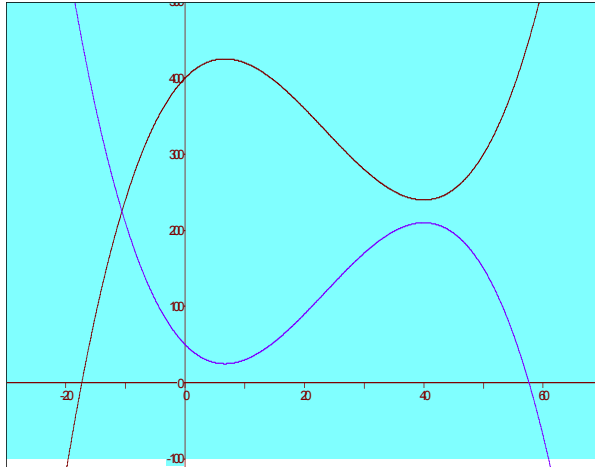
$y=150-8x+0.2x^2$

Spring 2010

© Erling Berge 2010

330

Third order polynomials Fig 5.6 in Hamilton



$$y=400+8x-0.7x^2+0.01x^3$$

$$y=50-8x+0.7x^2-0.01x^3$$

Spring 2010

© Erling Berge 2010

331

Choice of transformation

- Scatter plot or theory may provide advice
- Otherwise: transformation to symmetry gives the best option
- The regression reported in table 3.2 in Hamilton proved to be problematic
- Regression with transformed variables can reduce the problems

Spring 2010

© Erling Berge 2010

332

Choice of transformation in table 3.2 in Hamilton

Y = Water use 1981	$Y^* = Y^{0.3}$ provides approximate symmetry
X_1 = Income	$X_1^* = X_1^{0.3}$ provides approximate symmetry
X_2 = Water use 1980	$X_2^* = X_2^{0.3}$ provides approximate symmetry
X_3 = Education	Transformations are inappropriate
X_4 = Pensioner	Transformations do not work for dummies
X_5 = # people in 1981	$X_5^* = \ln(X_5)$ provides approximate symmetry
X_6 = Change in # people	$X_6 = X_5 - X_0$ (= # people in 1980)
X_7 = Relative change in #people	$X_7^* = \ln(X_5/X_0)$

Spring 2010

© Erling Berge 2010

333

Regression with transformed variables Tab 5.2 in Hamilton

Dependent Variable: (Wateruse81) ^{0.3}	B	Std. Err	t	Sig.
(Constant)	1,856	,385	4,822	,000
Income ^{0.3}	,516	,130	3,976	,000
Wateruse80 ^{0.3}	,626	,029	21,508	,000
Education in Years	-,036	,016	-2,257	,024
Retired?	,101	,119	,852	,395
Ln(# of people81)	,715	,110	6,469	,000
Ln(people81/people80)	,916	,263	3,485	,001

Spring 2010

© Erling Berge 2010

334

Table 3.2 (Hamilton p74)

Dependent Variable: Summer 1981 Water Use	B	Std. Error	t	Sig.	Beta
(Constant)	242.220	206.864	1.171	.242	
Income in Thousands	20.967	3.464	6.053	.000	.184
Summer 1980 Water Use	.492	.026	18.671	.000	.584
Education in Years	-41.866	13.220	-3.167	.002	-.087
Head of house retired?	189.184	95.021	1.991	.047	.058
# of People Resident, 1981	248.197	28.725	8.641	.000	.277
Increase in # of People	96.454	80.519	1.198	.232	.031

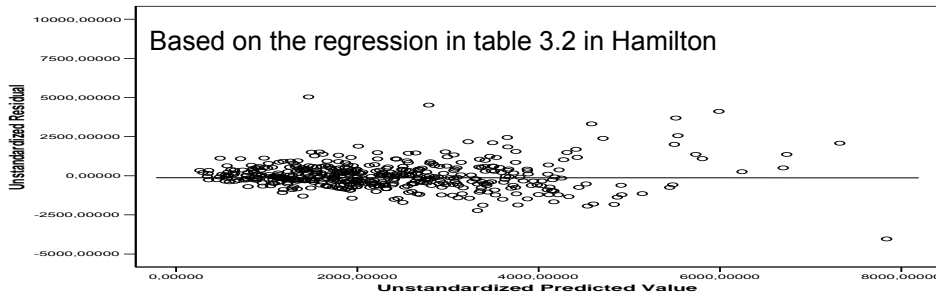
How do we interpret the coefficient of "Increase in # of People" ?

What leads to less water use after the crisis?

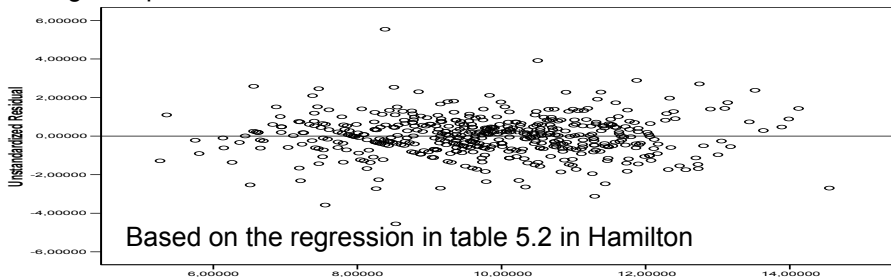
Spring 2010

© Erling Berge 2010

335



Residual against predicted Y



Spring 2010

© Erling Berge 2010

336

Other consequences of the transformations

- Two cases with large influence on the coefficient for income (large DFBTAS) do not have such influence (fig 4.11 and 5.9)
- One case with large influence on the coefficient for water use in 1980 do not have that large influence (fig 4.12 and 5.10)
- Transformation to symmetrical distributions will often solve many problems – but not always
- And it creates a new one: interpretation

Spring 2010

© Erling Berge 2010

337

Interpretation

- The model estimate now looks like this

$$y_i^{0.3} = 1.856 + 0.516x_{1i}^{0.3} + 0.626x_{2i}^{0.3} - 0.036x_{3i} + 0.101x_{4i} + 0.715\ln(x_{5i}) + 0.916\ln\left(\frac{x_{5i}}{x_{0i}}\right)$$

- The interpretation of the coefficients are not so straightforward any more. For example: the measurement units of the parameters have been changed
- The simplest way of interpreting is to use conditional effect plots

Spring 2010

© Erling Berge 2010

338

Conditional effect plot

- Should be used to study the relationship between the dependent variable and one x-variable with the rest of the x-variables given fixed values
- Typically we are interested in the relationship x-y when the other variables are given values that
 - Maximizes y
 - Are averages values of of the x-variables
 - Minimizes y

Spring 2010

© Erling Berge 2010

339

Example based on the regression in table 3.2 in Hamilton

Dependent Variable: Summer 1981 Water Use	Unstandardized Coefficients		t	Sig.
	B	Std. Error		
(Constant)	242,220	206,864	1,171	,242
Summer 1980 Water Use	,492	,026	18,671	,000
Income in Thousands	20,967	3,464	6,053	,000
Education in Years	-41,866	13,220	-3,167	,002
head of house retired?	189,184	95,021	1,991	,047
# of People Resident, 1981	248,197	28,725	8,641	,000
Increase in # of People	96,454	80,519	1,198	,232

Spring 2010

© Erling Berge 2010

340

To produce conditional effect plots it is useful to have a table of minimum, maximum and average variable values

	N	Minimum	Maximum	Mean
Summer 1981 water use	496	100	10100	2298,39
Summer 1980 water use	496	200	12700	2732,06
Income in thousands	496	2	100	23,08
Education in years	496	6	20	14,00
Head of household retired?	496	0	1	,29
# of people resident, 1981	496	1	10	3,07
Relative increase in # of people	496	-3	3	-,04
# People living in 1980	496	1	10	3,11

Spring 2010

© Erling Berge 2010

341

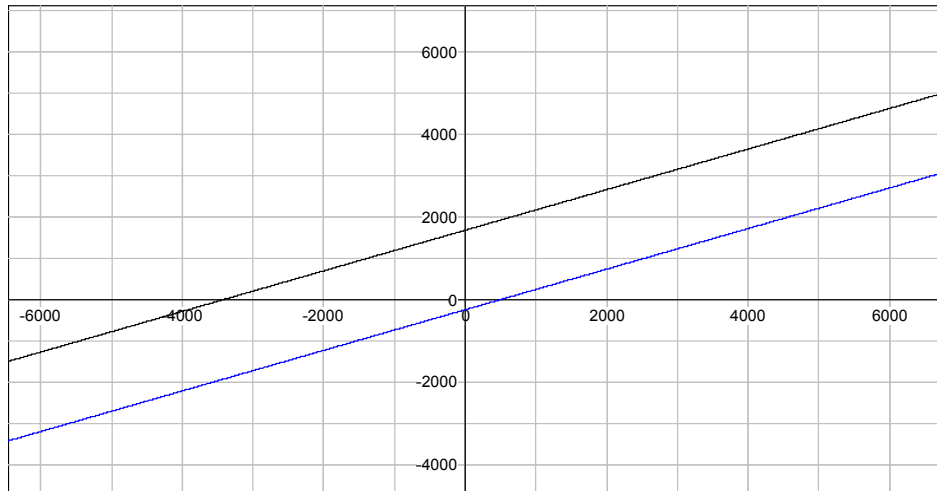
The equation

- Estimated $Y = 242,22 + 0,492X_1 + 20,967X_2 - 41,866X_3 + 189,184X_4 + 248,197X_5 + 96,454X_6$
- Maximizing the effect of X_1 on Y require maximum of X_2 , X_4 , X_5 , X_6 and minimum of X_3
- Average values of the effect of X_1 on Y is obtained by inserting average values of X_2 , X_3 , X_4 , X_5 , X_6
- Minimizing the effect of X_1 on Y require minimum of X_1 , X_2 , X_4 , X_5 , X_6 and maximum of X_3

Spring 2010

© Erling Berge 2010

342



$Y = 242.22 + 0.492X + 20.967 \times 10 - 41.866 \times 7 + 189.184 \times 1 + 248.197 \times 5 + 96.454 \times 1$
 $Y = 242.22 + 0.492X + 20.967 \times 1 - 41.866 \times 18 + 189.184 \times 0 + 248.197 \times 1 + 96.454 \times 0$

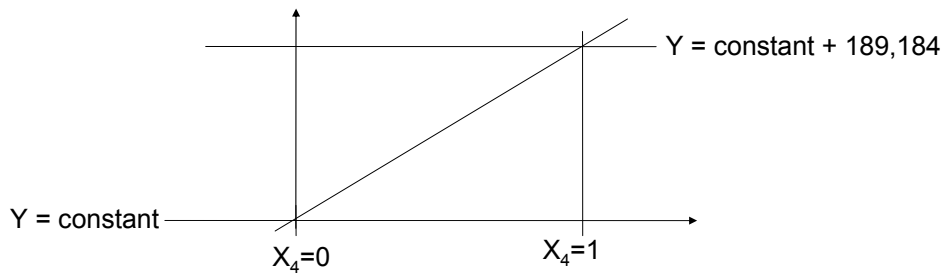
Spring 2010

© Erling Berge 2010

343

When x is dummy coded

- Estimated $Y = 242,22 + 0,492X_1 + 20,967X_2 - 41,866X_3 + 189,184X_4 + 248,197X_5 + 96,454X_6$
- Estimated $Y = \text{constant} + 189,184X_4$
 – X_4 can take the values of 0 or 1

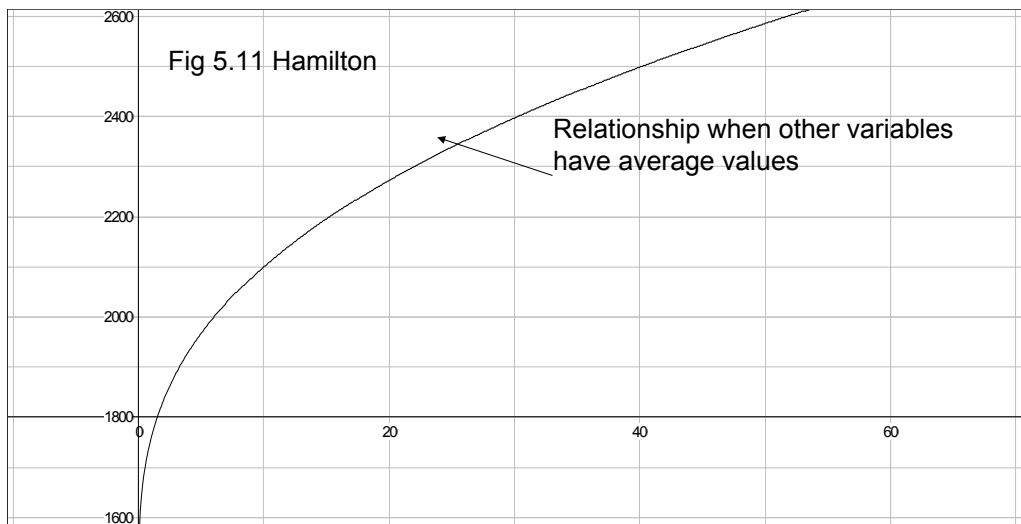


Spring 2010

© Erling Berge 2010

344

Water usage according to income controlled for the effect of other variables



$$y^{0.3} = 1.856 + 0.626(2732)^{0.3} - 0.036(14) + 0.101(0.294) + 0.715\ln(3.07) + 0.916(\ln(3.07) - \ln(3.11)) + 0.516(x)^{0.3}$$

Which plots might be of interest?

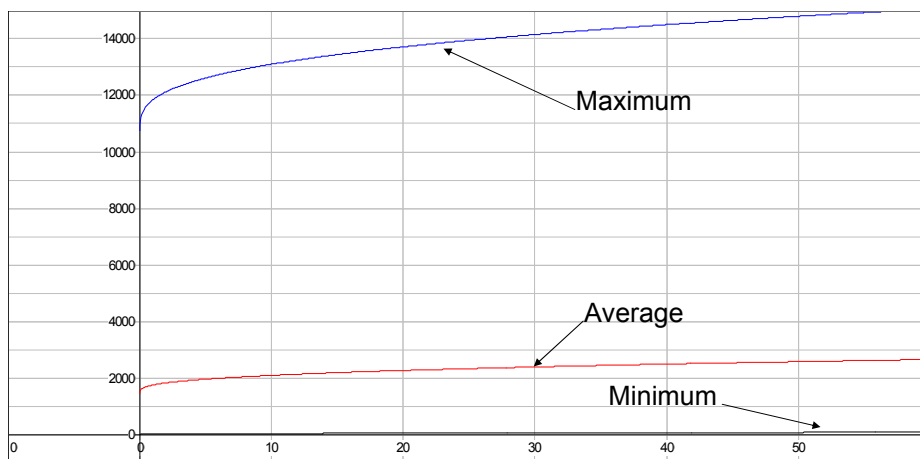
- The relationship between water usage and income controlled for the effect of other variables
 - Those minimizing water usage
 - Those maximizing water usage
 - Average values

$$1 \quad y^{0.3} = (1.856 + 0.626(200)^{0.3} - 0.036(20) + 0.101(0) + 0.715\ln(1) + 0.916(\ln(1) - \ln(10)) + 0.516(x)^{0.3})$$

$$2 \quad y^{0.3} = (1.856 + 0.626(12700)^{0.3} - 0.036(6) + 0.101(1) + 0.715\ln(10) + 0.916(\ln(10) - \ln(1)) + 0.516(x)^{0.3})$$

$$3 \quad y^{0.3} = (1.856 + 0.626(2732)^{0.3} - 0.036(14) + 0.101(0.29) + 0.715\ln(3.07) + 0.916(\ln(3.07) - \ln(3.11)) + 0.516(x)^{0.3})$$

Comparing three types of usage



Relationship between water usage and income Fig 5.12 in Hamilton
 Spring 2010 © Erling Berge 2010

347

The role of the constant in the plot

- The only difference between the three curves is the constant (konst)
 - In the maximum curve: (konst) = 14.046
 - In the minimum curve: (konst) = 4.204
 - In the average curve: (konst) = 8.507

$$y_i^{0.3} = (konst) + 0.516x_{1i}^{0.3}$$

- The effect of income varies with the value of (konst)
- When we transform the dependent variable all relationships become interaction effects

Spring 2010

© Erling Berge 2010

348

Comparing effects

- For some relationships the standardized regression coefficient can be used to compare effects, but it is sensitive for biased estimates of the standard error
- A more general method is to compare conditional effect plots where the scaling of the y-axis is kept constant

Spring 2010

© Erling Berge 2010

349

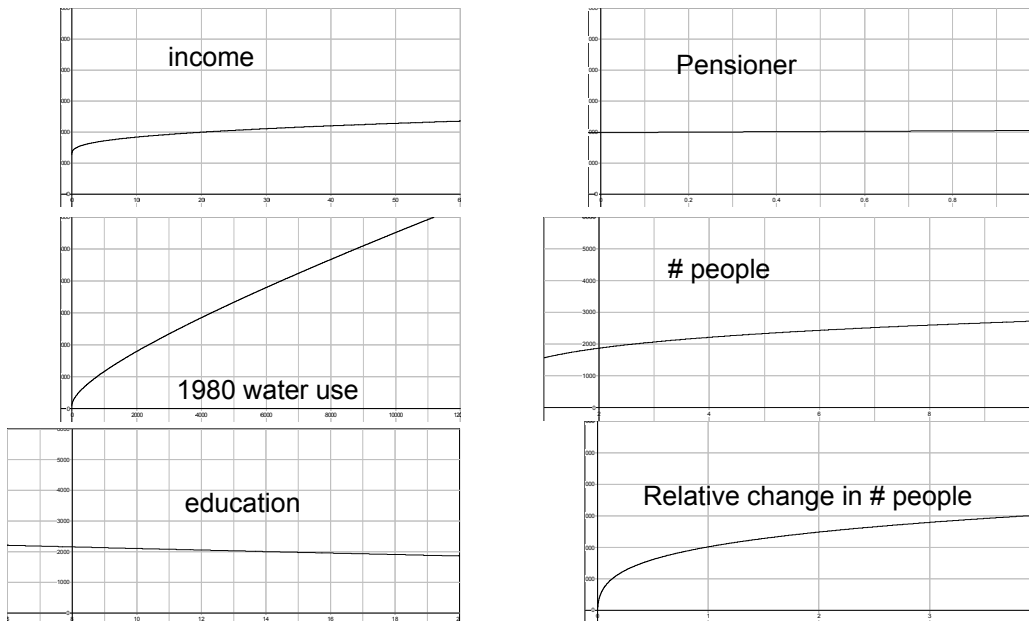


Fig 5.13 Hamilton

Spring 2010

© Erling Berge 2010

350

Non-linear models

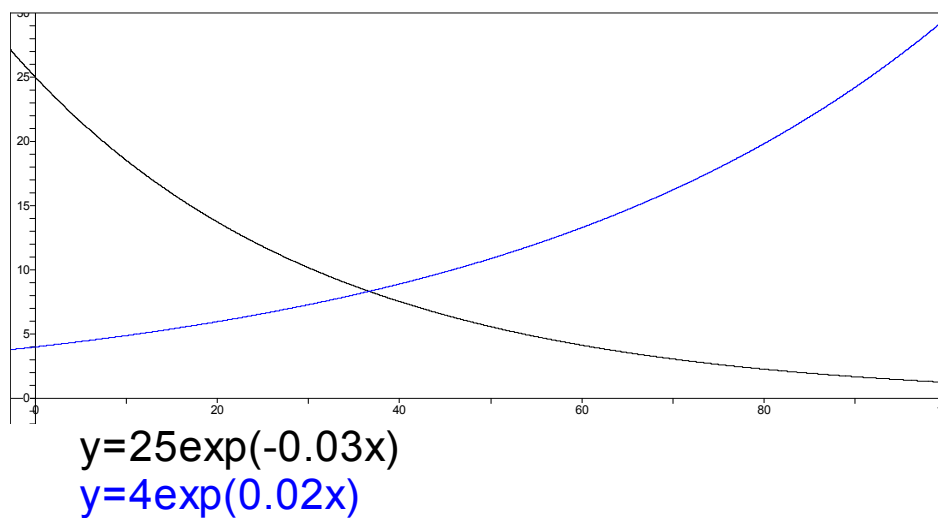
- If we do not have a model that is linear in the parameters other techniques than OLS are needed to estimate the parameters
- One may find two types of arguments for such models
 - Theory about the causal mechanism may say so
 - Inspection of the data may point towards one particular type of model
- We shall take a look at
 - Exponential models
 - Logistic models
 - Gompertz models

Spring 2010

© Erling Berge 2010

351

Exponential growth and decay Fig 5.14 in Hamilton

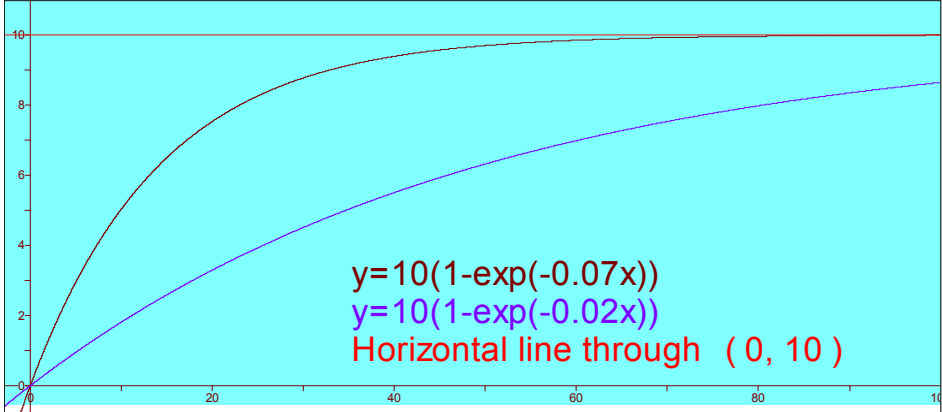


Spring 2010

© Erling Berge 2010

352

Negative exponential curves Fig 5.15 in Hamilton

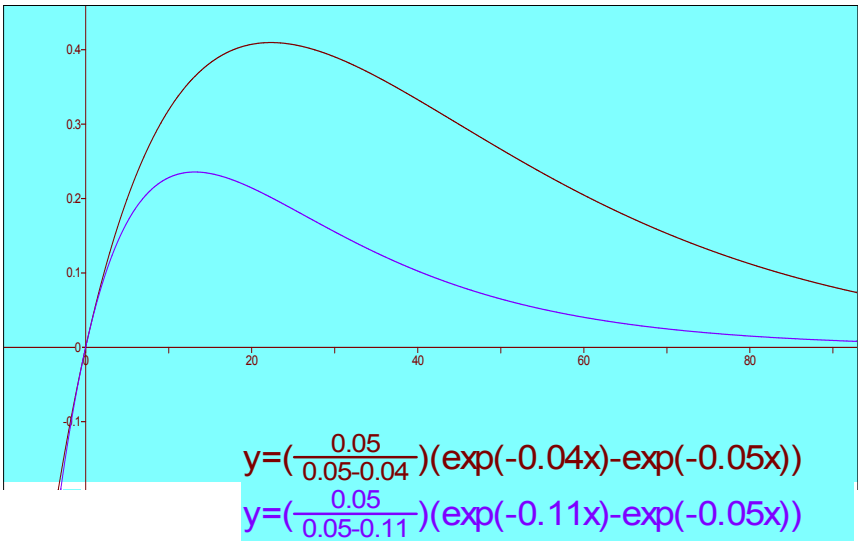


Spring 2010

© Erling Berge 2010

353

To-term exponential curves Fig 5.16 in Hamilton



Spring 2010

© Erling Berge 2010

354

Logistic models

- The logistic function is written
- As x grows towards infinity y will approach α
- When x declines towards minus infinity y will approach 0

$$y = \frac{\alpha}{1 + \gamma \exp(-\beta x)}$$

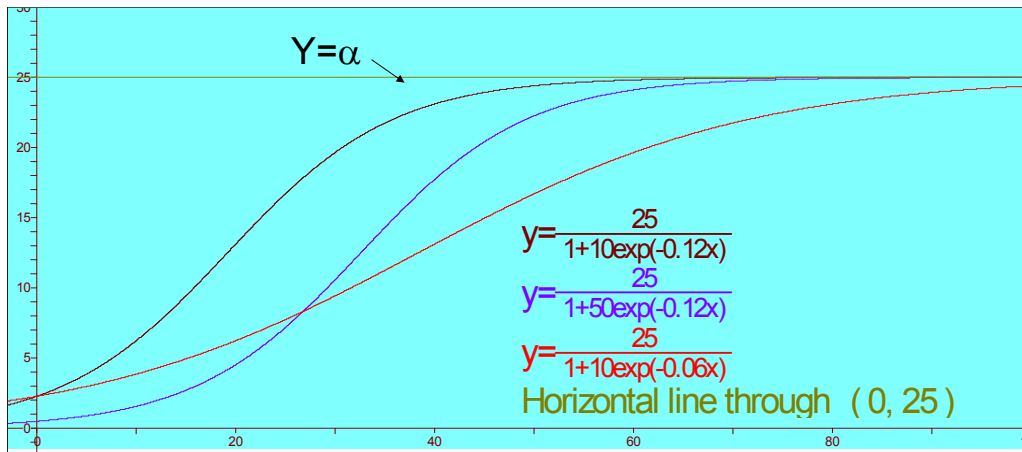
- Logistic models are appropriate for many phenomena
 - Growth of biological populations
 - Scattering of rumours
 - Distribution of illnesses

Spring 2010

© Erling Berge 2010

355

Logistic curves Fig 5.17 in Hamilton



- γ determines where growth starts
- β determines how fast the growth is

Spring 2010

© Erling Berge 2010

356

Logistic probability model

- If it is determined that $\alpha=\gamma=1$ y will vary between 0 and 1 as x goes from minus infinity to plus infinity
- Logistic curves can then be used to model probabilities

$$y_i = \frac{1}{1 + \exp(-\beta x_i)} + \varepsilon_i$$

Spring 2010

© Erling Berge 2010

357

Gompertz curves

- Gompertz curves are sigmoid curves like the logistic, but growth increase and growth reduction occur at different rates. Hence they are not symmetric

$$y = \alpha e^{-\gamma e^{-\beta x}} + \varepsilon$$

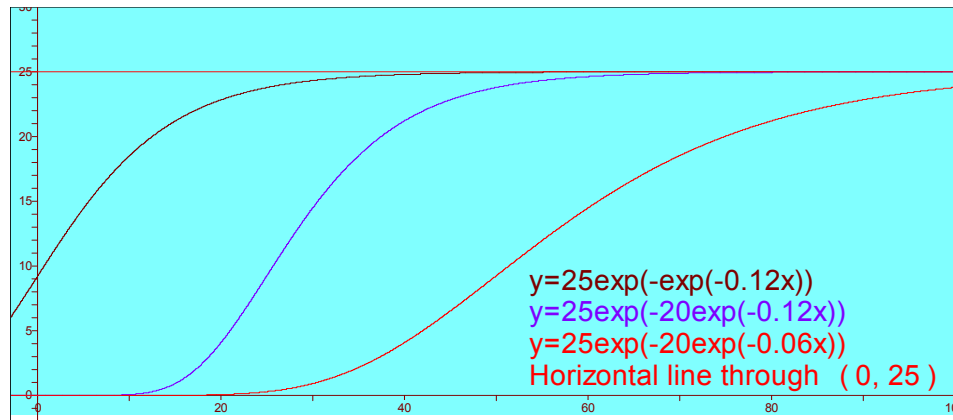
- Parameters α , γ , and β have the same interpretation as in the logistic model

Spring 2010

© Erling Berge 2010

358

Gompertz curves Fig 5.18 Hamilton



Spring 2010

© Erling Berge 2010

359

Estimation of non-linear models

- The criterion of fit is still minimum RSS
- It is uncommon to find analytical expressions for the parameters. One has to guess at a start value and go through several iterations to find which parameter value will give minimum RSS
- Good starting values are as a rule necessary, and everything from theory to inspection of data are used to find them

Spring 2010

© Erling Berge 2010

360

Per cent women with at least 1 child according to the woman's age and year of birth (England og Wales)

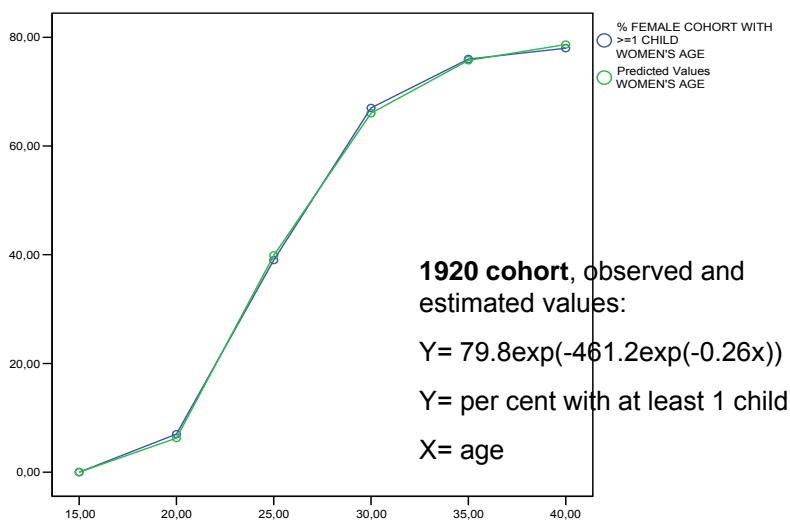
	1920	1930	1940	1945	1950	1955	1960	1965
15	0	0	0	0	0	0	0	0
20	7	9	13	17	19	18	13	11
25	39	48	59	60	53	45	39	-
30	67	75	82	82	75	68	-	-
35	76	83	87	88	83	-	-	-
40	78	86	89	90	-	-	-	-
45	-	86	89	-	-	-	-	-

Spring 2010

© Erling Berge 2010

361

Estimating Gompertz-models for cohorts (1)

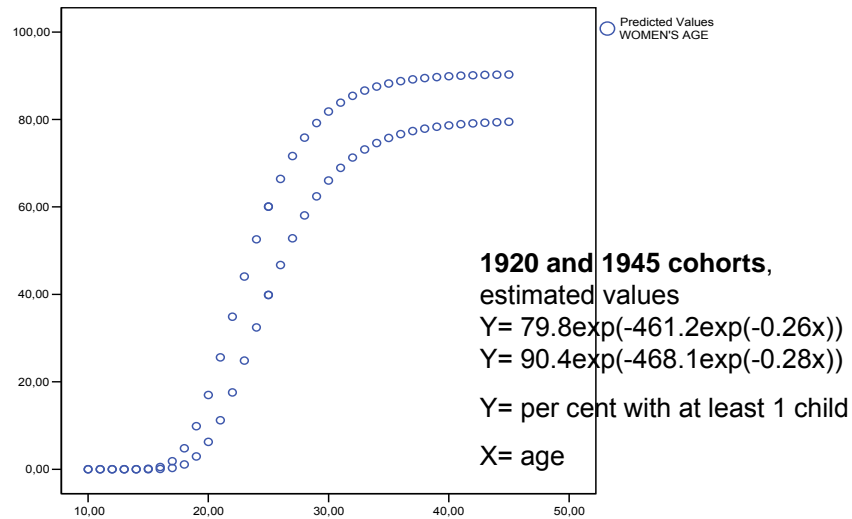


Spring 2010

© Erling Berge 2010

362

Estimating Gompertz-models for cohorts (2)



Spring 2010

© Erling Berge 2010

363

Model estimation and fit

- To evaluate a theoretically developed model
- To predict y within or outside the observed range of variation for x
- Substantial or comparative interpretation of the parameters of the model
 - On cohorts that are not finished with their births (thus predicting outside the observed range of x)
 - We can use the model to compare parameter values of different cohorts

Spring 2010

© Erling Berge 2010

364

Parameter interpretation Table 5.6 Hamilton

Cohort	α = upper limit	γ = ?	β = growth speed
1920	79.8	461.2	0.26
1930	86.5	538.0	0.27
1940	89.1	942.0	0.31
1945	90.4	468.1	0.28
1950	87.5	144.9	0.23
1955	88.9	60.3	0.18

Spring 2010

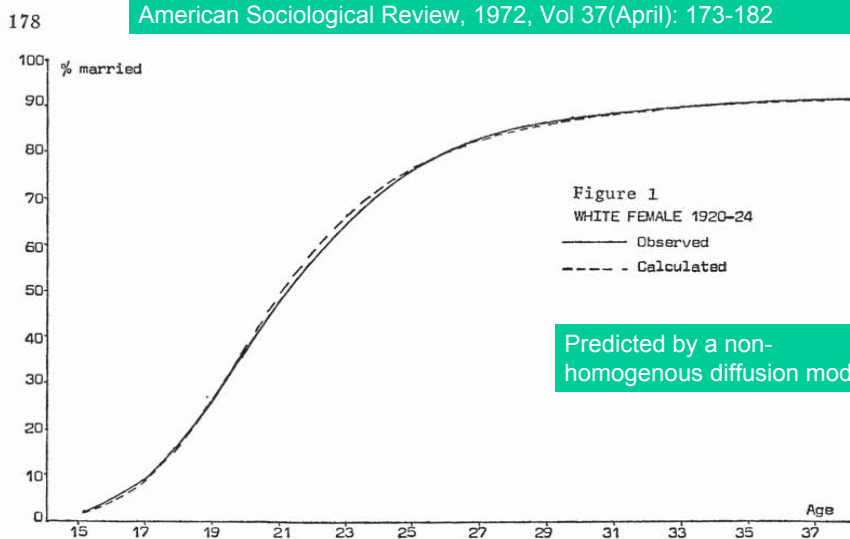
© Erling Berge 2010

365

The process of entry into first marriage

Gudmund Hernes

American Sociological Review, 1972, Vol 37(April): 173-182

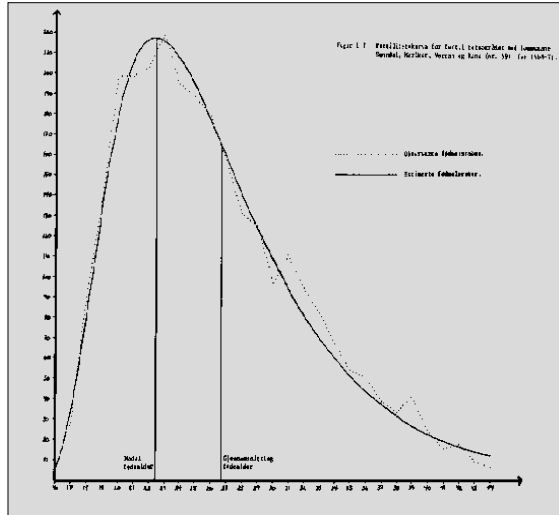


Spring 2010

© Erling Berge 2010

366

Birth rates in Sunndal, Meråker, Verran, and Rana 1968-71



- Estimated with a Hadwiger function
- Ref.: Berge, Erling. 1981. The Social Ecology of Human Fertility in Norway 1970. Ph.D. Dissertation. Boston: Boston University.

Spring 2010

© Erling Berge 2010

367

Conclusions of chapter 5 (1)

- Data analysis often starts with linear models. They are the simplest.
- Theory or exploratory data analysis (band regression, smoothing) can tell us if curvilinear or non-linear models are needed
- Transformation of variables give curvilinear regression. This can counteract several problems:
 - Curvilinear relationships
 - Case with large influence
 - Non-normal errors
 - Heteroscedasticity

Spring 2010

© Erling Berge 2010

368

Conclusions of chapter 5 (2)

- Non-linear regression use iterative procedures to find parameter estimates
- The procedures need initial values and are often sensitive for the initial values
- The interpretation of the parameters may be difficult. Graphs showing the relationship for different parameter values will provide valuable help for the interpretation

Spring 2010

© Erling Berge 2010

369

Ch 6 Robust Regression

- Has been developed to work well in situations where OLS breaks down. Where the OLS assumptions are satisfied robust regression are not as good as OLS, but not by very much
- Even if robust regression is better suited for those who do not want to put much effort into testing the assumptions, it is so far difficult to use
- Robust regression has focused on residuals with heavy tails (many cases with high influence on the regression)

Spring 2010

© Erling Berge 2010

370

Regression of mortality on air pollution

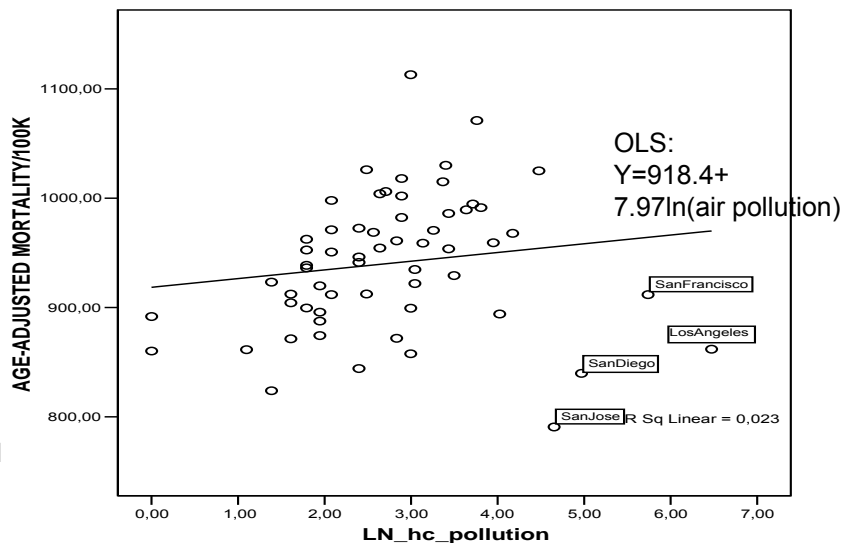


Figure 6.1
Hamilton

Spring 2010

© Erling Berge 2010

371

Robust regression of mortality on air pollution

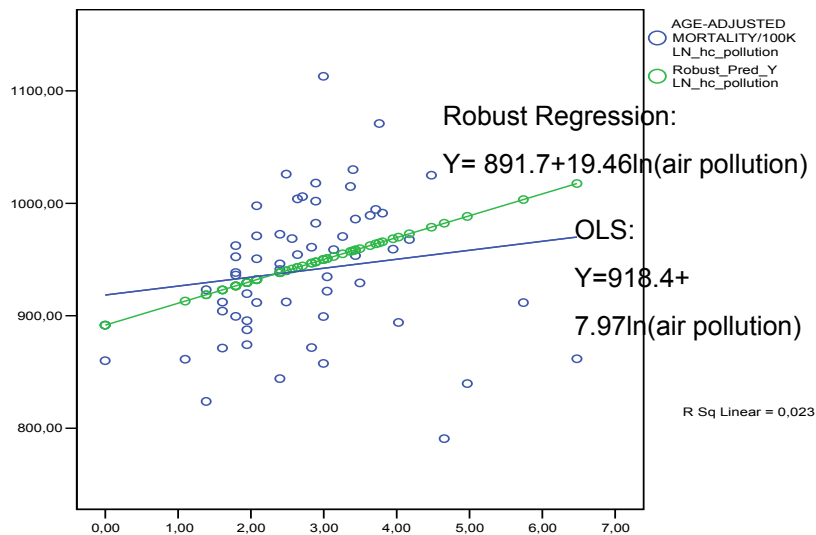


Figure 6.2
Hamilton

Spring 2010

© Erling Berge 2010

372

Robust regression and SPSS

- SPSS do not have a particular routine that performs robust regression
- It can possibly be done within the Generalized linear models procedure <but I have not tested it>
- It can be done by weighted OLS regression, but then it is required that we make the weight functions and go through the iterations one by one including computation of weights every time
- This procedure will be outlined below

Spring 2010

© Erling Berge 2010

373

ROBUST AND RESISTANT

- RESISTANT methods are not affected by small errors or changes in the sample data
- ROBUST methods are not affected by small deviations from the assumptions of the model
- Most resistant estimators are also robust in relation to the assumption about normally distributed residuals
-
- **OLS is neither ROBUST nor RESISTANT**

Spring 2010

© Erling Berge 2010

374

Outliers is a problem for OLS

Outliers affect the estimates of

- Parameters
- Standard errors (standard deviation of parameters)
- Coefficient of determination
- Test statistics
- And many other statistics

Robust regression tries to protect against this by giving less weight to such cases, not by excluding them

Spring 2010

© Erling Berge 2010

375

Protection against NON-NORMALE residuals

Robust methods can help when

- the tails in the distribution of the residuals are heavy, i.e. when it is too many outliers compared to the normal distribution
- Unusual X-values have leverage and may cause problems

But for other causes of non-normality robust methods will not help

Spring 2010

© Erling Berge 2010

376

Estimation methods for robust regression

- M-estimation (maximum likelihood) minimizes a weighted sum of the residuals. This can be approximated by the weighted least squares method (WLS)
- R-estimation (based on rank) minimizes a sum where a weighted rank is included. The method is more difficult to use than M-estimation
- L-estimation (based on quantiles) uses linear functions of the sample order statistics (quantiles)

Spring 2010

© Erling Berge 2010

377

IRLS- Iterated Reweighted Least Squares

M-estimation by means of IRLS needs

1. Start values from OLS. Save the residuals
2. Use OLS residuals to find weights. Larger residuals gives less weight
3. Find new parameter values and residuals with WLS
4. Go to step 2 and find new weights from the new residuals, go on to step 3 and 4, until changes in the parameters become small

Iteration: to repeat a sequence of operations

Spring 2010

© Erling Berge 2010

378

IRLS

- IRLS is in theory equivalent to M-estimation
- To use the method we need to compute
- Scaled residuals, u_i , and a
- Weight function, w_i , that gives least weight to the largest residuals

Spring 2010

© Erling Berge 2010

379

Scaling of residuals I

- Scaled residual u_i
 - s is the scale factor and e_i residual
- The scale factor in OLS is the estimate of the standard error of the residual: nb! s_e is not resistant
- A resistant alternative is based on MAD, "median absolute deviation"

$$u_i = \frac{e_i}{s}$$

$$s_e = \sqrt{\frac{RSS}{n-K}}$$

$$MAD = \text{median} | e_i - \text{median}(e_i) |$$

Spring 2010

© Erling Berge 2010

380

Scaling of residuals II

$$MAD = \text{median} | e_i - \text{median}(e_i) |$$

The scale factor (standard error of the distribution)

Using a resistant estimate will be

- $s = MAD / 0.6745 = 1.483MAD$

and the scaled residual

- $u_i = [e_i / s] = (0.6745 * e_i) / MAD$

In a normal distribution $s = MAD / 0.6745$ will estimate the standard error correctly like s_e

In case of non-normal errors $s = MAD / 0.6745$ will be better.

This is a resistant estimate, s_e is not resistant

Weight functions I

- Properties is measured in relation to OLS on normally distributed errors.
- The method should be “almost as good” as OLS on normally distributed errors and much better when the errors are non-normal
- Properties are determined by a “calibration constant” (c in the formulas)

Weight functions II

- **OLS-weights:** $w_i = 1$ for all i
- **Huber-weights:** weights down when the scaled residual is larger than c , $c=1,345$ gives 95% of the efficiency of OLS on normally distributed errors
- **Tukey's bi-weighted** estimates get 95% of the efficiency of OLS on normally distributed errors by gradually weighting down scaled errors until $|u_i| \leq c = 4.685$ and by dropping cases where the residual is larger

Spring 2010

© Erling Berge 2010

383

Huber-weights

$$w_i = 1 \quad \forall |u_i| \leq c$$

$$w_i = \frac{c}{|u_i|} \quad \forall |u_i| > c$$

$\forall =$ for alle

Spring 2010

© Erling Berge 2010

384

Tukey weights

$$w_i = \left[1 - \left(\frac{u_i}{c} \right)^2 \right]^2 \quad \forall |u_i| \leq c$$

$$w_i = 0 \quad \forall |u_i| > c$$

$\forall = \textit{for alle}$

- Tukey weighting in IRLS is sensitive for start values of the parameters (one may end up at local minima)

Spring 2010

© Erling Berge 2010

385

Standard errors and tests in IRLS

- The WLS program cannot estimate standard errors and test statistics correctly by IRLS
- A procedure that works is described by Hamilton on page 198-199

Spring 2010

© Erling Berge 2010

386

Use of Robust Estimation

- If OLS and Robust estimates are different it means that outliers have influence on the OLS results making them unreliable. Results cannot be trusted
- Robust predicted values will better portray the bulk of the data
- Robust residuals will be better at discovering which cases are unusual
- Weights from the robust regression will show which cases are outliers
- OLS and RR can support each other

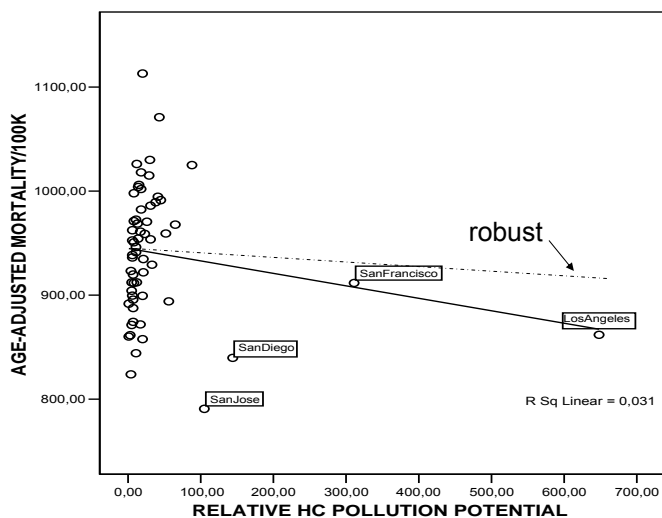
Spring 2010

© Erling Berge 2010

387

Fig 6.9 Hamilton: OLS and RR on untransformed data

Mortality regressed on air pollution
Effect of high leverage

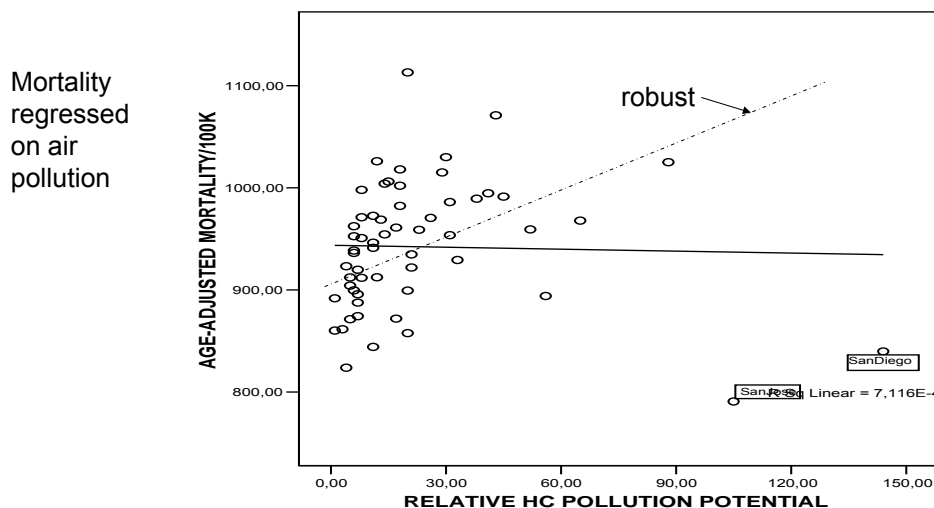


Spring 2010

© Erling Berge 2010

388

Fig 6.10 Hamilton: OLS and RR on untransformed data when two outliers are removed



Spring 2010

© Erling Berge 2010

389

RR do not protect against leverage

- RR with M-estimation protects against unusual y-values (outliers) but not necessarily against unusual x-values (leverage)
- Efforts to test and diagnose are still needed (heteroscedasticity is still a problem for IRLS)
- Studies of the data and transformation to symmetry will reduce the risk of problems appearing
- No method is “safe” if it is used without forethought and diagnostic studies of data

Spring 2010

© Erling Berge 2010

390

Robust Multippel Regresjon

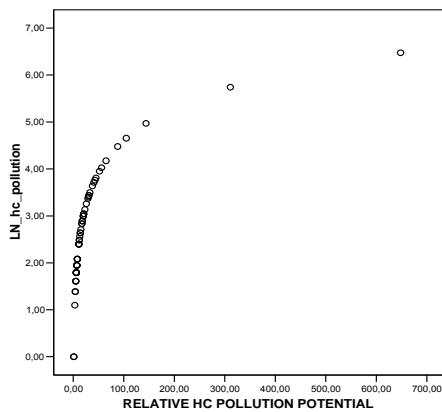
X ₁	RELATIVE HC POLLUTION POTENTIAL (natural log)
X ₂	AVG. YEARLY PRECIP. INCHES
X ₃	AVG. JANUARY TEMPERATURE, F
X ₄	MEDIAN EDUCATION OF POP 25+
X ₅	% NON-WHITE (square root)
X ₆	POPULATION PER HOUSEHOLD
X ₇	% 65 AND OVER
X ₈	% SOUND HOUSING UNITS
X ₉	PEOPLE PER SQUARE MILE (natural log)
X ₁₀	AVG. JULY TEMPERATURE, F
X ₁₁	% WHITE COLLAR EMPLOYMENT
X ₁₂	% FAMILIES WITH INCOME<\$3000 (negative reciprocal root)
X ₁₃	AVG RELATIVE HUMIDITY, %

Spring 2010

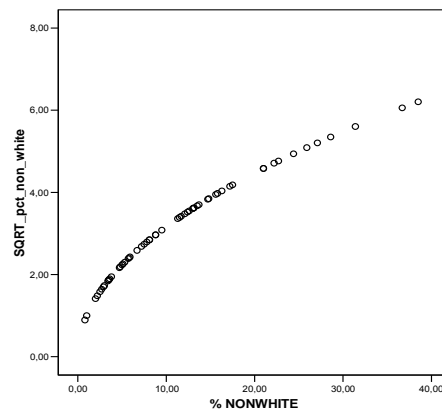
© Erling Berge 2010

391

Multiple OLS regression with transformed variables:
effect of transformation



In of air pollution



Square root of % non-white

Spring 2010

© Erling Berge 2010

392

OLS with backward elimination gives

Dependent Variable: AGE-ADJUSTED MORTALITY/100K	B	Std. Error	t	Sig.
(Constant)	986,261	82,674	11,929	,000
LN_hc_pollution	17,469	4,636	3,768	,000
AVG. YEARLY PRECIP. INCHES	2,352	,640	3,677	,001
AVG. JANUARY TEMPERATURE, F	-2,132	,504	-4,228	,000
MEDIAN EDUCATION OF POP 25+	-17,958	6,204	-2,895	,005
SQRT_pct_non_white	27,335	4,398	6,215	,000

- Robust regression gives predicted y:
- $Y = 1001.8 + 17.77x_{1i} + 2.32x_{2i} - 2.11x_{3i} - 19.1x_{4i} + 26.2x_{5i}$

Spring 2010

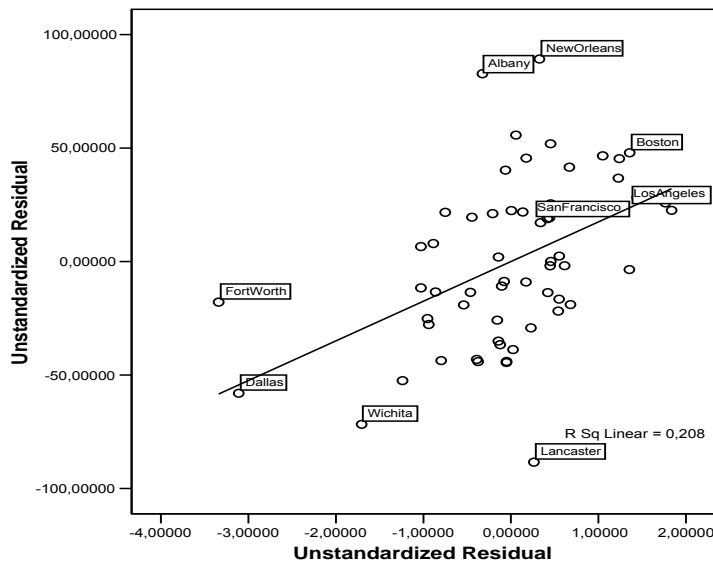
© Erling Berge 2010

393

Multiple OLS regression with transformed variables

Leverage plot of residual from mortality (y) and residual of ln_air_pollution (x)

Los Angeles and San Francisco are no longer outliers



Spring 2010

© Erling Berge 2010

394

Four estimates of the relationship mortality – air pollution

Effect of air pollution

	OLS	Robust
1 variable	7.97	19.46
5 variables	17.47	17.77

- Note that in RR the bivariate regression comes pretty close to the result of the multivariate regression
- In the five-variable model there are new cases with influence on the line of regression
- Removing the 5 cases that have the highest leverage parameter (h_i) do not give substantial changes in the coefficients

Spring 2010

© Erling Berge 2010

395

Robust Regression vs Bounded Influence Regression

- Robust Regression protect against the effect of outliers (unusual y-values) if these do not go together with unusual x-values
- Bounded Influence Regression is designed to protect against influence from unusual combinations of x-values

Spring 2010

© Erling Berge 2010

396

BI - Bounded Influence Regression

- BI-methods are made to limit the influence of high leverage cases (large h_i = high leverage)
- The simplest way of doing this is to modify the Huber-weights or Tukey-weights in the IRLS procedure for RR (robust regression) with a factor based on the leverage statistic

Spring 2010

© Erling Berge 2010

397

Bounded influence: modification of weights

- Expand the weight function with a weight based on the leverage statistic h_i
- $w_i^H = 1$ if $h_i \leq c^H$
- $w_i^H = (c^H / h_i)$ if $h_i > c^H$
- c^H is often set to the 90% percentile in the distribution of h_i
- Then the IRSL weight becomes $w_i w_i^H$ where w_i is either the Tukey- or Huber-weight that changes from iteration to iteration while w_i^H is constant

Spring 2010

© Erling Berge 2010

398

Bounded influence as a diagnostic tool

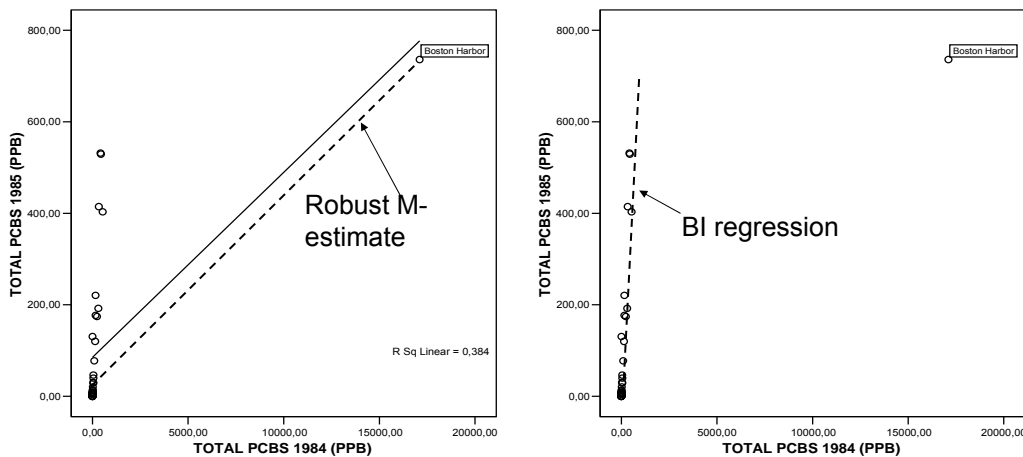
- Estimation of standard errors and test statistics becomes even more complicated than for the M-estimators mentioned above
- We can use BI estimates as a descriptive tool to check up on other estimates
- One (somewhat) extreme example: PCB pollution in river mouths in 1984 and 1985 (Hamilton table 6.4)

Spring 2010

© Erling Berge 2010

399

Fig 6.15 and 6.16 Hamilton



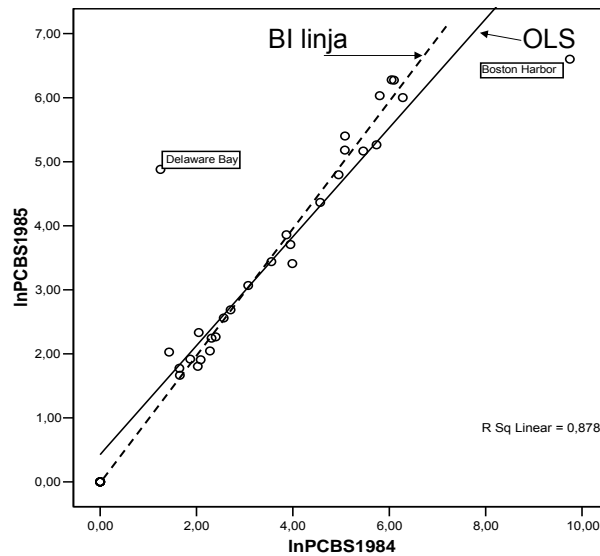
Spring 2010

© Erling Berge 2010

400

Fig 6.17 Hamilton

OLS and BI estimates with transformed variables give about the same result



Spring 2010

© Erling Berge 2010

401

Conclusions

- When data have many outliers robust methods will have better properties than OLS
 - They are more effective and give more accurate confidence intervals and tests of significance
- Robust regression can be used as a diagnostic tool
 - If OLS and RR agree we can have more confidence in the OLS results
 - If they disagree we will
 - Know that a problem exist
 - Have a model that fits the data better and identifies the outliers better
- Robust methods does not protect against problems that are due to curvilinear or non-linear models, heteroscedasticity, and autocorrelation

Spring 2010

© Erling Berge 2010

402

Logistic regression II

- Hamilton Ch 7 p217-242

Spring 2010

© Erling Berge 2010

403

Definitions I

- The probability that person no i shall have the value 1 on the variable Y_i will be written $\Pr(Y_i = 1)$.
- Then $\Pr(Y_i \neq 1) = 1 - \Pr(Y_i = 1)$
- The odds that person no i shall have the value 1 on the variable Y_i , here called O_i , is the ratio between two probabilities

$$O_i(y_i = 1) = \frac{\Pr(y_i = 1)}{1 - \Pr(y_i = 1)} = \frac{p_i}{1 - p_i}$$

Spring 2010

© Erling Berge 2010

404

Definitions II

- The LOGIT , L_i , for person no i (corresponding to $\Pr(Y_i=1)$) is the natural logarithm of the odds, O_i , that person no i has the value 1 on variable Y_i . This is written:

$$L_i = \ln(O_i) = \ln\{p_i/(1-p_i)\}$$
- The model assumes that L_i is a linear function of the explanatory variables x_{ji} ,
- i.e.:
- $L_i = \beta_0 + \sum_j \beta_j x_{ji}$, where $j=1,\dots,K-1$, and $i=1,\dots,n$

Spring 2010

© Erling Berge 2010

405

Logistic regression: assumptions

- The model is correctly specified
 - The logit is linear in its parameters
 - All relevant variables are included
 - No irrelevant variables are included
- x -variables are measured without error
- Observations are independent
- No perfect multicollinearity
- No perfect discrimination
- Sufficiently large sample

Spring 2010

© Erling Berge 2010

406

Assumptions that cannot be tested

- Model specification
 - All relevant variables are included
- x-variables are measured without error
- Observations are independent

Two will be tested automatically.

If the model can be estimated there is

- No perfect multicollinearity and
- No perfect discrimination

Spring 2010

© Erling Berge 2010

407

LOGISTIC REGRESSION

Statistical problems may be due to

- Too small a sample
- High degree of **multicollinearity**
 - Leading to large standard errors (imprecise estimates)
 - Multicollinearity is discovered and treated in the same way as in OLS regression
- High degree of **discrimination** (or separation)
 - Leading to large standard errors (imprecise estimates)
 - Will be discovered automatically by SPSS

Spring 2010

© Erling Berge 2010

408

Assumptions that can be tested

- Model specification
 - logit is linear in the parameters
 - no irrelevant variables are included
- Sufficiently large sample
 - What constitutes a sufficiently large sample is not always clear.
 - It depends on how the cases are distributed between 0 and 1 categories. If one of these is too small there will be problems estimating partial effects.
 - It also depends on the number of different patterns in the sample and how cases are distributed across these

Spring 2010

© Erling Berge 2010

409

Sample size in logistic regression

Large sample properties

- The good properties of ML estimates of binary logistic regression models are large sample properties that obtain as sample size goes towards infinity.
- What happens when you have too small a sample is largely unknown
- Long (1997) puts 100 cases as an absolute lower bound

Spring 2010

© Erling Berge 2010

410

Calculation of lower bounds

- A lower bound of 100 must be adjusted according to number of variables in the model and the distribution of cases on the dependent variable.
- Peduzzi et al. (1996) suggest:
- Let p be the smallest of the proportions of negative or positive cases in the population and k the number of covariates (the number of independent variables), then the minimum number of cases to include is:
- $N = 10 k / p$
- If the resulting number is less than 100 you should increase it to 100
- Or you may say that the maximum number of variables you can include in the model will be
- $k = N \cdot p / 10$

Spring 2010

© Erling Berge 2010

411

LOGISTIC REGRESSION: TESTING (1)

- Testing implies an assessment of whether statistical problems lead to departure from the assumptions

Two tests are useful

- (1) The **Likelihood ratio test** statistic: χ^2_H
 - Can be used analogous to the F-test
- (2) Wald test
 - The square root of this can be used analogous to the t-test but it follows a normal distribution

Spring 2010

© Erling Berge 2010

412

Logistic Regression in SPSS I

Case Processing Summary

Unweighted Cases ^a		N	Percent
Selected Cases	Included in Analysis	153	100.0
	Missing Cases	0	.0
	Total	153	100.0
Unselected Cases		0	.0
Total		153	100.0

a. If weight is in effect, see classification table for the total number of cases.

Dependent Variable Encoding

Original Value	Internal Value
OPEN	0
CLOSE	1

Logistic Regression in SPSS IIa

Iteration History^{a,b,c}

Iteration		-2 Log likelihood	Coefficients
			Constant
Step 1		209.212	-.275
0	2	209.212	-.276
	3	209.212	-.276

- a. Constant is included in the model.
- b. Initial -2 Log Likelihood: 209.212
- c. Estimation terminated at iteration number 3 because parameter estimates changed by less than .001.

Logistic Regression in SPSS IIb

Classification Table^{a,b}

Observed		Predicted		
		SCHOOLS SHOULD CLOSE		Percentage Correct
		OPEN	CLOSE	
Step 0	SCHOOLS SHOULD OPEN	87	0	100.0
	CLOSE	66	0	.0
Overall Percentage				56.9

a. Constant is included in the model.

b. The cut value is .500

Variables in the Equation

	B	S.E.	Wald	df	Sig.	Exp(B)	
Step 0	Constant	-.276	.163	2.864	1	.091	.759

Variables not in the Equation

	Score	df	Sig.	
Step 0	Variables lived	12.683	1	.000
	Overall Statistics	12.683	1	.000

Spring 2010

© Erling Berge 2010

415

Logistic Regression in SPSS IIIa

Iteration History^{a,b,c,d}

Iteration		-2 Log likelihood	Coefficients	
			Constant	lived
Step 1	1	195.684	.376	-.034
	2	195.269	.455	-.041
	3	195.267	.460	-.041
	4	195.267	.460	-.041

a. Method: Enter

b. Constant is included in the model.

c. Initial -2 Log Likelihood: 209.212

d. Estimation terminated at iteration number 4 because parameter estimates changed by less than .001.

Spring 2010

© Erling Berge 2010

416

Logistic Regression in SPSS IIIb

Omnibus Tests of Model Coefficients

		Chi-square	df	Sig.
Step 1	Step	13.944	1	.000
	Block	13.944	1	.000
	Model	13.944	1	.000

Model Summary

Step	-2 Log likelihood	Cox & Snell R Square	Nagelkerke R Square
1	195.267 ^a	.087	.117

a. Estimation terminated at iteration number 4 because parameter estimates changed by less than .001.

Spring 2010

© Erling Berge 2010

417

Logistic Regression in SPSS IIIc

Classification Table^a

			Predicted		
			SCHOOLS SHOULD CLOSE		Percentage Correct
Observed		OPEN	CLOSE		
Step 1	SCHOOLS SHOULD OPEN	59	28	67.8	
	CLOSE CLOSE	29	37	56.1	
	Overall Percentage			62.7	

a. The cut value is .500

Variables in the Equation

	B	S.E.	Wald	df	Sig.	Exp(B)
Step 1 ^a lived	-.041	.012	11.399	1	.001	.960
Constant	.460	.263	3.069	1	.080	1.584

a. Variable(s) entered on step 1: lived.

Spring 2010

© Erling Berge 2010

418

Conditional Effect Plot

- Set all x-variables except x_k to fixed values and enter these into the equation for the logit
- Plot $\text{Pr}(Y=1)$ as a function of x_k i.e.
- $P = 1/(1+\exp[-L]) = 1/(1+\exp[-\text{konst} - b_k x_k])$ for all reasonable values of x_k ,
 “konst” is the constant obtained by entering into the logit the fixed values of variables other than x_k

Spring 2010

© Erling Berge 2010

419

Excerpt from Hamilton Table 7.4

	B	S.E.	Wald	df	Sig.	Exp(B)	Minimum	Maximum	Mean
lived	-,040	,015	6,559	1	,010	,961	1,00	81,00	19,2680
educ	-,197	,093	4,509	1	,034	,821	6,00	20,00	12,9542
contam	1,299	,477	7,423	1	,006	3,664	,00	1,00	,2810
hsc	2,279	,490	21,591	1	,000	9,763	,00	1,00	,3072
nodad	-1,731	,725	5,696	1	,017	,177	,00	1,00	,1699
Constant	2,182	1,330	2,692	1	,101	8,866			

Logit:

$$L = 2.182 - 0.04*\text{lived} - 0.197*\text{educ} + 1.299*\text{contam} + 2.279*\text{hsc} - 1.731*\text{nodad}$$

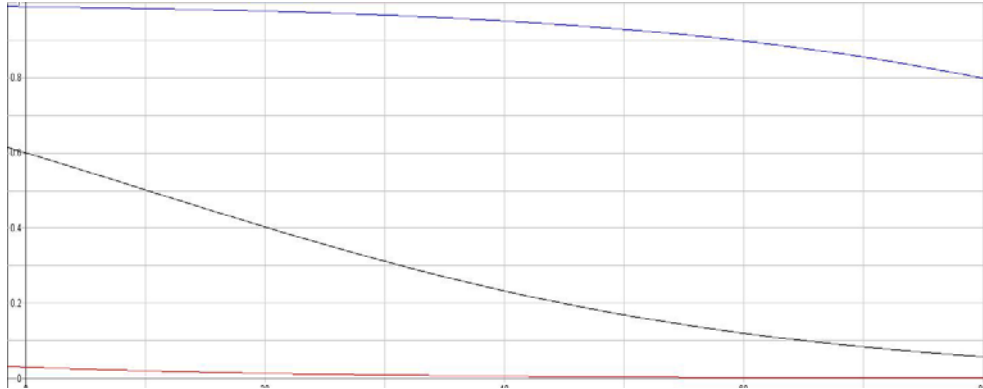
Here we let "lived" vary and set in reasonable values for other variables

Spring 2010

© Erling Berge 2010

420

Conditional effect plot from Hamilton table 7.4 (fig7.5):
effect of living for a long time in town



$$y = 1 / (1 + \exp(-2.182 - 0.04 \times 0.197 \times 12.95 + 1.299 \times 0.28 + 2.279 \times 0.31 - 1.731 \times 0.17))) \quad \text{Mean}$$

$$y = 1 / (1 + \exp(-2.182 - 0.04 \times 0.197 \times 6 + 1.299 \times 1 + 2.279 \times 1 - 1.731 \times 0))) \quad \text{Max}$$

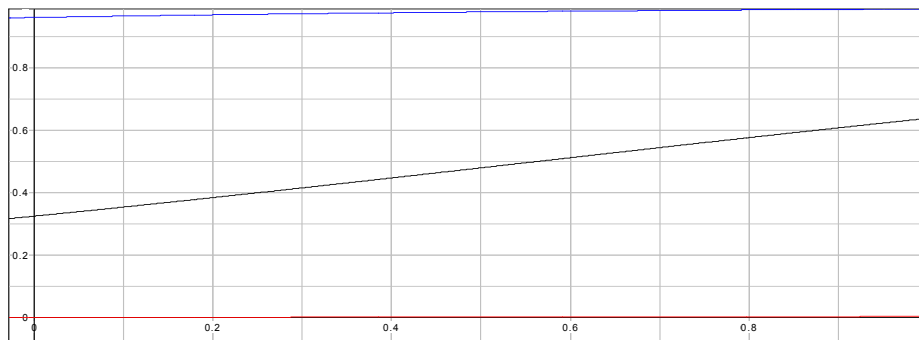
$$y = 1 / (1 + \exp(-2.182 - 0.04 \times 0.197 \times 20 + 1.299 \times 0 + 2.279 \times 0 - 1.731 \times 1))) \quad \text{Min}$$

Spring 2010

© Erling Berge 2010

421

Conditional effect plot from Hamilton table 7.4 (fig7.6):
effect of pollution on own land



$$y = 1 / (1 + \exp(-2.182 - 0.04 \times 19.27 - 0.197 \times 12.95 + 1.299 \times 1 + 2.279 \times 0.31 - 1.731 \times 0.17))) \quad \text{Mean}$$

$$y = 1 / (1 + \exp(-2.182 - 0.04 \times 1 - 0.197 \times 6 + 1.299 \times 1 + 2.279 \times 1 - 1.731 \times 0))) \quad \text{Max}$$

$$y = 1 / (1 + \exp(-2.182 - 0.04 \times 81 - 0.197 \times 20 + 1.299 \times 1 + 2.279 \times 0 - 1.731 \times 1))) \quad \text{Min}$$

Spring 2010

© Erling Berge 2010

422

Coefficients of determination

- Logistic regression does not provide measures comparable to the coefficient of determination in OLS regression
- Several measures analogous to R^2 have been proposed
- They are often called pseudo R^2
- Hamilton uses Aldrich and Nelson's pseudo $R^2 = \chi^2/(\chi^2+n)$
 where χ^2 = test statistic for the test of the whole model against a model with just a constant and n = the number of cases

Spring 2010

© Erling Berge 2010

423

Some pseudo R^2 in SPSS

- SPSS reports Cox and Snell, Nagelkerke, and in multinomial logistic regression also McFadden's proposal for R^2
- Aldrich and Nelson's pseudo R^2 can easily be computed by ourselves [pseudo $R^2 = \chi^2/(\chi^2+n)$]

Model Summary

Model Summary				Pseudo R-Square	
Step	-2 Log likelihood	Cox & Snell R Square	Nagelkerke R Square		
1	***	***	***	Cox and Snell	***
				Nagelkerke	***
				McFadden	***

Spring 2010

© Erling Berge 2010

424

Statistical problem: linearity of the logit

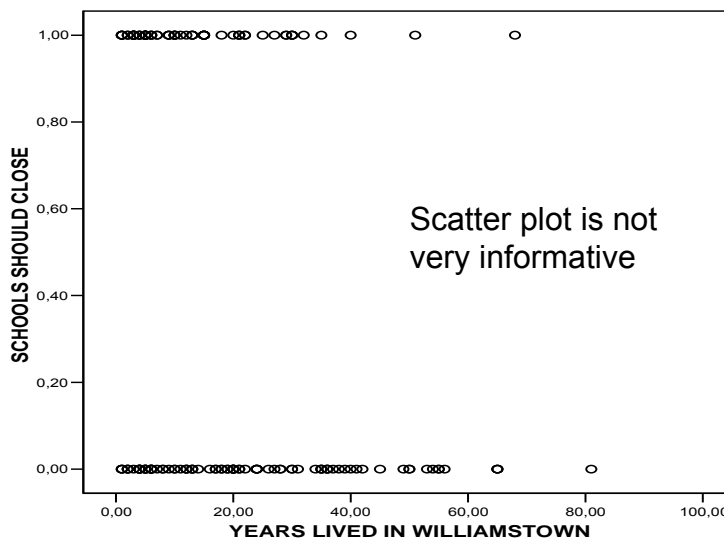
- Curvilinearity of the logit can give biased parameter estimates
- Scatter plot for $y - x$ is not informative since y only has 2 values
- To test if the logit is linear in an x -variable one may do as follows
 - Group the x variable
 - For every group find average of y and compute the logit for this value
 - Make a graph of the logits against the grouped x

Spring 2010

© Erling Berge 2010

425

Y="Closing school" vs. x= "Years lived in town"



Spring 2010

© Erling Berge 2010

426

Linearity in logit: example

Recall: $\text{Logit} = L_i = \ln(O_i) = \ln\{p_i/(1-p_i)\}$

SCHOOLS SHOULD CLOSE		YEARS LIVED IN WILLIAMSTOWN (Banded)						
		<= 3	4-6	7-11	12-22	23-33	34-44	45+
N	OPEN	7	14	7	22	11	13	13
N	CLOSE	13	14	10	17	8	2	2
Within group	Mean (=p)	,65	,50	,59	,44	,42	,13	,13
Logit	$\ln(p/(1-p))$	0,619	0	0,364	-0,241	-0,323	-1,901	-1,901

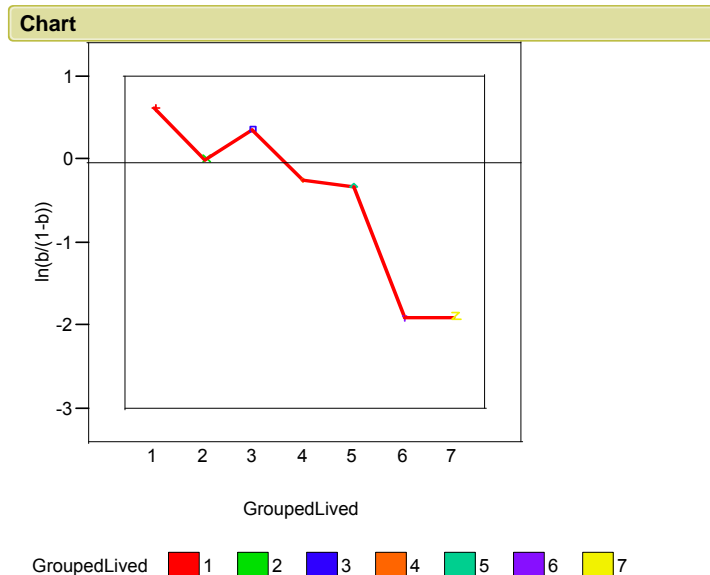
Spring 2010

© Erling Berge 2010

427

Is the logit linear in "years lived in town"?

Maybe!



Spring 2010

© Erling Berge 2010

428

In case of curvilinearity the odds ratio is non-constant

Assume the logit is curvilinear in education. Then the odds ratio for answering yes, adding one year of education, is:

$$\frac{e^{b_0 + b_a * Alder + b_k * Kvinne + b_{utd} * (E.utd+1) + b_{utd2} * (E.utd+1)^2}}{e^{b_0 + b_a * Alder + b_k * Kvinne + b_{utd} * E.utd + b_{utd2} * E.utd^2}} =$$

$$\frac{e^{b_{utd} + b_{utd2} * (E.utd^2 + 2E.utd + 1)}}{e^{b_{utd2} * E.utd^2}} = \frac{e^{b_{utd} + b_{utd2} * (2E.utd + 1)}}{e^0} = e^{b_{utd} + b_{utd2} * (2E.utd + 1)}$$

Spring 2010

© Erling Berge 2010

429

Statistical problems: influence

- Influence from outliers and unusual x-values are just as problematic in logistic regression as in OLS regression
- Transformation of x-variables to symmetry will minimize the influence of extreme variable values
- Large residuals are indicators of large influence

Spring 2010

© Erling Berge 2010

430

Influence: residuals

- There are several ways to standardize residuals
 - "Pearson residuals"
 - "Deviance residuals"
- Influence can be based on
 - Pearson residual
 - Deviance residual
 - Leverage (potential for influence): i.e. the statistic h_j

Spring 2010

© Erling Berge 2010

431

Diagnostic graphs

Outlier plots can be based on plots of estimated probability of $Y_i=1$ (estimated P_i) against

- Delta* B , ΔB_j , or
- Delta* Pearson Chisquare, $\Delta \chi^2_{P(j)}$, or
- Delta* Deviance Chisquare, $\Delta \chi^2_{D(j)}$

* "Delta" can be translated as "change in"

Spring 2010

© Erling Berge 2010

432

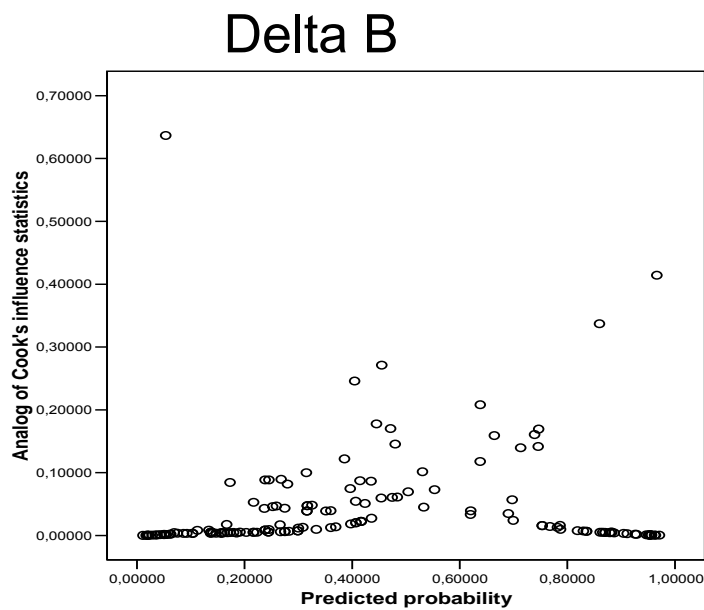
SPSS output

- **Cook's = delta B in Hamilton**
 - The logistic regression analogue of Cook's influence statistic. A measure of how much the residuals of all cases would change if a particular case were excluded from the calculation of the regression coefficients.
- **Leverage Value = h in Hamilton**
 - The relative influence of each observation on the model's fit.
- **DfBeta(s)** is not used by Hamilton in logistic regression
 - The difference in beta value is the change in the regression coefficient that results from the exclusion of a particular case. A value is computed for each term in the model, including the constant.

Spring 2010

© Erling Berge 2010

433



Spring 2010

© Erling Berge 2010

434

SPSS output from "Save" (1)

- **Unstandardized Residuals**

- The difference between an observed value and the value predicted by the model.

- **Logit Residual**

$$\tilde{e}_i = \frac{e_i}{\hat{\pi}_i(1 - \hat{\pi}_i)}; \text{ where } e_i = y_i - \hat{\pi}_i$$

π_i is the probability that $y_i = 1$; the "hat" means estimated value

Spring 2010

© Erling Berge 2010

435

SPSS output from "Save" (2)

- **Standardized = Pearson residual**

- The command "standardized" will make SPSS write a variable called ZRE_1 and labelled "Normalized residual"
- This is the same as the Pearson residual in Hamilton

- **Studentized = [SQRT(delta deviance chisquare)]**

- The command "Studentized" will make SPSS write a variable called SRE_1 and labelled "Standardized residual"
- This is the same as the square root of "delta Deviance chisquare" in Hamilton, i.e. "delta Deviance chisquare" = $(SRE_1)^2$

- **Deviance = Deviance residual**

- The command "Deviance" will make SPSS write a variable called DEV_1 and labelled "Deviance value"
- This is the same as the deviance residual in Hamilton

Spring 2010

© Erling Berge 2010

436

Computation of $\Delta\chi^2_{P(i)}$

- Based on the quantities provided by SPSS we can compute "delta Pearson chisquare"
- Where it says r_j in the formula we put in ZRE_1 and where it says h_j we put in LEV_1

$$\Delta\chi^2_{P(j)} = \frac{r_j^2}{(1-h_j)}$$

Spring 2010

© Erling Berge 2010

437

Computation of $\Delta\chi^2_{D(i)}$

Based on the quantities provided by SPSS we can compute "Delta Deviance Chisquare"

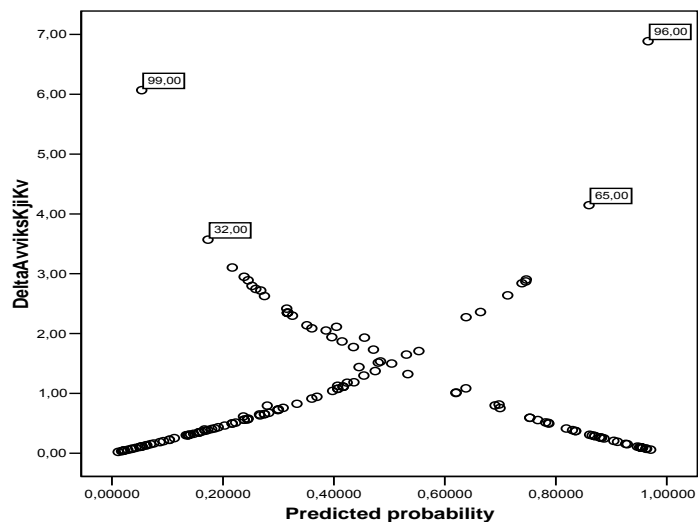
1. To find "delta deviance chisquare" we square SRE_1 $\Delta\chi^2_{D(j)} = SRE_1 * SRE_1$
2. Alternatively we put in $d_j=DEV_1$ and $h_j=LEV_1$ in the formula $\Delta\chi^2_{D(j)} = \frac{d_j^2}{(1-h_j)}$

Spring 2010

© Erling Berge 2010

438

DeltaDevianceChisquare (with/CaseNO)

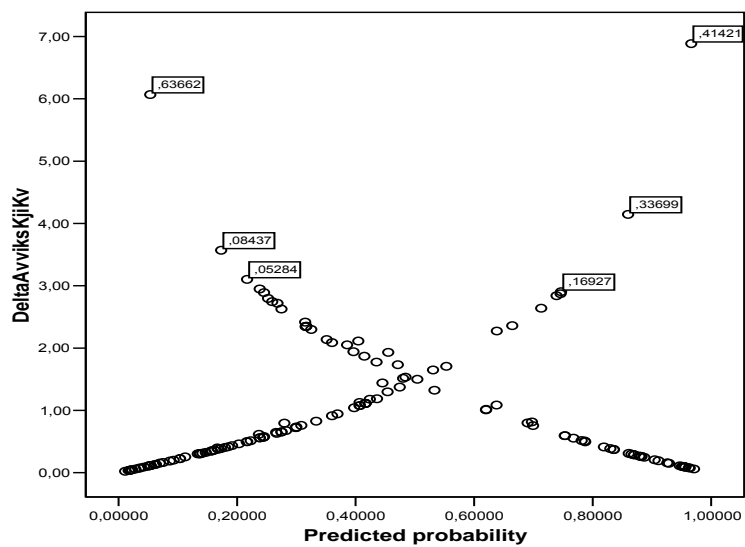


Spring 2010

© Erling Berge 2010

439

DeltaDevianceChisquare (with/delta B)

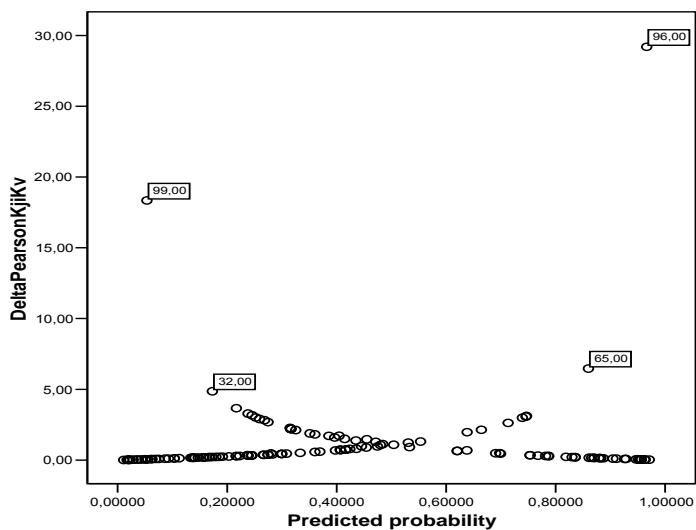


Spring 2010

© Erling Berge 2010

440

Delta Pearson Chisquare (with/CaseNO)

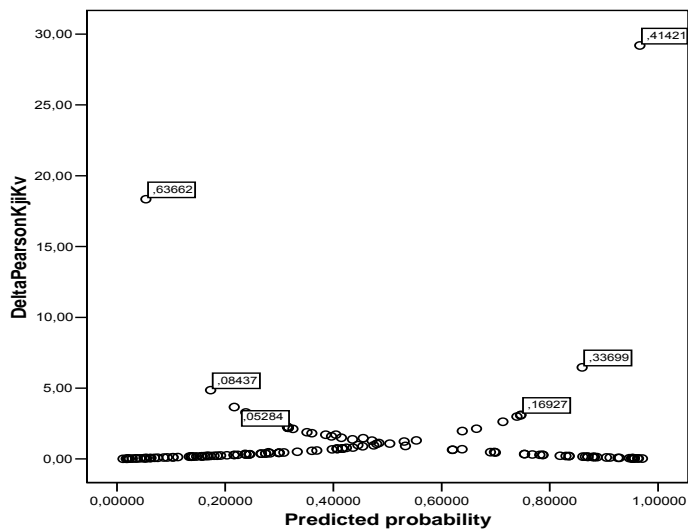


Spring 2010

© Erling Berge 2010

441

Delta Pearson Chisquare (with/ delta B)



Spring 2010

© Erling Berge 2010

442

Cases with large influence

Variables	CaseN o9 6	CaseN o6 5	CaseN o9 9	Variables	CaseN o9 6	CaseN o6 5	CaseN o9 9
Y=close	1,00	,00	,00	ZRE_1	4,21	-2,48	-5,36
lived	68,00	40,00	1,00	DEV_1	2,42	-1,98	-2,61
educ	12,00	12,00	12,00	DFB0_1	-,32	,01	-,36
contam	,00	1,00	1,00	DFB1_1	,01	,00	,00
hsc	,00	1,00	1,00	DFB2_1	,02	,01	,02
nodad	,00	,00	,00	DFB3_1	-,08	-,15	-,18
PRE_1	,05	,86	,97	DFB4_1	-,06	-,17	-,19
COO_1	,64	,34	,41	DFB5_1	-,08	,16	,14
RES_1	,95	-,86	-,97	DeltaPearsonKjiKv	18,34	6,47	29,20
SRE_1	2,46	-2,04	-2,62	DeltaAvviksKjiKv	6,07	4,14	6,89

Spring 2010

© Erling Berge 2010

443

From Cases to Patterns

- The figures shown previously are not identical to those you see in Hamilton
- Hamilton has corrected for the effect of identical patterns

Spring 2010

© Erling Berge 2010

444

Influence from a shared pattern of x-variables

- In a logistic regression with few variables many cases will have the same value on all x-variables. Every combination of x-variable values is called a pattern
- When many cases have the same pattern, every case may have a small influence, but collectively they may have unusually large influence on parameter estimates
- Influential patterns in x-values can give biased parameter estimates

Spring 2010

© Erling Berge 2010

445

Influence: Patterns in x-values

- Predicted value, and hence the residual will be the same for all cases with the same pattern
- Influence from pattern j can be found by means of
 - The frequency of the pattern
 - Pearson residual
 - Deviance residual
 - Leverage: i.e. the statistic h_j

Spring 2010

© Erling Berge 2010

446

Finding X-pattern by means of SPSS

- In the "Data" – menu find the command "Identify duplicate cases"
- Mark the x-variables that are used in the model and move them to "Define matching cases by"
- Cross for "Sequential count of matching cases in each group" and "Display frequencies for created variables"
- This produces two new variables. One, "MatchSequence", numbers cases sequentially 1, 2, ... where several patterns are identical. If the pattern is unique this variable has the value 0.
- The other variable, "Primary...", has the value 0 for duplicates and 1 for unique patterns

Spring 2010

© Erling Berge 2010

447

X-patterns in SPSS; Hamilton p238-242

	Frequency	Percent	Valid Percent	Cumulative Percent
Duplicate Case	21	13,7	13,7	13,7
Primary Case	132	86,3	86,3	100,0
Total	153	100,0	100,0	

Sequential count of matching cases	Frequency	Percent	Valid Percent	Cumulative Percent
0 [115 patterns with 1 case]	115	75,2	75,2	75,2
1 [17 patterns with 2 or 3 cases]	17	11,1	11,1	86,3
2 [17-4=13 patterns with 2 cases]	17	11,1	11,1	97,4
3 [4 patterns with 3 cases]	4	2,6	2,6	100,0
Total	153	100,0	100,0	

Spring 2010

© Erling Berge 2010

448

Hamilton table 7.6 Symbols

J	# unique patterns of x-values in the data ($J \leq n$)
m_j	# cases with the pattern j ($m_j \geq 1$)
\hat{p}_j	Predicted probability of $Y=1$ for case with pattern j
Y_j	Sum of y-values for cases with pattern j (= # cases with pattern j and $y=1$)
r_j	Pearson residual for pattern j
χ_P^2	Pearson Chisquare statistic
d_j	Deviance residual for pattern j
χ_D^2	Deviance Chisquare statistic
h_i	Leverage for case i
h_j	Leverage for pattern j

Spring 2010

© Erling Berge 2010

449

New values for $\Delta\chi^2_{P(i)}$ and $\Delta\chi^2_{D(i)}$

- By "Compute" one may calculate the Pearson residual (equation 7.19 in Hamilton) and delta Pearson chisquare (equation 7.24 in Hamilton) once more. This will provide the correct values
- The same applies for deviance residual (equation 7.21) and delta deviance chisquare (equation 7.25a)

Spring 2010

© Erling Berge 2010

450

Leverage and residuals (1)

- Leverage of a pattern is obtained as number of cases with the pattern times the leverage of a case with this pattern. The leverage of a case is the same as in OLS regression
- $h_j = m_j \cdot h_i$
- Pearson residual can be found from

$$r_j = \frac{Y_j - m_j \hat{P}_j}{\sqrt{m_j \hat{P}_j (1 - \hat{P}_j)}}$$

Spring 2010

© Erling Berge 2010

451

Leverage and residuals (2)

- Deviance residual can be found from

$$d_j = \pm \sqrt{\left\{ 2 \left[Y_j \ln \left(\frac{Y_j}{m_j \hat{P}_j} \right) + (m_j - Y_j) \ln \left(\frac{m_j - Y_j}{m_j (1 - \hat{P}_j)} \right) \right] \right\}}$$

Spring 2010

© Erling Berge 2010

452

Two Chi-square statistics

- Pearson Chi-square statistics

$$\chi_P^2 = \sum_{j=1}^J r_j^2$$

- Deviance Chi-square statistics

$$\chi_D^2 = \sum_{j=1}^J d_j^2$$

- Equations are the same for both cases and patterns

Spring 2010

© Erling Berge 2010

453

The Chisquare statistics

Both Chisquare statistics:

1. Pearson-Chisquare χ_P^2 and
 2. Deviance-Chisquare χ_D^2
- Can be read as a test of the null hypothesis of no difference between the estimated model and a “saturated model”, that is a model with as many parameters as there are cases/ patterns

Spring 2010

© Erling Berge 2010

454

Large values of measures of influence

- Measures of influence based on changes in (Δ) the statistic/ parameter value due to excluded cases with pattern j
 - ΔB_j “delta B” - analogue to Cook’s D
 - $\Delta\chi^2_{P(i)}$ “delta Pearson-Chisquare”
 - $\Delta\chi^2_{D(i)}$ “delta Deviance-Chisquare”

Spring 2010

© Erling Berge 2010

455

What is a large value of $\Delta\chi^2_{P(i)}$ and $\Delta\chi^2_{D(i)}$

- Both $\Delta\chi^2_{P(i)}$ and $\Delta\chi^2_{D(i)}$ measure how badly the model fits the pattern j. Large values indicates that the model would fit the data much better if all cases with this pattern were excluded
- Since both measures are distributed asymptotically as the chisquare distribution, values larger than 4 indicate that a pattern affects the estimated parameters “significantly”

Spring 2010

© Erling Berge 2010

456

ΔB_j "delta B"

- Measures the standardized change in the estimated parameters (b_k) that obtain when all cases with a given pattern j are excluded

$$\Delta B_j = \frac{r_j^2 h_j}{(1 - h_j)^2}$$

Larger values means larger influence

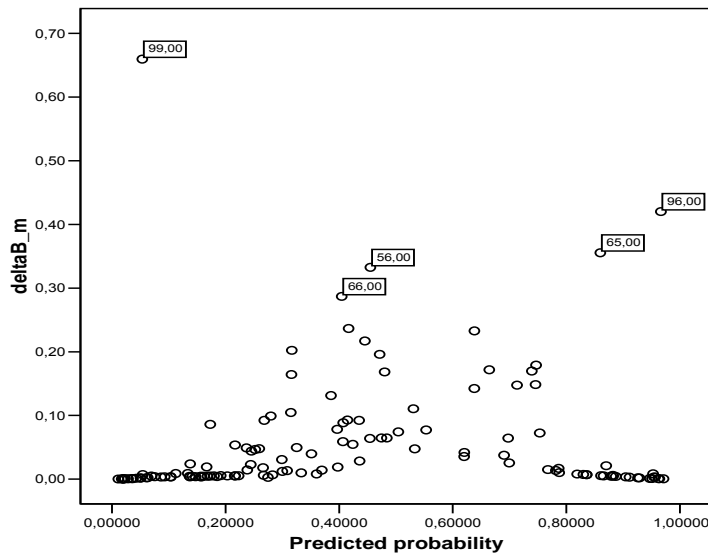
$\Delta B_j \geq 1$ must in any case be seen as "large influence"

Spring 2010

© Erling Berge 2010

457

delta B (with caseNO)



Spring 2010

© Erling Berge 2010

458

$\Delta\chi^2_{P(i)}$ "Delta Pearson Chisquare"

- Measures the reduction in Pearson χ^2 that obtains from excluding all cases with pattern j

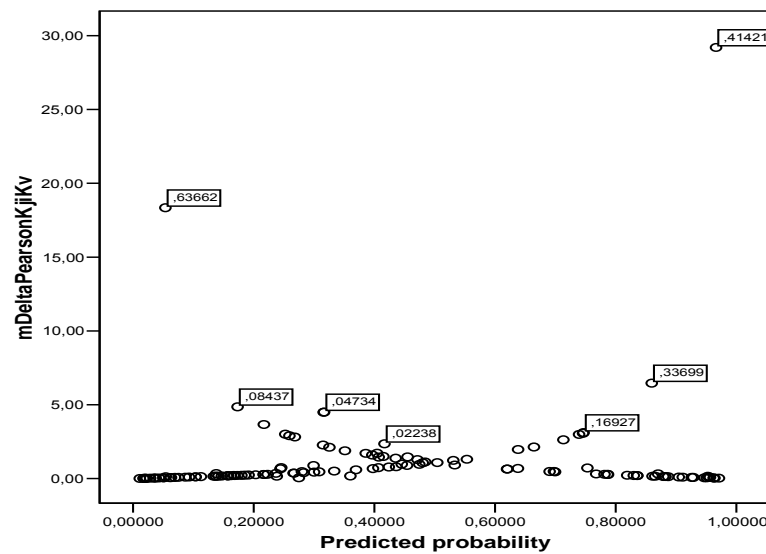
$$\Delta\chi^2_{P(j)} = \frac{r_j^2}{(1-h_j)}$$

Spring 2010

© Erling Berge 2010

459

Delta Pearson Chisquare (with delta B)



Spring 2010

© Erling Berge 2010

460

$\Delta\chi^2_{D(i)}$ “Delta Deviance Chisquare”

- Measures changes in deviance that obtains from excluding all cases with pattern j

$$\Delta\chi^2_{D(j)} = \frac{d_j^2}{(1-h_j)}$$

- This is equivalent to

$$\Delta\chi^2_{D(j)} = -2 \left[\mathcal{L}\mathcal{L}_K - \mathcal{L}\mathcal{L}_{K(j)} \right]$$

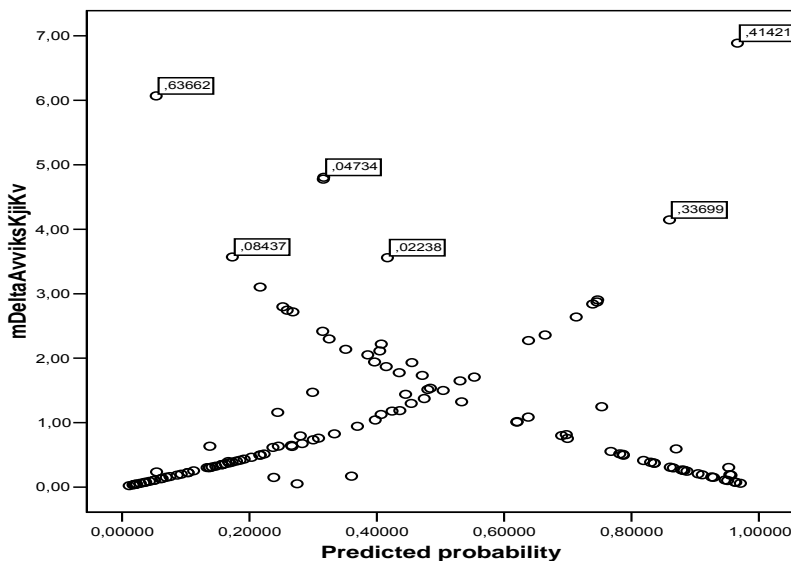
$\mathcal{L}\mathcal{L}_K$ is the LogLikelihood of a model with K parameters estimated on the whole sample and $\mathcal{L}\mathcal{L}_{K(j)}$ is from the estimate of the same model when all cases with pattern j are excluded

Spring 2010

© Erling Berge 2010

461

Delta Deviance Chisquare (with delta B)



Spring 2010

© Erling Berge 2010

462

Influence of excluded cases/patterns

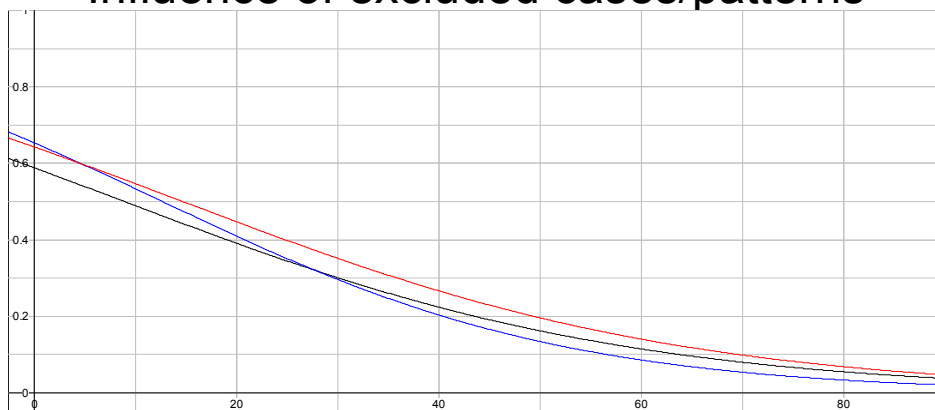
Variables in the model	Logit coefficient		
	Sample	Excluding case 99 $\Delta\chi^2P(i) = 18,34$	Excluding case 96 $\Delta\chi^2P(i) = 29,20$
lived	-,040	-,045	-,052
educ	-,197	-,224	-,214
contam	1,299	1,490	1,382
hsc	2,279	2,492	2,347
nodad	-1,731	-1,889	-1,658
Constant	2,182	2,575	2,530
2*LL(modell)	-142,652	-135,425	-136,124

Spring 2010

© Erling Berge 2010

463

Influence of excluded cases/patterns



$$y = 1 / (1 + \exp(- (2.18 - 0.04x - 0.2 \times 13 + 1.3 \times 0.28 + 2.28 \times 0.31 - 1.73 \times 0.17)))$$

$$y = 1 / (1 + \exp(- (2.53 - 0.05x - 0.21 \times 13 + 1.38 \times 0.28 + 2.35 \times 0.31 - 1.65 \times 0.17)))$$

$$y = 1 / (1 + \exp(- (2.58 - 0.04x - 0.22 \times 13 + 1.49 \times 0.28 + 2.49 \times 0.31 - 1.89 \times 0.17)))$$

Spring 2010

© Erling Berge 2010

464

Conclusions (1)

Ordinary OLS do not work well for dichotomous dependent variables since

- It is impossible to obtain normally distributed errors or homoscedasticity, and since
- The model predicts probabilities outside the interval [0-1]

The Logit model is for theoretical reasons better

- Likelihood ratio tests statistic can be used to test nested models analogous to the F-statistic
- In large samples the chisquare distributed Wald statistic [or the normally distributed $t = \sqrt{\text{Wald}}$] will be able to test single coefficients and provide confidence intervals
- There is no statistic similar to the coefficient of determination

Spring 2010

© Erling Berge 2010

465

Conclusions (2)

- Coefficient of estimated models can be interpreted by
 1. Log-odds (direct interpretation)
 2. Odds
 3. Odds ratio
 4. Probability (conditional effect plot)
- Non-linearity, case with influence, and multicollinearity leads to the same kinds of problems as in OLS regression (inaccurate or uncertain parameter values)
- Discrimination leads to problems of uncertain parameter values (large variance estimates)
- Diagnostic work is important

Spring 2010

© Erling Berge 2010

466

Causal analysis

Structural equation models

- Hamilton, Lawrence C. 2008. A Low-Tech Guide to Causal Modelling.
<http://pubpages.unh.edu/~lch/causal2.pdf>

Spring 2010

© Erling Berge 2010

467

Causal analysis

- Experiment
 - Randomized causal impacts ("treatment") provide precise causal conclusions about effects ("response") if there is significant differences in the mean response (effect)
 - Experiments can be impossible to achieve due to
 - Practical conditions
 - Economic constraints
 - Ethical judgements
- Instead one tries to obtain quasi-experiments
 - Using for example regression analysis

Spring 2010

© Erling Berge 2010

468

Model of causal effects Ref.:

- Research using observations utilize concepts from experimental design
 - “Treatment”, “Stimulus”
 - “Effect”, “Outcome”

Ref.:

Winship, Christopher, and Stephen L. Morgan 1999 “The Estimation of Causal Effects from Observational Data”, Annual Review of Sociology Vol 25: 659-707

Spring 2010

© Erling Berge 2010

469

Experiments allocate “cases” randomly to one of two groups:

- | | |
|--------------------------------------------|-------------------------------------------|
| • TREATMENT (T)
with observation | • CONTROLL (C)
with observation |
| – before treatment | – before non-treatment |
| – after treatment | – after non-treatment |

Spring 2010

© Erling Berge 2010

470

The counterfactual hypothesis for the study of causality

- Individual "i" can a priori be assumed selected for one of two groups
 - Treatment group, T, or control group, C.
- Treatment, t, as well as non-treatment, c, can a priori be given to individuals both in the T- and C-group
- In reality we are able to observe t only in the T-group and c in the C-group

Spring 2010

© Erling Berge 2010

471

Modelling of causal effects: The counterfactual hypothesis (1)

- There are for each individual "i" four possible outcomes
 - $Y_i(\mathbf{c}, \mathbf{C})$ or $Y_i(t, \mathbf{C})$; if allocated to a control group
 - $Y_i(\mathbf{c}, \mathbf{T})$ or $Y_i(\mathbf{t}, \mathbf{T})$; if allocated to a treatment group
 - Only $Y_i(\mathbf{c}, \text{given that "i" is a member of C})$ or
 - $Y_i(\mathbf{t}, \text{given that "i" is a member of T})$ can be observed for any particular individual

Spring 2010

© Erling Berge 2010

472

Modelling of causal effects:

The counterfactual hypothesis (2)

More formally one may write the possible outcomes for person no i :

	Treatment: t	Non-treat.: c
T-group	$Y_i^t \in T$	$Y_i^c \in T$
C-group	$Y_i^t \in C$	$Y_i^c \in C$

Spring 2010

© Erling Berge 2010

473

Modelling of causal effects:

The counterfactual hypothesis (3)

- Then the causal effect for individual i is
- $\delta_i = Y_i(t) - Y_i(c)$
- Only one of these two quantities can be observed for any given individual
- This leads to the “counterfactual hypothesis”

Spring 2010

© Erling Berge 2010

474

The counterfactual hypothesis: concluding

- “The main value of this counterfactual framework is that causal inference can be summarized by a single question: Given that the δ_i cannot be calculated for any individual and therefore that Y_i^t and Y_i^c can be observed only on mutually exclusive subsets of the population, what can be inferred about the distribution of the δ_i from an analysis of Y_i and T_i ?” (Winship and Morgan 1999:664)

Spring 2010

© Erling Berge 2010

475

Modelling of causal effects: from individual effects to population averages

- We can observe $Y_i(c | i \in C)$, but not $Y_i(t | i \in C)$
- The problem may be called a problem of missing data
- Instead of individual effects we can estimate average effects for the total population

Spring 2010

© Erling Berge 2010

476

Modelling of causal effects (1)

- Average effects can be estimated, but usually it involves difficulties
- One assumption is that the effect of the treatment will be the same for any given individual independent of which group the individual is allocated to
- This, however, is not self-evident

Spring 2010

© Erling Berge 2010

477

Modelling of causal effects (2)

The counterfactual hypothesis assumes:

- That changing the treatment group for one individual do not affect the outcome of other individuals (no interaction)
- That treatment in reality can be manipulated (e.g. sex can not be manipulated)

Spring 2010

© Erling Berge 2010

478

Modelling of causal effects (3)

- One problem is that in a sample the process of allocating person no i to a control or treatment group may affect the estimated average effect (*the problem of selection*)
- In some cases, however, the interesting quantity is the average effect for those who actually receive the treatment

Spring 2010

© Erling Berge 2010

479

Modelling of causal effects (4)

- It can be shown that there are two sources of bias for the estimates of the average effect
 1. An established difference between the C- and T- groups
 2. The treatment works in principle differently for those allocated to the T-group compared to those in the C-group
 - To counteract this one has to develop models of how people get into C- and T-groups (**selection models**)

Spring 2010

© Erling Berge 2010

480

Modelling of causal effects (5)

- A general class of methods that may be used to estimate causal effects are the **regression models**
- These are able to “control for” observable differences between the C- and T- groups, but not for unequal response to treatment

Spring 2010

© Erling Berge 2010

481

Causal modelling

- “path analysis” or “structural equations modelling” go back to the 60ies
- Jöreskog and Sörbom: LISREL
 - Use maximum likelihood to estimate model parameters maximising fit to the variance-covariance matrix
 - Commonly available in statistical packages
 - Covariance structural modelling
 - Structural equation modelling
 - Full information maximum likelihood estimation

Spring 2010

© Erling Berge 2010

482

Structural equation models: Low-Tech approach

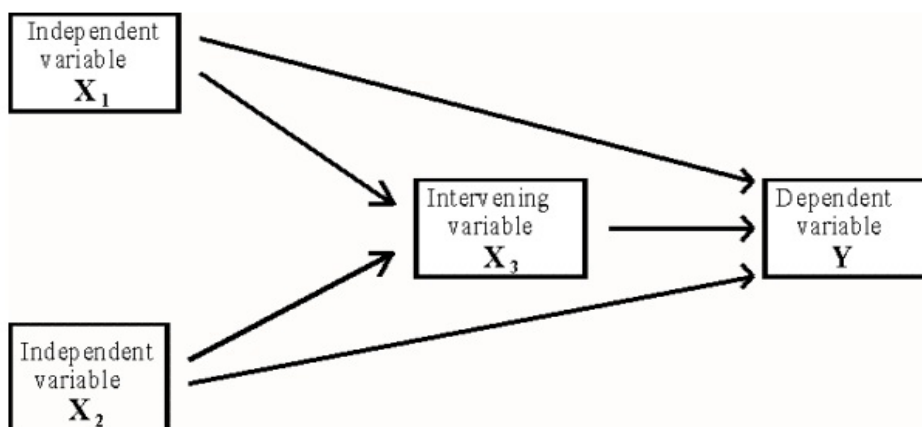
- Uses OLS to do simple versions of the structural equations models
- The key assumption is the causal ordering of variables. In survey data this ordering is supplied by theory
- The causal diagram visualize the order of causation:
 - Causality flows from left to right
 - Intervening variables give rise to indirect effects
 - “reverse causation” creates problems

Spring 2010

© Erling Berge 2010

483

Low-Tech causal modelling
Figure 1

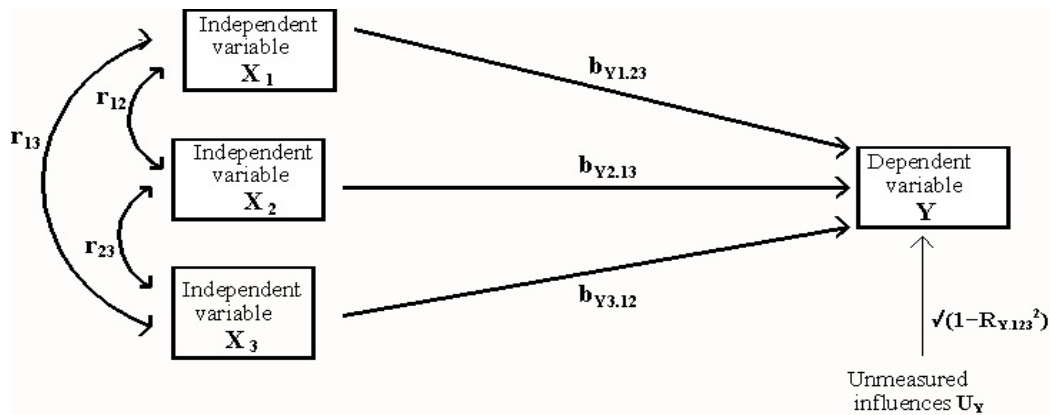


Spring 2010

© Erling Berge 2010

484

Multiple regression as a causal model Figure 2



Spring 2010

© Erling Berge 2010

485

Quantities in the diagram

r_{12}, r_{13}, r_{23}	Pearson correlations among x-variables
$b_{Y1.23}$, etc.	Usually a standardized regression coefficient (“beta weight”) taken from the regression of Y on X_1 , and “.” means controlled for X_2, X_3
$R_{Y.123}^2$	Coefficient of determination R^2 from the regression of Y on X_1, X_2, X_3
$\sqrt{1 - R_{Y.123}^2}$	Is an estimate of unmeasured influences called error term or disturbance

Spring 2010

© Erling Berge 2010

486

Multiple regression

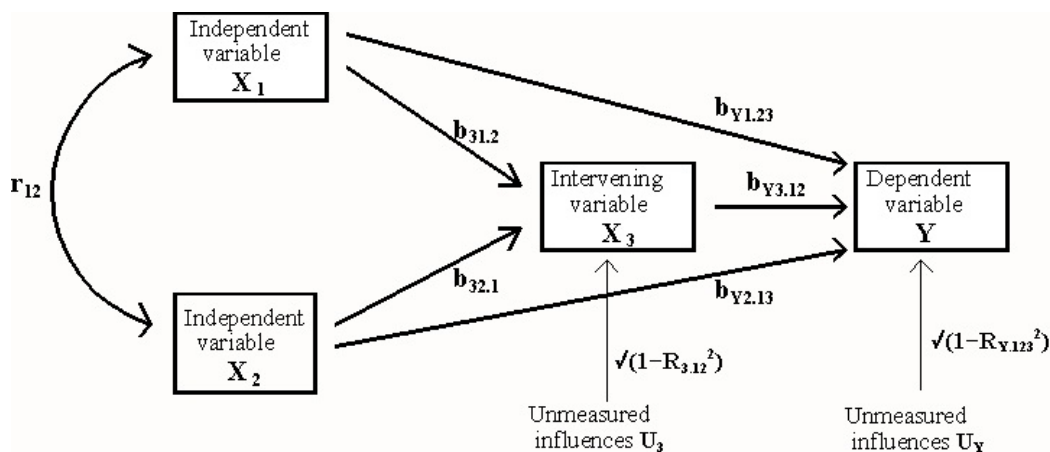
- All assumptions and all problems apply as before
 - Note in particular that error terms must be uncorrelated with included x-variables (all relevant variable are included)
- If some of the X-es are intervening in figure 2 the model is too simple, but it matters only if we are interested in causality

Spring 2010

© Erling Berge 2010

487

Path coefficients
Figure 3



Spring 2010

© Erling Berge 2010

488

New elements in figure 3

$b_{31.2}, b_{32.1}$	Standardized regression coefficients (“beta weight”) from the regression of X_3 on X_1 controlled for X_2 and from the regression of X_3 on X_2 controlled for X_1
$R_{3.12}^2$	Coefficient of determination (R^2) from the regression of X_3 on X_1 and X_2
$\sqrt{1-R_{3.12}^2}$	The error term from the regression of X_3 on X_1 and X_2

Spring 2010

© Erling Berge 2010

489

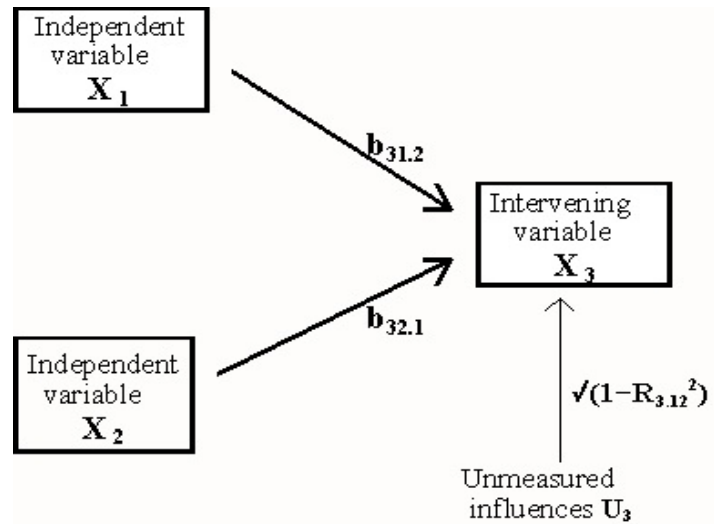
The structural model of figure 3

- $\hat{Y} = b_{Y1.23}X_1 + b_{Y2.13}X_2 + b_{Y3.12}X_3$
- $\hat{X}_3 = b_{31.2}X_1 + b_{32.1}X_2$
- In structural equations variables and coefficients are standardized
- That means that variables have an average of 0 and a standard deviation of 1 and that coefficients vary between -1 and +1

Spring 2010

© Erling Berge 2010

490

Figure 5: the regression of X_3 on X_1 and X_2 

Spring 2010

© Erling Berge 2010

491

Direct, Indirect and Total Effects

- *Indirect effects* equal the product of coefficients along any series of causal paths that link one variable to another
- *Total effects* equal the sum of all direct and indirect effects linking two variables

Spring 2010

© Erling Berge 2010

492

Indirect effects as products of path coefficients

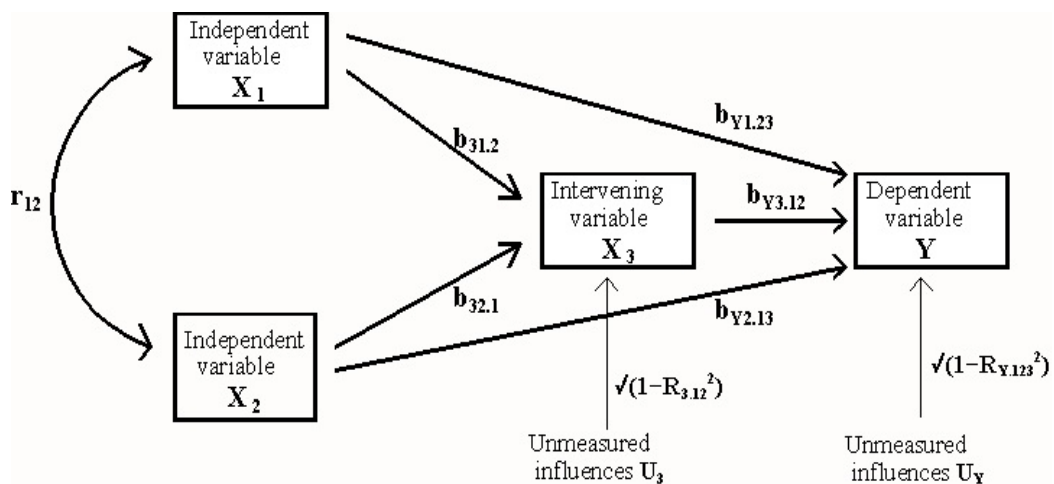
- $\hat{Y} = b_{Y1.23}X_1 + b_{Y2.13}X_2 + b_{Y3.12}X_3$
- $\hat{X}_3 = b_{31.2}X_1 + b_{32.1}X_2$
- Means that we have
- $\hat{Y} = b_{Y1.23}X_1 + b_{Y2.13}X_2 + b_{Y3.12}X_3$
- $= b_{Y1.23}X_1 + b_{Y2.13}X_2 + b_{Y3.12}(b_{31.2}X_1 + b_{32.1}X_2)$
- $= b_{Y1.23}X_1 + b_{Y2.13}X_2 + b_{Y3.12}b_{31.2}X_1 + b_{Y3.12}b_{32.1}X_2$
- $= (b_{Y1.23} + b_{Y3.12}b_{31.2})X_1 + (b_{Y2.13} + b_{Y3.12}b_{32.1})X_2$
- Compare compound coefficients to the diagram

Spring 2010

© Erling Berge 2010

493

Structural model



Spring 2010

© Erling Berge 2010

494

Path Coefficients

- X_1 to Y : $\mathbf{b}_{Y1.23}$ (regression coefficient of Y on X_1 , controlling for X_2 and X_3)
- X_2 to Y : $\mathbf{b}_{Y2.13}$ (regression coefficient of Y on X_2 , controlling for X_1 and X_3)
- X_3 to Y : $\mathbf{b}_{Y3.12}$ (regression coefficient of Y on X_3 , controlling for X_1 and X_2)
- X_1 to X_3 : $\mathbf{b}_{31.2}$ (regression coefficient of X_3 on X_1 , controlling for X_2)
- X_2 to X_3 : $\mathbf{b}_{32.1}$ (regression coefficient of X_3 on X_2 , controlling for X_1)

Spring 2010

© Erling Berge 2010

495

Direct effects

X_1 to Y : $\mathbf{b}_{Y1.23}$	regression coefficient of Y on X_1 , controlling for X_2 and X_3
X_2 to Y : $\mathbf{b}_{Y2.13}$	regression coefficient of Y on X_2 , controlling for X_1 and X_3
X_3 to Y : $\mathbf{b}_{Y3.12}$	regression coefficient of Y on X_3 , controlling for X_1 and X_2
X_1 to X_3 : $\mathbf{b}_{31.2}$	regression coefficient of X_3 on X_1 , controlling for X_2
X_2 to X_3 : $\mathbf{b}_{32.1}$	regression coefficient of X_3 on X_2 , controlling for X_1

Spring 2010

© Erling Berge 2010

496

Indirect and total effects

Indirect effects	
X_1 to Y , through X_3 :	$b_{31.2} \times b_{Y3.12}$
X_2 to Y , through X_3 :	$b_{32.1} \times b_{Y3.12}$
Total effects	
X_1 to Y :	$b_{Y1.23} + (b_{31.2} \times b_{Y3.12})$
X_2 to Y :	$b_{Y2.13} + (b_{32.1} \times b_{Y3.12})$

Spring 2010

© Erling Berge 2010

497

Additions to multiple regressions

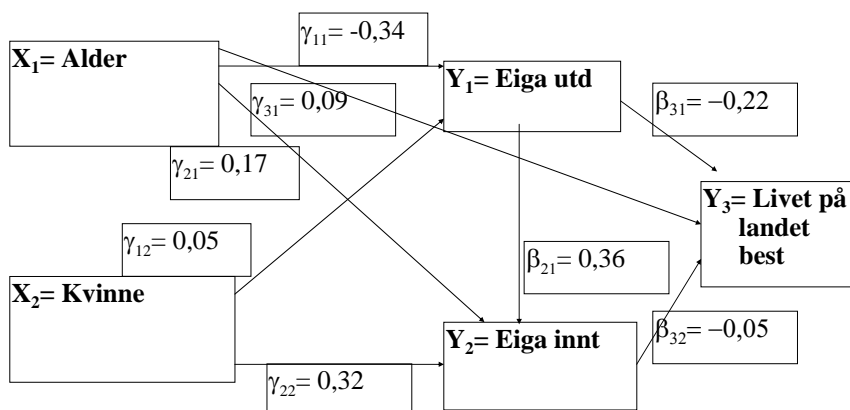
- We learn something new if the indirect effects are large enough to have substantial interest
- More than two steps of causation tends to become very weak
 - $0.3 \times 0.3 \times 0.3 = 0.027$
 - 0.3 standard deviation change in causal variables leads to a 0.027 standard deviation change in the dependent variable

Spring 2010

© Erling Berge 2010

498

Example of a model diagram with path coefficients



Figur 2.1

Note differences in symbols

Spring 2010

© Erling Berge 2010

499

Comment to the figure above

- The β coefficients go from one Y variable to another
- The γ coefficients go from one X variable a Y variable
- The coefficient indexing indicates which variables they link. The first index tells the dependent variable. The second index tells the independent variable
- The coefficients are standardized (OLS) regression coefficients (“beta weights”)

Spring 2010

© Erling Berge 2010

500

The structural model of the example

- $\hat{Y}_3 = \gamma_{31}X_1 + \gamma_{32}X_2 + \beta_{31}Y_1 + \beta_{32}Y_2$
- $\hat{Y}_2 = \gamma_{21}X_1 + \gamma_{22}X_2 + \beta_{21}Y_1$
- $\hat{Y}_1 = \gamma_{11}X_1 + \gamma_{12}X_2$
- $\hat{Y}_3 = 0.09X_1 - 0.22Y_1 - 0.05Y_2$
- $\hat{Y}_2 = 0.17X_1 + 0.32X_2 + 0.36Y_1$
- $\hat{Y}_1 = -0.34X_1 + 0.17X_2$

Spring 2010

© Erling Berge 2010

501

Direct and indirect effects on “Livet på landet best” from age

- Direct effect: $\gamma_{31} = 0.09$
- Indirect effect by way of “Eiga utd” and “Eiga innt”
- $\beta_{31} * \gamma_{11} + \beta_{32} * \beta_{21} * \gamma_{11} + \beta_{32} * \gamma_{21}$
- $(-0.22)*(-0.34)+(-0.05)*(0.36)*(-0.34)+(-0.05)*(0.17)$
- $0.22*0.34 + 0.05*0.36*0.34 - 0.05*0.17$
- $0.0748 + 0.00612 - 0.0085 = 0.07242$
- Total effect = $0.09 + 0.07242 = 0.16242$
- Increasing age by 1 st. dev. leads to an increase of 0.16 st.dev. in the strength of support for “Livet på landet best”

Spring 2010

© Erling Berge 2010

502

Variables and measurement

- All interval scale variables used in multiple regression (including non-linear transformed variables and interaction terms) can be included in structural equations models
- But interpretation becomes tricky when variables are complex. Conditional effect plots are very useful
- Robust, quantile, logit, and probit regression should not be used
- Categorical variables should not be used as intervening variables
- Scales or index variables can be used as usual in OLS regression

Spring 2010

© Erling Berge 2010

503

Concluding on structural equations modelling

- Including factors from factor analysis as explanatory variables make it possible to approximate a LISREL type analysis
- If assumptions are true LISREL will perform a much better and provides more comprehensive estimation, but too often assumptions are not true. Then the low-tech approach has access to the large toolkit of OLS regression for diagnostics and exploratory methods testing basic assumptions and discovering unusual data points
- Simple diagnostic work sometimes yields the most unexpected, interesting and replicable findings from our research

Spring 2010

© Erling Berge 2010

504

Principal components and factor analysis

Hamilton Ch 8 p249-282

Spring 2010

© Erling Berge 2010

505

Principal components and factor analysis

- Principal components and factor analysis are both methods for data reduction
- They seek underlying dimensions that are able to account for the pattern of variation among a set of observed variables
- Principal components analysis is a transformation of the observed data where the idea is to explain as much as possible of the observed variation with a minimum number of components

Spring 2010

© Erling Berge 2010

506

Factor analysis

- Estimates coefficients on - and variable values of - unobserved variables (Factors) to explain the co-variation among an observed set of variables
- The assumption is that a small set of the unobserved factors are able to explain most of the co-variation
- Hence factor analysis can be used for data reduction. Many variables can be replaced by a few factors

Spring 2010

© Erling Berge 2010

507

Factor analysis

- $Z_k = \ell_{k1}F_1 + \ell_{k2}F_2 + \dots + \ell_{kj}F_j + \dots + \ell_{kJ}F_J + u_k$
– $k = 1, 2, 3, \dots, K$
- Symbols
 - K observed variables, Z_k ; $k=1, 2, 3, \dots, K$
 - J unobserved factors, F_j ; $j=1, 2, 3, \dots, J$ where $J < K$
 - For each variable there is a unique error term, u_k , also called unique factors while the F factors are called common factors
 - For each factor there is a **standardized** regression coefficient, ℓ_{kj} , also called factor loading; k refers to variable no, j refers to factor no. An index denoting case no has been omitted here.

Spring 2010

© Erling Berge 2010

508

Correlation of factors

- Factors may be correlated or uncorrelated
 - Uncorrelated: they are then called **orthogonal**
 - Correlated: they are then called **oblique**
- Factors may be rotated
 - Oblique rotations create correlated factors
 - Orthogonal rotations create uncorrelated factors

Spring 2010

© Erling Berge 2010

509

Principal components

- Represents a simple transformation of variables. There are as many principal components as there are variables
- Principal components are uncorrelated
- $Z_k = \ell_{k1}F_1 + \ell_{k2}F_2 + \dots + \ell_{kj}F_j + \dots + \ell_{kK}F_K$
- If the last few principal components explain little variation we can retain $J < K$ components. Thus Principal Components also can be used to reduce data.
- $Z_k = \ell_{k1}F_1 + \ell_{k2}F_2 + \dots + \ell_{kj}F_j + \dots + \ell_{kJ}F_J + v_k$
 where $J < K$ and
 the residual v_k has small variance and consist of the discarded principal components

Spring 2010

© Erling Berge 2010

510

Principal components vs factor analysis

- Principal components analysis attempts to explain the observed variation of the variables
- Factor analysis attempts to explain their intercorrelations
- Use principal components to generate a composite variable that reproduce the maximum variance of observed variables
- Use factor analysis to model relationships between observed variables and unobserved latent variables and to obtain estimates of latent variable values
- The choice between the two is often blurred, to some degree it is a matter of taste

Spring 2010

© Erling Berge 2010

511

The number of principal components

- K variables yield K principal components
- If the first few components account for most of the variation, we can concentrate on them and discard the remaining
- The eigenvalues of the standardized correlation matrix provides a guide here
- Components are ranked according to eigenvalues
- A principal component with an eigenvalue $\lambda < 1$ accounts for less variance than a single variable
- Thus we discard components with eigenvalues below 1
- Another criterion for keeping components is that each component should have substantive meaning

Spring 2010

© Erling Berge 2010

512

Eigenvalues and explained variance

- In a covariance matrix the sum of eigenvalues equals the sum of variances.
- In a correlation matrix this = K (the number of variables) since each standardized variable has a variance of 1
- Thus the sum of eigenvalues of the principal components
 - $\lambda_1 + \lambda_2 + \lambda_3 + \dots + \lambda_K = K$ and
 - $\lambda_j / K =$ proportion of variance explained by component no j

Spring 2010

© Erling Berge 2010

513

Uniqueness and communality

- If K-J components are discarded and we have only J factors
- $Z_k = \ell_{k1}F_1 + \ell_{k2}F_2 + \dots + \ell_{kj}F_j + \dots + \ell_{kJ}F_J + v_k$
- And an error term v_k
- The variance of the error term is called the uniqueness of the variable
- Communality is the proportion of a variable's variance shared with the components
- Communality = $h_k^2 = 1 - \text{Uniqueness} = \sum_j \lambda_{kj}^2$, $j=1, \dots, J$; k = variable number

Spring 2010

© Erling Berge 2010

514

Rotation to simple structure

- The idea is to transform (rotate) the factors so that the loadings on each components make it easier to interpret the meaning of the component
- If the loading are close either to 1 or -1 on one factor and close to 0 on all others the structure is simpler to interpret: we rotate to “simple structure”. The rotated factors fit data equally well but are simpler to interpret
- Rotations may be
 - Orthogonal (rotation method typically: varimax)
 - Oblique (rotation method typically: oblimin, promax)

Spring 2010

© Erling Berge 2010

515

Why rotate?

- Underlying unobserved dimensions may in theory be seen as correlated
- Allowing correlated factors may provide even simpler structure than uncorrelated factors, thus easier to interpret
- All rotations fit data equally well
- Hence the one chosen depends on a series of choices done by the analyst
- Try different methods to see if results differ

Spring 2010

© Erling Berge 2010

516

SPSS output

- For rotated factor solutions with correlated factors SPSS provides two matrixes for interpretation
- The pattern matrix provides the direct regression of the variables on the factors. The coefficients tells about the direct contribution of a factor in explaining the variance of a variable. Due to the correlations of the factors there are also indirect contributions
- The structure matrix provides the correlations between the variables and the factors

Spring 2010

© Erling Berge 2010

517

Factor scores

- Both principal components and factor analysis may be used to compute composite scores called factor scores
- Recall that variables and factors are assumed to be related like
 - $Z_k = \ell_{k1}F_1 + \ell_{k2}F_2 + \dots + \ell_{kj}F_j + \dots + \ell_{kK}F_K$
- Then it is possible to find values c_{ij} making
 - $\hat{F}_j = c_{1j}Z_1 + c_{2j}Z_2 + \dots + c_{kj}Z_j + \dots + c_{Kj}Z_K$
- The coefficients c_{ij} are the factor score coefficients. They come from the regression of the factor F_j on the variables

Spring 2010

© Erling Berge 2010

518

Methods for extracting factors

- Principal factor analysis
 - The original correlation matrix \mathbf{R} is replaced by \mathbf{R}^* where the original 1-values of the diagonal has been replaced by estimates of the communality (the shared variance)
 - The factors extracted tries to explain the co-variance or correlations among the variables.
 - The unexplained variance is attributed to a unique factor (error term). The uniqueness may reflect measurement error or something that this variable measure that no other variable measure
 - The most common estimate of communality is R_k^2 the coefficient of determination from the regression of Z_k on all other variables

Spring 2010

© Erling Berge 2010

519

How many factor should we retain?

- In principal component analysis factors with eigenvalues above 1 is recommended
- In principal factor analysis factors with eigenvalues above 0 is recommended
- Procedure:
 - Extract initial factors or components
 - Rotate to simple structure
 - Decide on how many factors to retain
 - Obtain and use scores for the retained factors, ignoring discarded factors

Spring 2010

© Erling Berge 2010

520

Concluding (1)

- Principal components
 - transformation of the data, not model based.
Appropriate if goal is to compactly express most of the variance of k variables. Minor components (perhaps all except the first) may be discarded and viewed as a residual.
- Factor analysis
 - Estimates parameters of a measurement model with latent (unobserved) variables.

Spring 2010

© Erling Berge 2010

521

Concluding (2)

- Types of factor analysis
 - Principal factoring – principal components of a modified correlation matrix R^* in which communality estimates (R_k^2) replace “1” on the main diagonal
 - Principal factoring without iteration
 - Principal factoring with iteration
 - Maximum likelihood estimation – significance tests regarding number of factors and other hypotheses, **assuming multivariate normality**

Spring 2010

© Erling Berge 2010

522

Concluding (3)

- Rotation
 - If we retain more than one factor rotation simplifies structure and improves interpretability
 - Orthogonal rotation (varimax) maximum polarization given uncorrelated factors
 - Oblique rotation (oblimin, promax) further polarization by permitting interfactor correlations. The results may be more interpretable and more realistic than uncorrelated factors
- Scores
 - Factor scores can be calculated for use in graphs and further analysis, based on rotated or unrotated factors and principal components

Spring 2010

© Erling Berge 2010

523

Concluding (4)

- Factor analysis is based on correlations and hence as affected by non-linearities and influential cases as in regression
 - Use scatter plots to check for outliers and non-linearities
 - In maximum likelihood estimation this becomes even more important since it assumes multivariate normality making it even less robust than principal factors

Spring 2010

© Erling Berge 2010

524

Principal components of trust in Malawi

- Survey of 283 households in 18 villages in Malawi, 2007
- There are 8 related questions asked in one group
- Are there 1, 2 or more underlying dimensions shaping the attitudes expressed?
- Analysis of correlations (not co-variances)
- The questions:

Spring 2010

© Erling Berge 2010

525

M3 Would you say you trust all, most, some or just a few people in the following groups? (All=1 – None=5)

a	Your family members	<i>All</i>	<i>Most</i>	<i>Some</i>	<i>Only a few</i>	<i>None</i>	<i>Do not know</i>
b	Your relatives	<i>All</i>	<i>Most</i>	<i>Some</i>	<i>Only a few</i>	<i>None</i>	<i>Do not know</i>
c	Your village	<i>All</i>	<i>Most</i>	<i>Some</i>	<i>Only a few</i>	<i>None</i>	<i>Do not know</i>
d	People from outside the village	<i>All</i>	<i>Most</i>	<i>Some</i>	<i>Only a few</i>	<i>None</i>	<i>Do not know</i>
e	People of same ethnic group	<i>All</i>	<i>Most</i>	<i>Some</i>	<i>Only a few</i>	<i>None</i>	<i>Do not know</i>
f	People from outside ethnic group	<i>All</i>	<i>Most</i>	<i>Some</i>	<i>Only a few</i>	<i>None</i>	<i>Do not know</i>
g	People from same church/mosque	<i>All</i>	<i>Most</i>	<i>Some</i>	<i>Only a few</i>	<i>None</i>	<i>Do not know</i>
h	People <i>not</i> from same church/mosque	<i>All</i>	<i>Most</i>	<i>Some</i>	<i>Only a few</i>	<i>None</i>	<i>Do not know</i>

Spring 2010

© Erling Berge 2010

526

Trust in Malawi: descriptive

Descriptive Statistics

	Mean	Std. Deviation	Analysis N
M3.a. Trust in family members	1.60	.935	266
M3.b. Trust in relatives	2.12	1.136	266
M3.c. Trust in people in own village	2.69	1.090	266
M3.d. Trust in people outside the village	3.28	1.118	266
M3.e. Trust in people of same ethnic group	2.90	1.082	266
M3.f. Trust in people outside ethnic group	3.26	1.098	266
M3.g. Trust in people from same church/mosque	2.39	1.062	266
M3.h. Trust in people not from same church/mosque	3.02	1.197	266

Spring 2010

© Erling Berge 2010

527

Trust in Malawi: correlation of variables

Correlation Matrix

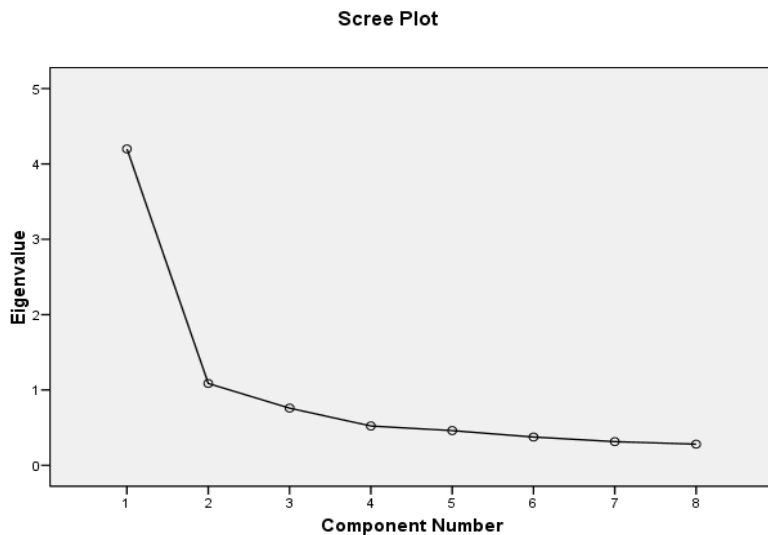
	M3.a. Trust in family members	M3.b. Trust in relatives	M3.c. Trust in people in own village	M3.d. Trust in people outside the village	M3.e. Trust in people of same ethnic group	M3.f. Trust in people outside ethnic group	M3.g. Trust in people from same church/mosque	M3.h. Trust in people not from same church/mosque
M3.a. Trust in family members	1.000	.500	.416	.236	.370	.316	.422	.305
M3.b. Trust in relatives	.500	1.000	.496	.315	.363	.353	.424	.292
M3.c. Trust in people in own village	.416	.496	1.000	.482	.588	.573	.465	.430
M3.d. Trust in people outside the village	.236	.315	.482	1.000	.526	.610	.233	.469
M3.e. Trust in people of same ethnic group	.370	.363	.588	.526	1.000	.702	.504	.643
M3.f. Trust in people outside ethnic group	.316	.353	.573	.610	.702	1.000	.430	.618
M3.g. Trust in people from same church/mosque	.422	.424	.465	.233	.504	.430	1.000	.536
M3.h. Trust in people not from same church/mosque	.305	.292	.430	.469	.643	.618	.536	1.000

Spring 2010

© Erling Berge 2010

528

Trust in Malawi: number of factors



Spring 2010

© Erling Berge 2010

529

Trust in Malawi: factor/ component matrix

Component Matrix ^a

	Component	
	1	2
M3.a. Trust in family members	.588	.586
M3.b. Trust in relatives	.624	.532
M3.c. Trust in people in own village	.776	.080
M3.d. Trust in people outside the village	.675	-.398
M3.e. Trust in people of same ethnic group	.832	-.221
M3.f. Trust in people outside ethnic group	.816	-.330
M3.g. Trust in people from same church/mosque	.690	.265
M3.h. Trust in people not from same church/mosque	.757	-.262

Extraction Method: Principal Component Analysis.

a. 2 components extracted.

Spring 2010

© Erling Berge 2010

530

Trust in Malawi: loadings on orthogonal factors

Rotated component matrix	Unrotated components		Orthogonal varimax	
	F1	F2	F1	F2
Variables				
M3.a. Trust in family members	.588	.586	.117	.821
M3.b. Trust in relatives	.624	.532	.178	.800
M3.c. Trust in people in own village	.776	.080	.572	.531
M3.d. Trust in people outside the village	.675	-.398	.779	.089
M3.e. Trust in people of same ethnic group	.832	-.221	.798	.324
M3.f. Trust in people outside ethnic group	.816	-.330	.850	.228
M3.g. Trust in people from same church/mosque	.690	.265	.391	.627
M3.h. Trust in people not from same church/mosque	.757	-.262	.762	.246

Spring 2010

© Erling Berge 2010

531

Trust in Malawi: communalities

Communalities

	Extraction
M3.a. Trust in family members	.689
M3.b. Trust in relatives	.671
M3.c. Trust in people in own village	.609
M3.d. Trust in people outside the village	.614
M3.e. Trust in people of same ethnic group	.741
M3.f. Trust in people outside ethnic group	.774
M3.g. Trust in people from same church/mosque	.546
M3.h. Trust in people not from same church/mosque	.641

Extraction Method: Principal Component Analysis.

Spring 2010

© Erling Berge 2010

532

Trust in Malawi: explained variance

Total Variance Explained

Component	Extraction Sums of Squared Loadings			Rotation Sums of Squared Loadings		
	Total	% of Variance	Cumulative %	Total	% of Variance	Cumulative %
1	4.199	52.487	52.487	3.071	38.387	38.387
2	1.087	13.582	66.069	2.215	27.681	66.069

Extraction Method: Principal Component Analysis.

Spring 2010

© Erling Berge 2010

533

Trust in Malawi: oblique factors pattern matrix

Rotated component matrix	varimax (orthogonal)		oblimin		promax	
	F1	F2	F1	F2	F1	F2
M3.a. Trust in family members	.117	.821	-.087	.868	-.145	.901
M3.b. Trust in relatives	.178	.800	-.014	.826	-.067	.855
M3.c. Trust in people in own village	.572	.531	.493	.414	.476	.409
M3.d. Trust in people outside the village	.779	.089	.838	-.133	.864	-.170
M3.e. Trust in people of same ethnic group	.798	.324	.797	.120	.806	.093
M3.f. Trust in people outside ethnic group	.850	.228	.881	-.001	.899	-.036
M3.g. Trust in people from same church/mosque	.391	.627	.268	.573	.237	.582
M3.h. Trust in people not from same church/mosque	.762	.246	.779	.045	.792	.016

Spring 2010

© Erling Berge 2010

534

Trust in Malawi: oblique factors structure matrix

Rotated component matrix	varimax		oblimin		promax	
	F1	F2	F1	F2	F1	F2
M3.a. Trust in family members	.117	.821	.327	.826	.351	.821
M3.b. Trust in relatives	.178	.800	.380	.819	.403	.817
M3.c. Trust in people in own village	.572	.531	.690	.649	.702	.671
M3.d. Trust in people outside the village	.779	.089	.775	.267	.771	.306
M3.e. Trust in people of same ethnic group	.798	.324	.854	.500	.857	.537
M3.f. Trust in people outside ethnic group	.850	.228	.880	.419	.880	.460
M3.g. Trust in people from same church/mosque	.391	.627	.541	.700	.557	.712
M3.h. Trust in people not from same church/mosque	.762	.246	.800	.416	.801	.452

Spring 2010

© Erling Berge 2010

535

Trust in Malawi: correlation of components

Component Correlation Matrix

Component	1	2
1	1.000	.477
2	.477	1.000

Extraction Method: Principal Component Analysis.

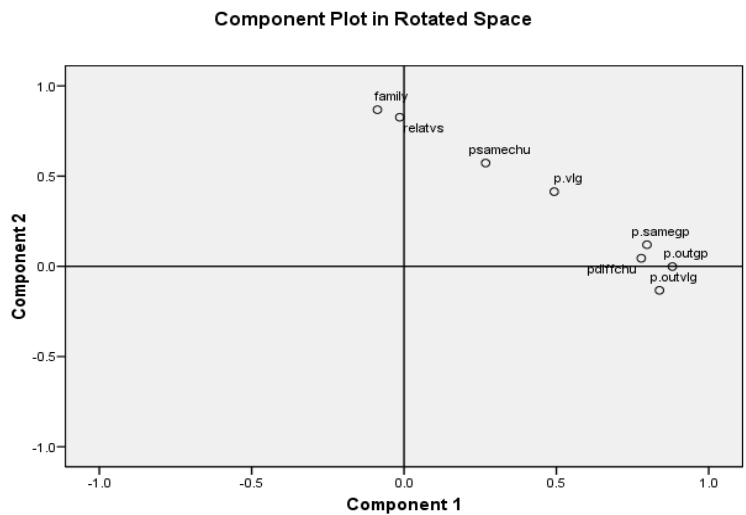
Rotation Method: Oblimin with Kaiser Normalization.

Spring 2010

© Erling Berge 2010

536

Trust in Malawi: variables in component plot

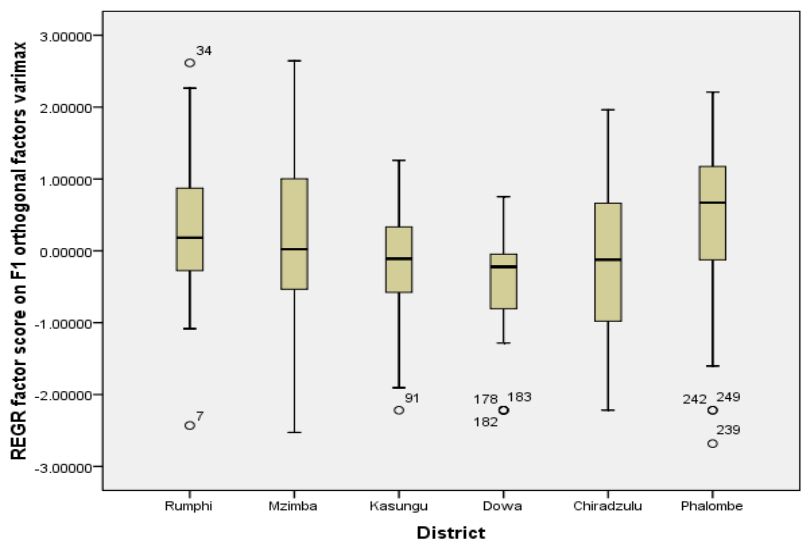


Spring 2010

© Erling Berge 2010

537

Trust in Malawi: Orthogonal Factor 1 by district

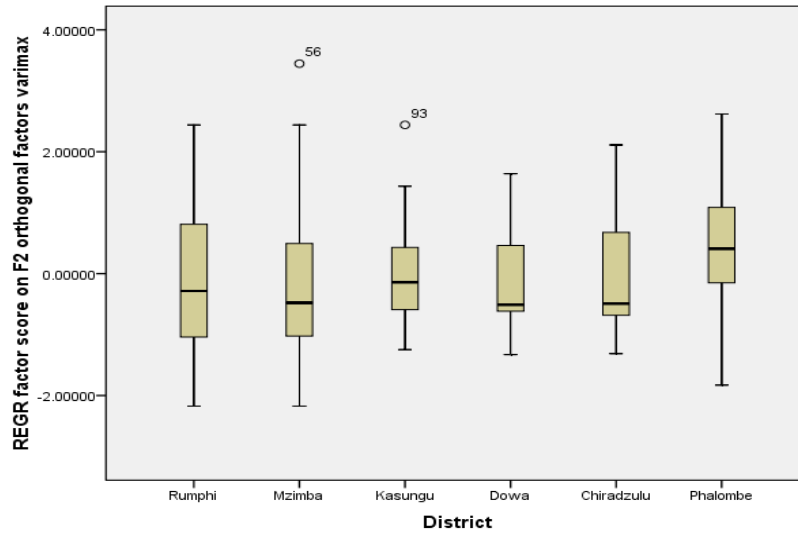


Spring 2010

© Erling Berge 2010

538

Trust in Malawi: Orthogonal Factor 2 by district



Spring 2010

© Erling Berge 2010

539

Trust in Malawi: Orthogonal factors by district

Case Processing Summary

	District	Cases					
		Valid		Missing		Total	
		N	Percent	N	Percent	N	Percent
REGR factor score on F1 orthogonal factors varimax	Rumphu	43	95.6%	2	4.4%	45	100.0%
	Mzimba	37	82.2%	8	17.8%	45	100.0%
	Kasungu	47	95.9%	2	4.1%	49	100.0%
	Dowa	49	98.0%	1	2.0%	50	100.0%
	Chiradzulu	46	93.9%	3	6.1%	49	100.0%
	Phalombe	44	97.8%	1	2.2%	45	100.0%

Case Processing Summary

	District	Cases					
		Valid		Missing		Total	
		N	Percent	N	Percent	N	Percent
REGR factor score on F2 orthogonal factors varimax	Rumphu	43	95.6%	2	4.4%	45	100.0%
	Mzimba	37	82.2%	8	17.8%	45	100.0%
	Kasungu	47	95.9%	2	4.1%	49	100.0%
	Dowa	49	98.0%	1	2.0%	50	100.0%
	Chiradzulu	46	93.9%	3	6.1%	49	100.0%
	Phalombe	44	97.8%	1	2.2%	45	100.0%

Spring 2010

© Erling Berge 2010

540

Missing data Biased samples

- Allison, Paul D 2002 "Missing Data", Sage University Paper: QASS 136, London, Sage,

Spring 2010

© Erling Berge 2010

541

There is a missing case in the sample

- If one person
 - Refuses to answer
 - Are not at home
 - Has moved away
 - Etc.
- The problem of missing data belong to the study of biased samples. In general biased samples is a more severe problem than the fact that we are missing answers for a few variables on some cases (see Breen 1996 "Regression Models: Censored, Sample Selected, or Truncated Data", QASS Paper 111, London, Sage)
- But the problems are related

Spring 2010

© Erling Berge 2010

542

There are missing answers for a few variables if

- Persons refuse to answer certain questions
- Persons forget, or do not notice some question, or the interviewer does it
- Persons do not know any answer to the question: "Do not know" are often a valid answer category. But the result is a missing answer
- The question is irrelevant (for the person)
- In administrative registers some documents may have been lost
- In research designs where variables with measurement problems may have been measured only for a minority of the sample

Spring 2010

© Erling Berge 2010

543

Missing data entail problems

- There are practical problems due to the fact that all statistical procedures assume complete data matrices
- It is an analytical problem since missing data as a rule produce biased parameter estimates
- It is important to distinguish between data missing for random causes and those missing from systematic causes

Spring 2010

© Erling Berge 2010

544

The simple solution: remove all cases with missing data

- Listwise/ casewise removal of missing data means to remove all cases missing data on one or more variables included in the model
- The method has good properties, but may in some cases remove most of the cases in the sample
- Alternatives like pairwise removal or replacement with average variable value has proved not to have good properties
- More recently developed methods like "maximum likelihood" and "multiple imputation" have better properties but are more demanding
- In general it pays to do good work in the data collection stage

Spring 2010

© Erling Berge 2010

545

Types of randomly missing

- **MCAR: missing completely at random**
 - Means that missing data for one person on the variable y is uncorrelated with the value on y and with the value on any other variable in the data set (however, internal case by case the value of missing may of course correlate with the value missing on other variables)
- **MAR: missing at random**
 - Means that missing data for person i on the variable y do not correlate with the value on y if one control for the variation of other variables in the model
 - More formally:
$$\Pr(Y_i = \text{missing} \mid Y_i, X_i) = \Pr(Y_i = \text{missing} \mid X_i)$$

Spring 2010

© Erling Berge 2010

546

Process resulting in missing

- Is ignorable if
 - The result is MAR and the parameters governing the process are unrelated to the parameters that are to be estimated
- Is non-ignorable if
 - The result is not MAR. Estimation of the model will then require a separate model of the missing process
 - See Breen 1996 "Regression Models: Censored, Sample Selected, or Truncated Data", QASS Paper 111, London, Sage
- Here the situation with MAR will be discussed

Spring 2010

© Erling Berge 2010

547

Conventional methods

Common methods in cases with MAR data:

- Listwise deletion
 - Pairwise deletion
 - Dummy variable correction
 - Imputation (guessing a value for the missing)
- Of the conventional methods listwise deletion is the best

Spring 2010

© Erling Berge 2010

548

Listwise deletion (1)

- Can always be used
- If data are MCAR we have a simple random subsample of the original sample
- Smaller n entails larger variance estimates
- In the case of MAR data and the missing values on an x-variable are independent of the value on y, listwise deletion will produce unbiased estimates

Spring 2010

© Erling Berge 2010

549

Listwise deletion (2)

- In logistic regression listwise deletion may cause problems only if missing is related both to dependent and independent variables
- If missing depends only on the values of the independent variable listwise deletion is better than replacement of missing values by maximum likelihood and multiple imputation

Spring 2010

© Erling Berge 2010

550

Pairwise deletion

- Means that all computations are based on all available information seen pairwise for all pairs of variables included in the analysis
- One consequence is that different parameters will be estimated on different samples (we see variation in n from statistic to statistic)
- Then all variance estimates are biased
- Common test statistics provides biased estimates (e.g. t -values and F -values)
- **DO NOT USE PAIRWISE DELETION !!**

Spring 2010

© Erling Berge 2010

551

Dummy variable correction

If data is missing for the independent variable x

- Let $x^*_i = x_i$ if x_i is not missing and
 $x^*_i = c$ (an arbitrary constant) if x_i is missing
- Define $D_i=1$ if x_i is missing, 0 otherwise
- Use x^*_i and D_i in the regression instead of x_i
- In nominal scale variables missing can get its own dummy

Investigations reveal that even if we have MCAR data parameter estimates will be biased

Do not use dummy variable correction!

Spring 2010

© Erling Berge 2010

552

Imputation

- The goal is to replace missing values with reasonable guesses about what the value might have been before one does an analysis as if this were real values; e.g.
 - Average of valid values
 - Regression estimates based on many variables and cases with valid observations
- Parameter estimates are consistent, but estimates of variances are biased (consistently too small), and the test statistics are too big
- Avoid if possible the simple forms of imputation

Spring 2010

© Erling Berge 2010

553

Concluding on conventional methods for missing data

- Conventional methods of correcting for missing data make problems of inference worse
- Be careful in the data collection so that the missing data are as few as possible
- Make an effort to collect data that may help in modelling the process resulting in missing
- If data are missing use listwise deletion if not maximum likelihood or multiple imputation is available

Spring 2010

© Erling Berge 2010

554

New methods for ignorable missing data (MAR data): Maximum Likelihood (ML)

- Conclusions
 - Based on the probability for observing just those values found in the sample
 - ML provides optimal parameter estimates in large samples in the case of MAR data
 - But ML require a model for the joint distribution of all variables in the sample that are missing data, and it is difficult to use for many types of models

Spring 2010

© Erling Berge 2010

555

ML-method: example (1)

- Observing y and x for 200 cases
- 150 distributed as shown
- For 19 cases with Y=1 x is missing and for 31 cases with Y=2 x is missing
- We want to find the probabilities p_{ij} in the population

	Y=1	Y=2
X=1	52	21
X=2	34	43

	Y=1	Y=2
X=1	p_{11}	p_{12}
X=2	p_{21}	p_{22}

Spring 2010

© Erling Berge 2010

556

ML-method: example (2)

- In a table with I rows and J columns, complete information on all cases and with n_{ij} cases in cell ij the Likelihood is

$$\mathcal{L} = \prod_{i,j} \left(p_{ij} \right)^{n_{ij}}$$

That is the product of all probabilities for every table cell taken to the power of the cell frequency

Spring 2010

© Erling Berge 2010

557

ML-method: example (3)

For a fourfold table the Likelihood will be

$$\mathcal{L} = \left(p_{11} \right)^{n_{11}} \left(p_{12} \right)^{n_{12}} \left(p_{21} \right)^{n_{21}} \left(p_{22} \right)^{n_{22}}$$

For the 150 cases in the table above where we have all observations the Likelihood will be

$$\mathcal{L} = \left(p_{11} \right)^{52} \left(p_{12} \right)^{21} \left(p_{21} \right)^{34} \left(p_{22} \right)^{43}$$

Spring 2010

© Erling Berge 2010

558

ML-method: example (4)

- For tables the ML estimator is $p_{ij} = n_{ij}/n$
- This provides good estimates for the table where we do not have missing data (listwise deletion)
- How can one use the information about y for the 50 cases where x is missing?
- Since MAR is assumed to be the case, the 50 extra cases with known y should follow the marginal distribution of y
- $\Pr(Y=1) = (p_{11} + p_{21})$ and $\Pr(Y=2) = (p_{12} + p_{22})$

Spring 2010

© Erling Berge 2010

559

ML-method: example (5)

- Taking into account all that is known about the 200 cases the Likelihood becomes

$$\mathcal{L} = (p_{11})^{52} (p_{12})^{21} (p_{21})^{34} (p_{22})^{43} (p_{11} + p_{21})^{19} (p_{11} + p_{21})^{31}$$

- The ML-estimators will now be

$$\hat{p}_{ij} = \hat{p}(x = i | y = j) \hat{p}(y = j)$$

Spring 2010

© Erling Berge 2010

560

ML-method: example (6)

- Taking into account the information we have about cases with missing data, parameter estimates change

Estimate of	Missing deleted	Missing included
p_{11}	0.346	0.317
p_{21}	0.227	0.208
P_{12}	0.140	0.156
p_{22}	0.287	0.319

Spring 2010

© Erling Berge 2010

561

The ML-method in practice

- For the general case with missing data there are two approaches
 - The expectation-maximization (EM) method, a two stage method where one starts out with the expected value of the missing data and use these to obtain parameter estimates that again will be used to provide better estimates of the missing values and so on ...
(this method provides biased estimates of standard errors)
 - Direct ML estimates are better but can be provided only for linear and log-linear models

Spring 2010

© Erling Berge 2010

562

New methods for ignorable missing data (MAR data): Multiple Imputation (MI)

- Conclusions
 - MI is based on a random component added to estimates of the missing data values
 - Has as good properties as the ML method and is easier to implement for all kinds of models
 - But it gives different results every time it is used

Spring 2010

© Erling Berge 2010

563

Multiple Imputation (1)

- MI have the same optimal properties as the ML method. It can be used on all kinds of data and with all kind of models. In principle it can be done with the ordinary analytical tools
- The use of MI can be rather convoluted. This makes it rather easy to commit errors. And even if it is done correctly one will never have the same result twice due to the random component in the imputed variable value

Spring 2010

© Erling Berge 2010

564

Multiple Imputation (2)

- Use of data from a simple imputation (with or without a random component) will underestimate the variance of parameters. Conventional techniques are unable to adjust for the fact that data have been generated by imputation
- The best way of doing imputation with a random component is to repeat the process many times and use the observed variation of parameter estimates to adjust the estimates of the parameter variances
- Allison, p.30-32, explains how this can be done

Spring 2010

© Erling Berge 2010

565

Multiple Imputation (3)

- MI requires a model that can be used to predict values of missing data. Usually there is an assumption of normally distributed variables and linear relationships. But models can be tailored to each problem
- MI can not handle interactions
- MI model should contain all variables of the analysis model
- (including the dependent variable)
- MI works only for interval scale variables. If nominal scale variables are used special programs are needed
- Testing of several coefficients in one test is complicated

Spring 2010

© Erling Berge 2010

566

When data are missing systematically

- Will usually require a model of how the missing cases came about
- ML and MI approaches can still be used, but with much stronger restrictions and the results are very sensitive for deviations from the assumptions

Spring 2010

© Erling Berge 2010

567

Summary

- If listwise deletion leaves enough data this is the simplest solution
- If listwise deletion do not work one should test out multiple imputation
- If there is a suspicion that data are not MAR one needs to create a model of the process creating missing. This can then be used together with ML or MI. Good results require that the model for missing is correct

Spring 2010

© Erling Berge 2010

568

Types of biased samples

- Censored
- Truncated
- Selected
- Such samples arise because society works “selectively”, and because we do not get complete answers to questions asked
- Which variables and how they are truncated determine the type of bias

Spring 2010

© Erling Berge 2010

569

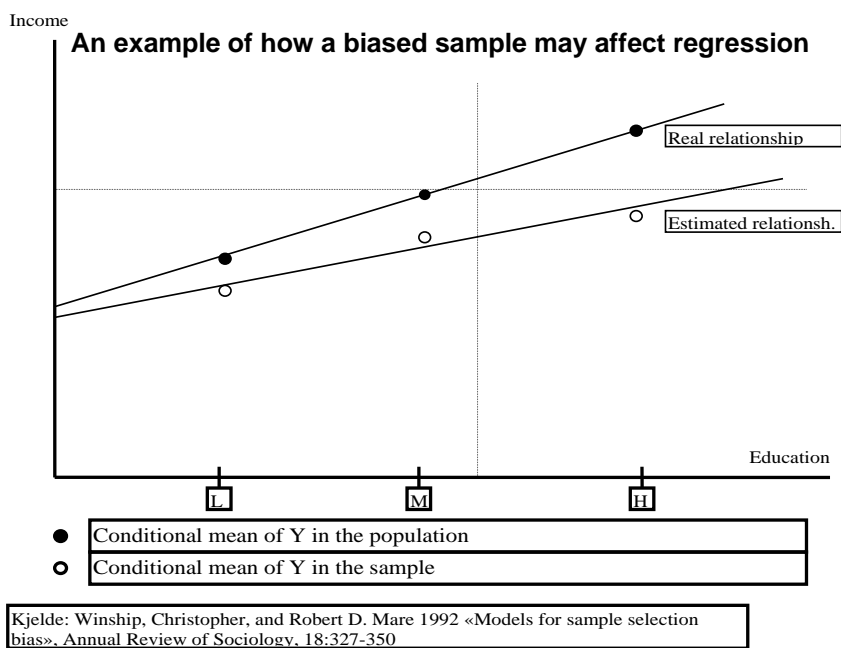
Causal analysis in biased samples

- Regression analysis
 - Will (as a rule) have severe problems if the sample is biased
- Hence
 - The process of selection needs to be included in the model or analysis

Spring 2010

© Erling Berge 2010

570



Spring 2010

© Erling Berge 2010

571

Comments to the figure

- Only persons with incomes below 15000 USD are included in the sample
- Result is erroneous estimate of the (real) impact of education
- Errors in reporting income creates a selected sample
- Large errors in the original sample leads to exclusion
- Large values on the independent variable leads to large (negative) errors
- The errors in the sample will be correlated with x

Spring 2010

© Erling Berge 2010

572

Truncation of variables

- A variable, X , is called truncated if we for $X < c$ or for $X > c$ do not know more than that $X < c$ or $X > c$
- This is known as left or right side truncation respectively
- We may have multiple truncation such as simultaneous left side and right side truncation

Spring 2010

© Erling Berge 2010

573

Biased samples and missing data I

- Censored samples (explicit selection on Y)
 - Y is unknown for cases where Y has value above or below c
 - X is known for all cases in the sample
- Selected samples (unsystematic selection)
 - Y is unknown for cases where f. e. $z=1$ and known if $z=0$
 - X is known for all cases in the sample

Spring 2010

© Erling Berge 2010

574

Selected or censored sample?

- The terminology is not very clear
- In general the distinction is a question of interpretation and theoretical meaning
 - If the missing observations on Y are caused by the measurement method or data collection method the sample is called censored
 - If the missing observations of Y are caused by the behaviour of the individuals the sample is called selected

Spring 2010

© Erling Berge 2010

575

Biased samples and missing data II

- Truncated samples (explicit selection on Y)
 - Y is unknown for cases where Y has value above or below c
 - X is known when Y is known
- Selection on the independent variable
 - Y is known for cases where X has a value above or below c
 - X is known when Y is known

Spring 2010

© Erling Berge 2010

576

Consequences of biased samples

- **Selection on the independent variable do not cause problems**
- Truncated, selected, and censored samples cause the residual to be correlated with the independent variables. Both external and internal validity is compromised

Spring 2010

© Erling Berge 2010

577

Causes of biased samples

- Data collection procedures and missing answers may lead to truncated, selected or censored samples
 - For example: "missing" on a dependent variable give a selected sample based on the variable Z, answer or no answer
- In every non-random sample there is a potential for erroneous conclusions due to biased sample

Spring 2010

© Erling Berge 2010

578

How to handle biased samples

- The analysis should at the outset acknowledge the problem and use models that are able to correct for bias in the sample unless there are good reasons to believe the problem is small
- The solution then is to
 - 1) construct a model that predicts selection
 - 2) use this model to construct a model that predicts y conditional on the person having been selected

Spring 2010

© Erling Berge 2010

579

A basic model for censored samples

$$E[Y | X] = \Pr[Y > c | X] * E[Y | Y > c \& X] + \Pr[Y \leq c | X] * E[Y | Y \leq c \& X]$$

Left side truncation of Y at c gives

$$E[Y | Y \leq c \& X] = c$$

It is always possible to transform Y so that $c=0$, hence the real regression, $E[Y | X]$, can be written

$$\bullet \quad E[Y | X] = \Pr[Y > 0 | X] * E[Y | Y > 0 \& X]$$

Spring 2010

© Erling Berge 2010

580

The model in a truncated sample

- $Y_i = E[Y_i | Y_i < a \ \& \ X_i] + e_i$

It can be shown that this is equivalent to

- $Y_i = E[Y_i | X_i] - \sigma \lambda'_i(m) + e_i$

where $\lambda'_i(m)$ is an estimate of the Hazard rate at the point

- $m = (a - E[Y_i | X_i]) / \sigma$

The parameters of $E[Y_i | X_i]$ are overestimated

The model can be estimated by the ML method

Spring 2010

© Erling Berge 2010

581

Two step estimation in censored samples

- The selection model, $\Pr[Y > c | X]$, can be modelled by probit regression on the censored sample
- The model of the outcome, $E[Y | Y > c \ \& \ X]$, can then be estimated on the censored sample
- The results are trustworthy only in large samples

Spring 2010

© Erling Berge 2010

582

Problems in the two step model

- Results are sensitive for assumptions about the distribution of the residual
 - Homoscedasticity: deviation for this assumption is more serious than in OLS since estimates in a censored model are neither consistent, nor efficient
 - Normal distributionBoth assumptions have to be properly tested
- There are also problems of identification of parameters due to multicollinearity between the hazard rate and the explanatory variables (see Breen 1996:16 (equation 2.7))

Spring 2010

© Erling Berge 2010

583

Two step estimation in OLS

is sensitive for

- Correlations between errors in the selection equation (u) and errors in the outcome equation (e)
- Correlations between variables in the selection and outcome equations
- Degree of censoring in the sample (how large a fraction of the cases have missing y values?)

Conclusion: use ML-estimation

Spring 2010

© Erling Berge 2010

584

Literature cited I

- Allison, Paul D. 2002. *Missing Data*. Sage University Paper: QASS 136. London: Sage.
- Breen 1996 "Regression Models: Censored, Sample Selected, or Truncated Data", QASS Paper 111, London, Sage
- Hamilton, Lawrence C. 1992. *Regression with Graphics: A Second Course in Applied Statistics*. Belmont: Duxbury Press.
- Hamilton, Lawrence C. 2008. *A Low-Tech Guide to Causal Modelling*. <http://pubpages.unh.edu/~lch/causal2.pdf>
- Hardy, Melissa A. 1993. *Regression with Dummy Variables*, Sage University Paper series on Quantitative Applications in the Social Sciences 07-093. Newbury Park, CA: Sage.
- Hosmer, David W., and Stanley Lemeshow. 1989. *Applied Logistic Regression*. New York: John Wiley & Sons.
- Long, J. Scott. 1997. *Regression Models for Categorical and Limited Dependent Variables*. London: Sage.

Spring 2010

© Erling Berge 2010

585

Literature cited II

- Menard, Scott. 1995. *Applied Logistic Regression Analysis*, Sage University Paper series on Quantitative Applications in the Social Sciences 07-106. Thousand Oaks, CA: Sage.
- Peduzzi, Peter, John Concato, Elizabeth Kemper, Theodore R. Holford, and Alvan R. Feinstein. 1996. A simulation study of the number of events per variable in logistic regression analysis. *Journal of Clinical Epidemiology* 49 (12):1373-1379.
- Weisberg, Sanford. 1985. *Applied Linear Regression*. Second edition. New York: John Wiley & Sons.
- Winship, Christopher, and Robert D. Mare 1992 «Models for sample selection bias», *Annual Review of Sociology*, 18:327-350
- Winship, Christopher, and Stephen L. Morgan 1999 "The Estimation of Causal Effects from Observational Data", *Annual Review of Sociology* Vol 25: 659-707

Spring 2010

© Erling Berge 2010

586