

Linear regression, collinearity, splines and extensions

Morten Frydenberg ©
Department of Biostatistics, Aarhus Univ, Denmark

General things for regression models:

Collinearity - correlated explanatory variables

Flexible modelling af response curves - Cubic splines

Normal regression models - extensions

Random coefficient model

Clustered data / data with several random components

Morten Frydenberg

Linear and Logistic regression - Note 3

1

Collinearity

Consider a subsample of the serum cholesterol data set and the **three** models:

model 0: regress logsc1 sex sbp dbp
model 1: regress logsc1 sex dbp
model 2: regress logsc1 sex sbp

variable	model0	model1	model2
sbp	.00126448 .00087992 0.1524		.0014988 .0005548 0.0075
dbp	.00056517 .00164485 0.7315	.00239702 .0010424 0.0226	
sex	.02080574 .02636149 0.4310	.02446746 .02631111 0.3536	.0197773 .02613048 0.4501
_cons	5.1444085 .09912234 0.0000	5.1555212 .09909537 0.0000	5.1615877 .08539118 0.0000
N	194	194	194

Estimate

Se

p

Each BP-measure is statistical significant, when the other is removed!

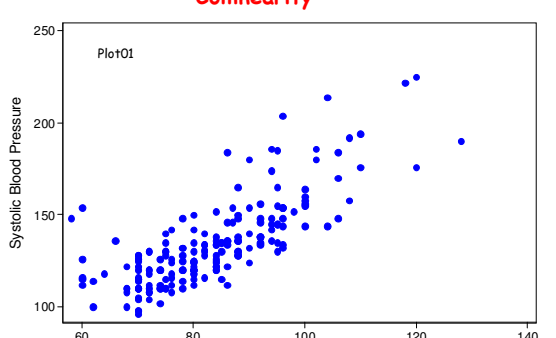
Legend: b/se/p

Morten Frydenberg

Linear and Logistic regression - Note 3

2

Collinearity



Plot01

SBP and DBP are **highly positively correlated** that will lead to **highly negatively correlated estimates!!!**

Morten Frydenberg

Linear and Logistic regression - Note 3

3

Collinearity

This can be seen by listing the **correlation between the estimates**.
In Stata by the command: `vce, cor`

```
regress logsc1 sbp dbp sex
vce,cor
```

	sbp	dbp	sex	_cons
sbp	1.0000			
dbp	-0.7750	1.0000		
sex	-0.0967	0.1135	1.0000	
_cons	-0.0780	-0.5044	-0.4665	1.0000

If two estimates are highly correlated, it indicates that it is very difficult to estimate the **"independent effect"** of each of the two variables.

Often it is even **nonsense** to try to do it!

Often it see better to try to **reformulate the problem**.

Morten Frydenberg

Linear and Logistic regression - Note 3

4

Collinearity

One way to work around the problem of colinearity is to **'ortogonalize'** it:

Create two new variable:
one measures the **blood pressure**
and another that measure the **difference** in systolic and diastolic blood pressure.

Some **candidates**:
(sbp+dbp)/2 and (sbp-dbp)
(sbp+dbp)/2 and (sbp/dbp)
 $\ln(\text{sbp} \cdot \text{dbp})/2$ and $\ln(\text{sbp}/\text{dbp})$

We will here consider the second pair.

Morten Frydenberg

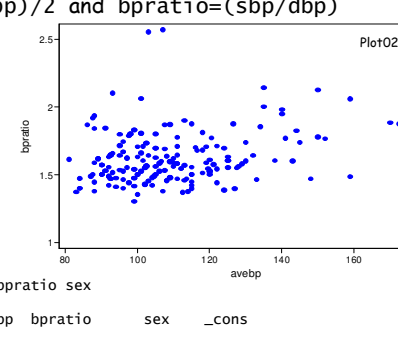
Linear and Logistic regression - Note 3

5

Collinearity

$\text{avebp} = (\text{sbp} + \text{dbp})/2$ and $\text{bpratio} = (\text{sbp}/\text{dbp})$

Only weakly associated



Plot02

```
regress logsc1 avebp bpratio sex
vce,cor
```

	avebp	bpratio	sex	_cons
avebp	1.0000			
bpratio	-0.2456	1.0000		
sex	0.0382	-0.1041	1.0000	
_cons	-0.4542	-0.6874	-0.2585	1.0000

Morten Frydenberg

Linear and Logistic regression - Note 3

6

Linear and Logistic Regression: Note 3

1

Collinearity

The serum cholesterol data set and the **three** models:
model 0: regress logsc1 sex avebp bpratio
model 1: regress logsc1 sex avebp
model 2: regress logsc1 sex bpratio

variable	model0	model1	model2
avebp	.00198973 .0007887 0.0125	.00206564 .00076285 0.0024	
bpratio	.02769662 .07067134 0.6956		.07148118 .06946246 0.3048
sex	.02060675 .02632924 0.4348	.02168128 .026128 0.4077	.01806662 .02667689 0.4991
_cons	5.1003417 .12936418 0.0000	5.1351912 .09374803 0.0000	5.2485724 .11685799 0.0000
N	194	194	194

Blood pressure seems to play a role,

The ratio between SBP and DBP might not.

Morten FrydenbergLinear and Logistic regression - Note 37

Collinearity

Look out for it:

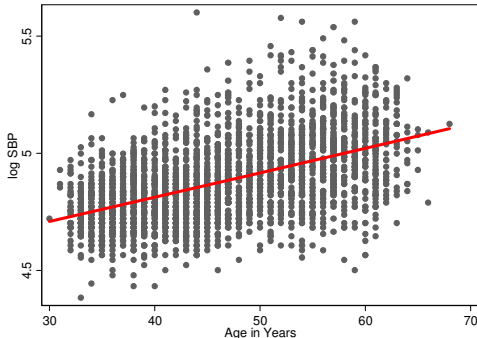
- systolic and diastolic blood pressure
- 24 hour blood pressure and 'clinical' blood pressure
- weight and height
- age and parity
- age and time since menopause
- BMI and skinfold measure
- age , birth cohort and calendar time
- volume and concentration
-

Remember you will need a **huge amount** of data to disentangle the effects of correlated explanatory variables

Morten FrydenbergLinear and Logistic regression - Note 38

Flexible modelling of response curves - cubic splines

Log SBP against age for 2650 women with fitted straight line.



Morten FrydenbergLinear and Logistic regression - Note 39

Flexible modelling of response curves - cubic splines

We want to model the relationship between SBP and age more flexible.

There several ways to do this, including fractional polynomial, splines and cubic splines.

We will here look at restricted cubic splines as they are implemented in Stata.

If one want used the restricted cubic splines you start by generating of set of new independent variables:

```
mkspline sage=age, cubic nk(6) disp
```

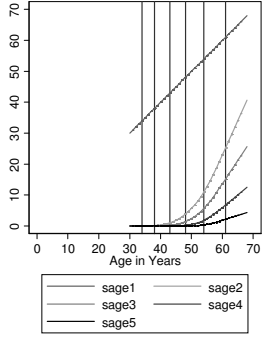
	knot1	knot2	knot3	knot4	knot5	knot6
age	34	38	43	48	54	61

Morten FrydenbergLinear and Logistic regression - Note 310

Flexible modelling of response curves - cubic splines

The mkspline command will generate 5 new variables named sage1 to sage5 which are function of age.

Where sage1=age.
sage2=0 if age<34
sage3=0 if age<38
sage4=0 if age<43
sage5=0 if age<48



Morten FrydenbergLinear and Logistic regression - Note 311

Flexible modelling of response curves - cubic splines

$knots: a_1, a_2, \dots, a_k$

$sage_i = age$

$$sage_{j+1} = (age - a_j)_+^3 - (age - a_{k-1})_+^3 \frac{a_k - a_j}{a_k - a_{k-1}} + (age - a_k)_+^3 \frac{a_{k-1} - a_j}{a_k - a_{k-1}}$$

Morten FrydenbergLinear and Logistic regression - Note 312

Flexible modelling of response curves - cubic splines

```
drop sage1
regress lsbp age sage?
```

	lsbp	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
age		.0067837	.0035322	1.92	0.055	-.0001425 .0137099
sage2		-.0005598	.0525269	-0.01	0.991	-.1035577 .1024381
sage3		-.0553357	.1336906	0.41	0.679	-.2068131 .3174845
sage4		-.1398205	.1547781	-0.90	0.366	-.4433189 .1636778
sage5		.0932052	.1207685	0.77	0.440	-.1436051 .3300155
_cons		4.527844	.1253021	36.14	0.000	4.282144 4.773544

```
testparm sage?
( 1) sage2 = 0
( 2) sage3 = 0
( 3) sage4 = 0
( 4) sage5 = 0
F( 4, 2644) = 3.81
Prob > F = 0.0043
```

Test of linearity
The hypothesis is rejected

The relationship is not linear, but how does look ?

Morten FrydenbergLinear and Logistic regression - Note 313

Flexible modelling of response curves - cubic splines

```
predict fit if e(sample)          /// fit values
predict fitsd if e(sample),stdp   /// standard error
generate low=fit-1.96*fitsd       /// lower ci-limit
generate hig=fit+1.96*fitsd       /// upper ci-limit
line fit low hig age             /// plot
```

Morten FrydenbergLinear and Logistic regression - Note 314

Flexible modelling of response curves - cubic splines

Compare with the straight line model:

Although, there is 'statistical significant' non-linearity, it has no practical implications- the straight line model is a valid approximation.

Morten FrydenbergLinear and Logistic regression - Note 315

Random coefficient models

Question
Is cerebral blood flow declining with age?

Data
Cross sectional data on age, sex and cerebral blood flow in grey matter from 7 studies:

study	sex		Total
	male	female	
1	7	0	7
2	4	6	10
3	6	6	12
4	8	7	15
5	5	4	9
6	17	0	17
7	6	0	6
8	1	1	2
Total	54	24	78

Morten FrydenbergLinear and Logistic regression - Note 316

All the data

Morten FrydenbergLinear and Logistic regression - Note 317

All the data - separate line for each study aender combination

Morten FrydenbergLinear and Logistic regression - Note 318

Seven simple linear regressions

We will here only consider the men.
Fitting a line for each of the seven studies:

$$CBF_{si} = \alpha_s + \beta_s \cdot (age_{si} - 50) + E_{si} \quad s = 1, \dots, 7, i = 1, \dots, n_s$$
$$E_{si} \sim N(0, \sigma_s^2)$$

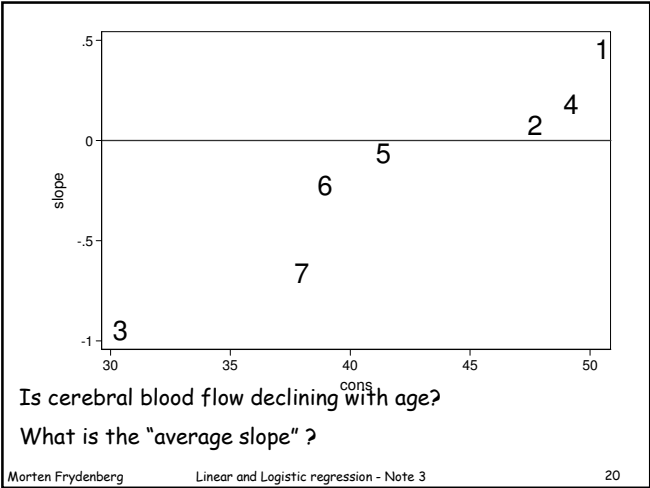
regress greymatter age50 if sex==0 & study==s

study	N	_cons	se(_cons)	age50	se(age50)	sd
1	7	50.51	13.63	0.465	0.564	4.070
2	4	47.71	6.49	0.082	0.682	11.428
3	6	30.42	18.11	-0.941	0.831	7.223
4	8	49.21	5.71	0.189	0.483	7.754
5	5	41.38	3.60	-0.055	0.433	6.701
6	17	38.94	1.96	-0.218	0.089	8.062
7	6	37.99	17.11	-0.654	1.095	14.420

Morten Frydenberg

Linear and Logistic regression - Note 3

19



Seven random slopes and intercepts

$$CBF_{si} = A_s + B_s \cdot (age_{si} - 50) + E_{si} \quad s = 1, \dots, 7, i = 1, \dots, n_s$$
$$A_s \sim N(\alpha, \sigma_A^2) \quad B_s \sim N(\beta, \sigma_B^2) \quad E_{si} \sim N(0, \sigma_E^2)$$

What is β ?

```
xtmixed greymatter age50 ||study: age50 if sex==0
```

greymatter	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
age50	-.089039	.11662	-0.76	0.445	-.31761 .1395321
_cons	44.44259	2.135614	20.81	0.000	40.25687 48.62832

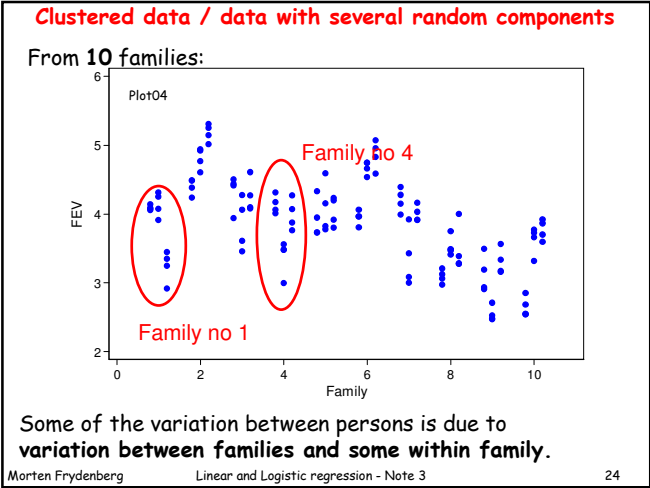
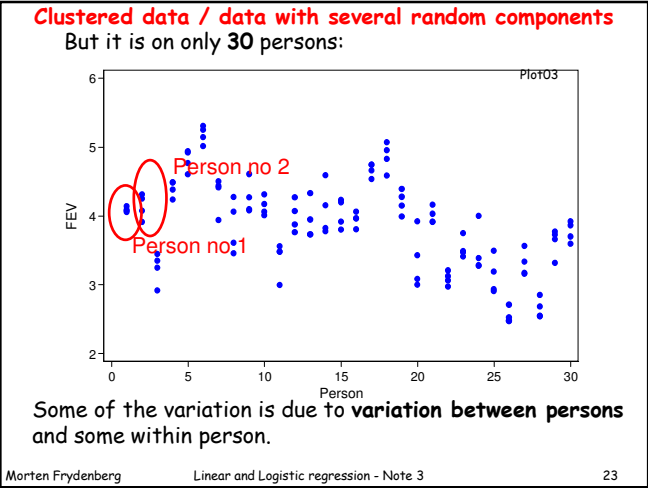
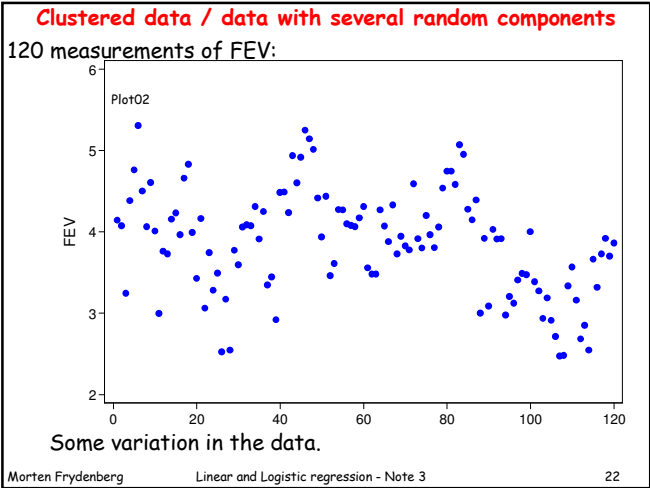
Random-effects Parameters	Estimate	Std. Err.	[95% Conf. Interval]
study: Independent			
sd(age50)	.1637849	.1691682	.0216319 1.240089
sd(_cons)	4.25174	2.182807	1.554415 11.62964
sd(Residual)	8.07755	.8410269	6.586479 9.906174

β : -0.089(-0.318;0.140) $H: \beta = 0$ $p = 45\%$

Morten Frydenberg

Linear and Logistic regression - Note 3

21



Clustered data / data with several random components

Structure of the data:

```
graph BT; Family --> FEV_fpd; Person --> FEV_fpd; Day --> FEV_fpd
```

Three sources of random variation:

Variation between **families**

Variation between **persons** (variation within family)

Variation between **days** (variation within person)

Morten Frydenberg

Linear and Logistic regression - Note 3

25

Clustered data / data with several random components

Factors of interest:

household **I**ncome Constant within **family**

Urbanization Constant within **family**

Age Constant within **person**; varies within family

Sex Constant within **person**; varies within family

Grass pollen Constant within **day**; varies within person

A model:

$$FEV = \beta_0 + \beta_I \cdot I + \beta_U \cdot U + \beta_A \cdot A + \beta_S \cdot S + \beta_G \cdot G$$

+random variation

Morten Frydenberg

Linear and Logistic regression - Note 3

26

Clustered data / data with several random components

$$FEV = \beta_0 + \beta_I \cdot I + \beta_U \cdot U + \beta_A \cdot A + \beta_S \cdot S + \beta_G \cdot G$$

+random variation

If the **three** levels/sources of **random** variation are **not** taken into account :

- The **precision** of the β_I and β_U are **highly overestimated**
- The **precision** of the β_A and β_S are **overestimated**
- The **estimates** of the β_I and β_U will be **biased** if the not all families are represented by the **same number of persons** and each person is measured the **same number of times**.
- The **estimates** of the β_A and β_S will be **biased** if the not all persons are measured the **same number of times**.

Morten Frydenberg

Linear and Logistic regression - Note 3

27

Clustered data / data with several random components

$$FEV = \beta_0 + \beta_I \cdot I + \beta_U \cdot U + \beta_A \cdot A + \beta_S \cdot S + \beta_G \cdot G$$

$+ F_f + P_{fp} + E_{fpd}$

F_f

P_{fp}

E_{fpd}

: Random family contribution

: Random person contribution

: Random day contribution

σ_F^2

σ_P^2

σ_E^2

variance

$$\text{var}(FEV_{fpd}) = \sigma_F^2 + \sigma_P^2 + \sigma_E^2$$

Variance components

Assumed to be normal distributed

Morten Frydenberg

Linear and Logistic regression - Note 3

28

Clustered data / data with several random components

Systematic part

$$FEV = \beta_0 + \beta_I \cdot I + \beta_U \cdot U + \beta_A \cdot A + \beta_S \cdot S + \beta_G \cdot G$$

$$+ F_f + P_{fp} + E_{fpd}$$

Random part

$\beta_0, \beta_I, \beta_U, \beta_A, \beta_S$ and β_G Quantify the **systematic** variation

σ_F^2, σ_P^2 and σ_E^2 Quantify the **random** variation

This is a:

- **Variance component** model
- **Mixed** model (both systematic and random variation)
- **Multilevel** model

The theory behind and the understanding of such models is well **established!!!**

Morten Frydenberg

Linear and Logistic regression - Note 3

29