

Working with logistic regression models

Morten Frydenberg ©
Section of Biostatistics, Aarhus Univ, Denmark

Further remarks on logistic regression

Test of fit: The Hosmer-Lemeshow test

Conditional logistic regression

Missing data

A small example - non completely random sample

Complete data analysis - bias

Missing at random vs missing **completely** at random

Introduction to techniques

Sampling weights

Imputation

Full modelling

Sensitivity analyses

Data with **several random components**: Binary outcome

Clustered binary data with **one random component**

1

Logistic regression models: Test of fit

A common, and to some extend informative, test of fit is the Hosmer-Lemeshow test.

Consider the model for obesity from Day 4

$$\text{logit}(\text{Pr}(\text{obese})) = \beta_0 + \beta_1 \cdot \text{woman} + \beta_2 \cdot (\text{age} - 45)$$

Logit estimates

Number of obs = 4690						
LR chi2(2) = 55.68						
Prob > chi2 = 0.0000						
Pseudo R2 = 0.0155						
obese	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
sex						
1	(base)					
2	.2743976	.0903385	3.04	0.002	.0973374	.4514579
age45	.0344723	.0051354	6.71	0.000	.0244072	.0445374
_cons	-2.147056	.0721981	-29.74	0.000	-2.288561	-2.00555

Significantly better than nothing - but is it good?

2

Logistic regression models: Test of fit

What about comparing the **estimated prevalence** with the **observed prevalence**?

In the Hosmer-Lemeshow test the data is **divided** into groups (traditionally 10) according to the **estimated probabilities**

and the **observed** and **expected** counts are compared in these groups by a chi-square test.

Most programs, that can fit a logistic regression model, can calculate this test.

In Stata it is done by (**after fitting the model**):

`estat gof, group(10) table`

The data is divided into **deciles** after the estimated probabilities.

3

Logistic regression models: Test of fit

OUTPUT

Logistic model for obese, goodness-of-fit test
(Table collapsed on quantiles of estimated probabilities)

Group	Prob	Obs_1	Exp_1	Obs_0	Exp_0	Total
1	0.0841	64	40.9	462	485.1	526
2	0.0953	43	45.5	453	450.5	496
3	0.1045	44	44.6	398	397.4	442
4	0.1112	42	50.3	422	413.7	464
5	0.1217	44	51.4	391	386.6	438
6	0.1332	52	63.0	441	430.0	493
7	0.1456	53	61.7	389	380.3	442
8	0.1592	62	69.8	392	384.2	454
9	0.1834	98	89.9	424	432.1	522
10	0.2407	99	83.8	314	329.2	413

number of observations = 4690
number of groups = 10
Hosmer-Lemeshow chi2(8) = 26.01
Prob > chi2 = 0.0010

One problem:
Too many in
the tails

Significant difference between observed and expected!

4

Logistic regression models: Test of fit

```
logit obese i.sex##age45
estat gof, group(10) table
Logistic model for obese, goodness-of-fit test
  (Table collapsed on quantiles of estimated probabilities)

+-----+
| Group |  Prob | Obs_1 | Exp_1 | Obs_0 | Exp_0 | Total |
+-----+
| 1 | 0.0796 | 36 | 35.9 | 466 | 466.1 | 502 |
| 2 | 0.1011 | 42 | 41.1 | 406 | 406.9 | 448 |
| 3 | 0.1053 | 49 | 49.6 | 429 | 428.4 | 478 |
| 4 | 0.1096 | 50 | 54.8 | 458 | 453.2 | 508 |
| 5 | 0.1124 | 52 | 54.2 | 436 | 433.8 | 488 |
| 6 | 0.1153 | 51 | 46.4 | 355 | 359.6 | 406 |
| 7 | 0.1182 | 52 | 53.9 | 410 | 408.1 | 462 |
| 8 | 0.1590 | 76 | 70.3 | 428 | 433.7 | 504 |
| 9 | 0.2133 | 96 | 91.8 | 391 | 395.2 | 487 |
| 10 | 0.3310 | 97 | 103.0 | 310 | 304.0 | 407 |
+-----+
  number of observations = 4690
  number of groups = 10
  Hosmer-Lemeshow ch2(8) = 2.43
  Prob > ch2 = 0.9650
```

The model 'fits' - when we look at it this way !!!!!!!

5

Conditional logistic regression When

Used in two situations:

1. Matched studies (binary response).
2. Unmatched studies with a confounder with many distinct values.

In 1. the models correspond to the way data was collected.

In 2. the method adjust for a 'mathematical' flaw in the unconditional method.

An example of situation 2:

The confounder is "kommune" having 275 distinct values.

6

Conditional logistic regression What

The logistic regression model (outcome disease yes/no):

$$\ln(\text{odds}) = \alpha + \sum_{i=1}^k (\beta_i \cdot x_i)$$

ln(odds) in reference ln(odds ratios)

Suppose the model above hold in each strata:

$$\ln(\text{odds}) = \alpha_s + \sum_{i=1}^k (\beta_i \cdot x_i)$$

ln(odds) in reference ln(odds ratios)
different in each strata the same in each strata

7

Conditional logistic regression What

$$\ln(\text{odds}) = \alpha_s + \sum_{i=1}^k (\beta_i \cdot x_i)$$

ln(odds) different in each strata

We are not interested in these !

In a matched study these are 'controlled'.

In a conditional logistic regression one 'condition on the odds in each strata', i.e. the case/control ratio.

In the conditional model the α 's disappear !

The β 's, the log OR's, are still in and can be estimated.

8

Conditional logistic regression How

A study of cancer in the oral cavity

Matched on **gender** and 10-year age groups

Ten strata (**genage**)

Here we focus on

textile-worker and

life time consumption of alcohol (three groups)

9

Conditional logistic regression How

logistic regression in *Stata*

`logit cancer textile i.alkcon i.genage, or`

`binreg cancer textile i.alkcon i.genage, or`

	cancer	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]
textile	1.652484	.6843458		1.21	0.225	.7338846 3.720889
alkcon						
0	1 (base)					
1	1.588614	.4485983		1.64	0.101	.9133833 2.763017
2	15.12845	4.890496		8.40	0.000	8.028433 28.50742
genage						
1	1.277731	1.598937		0.20	0.845	.1090655 14.84645
2	.6101724	.335794		-0.90	0.369	.2074977 1.794287
3						
4	1.196961	.7668028		0.58	0.779	.3410196 4.201272
5	.7482746	.4102097		-0.53	0.597	.2555206 2.191271
6	1.237034	.7746878		0.34	0.734	.3625102 4.221272
7	.7940664	.4255351		-0.43	0.667	.2779736 2.26835
8	1.734638	.9130996		1.05	0.295	.6182202 4.867148
9	1.02018	.6072521		0.05	0.957	.3257093 3.269977
10	1.745782	.9768952		1.00	0.319	.5830142 5.227581
_cons	.2301051	.1095992		-3.08	0.002	.0904687 .5852672 10

Conditional logistic regression in *Stata*

The syntax:

`clogit cancer textile i.alkcon, group(genage) or`

Part of the output:

	cancer	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]
textile	1.63708	.6717022		1.20	0.230	.732517 3.658661
alkcon						
0	1 (base)					
1	1.572508	.4390957		1.62	0.105	.909724 2.718168
2	14.30908	4.569879		8.33	0.000	7.651811 26.75835

11

Missing data - example 1

Consider the Frammingham study and imagine, that (due to a limited budget) only 500 measurements of SBP were allowed.

It was decided to take SBP measurements on 100 random participants in each of the age groups -40 and 60+ and 150 in each of the age groups 40-50 and 50-60.

That is we have missing SBP on 4190 of the 4,690 participants!

A short description of the design and the data:

agegrp	Freq.	N(sbp)	mean(sbp)	sd(sbp)
0-	1,325	100	122.18	15.4327
40-	1,684	150	130.85	22.2366
50-	1,346	150	140.93	22.4819
60-	335	100	149.51	26.9251
Total	4,690	500	135.87	24.0783

12

Missing data - example 1

agegrp	Freq.	N(sbp)	mean(sbp)	sd(sbp)
0-	1,325	100	122.18	15.4327
40-	1,684	150	130.85	22.2366
50-	1,346	150	140.93	22.4819
60-	335	100	149.51	26.9251
Total	4,690	500	135.87	24.0783

We note:

This is not a **completely** random sample
- the chance of being sample depends on age group!

The overall (total) average SBP is a biased estimate of the mean SBP among participants in the Framingham study!

I.e. an analysis of the 500 participants (a complete data analysis) will be biased.

13

Missing data - example 1

agegrp	Freq.	N(sbp)	mean(sbp)	sd(sbp)
0-	1,325	100	122.18	15.4327
40-	1,684	150	130.85	22.2366
50-	1,346	150	140.93	22.4819
60-	335	100	149.51	26.9251
Total	4,690	500	135.87	24.0783

We also note:

Within each age group the sample is **completely** random.
Within each age group the average SBP is an **unbiased** estimate of the mean SBP in the age group.

We know the size of each age group.

We can **calculate an unbiased** estimate of the total mean by weighing the group averages.

14

Missing data - example 1

agegrp	Freq.	N(sbp)	mean(sbp)	sd(sbp)
0-	1,325	100	122.18	15.4327
40-	1,684	150	130.85	22.2366
50-	1,346	150	140.93	22.4819
60-	335	100	149.51	26.9251
Total	4,690	500	135.87	24.0783

An unbiased estimate can be found as the **weighted average** of the group averages using the group sizes as weights:

$$\frac{122.18 \cdot 1325 + 130.85 \cdot 1684 + 140.93 \cdot 1346 + 149.51 \cdot 335}{4690} = 132.62$$

Conclusion: Although this is not a completely random sample, we have enough information in the data to find an unbiased estimate!!!!
(Assuming completely random sample **within** age group!)

15

Assuming that SBP is related to age:

Being missing is **not independent** of the **unobserved** SBP.

but

Being missing is **independent** of the unobserved SBP, **when we know the age group of the individual**.

The first statement means that the data is not **missing completely at random (MCAR)**.

The second statement corresponds to **missing at random (MAR)**, i.e. that given **all what we have observed** (including age group), then the missingness is (completely) random, i.e. independent of the unobserved data.

Mathematically Missing At Random implies that one (in theory) has enough information in the **observed data** to correct for the missing data - in principle.

16

Missing data: Standard terminology

Missing completely at random (MCAR).

The observed data is a (completely) random sample:
A complete data analysis will be unbiased

Missing at random (MAR)

Given all what we have observed, then the missingness is (completely) random (independent of the unobserved data):

The biased sampling might be adjusted for.

Missing not at random (MNAR)

None of the two above apply:

We will need further assumptions in order to analyse the data.

17

Missing at random

When the data is missing at random, then one can, in theory, make unbiased inference based on the observed data.

In the SBP example such an analysis could be to use the weighted average SBP instead of the biased unweighted average.

In general

If the sampled persons are not a completely random sample, but the i th person is sampled with a known probability, p_i , then we can obtain unbiased estimates by weighing the i th person with $1/p_i$.

The method is called Inverse Probability Weighing.

18

Inverse probability weighting

The SBP data:

Four different sampling probabilities and weights:

$$p_0 = 100/1325 = 0.0755 \quad w_0 = 1/p_0 = 13.25$$

$$p_1 = 150/1684 = 0.0891 \quad w_1 = 1/p_1 = 11.23$$

$$p_2 = 150/1346 = 0.1114 \quad w_2 = 1/p_2 = 8.97$$

$$p_3 = 100/335 = 0.2985 \quad w_3 = 1/p_3 = 3.35$$

That is information from each of the youngest should weight by 13.25 and information from the each of the oldest should weight by 3.35.

Sampling weights can be used in many Stata commands:

```
mean sbp [pw= sampw]
Mean estimation                               Number of obs = 500
-----+-----+-----+-----+-----+
      |     Mean   Std. Err.      [95% Conf. Interval]
-----+-----+-----+-----+-----+
      sbp | 132.6242  1.032943  130.5947  134.6536
-----+-----+-----+-----+-----+
```

19

Missing values – not by design

Most often the missing is not per design and both in the outcome and in the covariates:

<i>id</i>	<i>y</i>	x_1	x_2	x_3
1	o	o	o	o
2	o	m	o	o
3	m	o	o	o
4	m	m	o	o
5	o	o	o	o
6	o	m	m	o

o observed
m observed

Here we have only complete data on 2 persons, but partial information on 4 additional persons.

20

Missing values - not by design

If the missing is **completely at random**, then the analysis of the complete cases will be unbiased.

If this is not the case, then complete data analysis can give biased estimates.

If the data is **missing at random**, then it is **in theory** possible to make an unbiased analysis of all the data.

<i>id</i>	<i>y</i>	<i>x₁</i>	<i>x₂</i>	<i>x₃</i>
1	o	o	o	o
2	o	m	o	o
3	m	o	o	o
4	m	m	o	o
5	o	o	o	o
6	o	m	m	o

21

Imputation

One way to try solve the problem with missing is to **fill in** the data for the missing values and then make the analysis on the whole data set with the '**imputed**' values.

The imputation can be done in many ways.

One way is to fill in an "average" value.

This could be the total average of the observed values for the specific variable or the average in a **relevant subgroup**.

This method will not in general solve the bias problem.

And of course the **standard error** stated in the output, when you analyse the imputed data set, is **wrong**.

<i>id</i>	<i>y</i>	<i>x₁</i>	<i>x₂</i>	<i>x₃</i>
1	o	o	o	o
2	o	a ₁	o	o
3	a _y	o	o	o
4	a _y	a ₁	o	o
5	o	o	o	o
6	o	a ₁	a ₂	o

22

The missing SBP example

Imputation by **observed mean** in age group:

```
bysort agegrp: egen msbp=mean(sbp)
generate isbp=sbp
replace isbp=msbp if missing(sbp)
```

```
mean isbp
Mean estimation                               Number of obs = 4690
                                               
                                                | Mean   Std. Err. [95% Conf. Interval]
-----+-----+-----+-----+-----+-----+
      isbp | 132.6242  .1627486  132.3051  132.9432
```

Correct mean, but a much too small standard error - incorrectly **assuming 4690 independent observations**.

Correct analysis using sampling weights:

```
mean sbp [pw=sampw]
Mean estimation                               Number of obs = 500
                                               
                                                | Mean   Std. Err. [95% Conf. Interval]
-----+-----+-----+-----+-----+-----+
      sbp | 132.6242  1.032943  130.5947  134.6536
```

23

Imputation - random multiple

A fixed imputation will not take into account the random variation of the unobserved observation or the uncertainty of the parameters.

Imputation methods should add some random variation to the imputed data.

For that we need a **statistical model** for the missing data.

In **multiple imputations** one generates **several imputed** data sets.

For each imputed data set one fit the model of interest.

The point estimate, then the average across the imputed data sets.

One tricky thing is **calculation of the standard errors**.

<i>id</i>	<i>y</i>	<i>x₁</i>	<i>x₂</i>	<i>x₃</i>
1	o	o	o	o
2	o	m	o	o
3	m	o	o	o
4	m	m	o	o
5	o	o	o	o
6	o	m	m	o

24

Multiple imputations

Questions:

How to find **the models** from which to generate the missing data?

How should you handle missing data in this process?

How to find the uncertainty (**standard errors**) of the estimates?

Bookkeeping.

Most important: **Missing at random is required!**

<i>id</i>	<i>y</i>	<i>x₁</i>	<i>x₂</i>	<i>x₃</i>
1	0	0	0	0
2	0	m	0	0
3	m	0	0	0
4	m	m	0	0
5	0	0	0	0
6	0	m	m	0

25

The missing SBP example

```
use sbpdata,clear
mi set mlong
mi register imputed sbp
(4190 m=0 obs. now marked as incomplete)
```

```
mi impute regress sbp i.agegrp, add(20)
```

Univariate imputation		Imputations = 20		
Linear regression		added = 20	updated = 0	
Imputed: m=1 through m=20		observations per m		
variable	complete	incomplete	imputed	total
sbp	500	4190	4190	4690

(complete + incomplete = total; imputed is the minimum across m of the number of filled in observations.)

26

The missing SBP example

```
codebook, comp
```

variable	obs	unique	Mean	Min	Max	Label
sbp	84300	83383	132.3204	44.52609	270	Systolic Blood Pressure
id	88490	4690	2352.429	1	4699	
agegrp	88490	4	1.107481	0	3	
_mi_id	88490	4690	2357.795	1	4690	
_mi_miss	4690	2	.8933902	0	1	
_mi_m	88490	21	9.943496	0	20	

```
sum if _mi_m==1
```

Variable	obs	Mean	Std. Dev.	Min	Max
sbp	4190	131.2507	21.65931	59.92363	209.6556
id	4190	2352.611	1359.59	2	4699
agegrp	4190	1.105251	.8895275	0	3
_mi_id	4190	2358.483	1331.661	101	4690
_mi_miss	0				
_mi_m	4190	1	0	1	1

27

The missing SBP example

```
. table agegrp if _mi_m>0, c(count sbp mean sbp sd sbp)
```

agegrp	N(sbp)	mean(sbp)	sd(sbp)
0-	24,500	121.5843	22.32535
40-	30,680	131.1271	22.37045
50-	23,920	141.2539	22.4434
60-	4,700	150.2313	22.19089

20*1225=24500

20*235=4700

```
. table agegrp if _mi_m==0, c(count sbp mean sbp sd sbp)
```

agegrp	N(sbp)	mean(sbp)	sd(sbp)
0-	100	122.18	15.4327
40-	150	130.85	22.2366
50-	150	140.93	22.4819
60-	100	149.51	26.9251

28

The missing SBP example

```

mi estimate: mean sbp

Multiple-imputation estimates
Mean estimation
DF adjustment: Small sample
Within VCE type: ANALYTIC

Imputations      =      2
Number of obs   =    469
Average RVI    =  7.427
Complete DF     =    468
DF:             min   =    23.4
                           avg   =    23.4
                           max   =    23.4

-----  

      Mean |   Coef.   Std. Err.      t    P>|t|   [95% Conf. Interval]
-----  

      sbp | 132.6799 1.017506  130.40   0.000    130.5772  134.782
-----
```

Correct analysis using sampling weights

6

A more complicated example

```
use sbp2data,clear
codebook,comp

variable   Obs  Unique      Mean     Min     Max  Label
-----+-----+-----+-----+-----+-----+-----+
sex        4188      2  1.566141      1      2  Sex
sbp        4216     112 132.6945     80     270  Systolic Blood Pressure
dbp        4281      67  82.62766     40     148  Diastolic Blood Pressure
scl        4192     244 228.2011    115     568  Serum Cholesterol
age        4245      37  46.0636     30      66  Age in Years
bmi        4218     245 25.63148    16.2     57.6  Body Mass Index
id         4690    2349.172      1     4699
-----+-----+-----+-----+-----+-----+-----+
regress sbp age i.sex
Source |      SS          df          MS      Number of obs = 3,406
-----+-----+-----+-----+-----+-----+
          Model | 281261.425          2  140630.713  F(2, 3403) = 320.62
          Residual | 1492627.36      3,403  438.621029  Prob > F     = 0.0000
-----+-----+-----+-----+-----+-----+
          Total | 1773888.79      3,405  520.96587  R-squared     = 0.1586
                                      Adj R-squared = 0.1581
                                      Root MSE     = 20.943
-----+-----+-----+-----+-----+-----+
sbp |      Coef.      Std. Err.          t      P>|t|  [95% Conf. Interval]
-----+-----+-----+-----+-----+-----+-----+
          age |  1.072026  .0423621      25.31  0.000  .9889686  1.155084
          sex |
          Male |          0  (base)
          Female |  .2701054  .7247534      0.37  0.709  -1.150891  1.691101
          _cons |  83.39557  2.017962     41.33  0.000  79.43903  87.35211
-----+-----+-----+-----+-----+-----+-----+

```

— 30 —

A more complicated example

```
mi set mlong
mi register imputed sbp age sex dbp bmi scl
(2201 m=0 obs. now marked as incomplete)

mi describe

Style:  mlong
        last mi update 24nov2016 16:43:28, 0 seconds ago

Obs.:  complete      2,489
       incomplete    2,201  (M = 0 imputations)
       -----
       total         4,690

Vars.:  imputed:  6; sbp(474) age(445) sex(502) dbp(409) bmi(472) scl(498)
        passive:  0
        regular:  0
        system:   3; _mi_m _mi_id _mi_miss

        (there is one unregistered variable: id)
```

mi missable pattern, freq	
Missing-value patterns	
(1 means complete)	
	Pattern
Frequency	1 2 3 4 5
2,489	1 1 1 1 1
314	1 1 1 1 0
301	1 1 1 1 1
281	1 1 1 0 1
278	1 1 0 1 1
253	1 0 1 1 1
243	0 1 1 1 1
42	1 1 1 0 0
37	1 0 1 1 1
37	1 1 1 0 1
36	1 1 0 1 1
35	1 1 1 1 0
34	0 0 1 1 1
33	1 0 0 1 1
32	0 1 0 1 1
30	1 1 0 0 1
28	1 1 0 1 0
27	1 0 1 0 1
25	0 1 1 0 1
25	0 1 1 1 1
25	1 0 1 1 1

Variables are (1) dbp (2) age (3) bmi (4) sbp (5) scl (6) sex 32

A more complicated example

```
mi impute chained ///
(regress,include( i.sex age bmi      dbp scl))sbp ///
(regress,include( i.sex age bmi sbp      scl))dbp ///
(regress,include( i.sex age bmi      scl))bmi ///
(regress,include( i.sex age bmi sbp dbp scl))age ///
(regress,include( i.sex age bmi      scl))scl ///
(logit, include(      age bmi      ))sex ///
,add(100) noimputed
```

Conditional models:

```
dbp: regress dbp i.sex age bmi sbp scl
age: regress age i.sex bmi sbp dbp scl
bmi: regress bmi i.sex age scl
sbp: regress sbp i.sex age bmi dbp scl
scl: regress scl i.sex age bmi
sex: logit sex age bmi
```

Performing chained iterations ...

Multivariate imputation
Chained equations
Imputed: m=1 through m=100

Imputations = 100
added = 100
updated = 0

Initialization: monotone

Iterations = 1000
burn-in = 10

33

A more complicated example

variable	Observations per m			
	Complete	Incomplete	Imputed	Total
sbp	4216	474	474	4690
dbp	4281	409	409	4690
bmi	4218	472	472	4690
age	4245	445	445	4690
scl	4192	498	498	4690
sex	4188	502	502	4690

(complete + incomplete = total; imputed is the minimum across m of the number of filled-in observations.)

codebook, comp

Variable	Obs	Unique	Mean	Min	Max	Label
sbp	224316	47338	132.3185	47.08539	270	Systolic Blood Pressure
dbp	224381	40808	82.44368	40	148	Diastolic Blood Pressure
scl	224292	49896	227.064	34.90916	568	Serum Cholesterol
age	224345	44422	45.95276	12.26457	82.2043	Age in Years
bmi	224318	47293	25.52148	9.895696	57.6	Body Mass Index
id	224790	4690	2348.082	1	4699	
sex	224288	2	.5676273	0	1	RECODE of koen (sex)
_mi_m	224790	101	49.44637	0	100	
_mi_id	224790	4690	2329.055	1	4690	
_mi_miss	4690	2	.4692964	0	1	

34

A more complicated example

mi estimate: regress sbp age i.sex

Multiple-imputation estimates
Linear regression

Imputations = 100
Number of obs = 4,690
Average RVI = 0.1130
Largest FMI = 0.1394
Complete DF = 4687

DF adjustment: Small sample

DF: min = 2,256.34
avg = 2,715.41
max = 3,031.05

Model F test: Equal FMI
within VCE type: OLS

F(2, 3480.5) = 396.38
Prob > F = 0.0000

sbp	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
age	1.072957	.0376538	28.50	0.000	.9991277 - 1.146787
sex					
Male	0 (base)				
Female	.2033005	.6617939	0.31	0.759	-1.094488 - 1.501089
_cons	83.29757	1.802549	46.21	0.000	79.76314 - 86.83199

35

Clustered data / data with several random components Dichotomous outcome

A different outcome:

$$H_{fpd} = \begin{cases} 1 & \text{if the person has hayfewer} \\ 0 & \text{else} \end{cases}$$

A statistical model:

$$\text{logit}(H_{fpd} = 1) = \beta_0 + \beta_I \cdot I + \beta_U \cdot U + \beta_A \cdot A + \beta_S \cdot S + \beta_G \cdot G$$

$$+ F_j + P_{fp} + X_{ba}$$

Random part
This is not needed
due to the binomial
error

36

Clustered data / data with several random components
Dichotomous outcome

$$\text{logit}(H_{fpd} = 1) = \beta_0 + \beta_I \cdot I + \beta_U \cdot U + \beta_A \cdot A + \beta_S \cdot S + \beta_G \cdot G + F_f + P_{fp}$$

That is, an ordinary logistic regression + random components.

- A **generalized linear mixed model**
- A **multilevel model for dichotomous outcome**

Comments 1:

- It is **important** to include the **relevant random components** in the model.
- 'Multilevel models' is **essential** in medical/epidemiological research.

37

Clustered data / data with one random components
Dichotomous outcome

If the models only involve **one random component**, e.g. **variation between families** or between GP's, then methods exist which can **adjust the standard errors**.

Remember that if the **data contains clusters**, then the precision of the estimates are overestimated, that is, the reported **standard errors are too small**.

So-called **robust methods** or **sandwich estimates** of the standard errors will (try to) adjust for this problem.

Only a **few** programs have this option - Stata does!

39

Clustered data / data with several random components
Dichotomous outcome

Comments 2:

- The theory and insight into the models for non-normal data are **not yet fully developed**.
- The main problem being that it is very difficult to find **valid (unbiased) estimates**.
- Several software programs **falsely claim** to estimate the models.
- Some programs like Stata and NLwin can give you valid estimates if you take care and have **a lot of data**.

Advice:

Do not try to estimate this kind of models without consulting a specialist.

38