

## Lecture 5a: Logistic regression diagnostics (VER/MER 16.12)

<b>Index</b>	<b>Page</b>
Covariate patterns.....	2
Pearson residuals per covariate pattern.....	4
Goodness of fit tests.....	5
Residual analysis (covariate patterns).....	7
Leverage (h).....	9
Influential statistics .....	10
Dealing with influential observations.....	12
Summary logistic regression diagnostics.....	13
Predictive ability of a logistic model.....	14
Cross-validation.....	17
Overdispersion [L12a - L12b].....	18

### **Learning Objectives:**

- ★ Understand covariate patterns
- ★ Conduct residuals and influential analysis
- ★ Compute the predictive power and reliability

### **Tuesday:**

- ★ logistic regression exercises 2 and 3

### **Final Model:** Dataset:nocardia.dta

- All the examples based on VER Ex. 16.5 (model with **dcpct3**, dneo, dclox and dneo#dclox)

## Covariate patterns

- Unique combination of values of predictor variables
- Categorical predictors
  - ★ dcpct -> converted to categorical (dcpct3= 0-49, 50-99, 100)
  - ★ 11 covariate patterns

```
list cov dneo dclox dcpct3 obs
```

	cov	dneo	dclox	dcpct3	obs
1.	1	0	0	0	12
2.	2	0	0	50	2
3.	3	0	0	100	8
4.	4	0	1	50	1
5.	5	0	1	100	11
6.	6	1	0	0	11
7.	7	1	0	50	10
8.	8	1	0	100	38
9.	9	1	1	0	1
10.	10	1	1	50	5
11.	11	1	1	100	9

- Continuous predictors

- ★ production and SCC

- ★ 99 covariate patterns

```
. list cov prod bscc obs in 1/4
```

	cov	prod	bscc	obs
1.	1	11.40	182	1
2.	2	12.10	277	1
3.	3	15.90	167	1
4.	4	16.10	310	1
.....				
96.	96	30.60	167	1
97.	97	32.20	75	1
98.	98	34.00	206	1
99.	99	34.50	61	1

## Residuals in logistic regression

- One per observation (based on Hilbe, 2009<sup>1</sup>)
  - ★ (standard) residual and influential analyses
  - ★ mainly for visual assessment
  - ★ not very useful for assessing the model
  - ★ Stata command *-glm-*
    - ➔ Pearson
    - ➔ Anscombe (more normal distributed )
  
- One per covariate pattern
  - ★ goodness-of-fit tests
  - ★ residual analysis
    - ➔ Pearson residuals (standardized)
    - ➔ deviance residuals (not covered in this course)
  - ★ Stata command *-logit / logistic -*

### ● Example - residuals

				Residual	
Obs.	Cov. pattern	Disease	Pred. Value	1 per Obs.	1 per Cov. Pat.
1					
2					

<sup>1</sup> Hilbe J. Logistic Reg. Models. CRC Press: Boca Raton 2009

## Pearson residuals per covariate pattern

- Pearson residual (standardized)

- ★ 
$$r_j = \frac{y_j - m_j * p_j}{\sqrt{m_j * p_j * (1 - p_j)}}$$

- ➔  $y_j$  = nbr. pos. outcomes in  $j^{\text{th}}$  covariate pattern

- ➔  $m_j$  = nbr. obs. in the  $j^{\text{th}}$  covariate pattern

- ➔  $p_j$  = predicted prob. for the  $j^{\text{th}}$  covariate pattern

- ★ standardized residuals

- ➔ 
$$r_{sj} = \frac{r_j}{\sqrt{(1 - h_j)}} \quad ; \text{ makes variance to 1}$$

- $h_j$  = leverage (next slides)

- ★ 
$$\sum_{j=1}^J (r_j)^2 = \text{Pearson } \chi^2 \text{ statistic} \sim \chi^2 \text{ with } (J - k) \text{ df.}$$

- ➔  $k$  = number of parameters in the model

- ★ contribution of each cov. pattern to the  $\chi^2$  statistic

- ★ Stata command – *predict* after *logit/logistic* -

- ➔ eg. `logit casecon i.dneo`

- `predict pear, rresidual`

- `predict stdpear, rstandard`

## Goodness of fit tests

- Pearson  $\chi^2$  (1 per cov. Pattern) (see L3b)

- ★  $\chi^2$  distributions (J-k) df

- ➔ J = # covariate patterns; k = # of parameters in model

- ★ only if enough number of replications per group (eg cov. patterns)

- ➔ ~ guidelines ~  $X^2$ -statistic

- at least 1 expected count in each cell

- at least 80% expected values  $\geq 5$  counts

- ★ indicates the fit of the model (Ho = model fits the data)

- Example: nocardia

cov pat and Case		Nbr herds	Nbr pos	Pred. pr.	Pear. Res.	Pear. Res. <sup>2</sup>
- Control						
1	no yes	12	1	0.028	1.144	1.308949
2	no yes	2	0	0.102	-0.478	.2283039
3	no yes	8	1	0.182	-0.416	.1731429
4	no yes	1	1	0.152	2.360	5.569528
5	no yes	11	2	0.259	-0.584	.3405482
6	no yes	11	4	0.416	-0.353	.1245677
7	no yes	10	7	0.735	-0.254	.0643712
8	no yes	38	33	0.844	0.416	.1731366
9	no yes	1	0	0.082	-0.298	.0890559
10	no yes	5	1	0.258	-0.295	.0872724
11	no yes	9	4	0.403	0.252	.0634578

- ★ no indication of lack of fit

- ★ largest contribution is from cov. pattern with no replication (eg cov # 4)

- Pearson  $\chi^2$  after logit command

- ★ p-value is meaningful if enough replicates

- ★ Stata command (after logit)

- ➔ *estat gof*

```
. estat gof
Logistic model for casecont, goodness-of-fit test

      number of observations =          108
number of covariate patterns =           11
      Pearson chi2(5) =             8.22

              Prob > chi2 =             0.1444
```

- Hosmer-Lemeshow Test

- ★ the only useful test when there are few replicates per covariate pattern

- ★ group by:

- ➔ percentiles of predicted probabilities

- ➔ fixed points of predicted probabilities

- ★ compares predicted probabilities to observed probabilities in groups (g) of data

- ➔  $\chi^2$  with g-2 df, where “g” = number of groups

- ➔ low power if < 6 groups

- ★ Stata command (after logit)

- ➔ *estat gof, group(#) table*

## ● Example

```
. estat gof, g(10) table
```

Logistic model for casecont, goodness-of-fit test

(Table collapsed on quantiles of estimated probabilities)  
(There are only 7 distinct quantiles because of ties)

Group	Prob	Obs_1	Exp_1	Obs_0	Exp_0	Total
1	0.0284	1	0.3	11	11.7	12
2	0.1817	2	1.9	10	10.1	12
3	0.2589	3	4.1	13	11.9	16
4	0.4033	4	3.6	5	5.4	9
5	0.4161	4	4.6	7	6.4	11
6	0.7354	7	7.4	3	2.6	10
10	0.8439	33	32.1	5	5.9	38

```
number of observations =      108
number of groups      =         7
Hosmer-Lemeshow chi2(5) =      2.16
Prob > chi2           =      0.8262
```

## Residual analysis (covariate patterns)

### ● Pearson residuals (standardized)

★ identify large negative and positive residuals

★ visual assessment

➔ some guidelines about outlying obs.

### ● Example

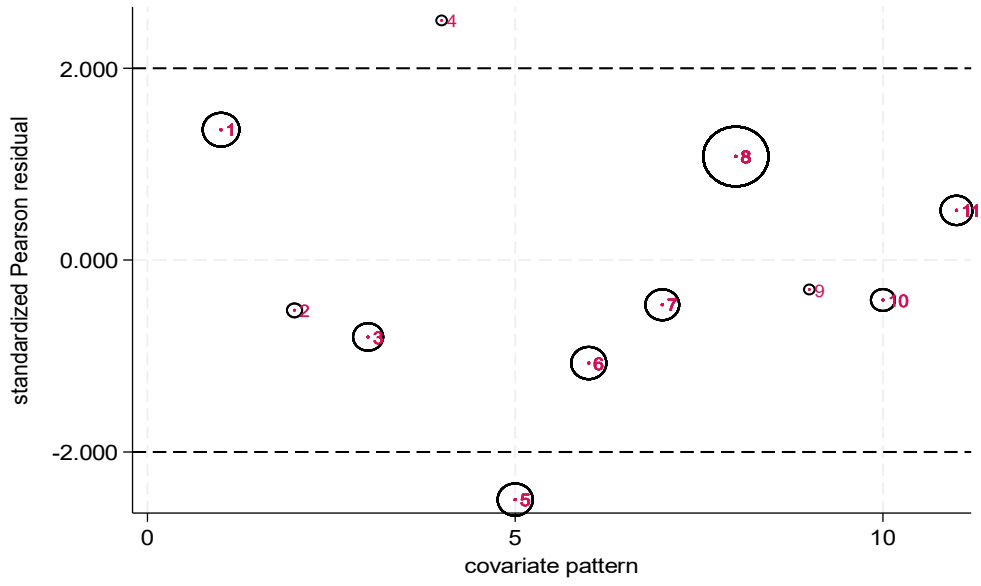
★ `logit casecont i.dneo###dclox i.dcpct3`

➔ `predict pear_std, rstandard`

➔ `predict cov, num`

★ `plot stdz. Pearson residuals (per cov. pattern) vs cov. pattern`

➔ size of the bubble ~ number of observations per cov. pattern



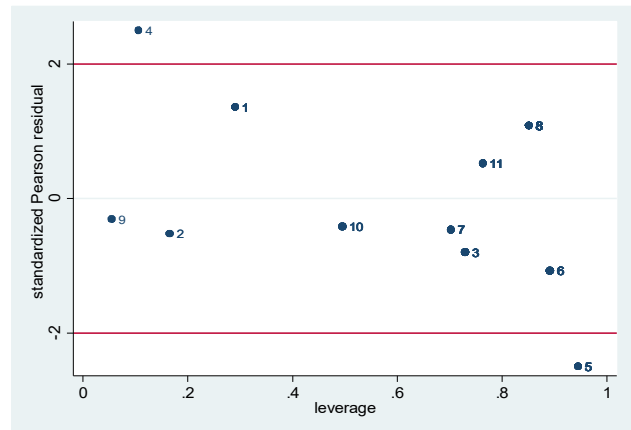
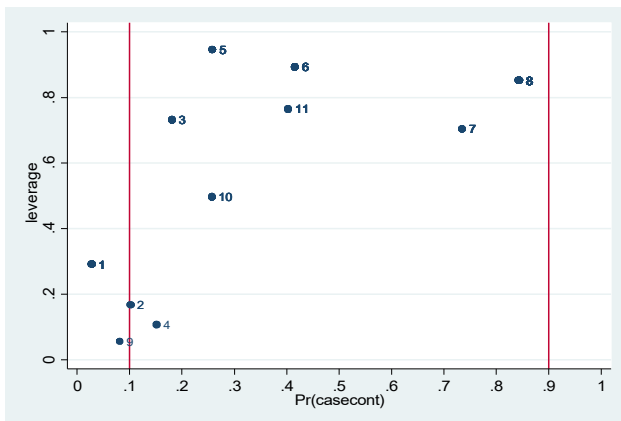
```
list cov cnt dcpct3 dneo dclox opr pv pear_std if wcov==1 & abs(pear_std)>2,noobs
```

cov	cnt	dcpct3	dneo	dclox	opr	pv	pear_std
4	1	50	no	yes	1.000	0.152	2.496
5	11	100	no	yes	0.182	0.259	-2.496

## Leverage ( $h$ )

- Potential impact of cov. pattern on the model
- Extent to which the  $j^{\text{th}}$  cov. pattern is separated for the others in terms of the explanatory variables
- One value per cov. pattern (Stata command `-logit / logistic -`)
- Leverage depends on x's and predicted probabilities
  - ★ if predicted probabilities  $<0.1$  and  $>0.9$  then leverage values will be low
  - ★ if predicted probabilities  $>0.1$  and  $<0.9$  then leverage values can be interpreted as distance
    - ➔ look for large leverage values within this range
    - ➔ look for points that fall some distance from the rest of the data

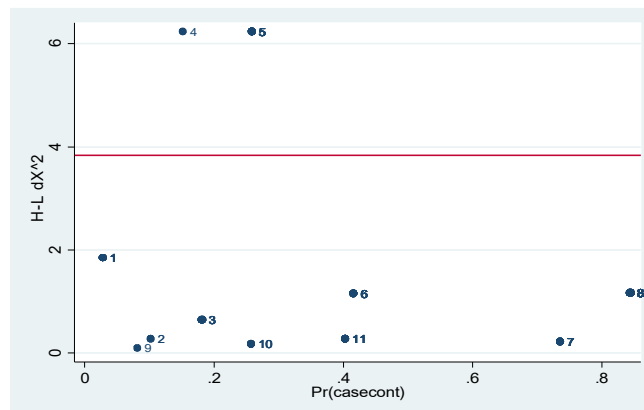
### ● Example



cov	cnt	dcpct3	dneo	dclox	opr	pv	pear_std	lev
4	1	50	no	yes	1.000	0.152	2.496	0.106
...								
11	9	100	yes	yes	0.444	0.403	0.518	0.764
8	38	100	yes	no	0.868	0.844	1.080	0.852
6	11	0	yes	no	0.364	0.416	-1.073	0.892
5	11	100	no	yes	0.182	0.259	-2.496	0.945

## Influential statistics

- ★  $\Delta \chi^2$  (  $\Delta \chi^2$  )
- ★ delta-beta (  $\Delta \beta$  ) (Cook's distance in normal regression)
- ★ one per covariate pattern (Stata command *-logit / logistic -* )
- Delta  $\chi^2$  (  $\Delta \chi^2$  ) (or  $r_{sj}^2$ )
  - ★ effect of a covariate pattern on fit of the model
    - ➔ identifies covariate patterns that do not fit well (outliers)
  - ★ plot delta values vs predicted probabilities
  - ★ plot delta values vs leverage
  - ★ delta-values  $\geq 3.84$  (95<sup>th</sup> percentile  $\chi^2$  distribution with 1df)
- Example - delta  $\chi^2$  (  $\Delta \chi^2$  ) vs Pr(Pr)



```
. list cov cnt dcpct3 dneo dclox pv dx2 lev pear if dx2>3.84 & wcov==1, noobs
```

cov	cnt	dcpct3	dneo	dclox	pv	dx2	lev	pear
5	11	100	no	yes	.2588893	6.232499	.9453593	-0.584
4	1	50	no	yes	.152218	6.232499	.1063734	2.360

- Delta-beta (  $\Delta\beta$  )
  - ★ analogous to Cook's distance
  - ★ measures influence of a cov. pattern on:
    - ➔ overall set of betas (Stata)
    - ➔ individual betas (SAS)
  - ★ depends on leverages and # of observations (  $m_j$  ) on the covariate pattern variable
    - ➔ Hosmer & Lemeshow<sup>2</sup> suggest that values >1 might be influential

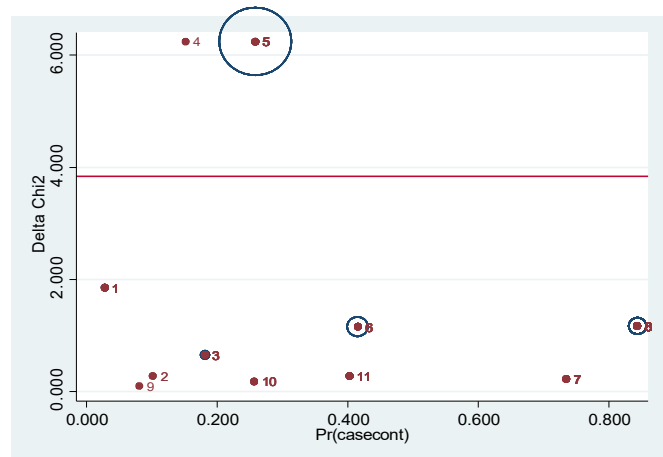
- Delta  $\chi^2$  and Delta-beta

- ★ values will depend on the predicted probabilities (similar to leverage)<sup>2</sup>

- Plots

- ★  $\Delta\beta$  vs pred. prob. and  $\Delta\beta$  vs leverage values

- ★  $\Delta\chi^2$  vs predicted probabilities with size proportional to  $\Delta\beta$



2 Hosmer and Lemeshow. Applied Log. Reg. 2<sup>nd</sup> Edition. pg-174-176

★ influential observations

```
. l cov cnt dcpct dneo dclox opr pv lev dx2 db if db > abs(1) & wcov==1, noobs
```

	cov	cnt	dcpct	dneo	dclox	opr	pv	lev	dx2	db
	3	8	100	no	no	0.125	0.182	0.730	0.642	1.739
	8	38	100	yes	no	0.868	0.844	0.852	1.167	6.693
	6	11	5	yes	no	0.364	0.416	0.892	1.152	9.504
	5	11	100	no	yes	0.182	0.259	0.945	6.232	107.831

★ delta-betas extreme for cov. 5 (same as VER), 6 (not in VER) and 8 (noted in VER to be due to a large group size)

### Dealing with influential observations

- Identify points with large residuals or large leverage values
- Evaluate their covariate patterns – why are they outliers?
- Delete from model and re-fit the model

★ does it change very much?

```
. estimates table final wocov5 wocov6 wocov8 , b(%5.3f) stats(N) star( .05 .01 .001)
```

Variable	final	wocov5	wocov6	wocov8
dneo				
yes	3.192***	3.248***	3.639***	2.518*
dclox				
yes	0.453	17.322	0.808	0.705
dneo#dclox				
yes#yes	-2.533*	-19.403	-3.018*	-2.053
dcpct3				
50	1.361	1.087	0.120	1.173
100	2.027**	2.132**	0.813	1.168
_cons	-3.531***	-3.581***	-2.672*	-2.972**
N	108	97	97	70

legend: \* p<.05; \*\* p<.01; \*\*\* p<.001

- ★ cp 5 – only cov. with dneo=0 and dclox=1 – part of the interaction
- ★ cp 6 – dneo=1 and dclox=0 with 0 dcpct
- ★ cp 8 – largest cov. pattern
- ★ cp 4 – largest contribution to the deviance and Pearson X<sup>2</sup> – however no influence (delta-beta = 0.74)

## Summary logistic regression diagnostics

- Covariate pattern residuals
  - ★ goodness-of-fit tests
    - ➔ inadequacies in the modelling of the predictors in the model
      - e.g non-linearity or missing interactions
    - ➔ can't detect missing predictors or clustering
  - ★ outlying observations (cov. patterns)
- Diagnostics
  - ★ consequences of the current model
    - ➔ eg. few replicates?
    - ➔ increase the nbr of obs. per cov. pattern (or reduce the nbr. of cov. patterns)
      - grouping (biological, frequency categories)
      - remove continuous predictors
        - eg. those that are borderline sig.
    - ➔ however... you will be working with a different model – not your final model!
  - ★ identify high influence cov. patterns (for instance  $\Delta\beta$ ) on the parameter estimates

## Predictive ability of a logistic model

$$\ln\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 * X_1 = X \beta \quad p = \frac{e^{(\beta X)}}{1 + e^{(\beta X)}} \quad \text{or} \quad p = \frac{1}{1 + e^{-\beta X}}$$

★ note: not a true probability if derived from a case-control study [see 14a]

## Sensitivity and Specificity

- Predict D+ if  $p \geq 0.5$

★ choose other cutpoint

Classified (predicted status)	true D+	true D-	Total
T+ = $p(D+) \geq 0.5$	40	8	48
T- = $p(D+) < 0.5$	14	46	60
Total	54	54	108

★ sensitivity (Se) =

★ specificity (Sp) =

★ positive predictive value (PPV) =

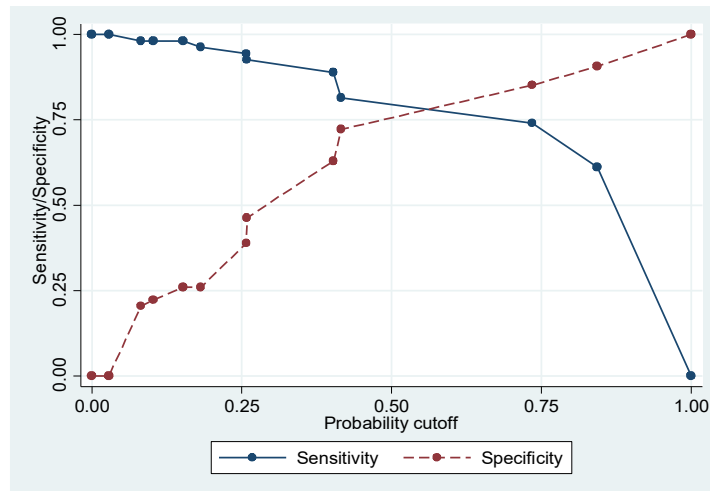
★ negative predictive value (NPV) =

★ overall correctly classified =

- Two-graph ROC (Se-Sp plot)

- ★ effect of changing the cut-point on Se and Sp.

- ➔ Se / Sp values corresponds the values  $\geq$  than the cutoff point where data is available (eg. for a cut-off  $\geq 0.5$ , the next value is 0.74)



- ROC curve

- ★ Se vs 1-Sp

- ★ assess discriminatory power of the model

- ➔ eg. probability that cases (or controls) are correctly classified by the model (at a given cutpoint)

- ★ quantify Area Under the Curve (AUC)

- ★ AUC interpretation<sup>1</sup>

AUC	Interpretation
0.5	No discrimination (better flip a coin!)
0.5 - 0.7	Poor discrimination
0.7 - 0.8	Acceptable discrimination
0.8 - 0.9	Excellent discrimination
>0.9	Outstanding discrimination

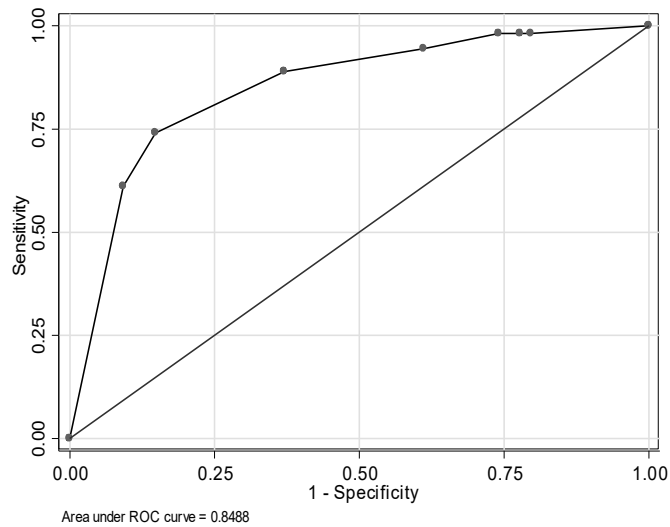
<sup>1</sup> Hosmer Lemeshow. Applied logistic regression (pg177)

● AUC = final model  
 . roctab casecont pv, sum detail graph

Detailed report of sensitivity and specificity

Cutpoint	Sensitivity	Specificity	Correctly classified	LR+	LR-
( >= .0284.. )	100.00%	0.00%	50.00%	1.0000	
( >= .0817.. )	98.15%	20.37%	59.26%	1.2326	0.0909
( >= .1024.. )	98.15%	22.22%	60.19%	1.2619	0.0833
( >= .152218 )	98.15%	25.93%	62.04%	1.3250	0.0714
( >= .1817.. )	96.30%	25.93%	61.11%	1.3000	0.1429
( >= .2577.. )	94.44%	38.89%	66.67%	1.5455	0.1429
( >= .2588.. )	92.59%	46.30%	69.44%	1.7241	0.1600
( >= .4032.. )	88.89%	62.96%	75.93%	2.4000	0.1765
( >= .4160.. )	81.48%	72.22%	76.85%	2.9333	0.2564
( >= .7353.. )	74.07%	85.19%	79.63%	5.0000	0.3043
( >= .8439.. )	61.11%	90.74%	75.93%	6.6000	0.4286
( > .8439.. )	0.00%	100.00%	50.00%		1.0000

Obs	ROC area	Std. err.	Asymptotic normal [95% conf. interval]	
108	0.8460	0.0377	0.77221	0.91984



★ probability of a case having higher predicted probability than a control is 0.84 for any randomly selected pair

## Cross-validation

- Similar to linear regression – leave-one-out analysis
  - ★ fit a model without observation “i”
  - ★ obtain the predicted probability “p” for the observation not included in the estimation
  - ★ then computed the “predicted AUC” and compare to the original final model

Example – nocardia : i.dneo###i.dclox i.dcpct3

### ★ Final model AUC

Obs	ROC Area	Std. Err.	-Asymptotic Normal-- [95% Conf. Interval]	
108	0.8488	0.0370	0.77621	0.92132

### ★ Cross-validated AUC

Obs	ROC area	Std. err.	Asymptotic normal [95% conf. interval]	
108	0.7805	0.0485	0.68544	0.87560

- ★ Cross-validated model provided acceptable discriminatory power
- ★ Discrepancy with original model – excellent discrimination

## Overdispersion [L12a - L12b]

- Assumption  $y_i \sim$  binomial distribution
  - ★ mean =  $n_i * p_i$
  - ★ variance =  $n_i * p_i * (1 - p_i)$
- Overdispersion = the data are more dispersed (larger variance) than would be expected
  - ★ apparent overdispersion – wrong model
    - ➔ missing important predictors
    - ➔ outliers
  - ★ real overdispersion – usually due to clustering
  - ★ too small S.E.