

# Chapter 7

## Limited dependent variables: classification

modeling binary, nominal, ordinal and  
count data

## Modeling binary variables

**Practical question:** a bank should decide about granting loans to new clients, i.e. forecast of the solvency

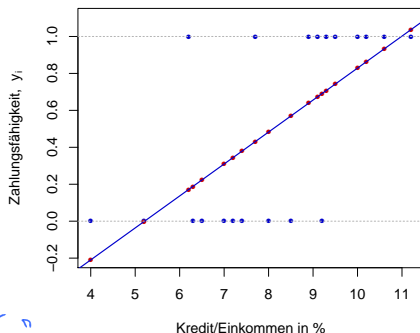
$$Y_i = \begin{cases} 0, & \text{the client } i \text{ is solvent} \\ 1, & \text{the client } i \text{ is insolvent} \end{cases}$$

- $X_{1i}$  – debt-to-income ratio (  $\times 100$ );
- $X_{2i}$  – years with the current employer;
- $X_{3i}$  – other debts (in 1000 Euro);
- $X_{4i}$  – age (in years).

**Question:** can we use a linear regression model for binary variables?  $\rightsquigarrow$   
linear probability model

## Linear Prob.-model

$$Y_i = \beta_0 + \beta_1 \cdot X_i + u_i$$



$E(Z) = \sum z_i \cdot P(Z=z_i)$  - expectation for discrete RV.

Note:

- (+) the forecast  $\hat{Y}_i$  can be seen as probability

$$E(Y_i|X_i) = 1 \cdot P(Y_i = 1|X_i) + 0 \cdot P(Y_i = 0|X_i) = p_i$$

- (-)  $\hat{Y}_i$  may lie outside of  $[0,1]$
- (-)  $R^2$  is useless as a goodness-of-fit measure
- (-) the residuals are not normally distributed
- (-)  $Var(Y_i|X_i) = p_i(1 - p_i) \neq const \leadsto$  heteroscedastic

expectation = probability of being insolvent

$$u_i = \begin{cases} 1 - \beta_0 - \beta_1 X \\ 0 - \beta_0 - \beta_1 X \end{cases}$$

$$F_2(z) = P(Z \leq z)$$

## Transition to Logit/Probit

Let  $Y_i$  be the observed binary variable and  $Y_i^*$  the corresponding unobserved metric variable. For  $Y_i^*$  it holds:

$$Y_i^* = \beta_0 + \beta_1 X_{1i} + \dots + \beta_k X_{ki} + u_i = \mathbf{X}_i' \boldsymbol{\beta} + u_i.$$

**Example:**  $Y_i^*$  is an unobserved solvency of the client  $i$  with

$$Y_i = 1 \text{ if } Y_i^* > 0 \text{ and } Y_i = 0 \text{ if } Y_i^* \leq 0.$$

one can always compare two good or two bad clients.

$$\begin{aligned} \underline{P(Y_i = 1 | \mathbf{X}_i)} &= P(Y_i^* > 0 | \mathbf{X}_i) = P(\mathbf{x}_i' \boldsymbol{\beta} + u_i > 0 | \mathbf{X}_i) \\ &= P(-u_i < \mathbf{X}_i' \boldsymbol{\beta} | \mathbf{X}_i) = \underline{F(\mathbf{X}_i' \boldsymbol{\beta})}, \end{aligned}$$

→ cdf of the residuals  $u_i$ .

where  $F(\cdot)$  is the cdf of the residuals.

- $F(z) = \frac{1}{1+e^{-z}}$  - the cdf of the logistic distribution  $\rightsquigarrow$  **logit**
- $F(z)$  - the cdf of the normal distribution  $\rightsquigarrow$  **probit**

$$\lim_{z \rightarrow -\infty} \frac{1}{1 + e^{-z}} = 0; \quad \lim_{z \rightarrow +\infty} \frac{1}{1 + e^{-z}} = 1$$

## Logistic regression

**Idea:** transformation with the **logistic** function

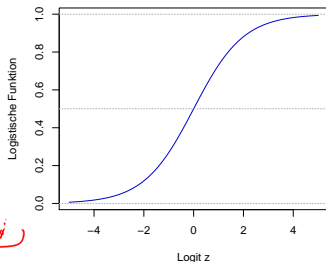
### Logistic model

$$P(Y_i = 1 | \mathbf{X}_i) = \frac{1}{1 + e^{-z_i}};$$

for **logits**  $z_i$  it holds

$$z_i = \beta_0 + \beta_1 X_{1i} + \dots + \beta_k X_{ki} + \cancel{\beta_{k+1} X_{k+1i}}$$

*assigned factor which determines your salary.*



**Note:** Alternatively we may use the CDF  $\Phi(z_i)$  of  $N(0, 1) \rightsquigarrow$  **probit**-model

$$y = \beta_0 + \beta_1 x_1 + \dots$$

$$+ \beta_j x_j + u_i \Rightarrow \text{cannot be visualized nicely}$$

**Note :** interpretation of the parameters is different than in case of LR

- The effect of a change in  $x_{ik}$  depends on  $x_i$  ab:

$$\frac{\partial P(Y_i=1)}{\partial x_{ik}} = \frac{\partial y_i}{\partial x_{ik}} = \frac{\partial F(x_i'\beta)}{\partial x_{ik}} = \frac{e^{x_i'\beta}}{(1 + e^{x_i'\beta})^2} \beta_k$$

- The direction of change in  $P(Y_i = 1)$  that arises due to a change in  $x_{ik}$  has the same sign as  $\beta_k$ .  
*Handwritten notes:*  $\beta_k > 0 \Rightarrow$  Prob. of involving  $\nearrow$  with  $x_{ik}$   
 $\beta_k < 0 \Rightarrow$  Prob. of involving  $\searrow$  with  $x_{ik}$ .
- If we have interactions, then it is even more difficult:

$$z_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i1} x_{i2} + \beta_3 x_{i2}$$

$$\frac{\partial F(x_i'\beta)}{\partial x_{i1}} = \frac{e^{z_i}}{(1 + e^{z_i})^2} \cdot (\beta_1 + \beta_2 x_{i2})$$

*Handwritten notes:*  $\beta_2 > 0$  or  $\beta_2 < 0 \Rightarrow$  even the signs are not informative.

- Odds** are helpful for interpretation (logit only!)

*Handwritten notes:* Odds = 2  $\Rightarrow$  Prob. of insurant is twice the prob. of being solvent  $\Rightarrow P(Y=1) = 2/3$   $P(Y=0) = 1/3$ .

$$\text{Odds} = \frac{P(Y=1|\mathbf{X})}{P(Y=0|\mathbf{X})} = e^z$$

	Logit ( $z$ )	Odds	$P(Y=1 \mathbf{X})$
$\beta > 0$	rises by $\beta$	rises by $e^\beta > 1$	rises
$\beta < 0$	falls by $\beta$	falls by $e^\beta < 1$	falls

## Estimation of the parameters: logit

The parameters are estimated using ML:

$$L = \prod_{i=1}^n \underbrace{\left( \frac{1}{1 + e^{-z_i}} \right)^{y_i}}_{P(Y_i=1)} \cdot \underbrace{\left( 1 - \frac{1}{1 + e^{-z_i}} \right)^{1-y_i}}_{P(Y_i=0)} \rightarrow \max, \text{ w.r.t. } \beta_0, \dots, \beta_k.$$

This simplifies to *to get rid of products and get sums.*

$$\frac{\partial \log L(\beta)}{\partial \beta} = \sum_{i=1}^n \left[ y_i - \frac{e^{z_i}}{1 + e^{z_i}} \right] \mathbf{x}_i = \sum_{i=1}^n \left[ y_i - \left( 1 - \frac{1}{1 + e^{z_i}} \right) \right] \mathbf{x}_i = \mathbf{0}$$

which leads to the 1st order conditions:

*system of nonlinear equations  
→ numerically*

$$\sum_{i=1}^n \hat{p}_i \mathbf{x}_i = \sum_{i=1}^n y_i \mathbf{x}_i$$

## Estimation of the parameters: probit

Again ML:

$$L(\boldsymbol{\beta}) = \prod_{i=1}^N \underbrace{\Phi(\mathbf{x}'_i \boldsymbol{\beta})}_{P(Y_i=1)}^{y_i} \cdot \underbrace{(1 - \Phi(\mathbf{x}'_i \boldsymbol{\beta}))}_{P(Y_i=0)}^{1-y_i} \longrightarrow \max, \text{ w.r.t. } \beta_0, \dots, \beta_K.$$

This simplifies to


$$\frac{\partial \log L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \sum_{i=1}^N \left[ y_i \frac{\phi(z_i)}{\Phi(z_i)} + (1 - y_i) \frac{-\phi(z_i)}{1 - \Phi(z_i)} \right] \mathbf{x}_i = \mathbf{0}.$$

It can be shown that (a general ML result):

$$\sqrt{N}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{d} \mathcal{N} \left( \mathbf{0}, E \left[ \frac{\partial^2 \log L}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} \right] \right)$$



## Note:

- In contrary to the LR the estimation is always numeric.
- Likelihood-Ratio tests can be used to check the significance of the parameters.
-  `glm(y ~ X,data=data, family=binomial(logit))`

**Example:** a data set with 700 observations *in LR!  $\hat{\sigma}^2(x|x)^{-1}$  here: Slide 439.*

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )	
(Intercept)	-1.434785	0.482326	-2.975	0.00293	**
debtinc	0.121391	0.019023	6.381	1.76e-10	***
employer	-0.161795	0.023742	-6.815	9.44e-12	***
debts	0.093460	0.045045	2.075	0.03801	*
age	-0.004397	0.014212	-0.309	0.75701	<i>&gt; 0.05 <math>\Rightarrow</math> irrelevant.</i>

*$\hat{\beta}_1 = 0.121 > 0 \Rightarrow$  if debtinc increases, then Prob. of involuntary increases.*

*$\hat{\beta}_2 = -0.16 < 0 \Rightarrow$  years with current employer increase (both chances of losing job)*

	debtinc	employer	debts	age	
$\hat{\beta}_i$	0.121*	-0.162*	0.093*	-0.004	<i><math>\Rightarrow</math> Prob of involuntary</i>
$e^{\hat{\beta}_i}$	1.129	0.851	1.098	0.996	<i><math>\Rightarrow</math> insignificant.</i>

(\*) - significant with  $\alpha = 0.05$

*$e^{\hat{\beta}_1} = 1.12 \Rightarrow$*

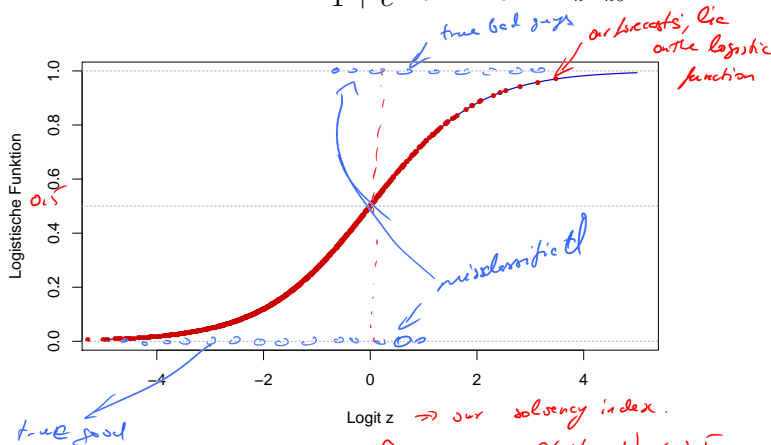
*2,26.*

$$\frac{P(Y=1 | \text{debtinc}+1) / P(Y=0 | \text{debtinc}+1)}{P(Y=1 | \text{debtinc}) / P(Y=0 | \text{debtinc})} = 1.13$$

*= 2*

## Forecasts:

$$P(\widehat{Y_0 = 1} | \mathbf{X_0}) = \hat{\pi}_0 = \frac{1}{1 + e^{-\hat{b}_0 - \hat{b}_1 X_{10} - \dots - \hat{b}_k X_{k0}}}.$$



$$z < 0 \Rightarrow \hat{Y} = 0 \Leftarrow P(Y=1) < 0.5$$

$$z > 0 \Rightarrow \hat{Y} = 1 \Leftarrow P(Y=1) > 0.5$$

# Goodness of the model

**Problem:** classical measures, such as  $R^2$ , cannot be used  $\rightsquigarrow$  pseudo- $R^2$ ; classification tables; graphical measures (ROC-curve)

$$P(Y=1) = \text{prob. of 1. without}$$

- pseudo- $R^2$ :

- Let  $LL_0$  be the Log-Likelihood of the null model ( $b_1 = \dots = b_k = 0$ )

log-lik without  
any features  
 $\Rightarrow$  prob. of getting  
the sample at hand.  
 $\left(\frac{N_1}{N}\right)^{N_1} \cdot \left(\frac{N_0}{N}\right)^{N_0}$

$$\begin{aligned} LL_0 &= \sum_{i=1}^N y_i \log(N_1/N) + \sum_{i=1}^N (1 - y_i) \log(1 - N_1/N) \\ &= N_1 \log(N_1/N) + N_0 \log(N_0/N) \end{aligned}$$

no features

$\swarrow$  our model.

- Let  $LL_v$  be the Log-Likelihood of the full model (with all variables)
- Let  $LL_s$  be the Log-Likelihood of the saturated model (model with perfect fit, here  $LL_s = 0$ )

Probab. = 1  $\Rightarrow$  log-prob = 0.

$LL_v \longrightarrow$  best model  $LL_s$   
 $LL_v \longrightarrow$  worst one  $LL_0$ .

- Deviance:  $D = -2 \cdot LL_v$  (close 0)
- Cox and Snell  $R^2$ :  $1 - (L_0/L_v)^{2/n}$
- McFaddens- $R^2$ :  $1 - LL_v/LL_0$  (starting from 0.4)

```

Null deviance: 804.36  on 699  degrees of freedom
Residual deviance: 626.49  on 695  degrees of freedom
> library("DescTools")
> PseudoR2(logit.model, which="all")
      McFadden      McFaddenAdj      CoxSnell      Nagelkerke
0.22113546      0.20870328      0.22438958      0.32849894
AldrichNelson VeallZimmermann      Efron
0.07170058      0.27698049      0.24041479
McKelveyZavoina      Tjur      AIC      BIC
0.38471667      0.24430962      636.49076003      659.24616170
logLik      logLik0      G2
-313.24538001      -402.18210244      177.87344485

```

→ not a very good model  $R^2 \approx 20-25\%$

# • Classification table

		predicted		
		$\hat{Y} = 1$	$\hat{Y} = 0$	
truth	1	$n_{11} = TP$	$n_{01} = FN$	$n_{.1} = P = n_{11} + n_{01}$
	0	$n_{10} = FP$	$n_{00} = TN$	$n_{.0} = N = n_{10} + n_{00}$
		$n_{1.}$	$n_{0.}$	

total number  
of truly bad  
guys



total number  
of truly  
good guys



Let  $\hat{y}_i = 1$  if  $P(\hat{Y}_i = 1 | \mathbf{X}_i) > 0.5$  and 0 else.

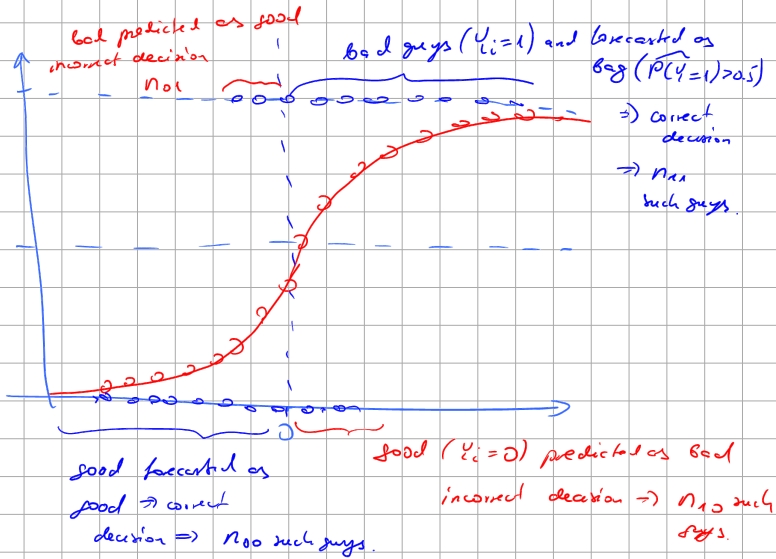
			$\hat{Y} = 1$	$\hat{Y} = 0$	
	$Y = 1$	$Y = 0$	72	111	
good guys also got no credit bank looks interest on the credit and good dia.			38	479	get money and do not pay back

$\leadsto (479+72)/700 = \underline{78,71\%}$  are correctly predicted.  $\Rightarrow$  in 78% of decisions  
no decision was correct!  
 $\Rightarrow$  losing race.

But: there are 73,86% solvent clients in the sample.  $\Rightarrow$  the above 73%  
look not so good  
now.

Question: is the threshold 0.5 a good choice?





Client case: boys 2 coins  $\left. \begin{array}{l} \rightarrow \text{Head Head} \\ \rightarrow \text{Head Tail} \\ \rightarrow \text{Tail Head} \\ \rightarrow \text{Tail Tail} \end{array} \right\} \Rightarrow \text{good}$

$\Rightarrow$  will give a correct classification in 25% of cases, because there are 74% of good clients.



## Goodness of the model and the choice of the threshold

- ROC (*receiver operating characteristics*), Lift and Gain curves are used to visualize and to quantify the goodness of the classification algorithms.

$$\begin{aligned}\text{sensitivity} &= \frac{n_{11}}{n_{.1}} = \frac{n_{11}}{n_{11} + n_{01}} \\ \text{specificity} &= \frac{n_{00}}{n_{.0}} = \frac{n_{00}}{n_{10} + n_{00}}\end{aligned}$$

**Sensitivity:** the fraction of correctly classified 1-values among all true 1-objects.

**Specificity:** the fraction of correctly classified 0-values among all true 0-object.

- Sensitivity =  $72 / (72 + 111) = 0.39$  - only 39% of insolvent clients are classified as insolvent
- Specificity =  $479 / (479 + 38) = 0.92$  - 92% of solvent clients are classified as solvent

*⇒ our bank is good in identifying good clients, but  
bad in identifying bad clients → really dangerous.*

$$\begin{aligned}\text{PPV or PV+} &= \frac{n_{11}}{n_{1.}} = \frac{n_{11}}{n_{11} + n_{10}} \\ \text{NPV or PV-} &= \frac{n_{00}}{n_{0.}} = \frac{n_{00}}{n_{01} + n_{00}}\end{aligned}$$

**PPV:** the fraction of correctly classified 1-values among all objects classified as 1.

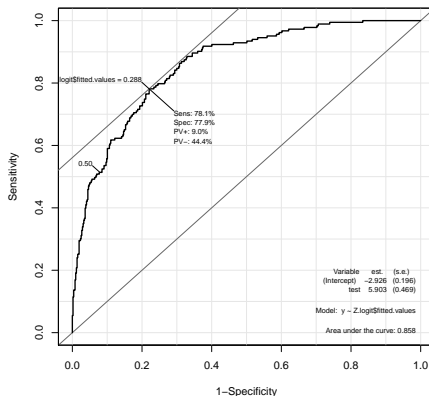
**NPV:** the fraction of correctly classified 0-values among all objects classified as 0.

(PPV-positive predicted value, NPV-negative predicted value)

- $\text{PPV} = 72 / (72 + 38) = 0.65$  - only 65% of all as insolvent classified clients are really insolvent
- $\text{NPV} = 479 / (479 + 111) = 0.81$  - 81% of all as solvent classified clients are really solvent

**ROC-curve:** sensitivity values as a function of specificity

- The steeper the function, the better the algorithm. **ROC-value** is the square under the curve.
- If the curve is close to the diagonal, then the algorithm is as good as random assignments.

**R:** roc-function from the pROC-package

Sensitivity

sens = 1

spec = 1

every bad as bad

every good as good

$\Rightarrow$  only correct decision.

"oracle" point

(0,0)

$$\text{sensitivity} = 0 = \frac{n_{11}}{n_{1.}}$$

$\Rightarrow \text{acc} = 0$

specificity = 1

$\Rightarrow$  no single bad as bad  
every good as good.

$\Rightarrow$  Strategy: every client is a good client  $\Rightarrow P(Y=1) > 1$ .

shortest distance

optimal threshold

sens = 1

spec = 0

every bad as bad

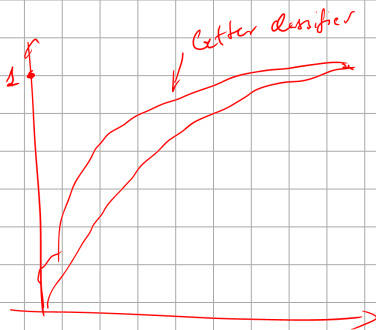
no single good as good

Strategy: every is a bad client

$$P(Y=1) \geq 0 \quad \text{instead of } 0.5$$

instead of 0.5

1 - specificity



compare the area under  
the curve.

Now let  $\hat{y}_i = 1$  if  $P(Y_i = 1|X_i) > 0.288$  and 0 else. *as before according to ROC.*

```
> confusionMatrix(as.numeric(Z.logit$fitted.values>0.288), y)
```

Confusion Matrix and Statistics

Reference

Prediction 0 1  
0 402 40  
1 115 143

Accuracy : 0.7786

95% CI : (0.746, 0.8088)

No Information Rate : 0.7386

P-Value [Acc > NIR] : 0.008201

Sensitivity : 0.7776

Specificity : 0.7814

Pos Pred Value : 0.9095

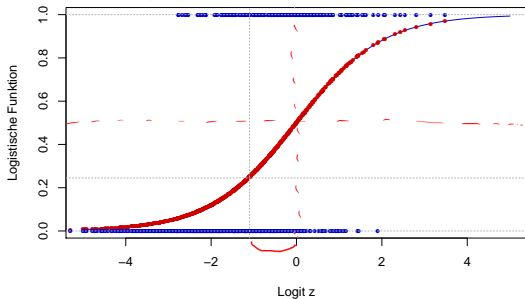
Neg Pred Value : 0.5543

*40 ⇒ bad as good (was 111) ⇒ good  
115 ⇒ good as bad (was 72) ⇒ bad*

*⇒ (402+143)/700. almost as before*

*⇒ much better ⇒ model is better in identifying bad clients.  
⇒ worse ⇒ worse in classifying good clients.*

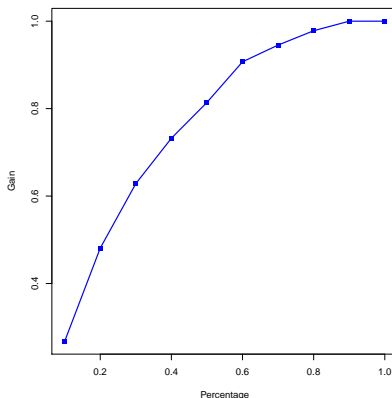
*The model is significantly better than doing a cut.*



## Gain-curve:

- If Gain equals 40% for 20%, this implies that if 20 % of the clients are classified as insolvent, then the algorithm will detect 40% of really insolvent clients .
- The diagonal shows a model-free classification: If 20 % of the clients are classified as insolvent, then the algorithm detects 20% of really insolvent clients.
- The steeper the curve, the better (with a single kink).

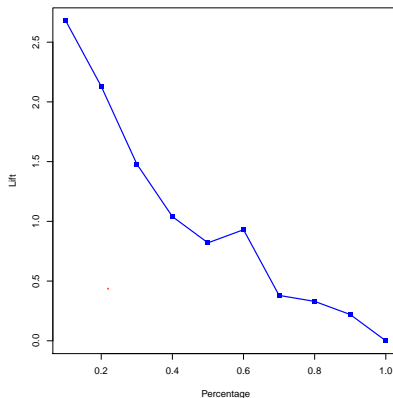
R: gains-package



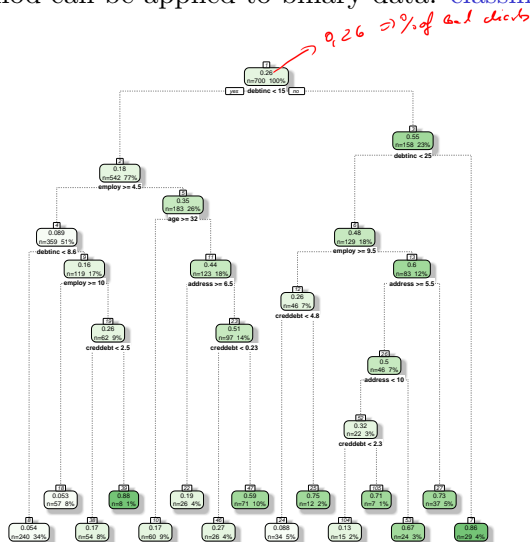


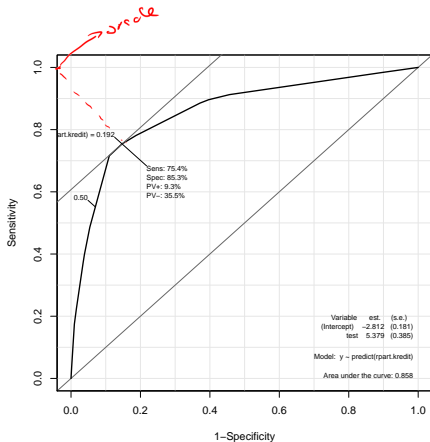
**Lift curve:** how much better is the predictive model compared to the model free classification?

- If lift is 2.1 for 20%, this implies that the model detection rate of insolvent clients is larger by 2.1 compared to model-free classification.



The CART method can be applied to binary data: **classification trees**





$p(Y=1) > 0,192 \rightarrow \text{Bad}$   
 $< 0,192 \rightarrow \text{Good}$

```
> confusionMatrix(as.numeric(predict(rpart.kredit)
>0.192), y)
```

Confusion Matrix and Statistics  
 Reference

Prediction	0	1
0	420	40
1	97	143

get the  
then logit.

Accuracy : 0.8043

95% CI : (0.7729, 0.8331)

No Information Rate : 0.7386

P-Value [Acc > NIR] : 2.822e-05

Kappa : 0.5395

Mcnemar's Test P-Value : 1.715e-06

Sensitivity : 0.8124

Specificity : 0.7814

Pos Pred Value : 0.9130

Neg Pred Value : 0.5958

Test of

$H_0: n_{01} = n_{10} = 0$

$\Rightarrow$  "boosting"

## Modelling nominal data

### Practical question:

- Choice of the political party depending of the characteristics of the voters;
- Choice of a product brand depending on the characteristics of the client;

### Example:

mode	–	“car”, “air”, “train”, oder “bus”
choice	–	decision
wait	–	waiting time, 0 for “car”
vcost	–	variable costs
travel	–	time
gcost	–	total costs
income	–	income
size	–	number of persons

	individual	mode	choice	wait	vcost	travel	gcost	income	size
1	1	air	no	69	59	100	70	35	1
2	1	train	no	34	31	372	71	35	1
3	1	bus	no	35	25	417	70	35	1
4	1	car	yes	0	10	180	30	35	1
5	2	air	no	64	58	68	68	30	2
6	2	train	no	44	31	354	84	30	2

*any car  
wife's car*

## Multinomial logit model

For the simple logit model it holds:

$$P(Y = 1|\mathbf{x}) = \frac{\exp(\mathbf{x}'\beta)}{1 + \exp(\mathbf{x}'\beta)} \quad \div \quad \frac{1}{1 + e^{-\mathbf{x}'\beta}}$$

*odd =  $e^z$*

$$\ln \left( \frac{P(Y = 0|\mathbf{x})}{P(Y = 1|\mathbf{x})} \right) = \mathbf{x}'\beta$$

For the  $k$  categories of  $Y$  we define:

*$\ln(\text{odd})$   
as above*

$$\ln \left( \frac{P(Y = r|\mathbf{x})}{P(Y = k|\mathbf{x})} \right) = \mathbf{x}'\beta_r, \quad r = 1, \dots, k-1$$

*category specific vector of parameters.*

*reference category*

with

$$P(Y = r|\mathbf{x}) = \frac{\exp(\mathbf{x}'\beta_r)}{1 + \sum_{s=1}^{k-1} \exp(\mathbf{x}'\beta_s)}, \quad r = 1, \dots, k-1$$

$$P(Y = k|\mathbf{x}) = \frac{1}{1 + \sum_{s=1}^{k-1} \exp(\mathbf{x}'\beta_s)}.$$

One category, i.e. the  $k$ -th, is the reference category.

**Note:**

- Estimation via ML assuming independence of the observations. This is a questionable assumption:
  - similar categories;
  - odds do not depend on other categories, etc.
  - Solution: Hausmann/McFadden test
- Goodness-of-fit, tests as for logit.

## Global and category specific variables

$$x'_r \beta_r \mapsto \underbrace{x'_{glob} \beta_r^*}_{\text{category specific.}} + \underbrace{x'_{spec,r} \alpha}_{\text{global}}$$

- Global variables (**income, number of persons**) do not depend on the categories and have individual parameters for each category:  $x'_{glob} \beta_r^*$ .  
The sign of the parameters cannot be interpreted.
- The category specific variables (**waiting time, costs**) depend on the categories and are evaluated relatively to the reference category.

$$(x_{spec,r} - x_{spec,k})' \alpha \quad \text{or} \quad x'_{spec,r} \alpha$$

The sign of the parameters can be interpreted.



Let *gcost* and *wait* be category specific and *income* and *size* are global variables. The reference category is *air*.

```
> library("mlogit")
> mlogit(choice~wait+gcost|income+size, ...)
```

Coefficients :

	Estimate	Std. Error	t-value	Pr(> t )	
train:(intercept)	-2.3115942	0.7525161	-3.0718	0.0021276	**
bus:(intercept)	-3.4504941	0.9064886	-3.8064	0.0001410	***
car:(intercept)	-7.8913907	0.9880615	-7.9867	1.332e-15	***
wait	-0.1013180	0.0112207	-9.0296	< 2.2e-16	***
gcost	-0.0197064	0.0053844	-3.6599	0.0002523	***
train:income	-0.0589804	0.0154532	-3.8167	0.0001352	***
bus:income	-0.0277037	0.0169812	-1.6314	0.1027991	
car:income	-0.0041153	0.0127301	-0.3233	0.7464866	
train:size	1.3289497	0.3141683	4.2301	2.336e-05	***
bus:size	1.0090796	0.3952899	2.5528	0.0106874	*
car:size	1.0392585	0.2665513	3.8989	9.663e-05	***
---					

Log-Likelihood: -176.77

McFadden R<sup>2</sup>: 0.37705

Likelihood ratio test : chisq = 213.98 (p.value = < 2.22e-16)

$P(Y = \text{train})$   $\rightarrow$  if income  $\uparrow$   
 $P(Y = \text{air})$   
 $\Rightarrow$  with higher income there is  
 a higher chance of  
 taking a plane than  
 a train.

more people  $\Rightarrow$   
 higher chances for  
 train, car, bus.

With the estimated parameters we can estimate the probabilities  $P(Y_i = r | \mathbf{x}_i)$  for all  $r$ .

	air	train	bus	car
[1,]	0.2368302	0.00000000	0.24496423	0.5182056
[2,]	0.2083323	0.27785076	0.00000000	0.5138170
[3,]	0.0000000	0.12686485	0.23058033	0.6425548
[4,]	0.1151004	0.05063597	0.02141839	0.8128452
[5,]	0.3405917	0.20694648	0.05624436	0.3962174
[6,]	0.1316850	0.36965292	0.26144217	0.2372200

largest prob  
→ car

train.

## Ordered response models

**Aim:**  $Y$  can take  $M$  different ordered (!) values (credit ratings, grades, income classes, etc)

up to now:  $u_i^k > 0 \Rightarrow y_i = k$   
 $u_i^k < 0 \Rightarrow y_i = 0$

Using a single latent variable we can specify



$$Y_i^* = \mathbf{x}_i' \boldsymbol{\beta} + u_i$$

$$Y_i = j \quad \text{if} \quad \gamma_{j-1} < Y_i^* \leq \gamma_j$$

for some unknown threshold values  $\gamma_j$  with  $\gamma_0 = -\infty$  and  $\gamma_M = \infty$ .



Assuming logistic cdf for  $u_i$  we obtain **ordered logit model** and assuming normality we obtain **ordered probit model**.

**Example:** a (simplified) rating of companies -  $Y = 1$  - lowest,  $Y = 2$  - average,  $Y = 3$  - highest

MARKET_VALUE	DIV_PER_SHR	TOTAL_DEBT	rating
-0.08911931	-0.08063048	-0.02276501	1
-0.09350059	-0.08148114	-0.14012151	3
-0.09452652	-0.08019333	-0.13869093	2
-0.09633656	-0.08090222	-0.13784721	2
-0.09254201	-0.08130392	-0.13822051	2
0.95192265	-0.06091170	13.33345544	3

Standardized.

without intercept.

in logit only  $\beta$ 's

Here:  $\beta$ 's and  $\gamma$ 's.

$$Y^* = \mathbf{x}'\boldsymbol{\beta} + u$$

$$Y = 1 \quad \text{if} \quad y^* \leq \gamma_{1|2}$$

$$= 2 \quad \text{if} \quad \gamma_{1|2} < y^* \leq \gamma_{2|3}$$

$$= 3 \quad \text{if} \quad \gamma_{2|3} < y^*$$

Assuming normal error terms we can state the corresponding prob's

$$\begin{aligned}
 P(Y_i \leq k | \mathbf{x}_i) &= P(Y_i^* \leq \gamma_{(k-1)|k} | \mathbf{x}_i) = \Phi(\gamma_{(k-1)|k} - \mathbf{x}_i' \boldsymbol{\beta}) \\
 P(Y_i = 1 | \mathbf{x}_i) &= P(Y_i^* \leq \gamma_{1|2} | \mathbf{x}_i) = \Phi(\gamma_{1|2} - \mathbf{x}_i' \boldsymbol{\beta}) \\
 P(Y_i = 3 | \mathbf{x}_i) &= P(Y_i^* > \gamma_{2|3} | \mathbf{x}_i) = 1 - \Phi(\gamma_{2|3} - \mathbf{x}_i' \boldsymbol{\beta}) \\
 P(Y_i = 2 | \mathbf{x}_i) &= P(\gamma_{1|2} < Y_i^* \leq \gamma_{2|3} | \mathbf{x}_i) = \Phi(\gamma_{2|3} - \mathbf{x}_i' \boldsymbol{\beta}) - \Phi(\gamma_{1|2} - \mathbf{x}_i' \boldsymbol{\beta})
 \end{aligned}$$

*Handwritten notes:*  $\mathbf{x}_i' \boldsymbol{\beta} + u \leq \gamma \Rightarrow u \leq \gamma - \mathbf{x}_i' \boldsymbol{\beta}$  (circled in red, with arrows pointing to the equations above)

The log-likelihood function is then given by

$$LL(\boldsymbol{\beta} | \mathbf{X}) = \sum_{i=1}^N \log P(Y_i = y_i | \mathbf{x}_i) \rightarrow \max, \text{ w.r.t. } \boldsymbol{\beta}, \text{ s.t.}$$

The inferences follow in a similar fashion as for the simple logit.

*For logit:*  $\prod \left( \frac{1}{1 + e^{-x}} \right) \cdot \left( 1 - \frac{1}{1 + e^{-x}} \right)$

## Example:

Coefficients:

	Value	Std. Error	t value
MARKET_VALUE	2.14118	0.6803	3.1474
DIV_PER_SHR	-0.70836	0.3189	-2.2212
TOTAL_DEBT	0.05553	0.1367	0.4063

Intercepts:

	Value	Std. Error	t value
1 2	-0.0309	0.0789	-0.3916
2 3	1.0181	0.0879	11.5797

Residual Deviance: 1480.408

AIC: 1490.408

```
> ordlog.pred = predict(ordlog, type="probs")
```

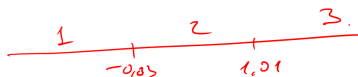
```
> ordlog.pred
```

	1	2	3
5.259901e-01	2.340771e-01	0.2399328	
5.298021e-01	2.330434e-01	0.2371545	
5.305567e-01	2.328364e-01	0.2366069	
5.313852e-01	2.326083e-01	0.2360065	
5.292958e-01	2.331818e-01	0.2375224	
5.453970e-02	8.685558e-02	0.8586047	

$\beta_1 = 2.14 > 0 \Rightarrow$  increasing market share increases the probability to move to a better category.

$\beta_2 = -0.7 \Rightarrow$  high dividend indicate a lower quality company.

$\rightarrow$  small  $\Rightarrow$  insignificant



1st category

$\rightarrow$  the 3rd category.

## Modelling for count data

### Practical questions:

- the number of claims by an insurance company per time period;
- the number of consultations by a doctor per year ;
- the number of insolvent companies per time period;
- occurrences of a seldom disease per season;
- ....

**Note:** the modelling is particularly important for small values of the target variable (rare events) and the distribution is heavily skewed.

# Poisson distribution

The Poisson distribution is frequently used to model rare events

*# of events*

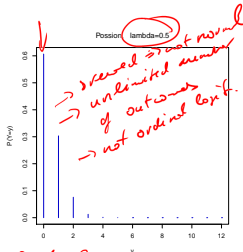
$$P(Y = y) = \begin{cases} \frac{\lambda^y}{y!} e^{-\lambda}, & \text{for } y = 0, 1, 2, \dots \\ 0, & \text{else,} \end{cases}$$

with the **intensity parameter**  $\lambda$ . It fulfils the *equidispersion*-condition:

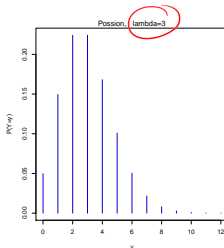
$$E(Y) = Var(Y) = \lambda$$

*expected # of events per period of time*

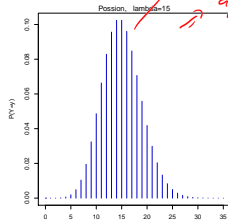
*in 60% of cases you don't go to a doctor*



*0, 0, 1, 0, 0, 1, 2, 0, 0, ...*



*2, 3, 4, 3, 3, 0, ...*



*0, 14, 15, 20, 10, 15*

*live normal  
→ apply LR*



## Poisson regression model

Let  $Y_i, \mathbf{x}_i$  be independent realisations, while  $Y_i$  follows Poisson distribution with

$$E(Y_i | \mathbf{x}_i) = h(\mathbf{x}_i' \boldsymbol{\beta}) = \exp(\mathbf{x}_i' \boldsymbol{\beta}) = \lambda_i.$$

*expected number of visits for person i depends on individual characteristics in  $\mathbf{x}_i$ .*

*> 0*

- The interpretation of the parameters follows as for the logit model.
- The parameters are estimated via ML:

$$LL(\boldsymbol{\beta}) = \sum_{i=1}^n y_i \ln(h(\mathbf{x}_i' \boldsymbol{\beta})) - h(\mathbf{x}_i' \boldsymbol{\beta}) - \ln(y_i!) \longrightarrow \max, \text{ w.r.t. } \boldsymbol{\beta}$$

$$\ln \left( \frac{(h(\mathbf{x}_i' \boldsymbol{\beta}))^{y_i} \cdot e^{-h(\mathbf{x}_i' \boldsymbol{\beta})}}{y_i!} \right)$$

## Goodness of the model

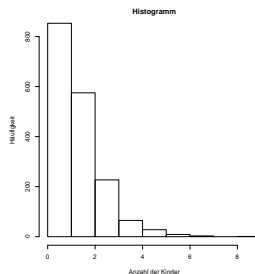
To measure the goodness of the model we use **deviance**, i.e. the difference between the log-likelihood for the actual observations (**perfect/saturiertes model**) and the log-likelihood for the predicted values:

$$D = -2 \sum_{i=1}^n [LL_i(\hat{Y}_i) - LL_i(Y_i)] = 2 \sum_{i=1}^n \left[ Y_i \ln(Y_i / \hat{\lambda}_i) \right] \sim \chi_{n-p}^2$$

## Example: number of children

- child - number of children
- age - age of the woman
- dur - years at school/college
- nation - nationality, 0 = german , 1 = else
- god - trust in God: 1 = strong, ..., 6 = never thought about it
- univ - university degree: 0 = no, 1 = yes

```
mean(children$child)
[1] 1.57297
> var(children$child)
[1] 1.552769
```



```
glm(formula = child ~ age + I(age^2) + I(age^3) + I(age^4) +  
    dur + I(dur^2) + nation + god + univ, family = poisson(link = log),  
    data = children)
```

Poisson regression.

$\exp(x'\beta)$

Deviance Residuals:

Min	1Q	Median	3Q	Max
-2.1514	-0.7559	0.0102	0.4832	3.6715

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-1.228e+01	1.484e+00	-8.277	< 2e-16 ***
age	9.359e-01	1.239e-01	7.553	4.26e-14 ***
I(age^2)	-2.490e-02	3.786e-03	-6.577	4.80e-11 ***
I(age^3)	2.842e-04	4.915e-05	5.781	7.42e-09 ***
I(age^4)	-1.180e-06	2.297e-07	-5.137	2.80e-07 ***
dur	1.118e-01	6.652e-02	1.680	0.092904 .
I(dur^2)	-8.328e-03	2.997e-03	-2.779	0.005454 **
nation1	5.686e-02	1.386e-01	0.410	0.681599
god2	-1.025e-01	5.903e-02	-1.736	0.082599 .
god3	-1.448e-01	6.780e-02	-2.136	0.032683 *
god4	-1.279e-01	7.088e-02	-1.805	0.071128 .
god5	-3.621e-02	6.695e-02	-0.541	0.588569
god6	-9.241e-02	7.505e-02	-1.231	0.218239
univ1	6.372e-01	1.729e-01	3.686	0.000228 ***

$$E(\# \text{ of kids}) = \exp\{\beta_0 + \beta_1 \cdot \text{age} + \beta_2 \cdot \text{age}^2 + \beta_3 \cdot \text{dur} + \beta_4 \cdot \text{dur}^2 + \beta_5 \cdot I(\text{god} = 2) + \beta_6 \cdot I(\text{god} = 3) + \dots + \beta_{10} \cdot I(\text{god} = 6) + \dots\}$$

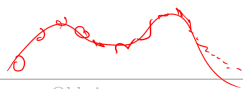
test using GLH

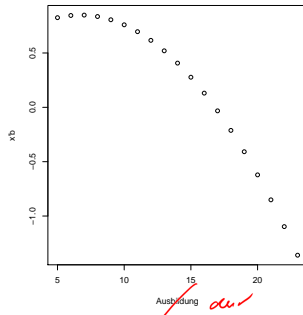
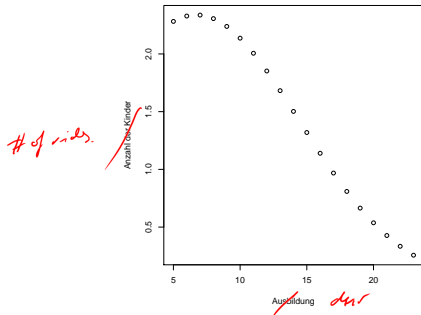
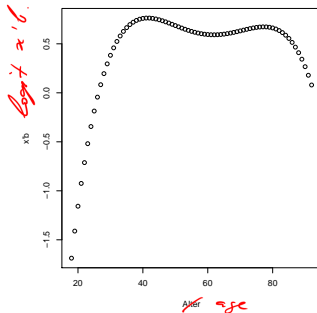
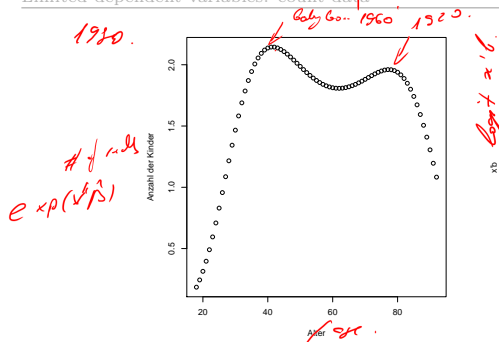
(Dispersion parameter for poisson family taken to be 1)

Null deviance: 2067.4 on 1760 degrees of freedom  
Residual deviance: 1718.6 on 1747 degrees of freedom  
AIC: 5196.8

$$y = \dots \beta \cdot \text{God} + \dots$$

will





**Note:** for the Poisson distribution it should hold  $E(Y_i) = Var(Y_i) = \lambda_i$ .

If this assumption is not fulfilled then we have  
*overdispersion/underdispersion*.

**Solution:** as an alternative we can use **Quasi-Poisson-** or the **negative binomial distribution (negbin)**. Both distributions allow for different expectations and variances.

For **negbin** it holds:

$$P(Y_i|\mathbf{x}_i) = \frac{\Gamma(Y_i + \nu)}{\Gamma(\nu)\Gamma(Y_i + 1)} \cdot \left(\frac{\lambda_i}{\lambda_i + \nu}\right)^{Y_i} \cdot \left(\frac{\nu}{\lambda_i + \nu}\right)^{\nu}$$

with  $E(Y_i) = \lambda_i = \exp(\mathbf{x}_i'\boldsymbol{\beta})$  and  $Var(Y_i) = \lambda_i + \lambda_i^2/\nu$ .

```
glm(formula = child ~ age + I(age^2) + I(age^3) + I(age^4) +
     dur + I(dur^2) + nation + god + univ, family = negative.binomial(theta = 1,
     link = log), data = children)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-1.56820	-0.50984	-0.01054	0.29990	1.90633

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-1.338e+01	1.267e+00	-10.555	< 2e-16 ***
age	1.022e+00	1.075e-01	9.502	< 2e-16 ***
I(age^2)	-2.730e-02	3.342e-03	-8.169	5.90e-16 ***
I(age^3)	3.126e-04	4.395e-05	7.113	1.65e-12 ***
I(age^4)	-1.302e-06	2.074e-07	-6.277	4.34e-10 ***
dur	1.269e-01	5.990e-02	2.118	0.034294 *
I(dur^2)	-9.577e-03	2.637e-03	-3.632	0.000289 ***
nation1	8.309e-02	1.349e-01	0.616	0.538128
god2	-1.186e-01	5.849e-02	-2.028	0.042743 *
god3	-1.681e-01	6.642e-02	-2.530	0.011483 *
god4	-1.563e-01	6.923e-02	-2.258	0.024075 *
god5	-3.273e-02	6.602e-02	-0.496	0.620135
god6	-1.205e-01	7.384e-02	-1.632	0.102848
univ1	7.749e-01	1.581e-01	4.900	1.04e-06 ***

---

(Dispersion parameter for Negative Binomial(1) family taken to be 0.3516262)

Null deviance: 1023.1 on 1760 degrees of freedom  
 Residual deviance: 852.3 on 1747 degrees of freedom  
 AIC: 5911.9

related to  $\lambda \Rightarrow$  allow  
 $E$  and  $Var$  to  
 be different.