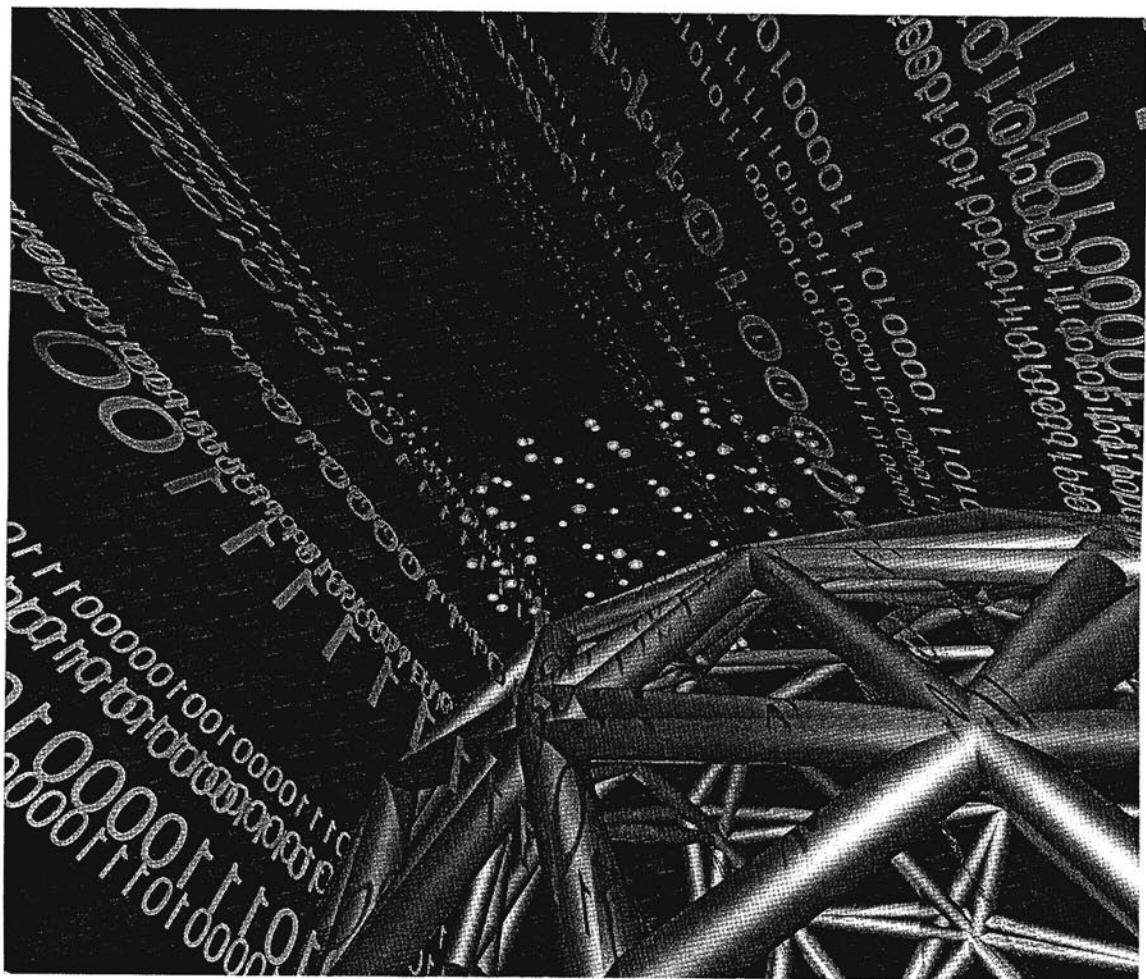


Celostátní seminář – sborník přednášek

Analýza dat 2006/II Statistické metody pro praxi

Lázně Bohdaneč 24. - 27. 10. 2006



*Pokročilé statistické metody
pro řízení jakosti, technologickou, výzkumnou a zkušební praxi*



TriloByte Statistical Software

CQR Výzkumné centrum pro jakost a
spolehlivost



Editor: Karel Kupka

Celostátní seminář – sborník přednášek
Analýza dat 2006/II
*(Pokročilé statistické metody
pro řízení jakosti, technologickou, výzkumnou a zkušební praxi)*

Lázně Bohdaneč 24. - 27. 10. 2006

Vydal (c) 2007



TriloByte Statistical Software, Jiráskova 21
530 02 Pardubice, Czech Republic
info@trilobyte.cz, <http://www.trilobyte.cz>

ISBN 978-80-239-8998-4

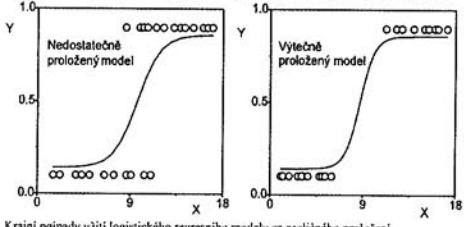
Základy logistické regrese

Milan Meloun
Univerzita Pardubice

1. Zaměření metody LR

Navržena v 60tých letech jako alternativní postup k MNČ, když je závisle proměnná binární (medicina), kde y představuje přítomnost (1) nebo nepřítomnosti (0) choroby.

Jde o klasifikaci, když není splněna normalita vícerozměrného modelu.



Krajní případ užití logistického regresního modelu za rozlišeného proložení.

PŘÍKLAD 8.2 Významnost sledovaných znaků ovlivňujících ústup leukemie
Lee (1980) publikoval data o leukemii pacientů a o ovlivnění jejího ústupu. Cílem je nalézt znaky, které jsou v navrženém logistickém regresním modelu statisticky významné k ovlivnění ústupu leukemie.

o Data: Pacientů (řádků) $n = 27$ a znaků (sloupců) $m = 7$.

Závisle proměnnou je REMISS: zda se objeví (1) či neobjeví (0) ústup leukemie.

Nezávisle proměnnými jsou u pacientů naměřené hodnoty 6 znaků:

CELL značí celulritu, buněčnost sraženiny kostní dřeně,
SMEAR značí skvrnu diferenčního procenta napadení,
INFIL značí procento infiltrátu kostní dřeně buňkou leukemie,
LI je procento označeného indexu leukemicích buněk kostní dřeně,
BLAST je absolutní počet napadení v periferní krvi,
TEMP značí nejvyšší naměřenou teplotu před začátkem léčby.

Index	REMISS	CELL	SMEAR	INFIL	LI	BLAST	TEMP
1	1	0.8	0.83	0.66	1.9	1.1	0.996
..
27	0	1	0.73	0.73	0.7	0.398	0.986

o Auteři: Byly užity programy NCSS2000 [67], MINITAB [86] a STATISTICA [102].

Podmínky výpočtu:

Závisle proměnná: REMISS

Model obsahuje 6 nezávisle proměnných: CELL | SMEAR | INFIL | LI | BLAST | TEMP.

Objektů (řádků): 27

Znaků (sloupců): 7

Volba proměnných a výstavba modelu:

Třída závisle proměnnou může nabývat hodnot, 0 a 1.

Počet je součet četnosti závisle proměnné pro každou třídu.

Řádky přináší počet objektů v každé třídě tak, jak byl vypočten z hodnot nezávisle proměnných.

Prior a priori pravděpodobnost v každé třídě zadána předem uživatelem.

Aktuálně vs. Predikce, R^2 značí hodnotu R^2 , kterou obdržíme z regrese, kde závisle proměnná ve třídě bude lineární funkcií predikované pravděpodobnost této třídy.

% Správně klasifikováno přináší procento objektů v této třídě, které byly správně klasifikovány logistickým regresním modelem.

REMISS		Aktuálně		% Správně	
Třída	Počet	Řádky	Prior	vs. Predikce, R^2 klasifikováno	
0	18	18	0.66667	0.38775	83.333
1	9	9	0.33333	0.38775	55.556
Total	27	27			74.074

Iterační přiblížení odhadů parametrů:

zvyšuje logit metodou maximální věrohodnosti:

0. iter.	1. iter.	2. iter.	3. iter.	4. iter.	5. iter.	6. iter.	7. iter.	
Úsek	0.0000	-18.7941	-43.8181	-64.4856	-66.9072	-59.9205	-58.1082	-58.0387
LI	0.0000	0.7895	0.4785	-1.7390	-8.4324	-19.9753	-24.3896	-24.6605
CELL	0.0000	5.8190	7.8289	6.5724	-0.7276	-14.0188	-18.9941	-19.2925
TEMP	0.0000	-5.9853	-8.4942	-7.5008	0.0725	14.0932	19.2909	19.6001
SMEAR	0.0000	-2.1173	-2.8763	-3.3905	-3.6410	-3.7879	-3.8874	-3.8959
INFIL	0.0000	-0.0030	-0.0152	-0.0980	-0.1186	-0.1456	-0.1513	-0.1511
BLAST	0.0000	20.6740	47.3696	71.2305	80.4766	84.6743	87.2300	87.4331
Logit	-18.7150	-12.3083	-11.3620	-11.0574	-10.9449	-10.8812	-10.8753	-10.8753

Logit představuje rozhodčí kritérium, pomocí kterého se rozhodne, zda důtyčný parametr logistický model zlepší nebo zhorší.

Průběh postupného krokového zavádění parametrů do logistického regresního modelu

R^2 modelu je odhad koeficientu determinace R^2 pro logistický model proložený daty. Změna v R^2 udává hodnotu, která se přidá k celkovému R^2 , když se tento parametr přidá do modelu. Změny R^2 mají vesměs kladné znaménko, zavedení dalšího parametru do modelu způsobí zvýšení hodnoty R^2 čili zlepšení modelu až na konečnou 36.719 %.

R^2 je však pouze přibližná hodnota dle vzorce
 $R^2 = \chi^2(df)/[\chi^2(df) + n - p - 1]$,

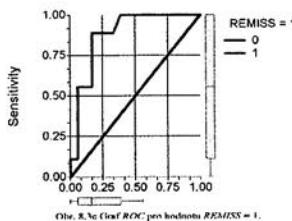
kde df značí stupně volnosti u χ^2 -testu a jsou rovny počtu nezávisle proměnných.

Vysvětlení: Z počátečního logitu $\ln L_{(0)} = -17.18588$
až do terminace $\ln L_{(1)} = -10.87533$ pro REMISS = 0.

Krok výstavby	Parametr zaveden	Logit	Dosažená R^2	Změna v R^2
1	Úsek	-17.18588	0.00000	0.00000
2	LI	-13.03648	0.24144	0.24144
3	CELL	-12.17036	0.29184	0.05040
4	TEMP	-10.97669	0.36130	0.06946
5	SMEAR	-10.92900	0.36407	0.00277
6	INFIL	-10.87752	0.36707	0.00300
7	BLAST	-10.87533	0.36719	0.00013

Graf prahové operační charakteristiky ROC
vystihující správnost diagnostického testu,
zda logistickým modelem vypočtené Ano nebo Ne je správné.

Na ose y se vynáší *sensitivity* a na ose x hodnota „1 minus specificity“. Křivky pomohou nalézt hodnotu dělicího bodu P_c ke klasifikaci objektů.



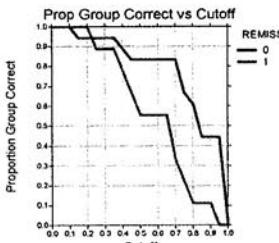
Když leží dělicí bod P_c v levém horním rohu grafu, je dosaženo nejvhodnější hodnoty a dochází k nejlepší klasifikaci objektů.

Kvalita nalezeného logistického modelu se posuzuje také dle plochy AUC pod křivkou ROC. Čím více se AUC blíží jedné nebo 100 %, tím je klasifikace objektů lepší.

Graf podílu správně zařazených objektů v závislosti na P_c

je velice užitečný graf k určení nejlepší hodnoty prahového dělicího bodu P_c .

Na ose y je procento správně zařazených objektů a na ose x hodnoty pravděpodobnosti dělicího bodu P_c v jednotkách vyčíslované pravděpodobnosti.



Obr. 8.3d Graf závislosti obou křivek podílu správně zařazených objektů na P_c (cutoff).

Nalezený logistický regresní model:

Dle statistické významnosti odhadů parametrů byl stanoven logistický regresní model pro $REMISS = 0$:

$$\begin{aligned} & -58.04 - 24.66 \text{ CELL} + 19.60 \text{ INFIL} - 3.90 \text{ LI} \\ & - 19.29 \text{ SMEAR} + 87.43 \text{ TEMP} + 0.15 \text{ BLAST}. \end{aligned}$$

○ Závěr:

1. Byl nalezen logistický regresní model znaků *CELL*, *SMEAR*, *INFIL*, *LI*, *BLAST*, a *TEMP*, které významně ovlivňují znak ústupu leukemie *REMISS*.

2. Znaky *LI* a *BLAST* nejsou statisticky významné.

3. Z ROC byl odhadnut prahový dělicí bod pravděpodobnosti, dle kterého se objekty spolehlivě zařadí do dvou tříd ústupu a neústupu leukemic.

Logitová transformace vychází z poměru šancí či naděje.

Dle typu závisle proměnné y se rozlišují:

Binární logistická regrese: binární závisle proměnná nabývá pouze dvou hodnot, například přítomnost-absence, muž-žena. Vektor nezávisle proměnných x obsahuje jednu či více spojitých proměnných (*prediktory*) nebo diskrétních, kategorických (*faktory*).

Ordinální logistická regrese: ordinální závisle proměnná nabývá tří a více možných stavů, např. silný nesouhlas, nesouhlas, souhlas, silný souhlas. Vektor x nezávisle proměnných obsahuje jak *prediktory* tak i *faktory*.

Nominální logistická regrese: nominální závisle proměnná o více než třech úrovních, např. mezi kterými je definována pouze odlišnost. Vektor x může obsahovat jak *prediktory*, tak i *faktory*.

Logistická regrese LR se liší od lineární regrese:

predikuje pravděpodobnost události, která se buď stala (1) nebo nestala (0).

Logitová transformace vede na sigmoidální vztah mezi závisle proměnnou y a vektorem nezávisle proměnných x.

Při velmi nízkých hodnotách x se pravděpodobnost proměnné y blíží k nule. Při vysokých hodnotách x se blíží k jedné.

Logistická regrese používá *kategorickou závisle proměnnou* zatímco lineární regrese užívá pouze *spojitou vysvětlovanou proměnnou*.

2. Logistický regresní model

V LR potřebujeme vědět, zda se událost stala (1) nebo nestala (0).

Jde o dichotomickou hodnotu 0 - 1 závisle proměnné y, ze které se predikuje odhad pravděpodobnosti, že se událost stala (1) či nestala (0).

Je-li predikovaná pravděpodobnost větší než 0,50, pak se událost stala (1), je-li menší než 0,50, pak se nestala (0).

Postup LR porovnává pravděpodobnost události odehrané $L_{(1)}$ vůči pravděpodobnosti události neodehrané $L_{(0)} = 1 - L_{(1)}$.

Využijeme *pravděpodobnostní poměr* $L_{(1)}/L_{(0)}$, ve kterém pravděpodobnost $L_{(1)}$ je vyjádřena logistickou funkcí

$$L_{(1)} = \frac{1}{1 + e^{C - z}}$$

vděpodobnostní poměr (zvaný "poměr šancí") je vyjádřen

$$\frac{L_{(1)}}{L_{(0)}} = e^{a_0 + a_1 x_1 + a_2 x_2 + \dots + a_p x_p}$$

odhadované koeficienty $a_0, a_1, a_2, \dots, a_p$ jsou míry změny poměru obou pravděpodobnosti $L_{(1)}/L_{(0)}$.

něr je lineární funkčí diskriminační funkce o p nezávisle proměnných

$$Z = a_0 + a_1 x_1 + a_2 x_2 + \dots + a_p x_p$$

zlogaritmování a úpravč vyjde

$$C - Z = \ln \left(\frac{L_{(1)}}{L_{(0)}} \right)$$

je C je absolutní člen a_0 .

Dle klasifikačního postupu je

$$\begin{aligned} L_{(0)} &= P(G = 1 \mid x) \\ L_{(1)} &= P(G = 0 \mid x) = 1 - P(G = 1 \mid x). \end{aligned}$$

a po úpravách bude

$$\ln \left(\frac{L_{(1)}}{L_{(0)}} \right) = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_p x_p$$

kde $b_0 = -C + a_0$, $b_i = a_i$ pro $i = 1, \dots, p$.

Například: ve sportu řekneme, že tým má šanci 3:1. Tvrzení říká, že favorizovaný tým má pravděpodobnost vítězství $\frac{3}{3+1} = \frac{3}{4} = 0.75$.

$$\text{Platí tedy pravděpodobnostní poměr } \frac{L_{(1)}}{L_{(0)}} = \frac{0.75}{1 - 0.75} = \frac{3}{1}.$$

posteriorní pravděpodobnost $P(G = j \mid x)$ zařazení do j -té kategorie: logistický model lze rozšířit na případ K tříd, a předpokládat, že posteriorní pravděpodobnost $P(G = j \mid x)$ zařazení do j -té kategorie bude

$$\ln \frac{P(G = 1 \mid x)}{P(G = K \mid x)} = b_{1,0} + b_1^T x$$

$$\ln \frac{P(G = 2 \mid x)}{P(G = K \mid x)} = b_{2,0} + b_2^T x$$

...

$$\ln \frac{P(G = K-1 \mid x)}{P(G = K \mid x)} = b_{K-1,0} + b_K^T x$$

Po zpětné transformaci vychází

$$P(G = j \mid x) = \frac{\exp(b_{j,0} + b_j^T x)}{1 + \sum_{l=1}^{K-1} \exp(b_{l,0} + b_l^T x)}$$

$$a \quad P(G = K \mid x) = \frac{1}{1 + \sum_{l=1}^{K-1} \exp(b_{l,0} + b_l^T x)}$$

Označíme pravděpodobnost

$$P(G = K \mid x) = p_k(x, b),$$

aby se zvýraznilo, že jde o funkci regresních parametrů

$$b = [b_{1,0}, b_1, b_{2,0}, b_2, \dots, b_{K-1,0}, b_{K-1}]$$

(Pro $K = 2$ přechází tento model na standardní logistický model pro binární vysvětlovanou proměnnou $y = G$).

Odhady parametrů:

Metoda odhadu parametrů:

Pro odhad parametrů logistických modelů se používá metoda maximální věrohodnosti.

Přítomnost v první třídě $y = 1$ je pro $G = 1$.

Nepřítomnost v první třídě $y = 0$ je pro $G = 2$ čili přítomnost ve druhé třídě.

Výchozí data: vektor y rozměru $n \times 1$ a matici X rozměru $n \times m$.

Pro i -tý objekt má y_i hodnotu buď 0, nebo 1 a x_i^T je i -tý řádek matici X .

$$\begin{aligned} \text{Označme} \quad p(x, b) &= p_1(x, b) \quad a \\ 1 - p(x, b) &= p_2(x, b) \end{aligned}$$

a za předpokladu binomického rozdělení y lze zapsat logaritmus věrohodnostní funkce ve tvaru

$$\begin{aligned} \ln L(b) &= \sum_{i=1}^n (y_i \ln p(x_i, b) + (1 - y_i) \ln(1 - p(x_i, b))) \\ &= \sum_{i=1}^n (y_i b^T x_i - \ln(1 + \exp(b^T x_i))) \end{aligned}$$

kde $b^T = \{b_0, b_1\}$ a předpokládá se, že první sloupec matici X obsahuje pouze jedničky (absolutní člen).

Pro maximalizaci $\ln L(b)$ se využívá multy prvních derivací

$$J = \frac{d \ln L(b)}{db} = \sum_{i=1}^n x_i (y_i - p(x_i, b)) = 0$$

Jde o soustavu $m + 1$ nelineárních rovnic vzhledem k b . Řešení soustavy nelineárních rovnic využívá Newtonův-Raphsonovův algoritmus, který vyžaduje matici druhých derivací (hessiánu)

$$H = \frac{d^2 L(b)}{db db^T} = - \sum_{i=1}^n x_i x_i^T p(x_i, b) (1 - p(x_i, b))$$

Newtonova-Raphsonova metoda je iterativní, takže výsledkem j -té iterace je zpřesněný odhad

$$b_{(j+1)} = b_{(j)} - H_{(j)}^{-1} J_j$$

kde pro vektor pravděpodobnosti p rozměru $n \times 1$ s prvky $p(x_i, b_{(j)})$, a diagonální matici vah W rozměru $n \times m$ s prvky lze psát

$$J_{(j)} = X^T (y - p) \quad a \quad H = -X^T W X$$

Interpretace regresních koeficientů

Žádné předpoklady o x neexistují a x mohou být jak diskrétní (*faktory*) tak i spojité veličiny (*prediktory*).

Předpoklad říká, že logit $\ln(L_{(1)}/L_{(0)})$ je lineární funkcií nezávisle proměnných.

Pro $\ln(L_{(1)}/L_{(0)})$ se užívá termín logit

nebo-li logit transformace pravděpodobnosti.

Logistický model se nazývá vícenásobný logistický regresní model
(krátce *logit*) a koeficienty b_i jsou interpretovány jako regresní koeficienty.

Logit lze ale také upravit: dosazením za $L_{(1)} = (1 - L_{(0)})$ dostaneme

$$L_{(0)} = \frac{1}{1 + \exp[-(b_0 + b_1 x_1 + b_2 x_2 + \dots + b_p x_p)]}$$

Obecně:

Kladné znaménko koeficientu b_i zvyšuje pravděpodobnost $L_{(1)}$, a záporné znaménko tuto pravděpodobnost snižuje.

Diskuse koeficientu b_i :

- 1) Je-li b_i kladné, funkce exp je větší než 1 a pravděpodobnostní poměr $L_{(1)}/L_{(0)}$ se bude zvyšovat.

Zvýšení se objeví, když predikovaná pravděpodobnost odchrané události $L_{(1)}$ se zvýší a predikovaná pravděpodobnost neodehrané události $L_{(0)}$ se sníží.

Proto má model vyšší predikovanou pravděpodobnost odchrané události $L_{(1)}$.

- 2) Je-li b_i záporné, je funkce exp menší než 1 a pravděpodobnost se sníží.

- 3) Pro koeficient rovný nule vede funkce exp k hodnotě 1 čili k žádné změně pravděpodobnosti.

Test významnosti regresních koeficientů

Logistická regrese umožňuje testovat významnost koeficientů čili ověřit, že regresní koeficient se liší od nuly.

Nula zde značí, že pravděpodobnostní poměr $L_{(1)}/L_{(0)}$ se nemění a pravděpodobnost tím pádem není ovlivněna.

Studentův *t*-test k vyšetření statistické významnosti jednotlivých regresních koeficientů.

Pro velké výběry lze užít Waldovo testační kritérium $W_{a,i} = (b_i/s(b_i))^2$, které vyčíslouje statistickou významnost pro odhady regresních koeficientů stejně jako ve vícenásobné regresi.

Waldova statistika $W_{a,i}$ má χ^2 -rozdělení s 1 stupněm volnosti a představuje čtverec poměru odhadu regresního koeficientu a jeho směrodatné odchyly $W_{a,i} = (b_i/s(b_i))^2$.

Pro kategorické proměnné má $W_{a,i}$ počet stupňů volnosti roven o 1 méně než je počet kategorií.

Waldova statistika W_a má ale jednu nežádoucí vlastnost. Když je absolutní hodnota regresního koeficientu b_i veliká a odhad je ho směrodatně odchyly $s(b_i)$ je také veliký, je výsledkem příliš malá hodnota testačního kritéria $W_{a,p}$, která vede k selhání zamítnutí nulové hypotézy, že regresní koeficient je nulový. Proto, je-li regresní koeficient veliký, neužijeme Waldova kritéria.

Parciální korelace

Je obtížné určit příspěvek jednotlivých proměnných.

Příspěvek každé proměnné závisí také na ostatních proměnných v logistickém modelu.

K vyšetření parciální korelace mezi závisle proměnnou a každou nezávisle proměnnou se užívá korelační koeficient R_i , (v intervalu od -1 do +1).

- 1) Kladné hodnoty R_i : když roste hodnota R_i , zvyšuje se pravděpodobnost objektu "v události" $L_{(1)}$.
- 2) Záporné hodnoty R_i : naopak snižuje se pravděpodobnost objektů "v události" $L_{(1)}$.
- 3) Malé hodnoty R_i : proměnná má malý vliv na model.

Kategorické proměnné

Jednou z důležitějších výhod logistického modelu je možnost užívat i kategorické nezávislé proměnné x , zvané faktory.

Za faktor lze použít numerickou, textovou nebo datumovou hodnotu, zvanou úroveň nebo referenční hladina.

Interpretace odhadovaných regresních koeficientů je relativní vůči této hladině.

Nejjednodušší situací je jediný faktor x se dvěma možnými hodnotami, například, deprese u 143 žen a 101 mužů je ovlivněna pohlavím, kde faktor pohlaví má dvě úrovně: pro muže je $x = 0$ a pro ženy je $x = 1$.

objektem žena, pak pravděpodobnostní poměr, že je žena v depresi, je například 40/143. Podobně je tento poměr u mužů například 10/101. Pravděpodobnostní poměr, že jedinec je v depresi bude

$$\frac{L_{(1)}}{L_{(0)}} = \frac{40/143}{10/101} = 2.825.$$

oto šance žen čili pravděpodobnostní poměr žen nacházet se v depresi je 2.825krát větší než šance mužů.

dobně můžeme vyčíslit také šanci "nebýt v depresi":

$$\frac{L_{(0)}}{L_{(1)}} = \frac{143/40}{101/10} = 0.354.$$

nu šance se hodně využívá v biomediálních aplikacích. Je mírou spojení binární proměnné, jako je faktor riziku výskytu dané události, například nemoci.

Kategorická proměnná čili faktor má dvě úrovně, tj. $x = 0$ značí muže a $x = 1$ značí ženy a logistickou rovnici pak bude

$$L_{(1)} = \frac{1}{1 + e^{-a - b x}}$$

odhad parametru $a = -2.313$ a odhad $b = 1.039$.

Odhad b představuje přirozený logaritmus pravděpodobnostního poměru žen a mužů, 1.039 = $\ln 2.825$, a proto pravděpodobnostní poměr $e^b = e^{1.039} = 2.825$.

Odhad a je přirozený logaritmus pravděpodobnostního poměru mužů ($x = 0$) nebo $-2.313 = \ln 10/101$. Existuje-li pouze jedna dichotomní proměnná, není potřebné provádět logistickou regresní analýzu.

Přibližný interval spolehlivosti pro pravděpodobnostní poměr jako pro binární proměnnou se vypočte užitím odhadu směrnice b a odhadu je ji směrodatné odchyly.

Například 95% interval spolehlivosti pro pravděpodobnostní poměr se vyčíslí jako $\exp(b \pm 1.96 s(b))$.

Volba proměnných

apříklad úloha logistické předpovědi infarktu:

ata jsou z dlouhodobého sledování z počátku zdravých pacientů, u kterých byla dlouhodobě provedena opakována měření. Několik jedinců bylo postiženo infarktem, několik ne.

y sledován výběr nezávisle proměnných, které by mohly odhalit bližící se infarkt.

výběr účinných nezávisle proměnných byl předem lékaři vytypován.

Castěji však uživatel předem neví nic o nezávisle proměnných.

Proměnné x jsou nejprve vyšetřovány, která je nejvíce spjata z dichotomní závisle proměnnou.

Studentův t-test významnosti jednotlivých parametrů: užívá se dostatečně vysoká hladina významnosti, například $\alpha = 0.15$, aby užitečná nezávisle proměnná nemohla být odstraněna.

Vyšetření zredukuje počet nezávisle proměnných na 10 či ještě méně.

Pak nastoupí **kroková logistická regresní analýza**: jde o test, zda proměnná x zlepší prediktivní schopnost modelu. Postupy a jejich kritéria jsou užita k rozhodování, kolik proměnných x , a které je třeba užít.

Testy v dopředné krokové analýze jsou postaveny na χ^2 -statistiky: velká hodnota χ^2 nebo malá spočtená hladina významnosti P ukazují, že nezávisle proměnná by měla být zařazena do proměnných.

Nalezená velká hodnota χ^2 ukazuje, že proměnné jsou užitečné.

4. Těsnost proložení logistickým modelem

Před analýzou je třeba posoudit, zda nejsou odlehle hodnoty. Rozptylové diagramy snadno odhalí odlehle body.

Proměnné nemusí být normálně rozděleny.

Rozregresní diagnostika s analýzou vlivných bodů odhalí O a E .

Logistická křivka má esovitý tvar a vystihuje logistický model, který je vzhledem ke koeficientům b nelinéarní.

Po linearizační transformaci budou koeficienty představovat směrnice u proměnných lincárního regresního modelu.

Míra těsnosti proložení navrženého modelu dat je hodnota pravděpodobnosti $L_{(1)}$, že se událost uskuteční.

Místo veličiny $L_{(1)}$ se používá tzv. **odchylka, deviance** $D = -2 \ln L_{(1)}$ čili $D = -2LL$.

D představuje míru těsnosti proložení dat logistickým regresním modelem:

1) **Dobrý model** vede k vysoké pravděpodobnosti objektů v události $L_{(1)}$, což přetrasformáno do veličiny $-2 \ln L_{(1)}$ poskytnec malou hodnotu blízkou nule.

2) **Minimální hodnotou** pro $-2 \ln L_{(1)}$ je **nula**, při které je dosaženo naprostou perfektní těsnosti proložení.

Rozdíl v odchylce je definován vztahem $G = D(\text{model bez proměnné}) - D(\text{model s proměnnou})$

čili

$$G = -2 \ln \frac{\text{pravděpodobnost modelu bez proměnné}}{\text{pravděpodobnost modelu s proměnnou}}$$

Veličina G proto odpovídá věrohodnostnímu poměru.

Těsnost proložení: spočívá porovnání experimentálních hodnot E s vypočtenými V :

Pearsonův test dobré shody χ^2 se užije, když model platí:

- Velká hodnota χ^2 indikuje špatné proložení modelu.
- Malé hodnoty vypočtené hladiny významnosti P indikují špatné proložení modelu.

Nejužívanější způsoby posouzení těsnosti proložení:

Nejužívanější způsoby posouzení těsnosti proložení:

a) **Klasický Pearsonův přístup** začíná s identifikováním různých kombinací hodnot proměnných v regresním modelu, tj. vzorů.

Například dvě dichotomní proměnné, (pohlaví a zaměstnání) vedou na 4 kombinace: muž zaměstnán, muž nezaměstnán, žena zaměstnána, žena nezaměstnána.

- Pro každou kombinaci vyčíslíme počet E experimentálních hodnot jednotlivců (objektů) ve třídě I a II.

- Podobně pro každého jednotlivce vypočteme pravděpodobnost, že se nachází ve třídě I a ve třídě II logistickou regresní analýzou.

- Suma těchto pravděpodobností pro daný vzor se označí V .

- Testovací statistika testu dobré shody χ^2 se vyčíslí jako

$$\chi^2_{\text{exp}} = \sum_{i=1}^n 2E \left(\ln \frac{E}{V} \right)$$

kde suma se provede přes všechny odlišné vzory.

- Rezidua se sledují právě pro tyto odlišné vzory.

b) **Hosmerův-Lemeshowův test dobré shody** byl navržen v 1982. Pearsonův χ^2 -test dobré shody k redukci v logaritmech hodnoty pravděpodobnosti je měrou sledování zlepšení těsnosti zavedením jedné či více nezávisle proměnných.

Základní model, který je podobný výpočtu sumy čtverců při použití pouze průměrů, poskytuje nulovou linii k porovnání.

Vedle χ^2 -testu existuje několik R^2 -podobných měr k posouzení těsnosti proložení, obdoba koeficientu determinace ve vicenásobné regresi.

"Pseudo R^2 " v logistické regrese pro logitový model se vypočte dle

$$R^2_{\text{logit}} = \frac{2 \ln L_{\text{mul}} - (-2 \ln L_{\text{model}})}{-2 \ln L_{\text{mul}}} = -\frac{D_{\text{model}} + D_{\text{mul}}}{D_{\text{mul}}}$$

c) **Metoda klasifikačních matic**, vyvinutých v diskriminační analýze slouží k vyhodnocení prediktivní schopnosti v pojmech zařazení do třídy.

Pravděpodobnost zařazení do třídy I je vypočtena pro každého jednotlivce (objekt) ve výběru a výsledný počet je uspořádán vzestupně.

Pravděpodobnosti jsou pak rozděleny do 10 skupin (decily).

Pro každý naměřený počet jednotlivců ve třídě I je vyčíslen počet E . Užití logické regrese jsou pro jedince v každém deciliu vypočteny počty V . Počty se vyčíslí Pearsonova χ^2 -statistika testu dobré shody

$$\chi^2_{\text{exp}} = \sum_{i=1}^n \frac{(E - V)^2}{V}$$

kde sumace se provede přes obě třídy a 10 deciliů.

Velká hodnota χ^2 nebo malá hodnota P indikují, že proložení není dobré.

5. Kvalita vyhodnocení logistickou regresí

Třídíme objekty do tříd, musíme nalézt prahový bod pravděpodobnosti P_c : objekt je "v události", když pravděpodobnost události větší nebo rovná hodnotě P_c .

Graf prahové operační charakteristiky ROC k detekci signálu, když signál nebylo vždy možné správně přejmout.

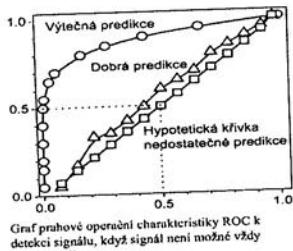
Na ose y je procento správně zařazených objektů "v události" nazvané pozitivní podíl (a v lékařském výzkumu nazývané citlivost).

Na ose x je procento nesprávně zařazených objektů nazvané falešný podíl nebo v lékařském výzkumu "1 minus specificita" (v lékařském výzkumu nazývané senzitivita zařazených krys pro správně zařazené krysy a specificita krys pro falešně zařazené krysy).

1) Horní křivka v grafu ROC představuje výtečnou predikci: i pro hodnoty podílu nesprávně zařazených objektů se získá vysoké procento správně zařazených objektů, které skutečně jsou "v události".

2) Střední křivka je skutečná křivka při uvažování malého počtu nezávislých proměnných, - třeba dvou. Vysoké procento (80 %) objektů správně zařazených v události je v poměru k 65 % chyběně zařazených v události nepřijatelně hladině.

3) Dolní hypotetická křivka, (přímka) odpovídá nahodilým výsledkům například házení mincí. Blízkost střední křivky k dolní ukazuje, že potřeba budou volit jinou, anebo přidat ještě další nezávisle proměnné abychom získali lepší model, i když je ale tento model statisticky významný na spočtené hladině $P = 0.009$.



Graf prahové operační charakteristiky ROC k detekci signálu, když signál není možné vždy správně pojmitout.

Vybereme prahový bod na dolní části křivky grafu ROC a nechceme mít příliš mnoho objektů, zařazených jako "v události", bude se nazývat přísný práh.

Nevýhoda: je ztráta mnoha objektů, které jsou "v události".

) Vybereme prahový bod na horní části křivky grafu ROC a chceme mít hodně objektů zařazených jako "v události", bude se nazývat nedbalý práh.

Nevýhoda: sice velmi málo objektů "v události" bude ztraceno ale mnoho objektů "v neudálosti" bude chybět označeno jako "v události".

Křivky v grafu ROC musí procházet body (0, 0) a (1, 1).

Maximální plocha pod křivkou je jen číslo 100%.

Numerická hodnota velikosti plochy bude blízká 1, když predikce modelu bude výtečná.

Když bude plocha blízká hodnotě 0,5, bude predikce modelu špatná.

Křivka ROC je proto užitečná při rozhodování, který ze dvou logistických modelů vybrat; lepší model dosáhne větší plochy pod křivkou ROC ale také větší výšky prahového bodu na křivce ROC.

Většina programů vybírá logistický model podle kritéria největší plochy pod křivkou ROC.

Aplikace logistické regrese

Modelu vícenásobné logistické regrese se často užívá k odhadu pravděpodobnosti jisté události, která se přihodí danému objektu.

Určení logistického regresního modelu je třeba výběru dat, ve kterém každý objekt, jedinec byl sledován v uvedeném časovém období a hodnoty závažných proměnných byly od začátku pečlivě zaznamenávány.

Výběr dat může být uskutečněn dvojím způsobem:

1. Výběr cross-validation: je získán náhodným způsobem a pozorování provedeno v uvedeném časovém období. Z tohoto výběru se vyčlení dva podvýběry: **první podvýběr**, který obsahuje hodně zkušeností o události, a **druhý podvýběr**, který obsahuje zbylé údaje.

Na datech prvního podvýběru se vyčíslí logistický regresní model, který pak může být aplikován na člena druhého podvýběru.

Předpokládá se, že **původní výběr** je v ustáleném stavu, tzn. neobjevily se žádné podstatné změny, které by pozměnily vztah mezi nezávisle proměnnými a výskytem události.

2. Případ řídícího výběru: spočívá v získání dvou náhodných výběrů: **první výběr**, ve kterém se událost objeví, a **druhý výběr**, ve kterém se událost neobjeví.

Hodnoty predikovaných proměnných se musí získat retrospektivním způsobem, z minulých záznamů nebo ze vzpomínek.

Konstanta a musí být nastavena tak, aby vyjadřovala pravý poměr objektu v události.

Existují důležité požadavky, které je třeba respektovat:

1. Model předpokládá, že logaritmus pravděpodobnostního poměru je lineárně závislý na nezávislých proměnných. Nesplnění by mělo být předem prověřeno buď užitím měr řešnosti proložení, nebo jinými způsoby. To může vyžadovat transformaci dat.

2. Výpočty jsou často časově náročné, a proto by měl uživatel rozumne redukovat počet proměnných.

3. Logistická regrese se neměla užívat k vyhodnocení faktorů riziku v dlouhodobých studiích, ve kterých jsou jo dnotlivé studie rozličné délky.

4. Regresní koeficienty pro nezávisle proměnnou v logistickém regresním modelu závisí na ostatních proměnných, zařazených do logistického modelu. Koeficienty pro stejnou nezávisle proměnnou, když se použijí různé výběry proměnných, mohou být zcela odlišné.

5. Je-li užita sehnána analýza, kterákoliv proměnná pro sehnání nemůže být použita jako nezávisle proměnná.

6. Jsou okolnosti, kde metoda maximální věrohodnosti odhadovaných regresních koeficientů neposkytuje odhady, tj. nekonverguje.

Příklad 4.26 Volba proměnných k popisu leukemie

Lee (1980) publikoval data o leukemii pacientů. Závisle proměnnou je binární proměnná *REMISS*, zda se objeví ústup leukemie y (1) či neobjeví (0). Nezávisle proměnnými x jsou:

CELL celulrita, buněčnost sraženiny kostní dřeně,

SMEAR skvrna diferenčního procenta napadení,

INFIL procento infiltrátu kostní dřeně buňkou leukemie,

LJ procento označeného indexu leukemických buněk kostní dřeně,

BLAST absolutní počet napadení v periferní krví,

TEMP nejvyšší teplota před začátkem léčby.

Otázkou je, které nezávisle proměnné jsou statisticky významné v logistickém regresním modelu.

Řešení: Byl užit program NCSS2000.

1. Odhad regresních koeficientů.

Proměnná	Regresní koeficient	Směrodatná odchylnka	χ^2 pro $\beta = 0$	Spočtená hladina P	Poslední R^2
Úsek	58.0387	71.23627	0.66	0.415224	0.032124
CELL	24.66053	47.83722	0.27	0.606197	0.013113
SMEAR	19.29247	57.94952	0.11	0.739196	0.005511
INFIL	-19.60012	61.68098	0.10	0.750662	0.005023
LJ	3.895928	2.3371	2.78	0.095516	0.121993
BLAST	0.1510942	2.278567	0.00	0.947130	0.000220
TEMP	-87.43308	67.57322	1.67	0.195699	0.077243

χ^2 udává Pearsonovo testování kritérium χ^2 pro 1 stupně volnosti k testu H_0 : $\beta_i = 0$ vs. H_A : $\beta_i \neq 0$. Vyčíslí se Waldovo kritérium $W_{\alpha,1}^2 = [b_i / s(b_i)]^2$.

Test významnosti b_i : je-li spočtená hladina P menší než předvolená $\alpha = 0.05$, je parametr b_i statisticky významný. Všechny prediktory se jeví jako statisticky nevýznamné.

Poslední R^2 udává hodnotu, která se přísluší k celkové R^2 , když se tato nezávisle proměnná přidá do logistického regresního modelu. Vypočte se dle vzorce

$$R^2 = \chi^2(df) / [\chi^2(df) + n - p - 1]$$

2. Nalezený model v transformované formě.

Nalezený logistický regresní model: $58.0387 + 24.66053 * CELL + 19.29247 * SMEAR - 19.60012 * INFIL + 3.895928 * LJ + 0.1510942 * BLAST - 87.43308 * TEMP$.

3. Přehled modelu.

R ² modelu	df	Odhylka D	Spočtená hladina významnosti P
0.386900	6	12.62	0.049463

Odhylka D testuje, zda všechny regresní koeficienty β , kromě úseku β_0 , jsou rovny nule. Protože je spočtená P menší než $\alpha = 0.05$, je regresní model statisticky významný.

4. Klasifikační tabulka.

		Nalezeno predikci logistickým modelem		
Dáno závisle proměnnou		Nc	Ano	Celkově
Ne	Četnost	15	3	18
	Řádkové procento	83.33	16.67	100.00
Ano	Sloupcové procento	78.95	37.50	66.67
	Četnost	4	5	9
Celkově	Řádkové procento	44.44	55.56	100.00
	Sloupcové procento	21.05	62.50	33.33
Procento správně klasifikovaných = 74.07				

Tabulka přináší četnosti, řádková procenta a sloupcová procenta predikovaných objektů a nakonec je procento správně klasifikovaných objektů. Jde o procento z celkového počtu, které padne na diagonálu tabulky.

5. Predikovaná klasifikace.

Rádek	Dána třída	Nalezená třída	Logistické skóre	Reziduum
2	Ano (1)	Ne (0)	0.799341	
3	Ano (1)	Ne (0)	0.434904	
4	Ne (0)	Ne (0)	0.155573	
5	Ano (1)	Ano (1)	0.094662	
6	Ne (0)	Ne (0)	0.094699	0.362981
7	Ano (1)	Ne (0)	0.374582	0.621413
8	Ne (0)	Ano (1)	0.600523	-0.600523
9	Ne (0)	Ne (0)	0.165277	-0.165277
10	Ne (0)	Ne (0)	0.060958	-0.060958
11	Ne (0)	Ne (0)	0.027695	-0.027695
12	Ne (0)	Ne (0)	0.027695	-0.027695
13	Ne (0)	Ne (0)	0.009946	-0.009946
14	Ne (0)	Ne (0)	0.000001	-0.000001
15	Ne (0)	Ano (1)	0.372588	-0.372588
16	Ano (1)	Ano (1)	0.725752	0.274248
17	Ne (0)	Ne (0)	0.000055	-0.000055
18	Ne (0)	Ne (0)	0.228173	-0.228173
19	Ne (0)	Ne (0)	0.000001	-0.000001
20	Ano (1)	Ano (1)	0.673863	0.322137
21	Ne (0)	Ne (0)	0.015910	-0.015910
22	Ne (0)	Ne (0)	0.007445	-0.007445
23	Ano (1)	Ne (0)	0.247686	0.752316
24	Ne (0)	Ano (1)	0.851096	-0.851096
25	Ano (1)	Ano (1)	0.938464	0.061536
26	Ano (1)	Ne (0)	0.461177	0.538823
27	Ne (0)	Ne (0)	0.279469	-0.279469

Jsou zde zobrazeny pouze chyběně zařazené řádky.

Dána třída určuje zadanou skutečnou třídu. *Nalezená třída* představuje nalezenou třídu na základě logistického regresního modelu. Logistické skóre je odhad pravděpodobnosti, že objekt patří do třídy Ne. Reziduum představuje jíto rozdíl mezi Logistikou skóre a indexem skutečné třídy. Index třídy Ne je 0 a index třídy Ano je 1.