



UNIVERZITET U NOVOM SADU
PRIRODNO – MATEMATICKI FAKULTET
DEPARTMAN
ZA MATEMATIKU I INFORMATIKU



**RAZLICITI PRISTUPI KREDITNOM
SKORING SISTEMU**

- master rad -

Profesor: dr. Zorana Lužanin

Autor: Jelena Burgijašev

Novi Sad, Septembar 2009.

Sadržaj:

1. Uvod	2
2. Osnove kreditnog scoring sistema	3
3. Matematicki aparat	5
3.1 Linearno programiranje	5
3.2 Problem separacije i klasifikacije	7
3.2.1 Separacija pomocu linearne površi	7
3.2.2 Separacija pomocu kvadratne površi	9
3.2.3 Model klasifikacije uz pomoc funkcije korisnosti	10
3.2.4 Klasifikacija pomocu linearne funkcije korisnosti	11
3.2.5 Klasifikacija uz pomoc kvadratne funkcije korisnosti	12
3.3 Logit analiza	13
3.3.1 Selektioni kriterijumi	15
3.3.2 Metod maksimalne verovatnoće	18
3.3.3 Newton-Raphson algoritam (NRA)	19
3.3.4 Testiranje hipoteza	20
3.3.5 Konstruisanje intervala poverenja	21
3.3.6 Primena softvera	22
3.4 Probit analiza	27
3.4.1 Primena softvera	28
3.5 Tobit model	33
4. Modeli kreditnog scoringa sistema	34
4.1 Primena problema separacije	34
4.2 Dve faze kreditnog ocenjivanja u procesu odobravanja kredita	36
4.3 Primena probit i tobit modela	43
4.4 Primena logit modela	51
4.5 Primer genetskog programiranja	55
4.5.1 Uporedivanje genetskog programiranja i probit analize na celom uzorku	56
4.5.2 Uporedivanje genetskog programiranja i probit analize na poduzorku	56
4.6 Primena analize obavljanja podataka (DEA)	57
4.6.1 Odabir finansijskih koeficijenata	57
4.6.2 Racunanje skorova uz pomoc DEA	58
4.6.3 Provera valjanosti DEA skorova	59
4.6.4 Metod kreditnog rejtinga	60
5. Zaključak	62
6. Literatura	64

1. Uvod

Kako se kreditno tržište sve više razvija, sve je veća potreba za dobro razvijenim sistemom za zaštitu od rizika i gubitaka, i u tu svrhu je uveden kreditni skoring sistem. U ovom radu su prezentovani razliciti pristupi kreditnom skoring sistemu koji su razvijeni primenom matematičkog aparata. Ovaj model se najčešće sreće u bankama kada odobravaju kredit, pri proceni kreditne sposobnosti klijenta i rizika koji taj klijent nosi. U zavisnosti od kreditne politike banke postoje razni pristupi ovom problemu, kao što su na primer minimizacija rizika, maksimizacija profita ili nešto treće.

Za razvijanje modela najčešće se koriste statističke i ekonometrijske tehnike, kao što su na primer logit i probit analiza. Metode separacije nam služe za klasifikaciju klijenata u dve grupe: dobre i loše, tj. one koji redovno izmiruju svoje obaveze i one koji imaju problema u poslovanju pa samim tim kasne u otplati kredita. Cilj je što tacnije klasifikovati klijente, pošto se time u isto vreme minimiziraju troškovi koji nastaju pri pogrešnoj klasifikaciji. Prilikom separacije pomocu linearne, odnosno kvadratne površi primenjuje se problem linearног programiranja za minimizaciju greške klasifikacije. U jednom primeru cemo uporedivanjem genetskog programiranja i probit analize pokazati da se uporedno sa razvojem kreditnog tržišta razvijaju i precizniji modeli.

Kroz date primere se vidi da se u razlicitim zemljama primenjuju razliciti modeli i pristupi kreditnom skoringu, koji u nekim slučajevima daju i nelogične rezultate sa naše tacke gledišta, kao što cemo videti na primeru Vijetnamskog bankarskog tržišta. Vecina modela se zasniva na istorijskim podacima, i zbog toga je bitno da uzorak bude što veci da bi dobili što bolji i realniji model. Kao problem može da se javi nedovoljna kolicina podataka, jer je kreditno tržište još uvek u fazi razvoja, kao i odredena pristrasnost modela pošto se posmatraju samo klijenti kojima je odobren kredit. Ovo su problemi kojima se treba baviti u buducnosti, da bi se razvio što pouzdaniji i realniji model.

Pojam kreditnog skorинга se uglavnom vezuje za analizu kreditne sposobnosti pojedinacnog klijenta, međutim sve je veća njegova primena i pri ocenjivanju rizika koji nose firme. Ovaj tip kreditnog skorингa je još uvek u fazi razvoja i još uvek nije toliko zastupljen u praksi, iako cemo kroz primer videti da on daje dosta precizne rezultate.

2. Osnove kreditnog scoring sistema

Kada jedan bankar pita drugog „Koji je skor?“, vecina ljudi ce pomisliti da ovaj prvi nije gledao sinocnu utakmicu. Medutim, oni samo rade svoj posao, tacnije, raspravljuju o tome koliki je kreditni skor klijenta koji podnosi zahtev za kredit. **Kreditni scoring** je statisticki metod koji se koristi da se oceni koliki je rizik da ce novi klijent koji podnosi zahtev za kredit ili vec postojeci klijent kasniti u otplati kredita ili da uopšte nece biti u stanju da isplati kredit. Drugim recima, kreditni scoring je metod za odredivanje kreditnog rizika koji nosi klijent. Koristeci istorijske podatke i statisticke tehnike kreditni scoring pokušava da izoluje efekte odredenih karakteristika koje klijenta vode u situaciju da ne može da otplacuje kredit. Rezultat ovog metoda je skor koji banka koristi da rangira klijente na osnovu rizika koji oni nose. Na osnovu toga koliki rizik banka želi da prihvati, odreduje se granicni skor, tako da ce klijentima koji imaju veci skor od granicnog biti odobren kredit, a oni koji imaju manji skor ce biti odbijeni. Da bi napravili model, analiticari analiziraju istorijske podatke ranije odobrenih kredita, tj. analiziraju kako su se raniji klijenti ponašali pri otplati kredita, da bi odredili karakteristike koje su korisne pri oceni da li je klijent sposoban da redovno otplacuje kredit. Samim tim su i klijenti podeljeni na dobre i loše. Klijenta definišemo da je *loš* ukoliko nije platio tri uzastopne rate, tj. u kašnjenju je dužem od devedeset dana. Problem se javlja kada se pravi model za relativno nove proizvode zbog nedovoljnog obima istorijskih podataka na osnovu kojih bi se razvio model. U tom slučaju se model može praviti na malom uzorku, ili na nekom slicnom proizvodu, ali tada rezultati nece biti dovoljno dosledni. Dobro napravljen model bi trebao da dodeli veci skor klijentima koji ce redovno otplacivati kredit, a niži skor klijentima koji ce kasniti u otplacivanju kredita. Medutim, nijedan model nije savršen, pa se dešava da neki loši klijenti dobiju veci skor od nekih dobrih klijenata.

Razlikujemo dva modela kreditnog scoring sistema:

1. model koji pomaže analiticaru da doneše odluku da li da odobri kredit klijentu ili da ga odbije
2. model koji pomaže u analiziranju ponašanja vec postojećih klijenata

Mi cemo se u ovom radu fokusirati na prvi model kreditnog scoring sistema.

Pri oceni kreditne sposobnosti klijenta koriste se dva izvora informacija: prvi je aplikacioni formular koji klijent popunjava kada podnosi zahtev za kredit, a drugi je kreditni biro. **Kreditni biro** je institucija koja poseduje veliku bazu podataka o svim klijentima, kako licnih podataka tako i informacije o urednosti otplate kredita. Ova informacija se naziva **kreditna istorija**. Kada klijent podnosi zahtev za kredit, kreditni analiticar proverava sa kreditnim birom kakva je klijentova kreditna istorija. Ako je kreditna istorija «siromašna», analiticar može odbiti da odobri kredit. U situaciji kada analiticar i odluci da odobri kredit, sigurno ce odrediti i vecu kamatnu stopu.

Najzastupljeniji statisticki modeli koji se koriste da bi se napravio kreditni scoring sistem su linearni model, logit model, probit model i model diskriminantne analize. Prva tri modela su standardne statisticke tehnike koje na osnovu istorijskih podataka i karakteristika klijenta ocenjuju verovatnocu da ce taj klijent biti loš.

Razlika medu ovim modelima je ta što linearni model prepostavlja linearnu zavisnost izmedu verovatnoce da ce klijent biti loš i karakteristika koje se koriste u modelu, logit model prepostavlja da ova verovatnoca ima logaritamsku raspodelu, a probit model da ima normalnu raspodelu. Diskriminantna analiza se od ovih statistickih modela razlikuje po tome što umesto da ocenjuje verovatnocu da ce klijent biti loš, ona deli klijente u dve klase, manje i više rizicne. Na jednom primeru cemo prikazati primenu genetskog programiranja u kreditnom scoring modeliranju kao i primenu analize obavijanja podataka, ali necemo ulaziti u detalje, pošto to izlazi iz opsega ovog rada. Napomenimo da za kreditni scoring model mogu da se koriste i razni neparametarki i modeli veštacke inteligencije.

Koliko god model bio dobar, on nece sa sigurnošcu predvideti kako ce se klijent ponašati pri otplati kredita, ali bi ipak trebalo da dâ dovoljno tacnu ocenu verovatnoce da klijent sa određenim karakteristikama nece biti u mogucnosti da otplati kredit. Takođe treba uzeti u obzir da postoji i odredena pristrasnost u uzorku, pošto se uzimaju u obzir samo klijenti kojima je odobren kredit. Ovo se javlja zbog toga što se ne zna kako bi se odbijeni klijenti ponašali da im je odobren kredit.

Prednosti kreditnog scoringa koje su dovele do njegove sve vece primene su kao prvo ušeda vremena koje je potrebno da se doneše odluka u procesu odobravanja kredita, što je vrlo bitna karakteristika u trenutnoj ekspanziji kreditnog tržišta. Takođe, ovaj sistem donosi odluku na osnovu istih kriterijuma za sve klijente, nezavisno od pola, rase ili nekog drugog kriterijuma koji može dovesti do odredene diskriminacije. Da bi se neka od ovih karakteristika koristila u modelu, moraju da se imaju dovoljno jaki argumenti zašto je ona toliko bitna za model.

Prvo cemo se upoznati sa matematickim aparatom koji koristimo da bismo razvili model i došli do željenih rezultata, a zatim cemo kroz primere objasniti ove matematicke modele.

3. Matematicki aparat

3.1 Linearno programiranje

Problem linearog programiranja služi za modeliranje tzv. uslovne optimizacije u kojoj treba naci optimalno rešenje, tj. ono rešenje za koje se postiže najbolja vrednost nekog cilja na skupu svih mogucih alternativnih rešenja problema, pri cemu svako rešenje iz ovog skupa zadovoljava zadate uslove (ogranicenja).

Opšti zadatak linearog programiranja se može iskazati u formalizovanom razvijenom obliku kao

$$\begin{cases} \max \\ \min \end{cases} f(x) = c_1x_1 + c_2x_2 + \dots + c_nx_n$$

uz ogranicenja

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n &\stackrel{<}{\underset{>}{=}} b_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n &\stackrel{<}{\underset{>}{=}} b_2 \quad (*) \\ &\vdots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n &\stackrel{<}{\underset{>}{=}} b_m \\ x_1 \geq 0, x_2 \geq 0, \dots, x_n \geq 0 & \quad (**) \end{aligned}$$

ili krace

$$\begin{cases} \max \\ \min \end{cases} f(x) = \sum_{i=1}^n c_i x_i$$

$$\sum_{j=1}^n a_{ij} x_j \stackrel{<}{\underset{>}{=}} b_i, \quad i = 1, \dots, m$$

$$x_j \geq 0, \quad j = 1, \dots, n$$

gde su a_{ij} , c_i , b_j , $i = 1, 2, \dots, m$, $j = 1, 2, \dots, n$ zadati realni brojevi. Funkcija $f(x)$ se naziva **funkcija cilja**, dok se m ogranicenja iz skupa (*) nazivaju jednostavno **ogranicenja**. Ogranicenja (**) se nazivaju **uslovi nenegativnosti**.

Dopustiva oblast skupa (*) - (**) se naziva **dopustiva oblast** problema linearog programiranja, a svaka tacka dopustive oblasti predstavlja **dopustivo rešenje** ovog problema. Ono dopustivo rešenje u kome funkcija dostiže svoj ekstrem se naziva **optimalno rešenje** problema. Matrica $A = [a_{ij}]_{m \times n}$ se cesto naziva matricom ogranicenja, a b_i , $i = 1, 2, \dots, m$ slobodnim clanom i -tog ogranicenja iz sistema (*).

Ako su sva ogranicenja u (*) istog tipa i $c = (c_1, c_2, \dots, c_n)$, $b = (b_1, b_2, \dots, b_m)$ i $x = (x_1, x_2, \dots, x_n)$ tada se problem linearog programiranja može prikazati i u matricno-vektorskoj formi. Ako su sva ogranicenja tipa jednakosti, tada imamo **standardni oblik** linearog programiranja koji ima sledecu matricno-vektorskiju formu:

$$\begin{cases} \max \\ \min \end{cases} c^T x$$

$$Ax = b \quad (***)$$

$$x \geq 0$$

Prepostavljamo da ovaj problem zadovoljava sledeće dodatne uslove:

1. svi slobodni clanovi $b_i, i = 1,2,\dots,m$ su nenegativni, tj. $b \geq 0$
2. $m < n$ i $\text{rang } A = m$, tj. sve jednacine ogranicenja su linearno nezavisne

Posmatramo sistem linearnih jednacina:

$$Ax = b \quad (\&)$$

Definicija *Bazno rešenje* sistema ($\&$) predstavlja ono rešenje ovog sistema za koje postoji neka $m \times m$ podmatrica A_B matrice A koja je regularna, tj. $\det A_B \neq 0$ i takva da je u ovom rešenju svaka promenljiva cija se kolona koeficijenata ne nalazi u A_B jednaka nuli. Promenljive koje odgovaraju kolonama podmatrice A_B se nazivaju **bazne promenljive**, a preostale promenljive su nebazne. Skup svih baznih promenljivih cini **bazu** baznog rešenja, dok je A_B **matrica baze**.

Primetimo da rešenje ovakvog sistema uvek postoji pošto je A_B regularna podmatrica. Kako je $\text{rang } A = m$, postoji bar jedna regularna $m \times m$ podmatrica od A , pa time i bar jedno bazno rešenje sistema ($\&$).

Definicija Ako neko bazno rešenje problema ($\&$) zadovoljava i uslove nenegativnosti onda ono predstavlja **bazno dopustivo rešenje** problema (***)

Definicija Dva bazna dopustiva rešenja problema ($\&$) su **susedna** ako im se baze razlikuju samo u jednoj promenljivoj.

Jedan od razloga tako rasprostranjene primene modela linearog programiranja u praksi je i postojanje efikasnih procedura za njihovo rešavanje koje sa uspehom rešavaju i probleme velikih dimenzija. Medu njima najpoznatija je **Simpleks metoda**. Ona predstavlja algebarsku proceduru koja nalazi niz baznih dopustivih rešenja x_k problema ($\&$) sa odgovarajucim bazama $B_k, k = 0,1,2,\dots$ tako da je vrednost funkcije cilja u x_{k+1} veca od njene vrednosti u x_k , sve dok trenutno generisano bazno dopustivo rešenje ne zadovolji kriterijum optimalnosti.

Kriterijum optimalnosti Ako je $c_j^k \leq 0$, za svako $j = 1,2,\dots,n$ tada je x_k optimalno rešenje problema linearog programiranja, a f_k maksimalna vrednost funkcije cilja.

Drugim recima, x_k je optimalno ako se nikakvim povecanjem vrednosti nebaznih promenljivih ovog rešenja ne može povecati vrednost funkcije cilja. Necemo dublje ulaziti u postupak Simpleks metoda.

3.2 Problem separacije i klasifikacije

Problem separacije je problem nalaženja kriterijuma za razdvajanje elemenata u dva disjunktna skupa objekata. Osnovni problem separacije se može definisati na sledeći nacin: neka su data dva skupa objekata A i B , pri cemu skup A sadrži m objekata, a skup B sadrži k objekata i pri tome je svaki objekat jedan n -dimenzionalni vektor. Problem separacije se tada svodi na pronalaženje površi u n -dimenzionalnom prostoru tako da se sve tacke koje predstavljaju objekte iz skupa A nalaze sa jedne strane ove površi, dok se sve tacke koje predstavljaju objekte iz skupa B nalaze sa druge strane površi.

3.2.1 Separacija pomocu linearne površi

Definišemo *skup objekata* kao nepraznu matricu realnih brojeva, gde svaka vrsta definiše pojedinacni objekat. Objekat se sastoji od n realnih brojeva koji se nazivaju *opservacije*. Ako matricu objekata označimo sa A , pojedinacne objekte kao vektor vrste a_i , tada je j -ta opservacija i -tog objekta označena sa a_{ij} . Osnovni problem **linearne separacije** glasi: za data dva skupa objekata koja su definisana matricama $A \in \mathbf{R}^{m \times n}$ i $B \in \mathbf{R}^{k \times n}$, odrediti hiperravan u prostoru \mathbf{R}^n tako da ako m vrsta matrice A i k vrsta matrice B posmatramo kao tacke u ovom prostoru, tada se one moraju nalaziti sa razlicitih strana ove ravni. Neka su a_i , $i=1,\dots,m$ objekti koji pripadaju matrici A , a b_j , $j=1,\dots,k$ objekti koji pripadaju matrici B . Ako postoji nenula vektori $c \in \mathbf{R}^n$ i $c_0 \in \mathbf{R}^n$ takvi da je

$$\begin{aligned} c^T a_i &> c_0, \quad i=1,\dots,m \\ c^T b_j &< c_0, \quad j=1,\dots,k \end{aligned}$$

tada se $H(c, c_0) = \{x \in \mathbf{R}^n \mid c^T x = c_0\}$ naziva **hiperravan**. Ako $x_0 \in H(c, c_0)$ tada važi da je $c^T x_0 = c_0$ pa je

$$H(c, c_0) = \{x \in \mathbf{R}^n \mid c^T x = c^T x_0 = c_0\} = \{x \in \mathbf{R}^n \mid c^T (x - x_0) = c_0\}$$

Izraz $c^T (x - x_0) = 0$ predstavlja ortogonalnost vektora c i $x - x_0$ i zbog toga se vektor c naziva normalni vektor hiperravnih $H(c, c_0)$.

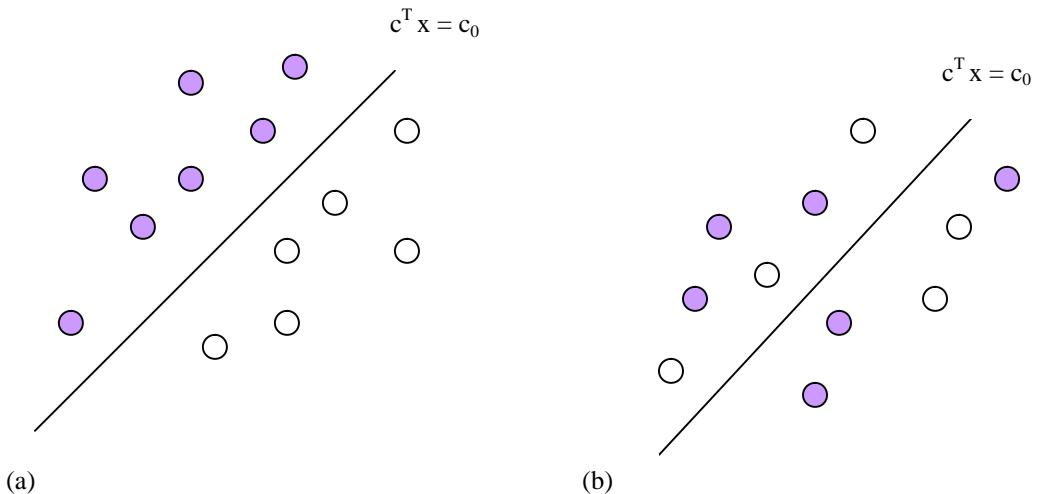
Ako uvedemo granicu debljine δ izmedu klasa, onda važi:

$$\begin{aligned} c^T a_i &\geq c_0 + \delta, \quad i = 1, \dots, m \\ c^T b_j &\leq c_0 - \delta, \quad j = 1, \dots, k \end{aligned}$$

U opštem slučaju ove nejednakosti nisu zadovoljene, ali mi težimo da one budu bar približno zadovoljene.

Poluprostore definišemo na sledeći nacin:

$$\begin{aligned} H_+(c, c_0) &= \{x \in \mathbf{R}^n \mid c^T x > c_0\} \\ H_-(c, c_0) &= \{x \in \mathbf{R}^n \mid c^T x < c_0\} \end{aligned}$$

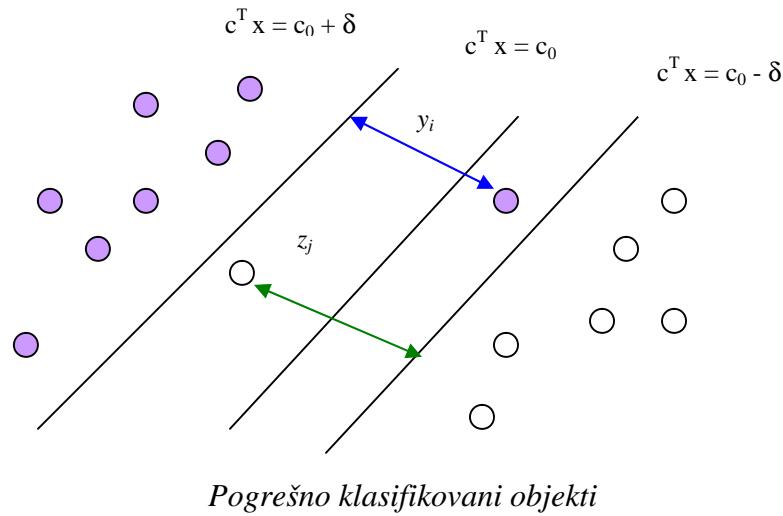


- (a) Hiperravan potpuno razdvaja objekte dve razlicite klase;
 (b) Hiperravan koja razdvaja objekte razlicitih klasa ne postoji

Klasifikacija može da se definiše kao funkcija klasifikacije koja svakom objektu dodeljuje broj, odnosno svaki objekat «smešta» u jednu od klasa, u ovom slučaju u jedan poluprostor. Za objekat A_i za koji je $a_i \in H_-(c, c_0)$ kažemo da je *pogrešno klasifikovan objekat prvog reda*. Za objekat B_j za koji je $b_j \in H_+(c, c_0)$ kažemo da je *pogrešno klasifikovan objekat drugog reda*. Sa y_i označavamo rastojanje $a_i \in H_-(c, c_0)$ od hiperravnih $H(c, c_0 + \mathbf{d})$, a sa z_j rastojanje $b_j \in H_+(c, c_0)$ od hiperravnih $H(c, c_0 - \mathbf{d})$. Klasifikaciju vršimo tako što minimiziramo težinsku sumu rastojanja pogrešno klasifikovanih objekata, tj. rešavamo sledeći problem linearne programiranja:

$$\begin{aligned} \min_{y,z} & \mathbf{a} \sum_{i=1}^m y_i + \mathbf{b} \sum_{j=1}^k z_j \\ c^T a_i + y_i & \leq c_0 + \mathbf{d}, \quad i=1,\dots,m \\ c^T b_j - z_j & \leq c_0 - \mathbf{d}, \quad j=1,\dots,k \\ y_i, z_j & \geq 0 \end{aligned}$$

gde su α , β težinski koeficijenti. Ovaj problem linearne programiranje ima optimalno rešenje $(c^*, c_0^*, y_1^*, \dots, y_m^*, z_1^*, \dots, z_k^*)$ i pri tome je $H(c^*, c_0^*)$ hiperravan sa najmanjom greškom klasifikacije.



3.2.2 Separacija pomocu kvadratne površi

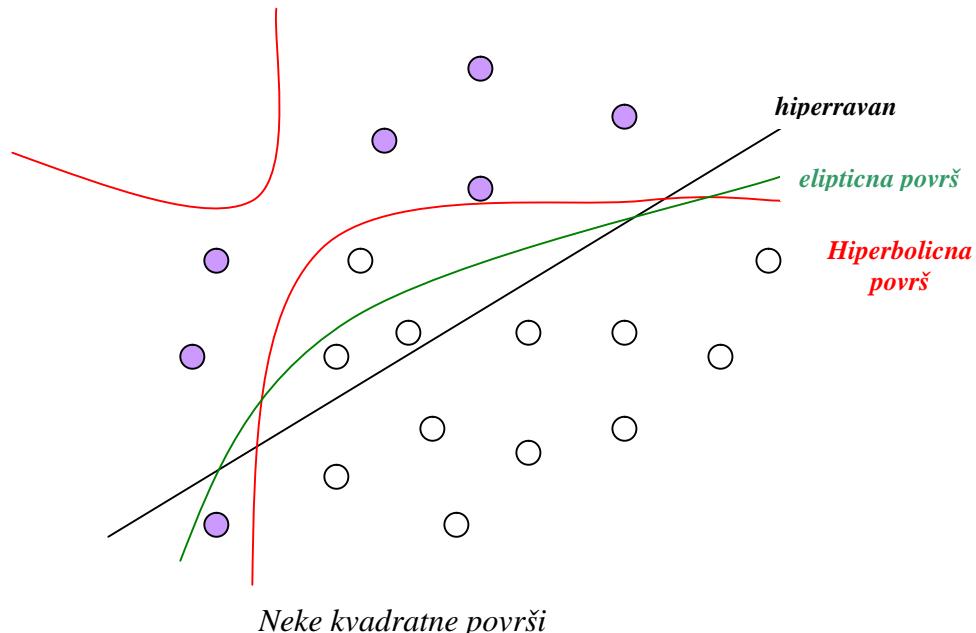
Ponekad nije moguce da se dva skupa objekata razdvoje pomocu ravni. Zbog toga se za separaciju koriste površi koje su nelinearne, takva je na primer kvadratna separacija. Neka je $x=(x_1, \dots, x_n)^T$ vektor koji predstavlja tacku u n -dimenzionalnom prostoru. Problem **kvadratne separacije** skupova A i B sastoji se u određivanju jedinstvene kvadratne površi:

$$x^T D x + x^T c - c_0 = 0$$

gde su $D \in \mathbf{R}^{n \times n}$, $c \in \mathbf{R}^n$, $c_0 \in \mathbf{R}$, tako da je

$$\begin{aligned} a_i^T D a_i + a_i^T c - c_0 &> 0, \quad i = 1, \dots, m \\ b_j^T D b_j + b_j^T c - c_0 &< 0, \quad j = 1, \dots, k \end{aligned}$$

Dva skupa A i B su **kvadratno separabilna** ako i samo ako postoje D , c i c_0 takvi da su gornje jednakosti zadovoljene.



Sledeći problem linearog programiranja minimizira grešku klasifikacije:

$$\begin{aligned} \min_{y,z} & \mathbf{a} \sum_{i=1}^m y_i + \mathbf{b} \sum_{j=1}^k z_j \\ & a_i^T D a_i + c^T a_i + y_i \geq c_0 + \mathbf{d}, i=1, \dots, m \\ & b_j^T D b_j + c^T b_j - z_j \leq c_0 - \mathbf{d}, j=1, \dots, k \\ & y_i, z_j \geq 0 \end{aligned}$$

gde su $D \in \mathbf{R}^{n \times n}$, $c \in \mathbf{R}^n$, $c_0 \in \mathbf{R}$, $y_i, z_j \in \mathbf{R}$ promenljive koje treba odrediti. Promenljive y_i, z_j predstavljaju rastojanje pogrešno klasifikovanih objekata od kvadratne površi.

3.2.3 Model klasifikacije uz pomoc funkcije korisnosti

Prvo se podsetimo pojma relacija preferencije i funkcije korisnosti.

Definicija Binarna relacija \succ na skupu X koja ima osobine:

- 1) refleksivnost: $\forall x \in X \quad x \succ x$
- 2) kompletност: $\forall x, y \in X \quad x \succ y \vee y \succ x$
- 3) tranzitivnost: $\forall x, y, z \in X \quad x \succ y \wedge y \succ z \Rightarrow x \succ z$

naziva se **relacija preferencije** na skupu X .

Definicija Neka je \succ relacija preferencije na skupu X . Funkcija $u: X \rightarrow \mathbf{R}$ za koju važi:

$$x \succ y \Leftrightarrow u(x) \geq u(y)$$

naziva se **funkcija korisnosti** relacije preferencije \succ .

Pošto je funkcija klasifikacije prekidna funkcija, za razdvajanje objekata koji pripadaju razlicitim klasama koristi se neprekidna funkcija korisnosti. Najbolja funkcija korisnosti neke klase se odreduje minimiziranjem greške klasifikacije. Ako tražimo funkciju korisnosti koja je linearna kombinacija nekih drugih funkcija onda problem možemo formulisati kao problem linearog programiranja.

Sada posmatramo problem klasifikacije u J klasa. Prepostavimo da postoji skup objekata $I = \{1, \dots, m\}$ sa poznatim klasifikacijama i svaki objekat je predstavljen tackom iz \mathbf{R}^n . Ovaj skup tacaka $X = \{x^i / i \in I\}$ se naziva *training skup*. Prepostavimo da razlaganje skupa I , $\{I_j\}_{j=1}^J$ ($I_{j_1} \cap I_{j_2} = \emptyset, j_1 \neq j_2, I = \bigcup_{j=1}^J I_j$) definiše klasifikaciju datih m objekata. Kažemo da tacka x^i pripada klasi k ako $i \in I_k$.

Ako je data funkcija korisnosti $u(x)$ koja je definisana na skupu $X = \{x^i / i \in I\}$ i skup pragova $U_p = \{u_0, u_1, \dots, u_{j-1}, u_j\}$ tada je $u_0 = -\infty < u_1 < \dots < u_{j-1} < u_j = \infty$ i ako važi

$$i_j = \{i \in I / u_{j-1} < u(x^i) < u_j\}$$

tada funkcija $u(x)$ i skup U_p definišu razlaganje skupa I , $\{I_j\}_{j=1}^J$ i na taj nacin klasificuju training skup. Medutim, cešci je slučaj da nam je poznata klasifikacija i , $\{I_j\}_{j=1}^J$, a da je potrebno pronaci funkciju korisnosti i skup pragova koji definišu ovu klasifikaciju. Neka $u^g(x)$ predstavlja klasu funkcija korisnosti koje su odredene parametrom $\gamma \in \mathbf{R}^n$. Pretpostavimo da je training skup $X = \{x^i / i \in I\}$ podeljen klasama $\{I_j\}_{j=1}^J$. Sa $F(\mathbf{s}^+, \mathbf{s}^-)$ označicemo ukupnu kaznu za sve pogrešno klasifikovane tacke. Za training skup X , sledeći problem optimizacije pronalazi najbolju funkciju korisnosti $u^g(x)$ u prethodno datoj klasi funkcija:

$$\begin{aligned} & \min_{\mathbf{g}, u, \mathbf{s}^+, \mathbf{s}^-} F(\mathbf{s}^+, \mathbf{s}^-) \\ & u_{j-1} - \mathbf{s}_i^- + \mathbf{d} \leq u^g(x^i) \leq u_j + \mathbf{s}_i^+, i \in I_j, j = 1, \dots, J \\ & u_j - u_{j-1} \leq s, j = 2, \dots, J-1 \\ & \mathbf{s}_i^+ \geq 0, \mathbf{s}_i^- \geq 0 \\ & \mathbf{g} \in \mathbf{R}^n \end{aligned}$$

Funkcija $F(\mathbf{s}^+, \mathbf{s}^-)$ je neopadajuća u odnosu na greške klasifikacije $\mathbf{s}_i^+, \mathbf{s}_i^-$, $i \in I$. Veliko odstupanje od savršene klasifikacije dovodi do velikih kazni. Posmatramo linearnu kaznenu funkciju

$$F(\mathbf{s}^+, \mathbf{s}^-) = \sum_{i=1}^m (\mathbf{a}_i^- \mathbf{s}_i^- + \mathbf{a}_i^+ \mathbf{s}_i^+)$$

gde su $\mathbf{a}_i^-, \mathbf{a}_i^+ \geq 0$ koeficijenti kazne za tacki i . Ako objekat x^i pripada klasi I_j gornji problem matematičkog programiranja implicira da je:

$$\begin{aligned} \mathbf{s}_i^+ &= \max\{0, u^g(x^i) - u_j\} \\ \mathbf{s}_i^- &= \max\{0, u_{j-1} + \mathbf{d} - u^g(x^i)\} \end{aligned}$$

Ako funkcija korisnosti $u^g(x)$ u tacki $x^i \in I_j$ premašuje gornji prag u_j , tada je greška \mathbf{s}_i^+ jednaka razlici izmedu $u^g(x^i)$ i gornjeg praga u_j . U suprotnom $\mathbf{s}_i^+ = 0$. Slicno, ako je vrednost funkcije korisnosti $u^g(x^i)$ ispod $u_{j-1} + \mathbf{d}$ tada je kazna \mathbf{s}_i^- jednaka razlici izmedu $u_{j-1} + \mathbf{d}$ i $u^g(x^i)$, inace je $\mathbf{s}_i^- = 0$.

3.2.4 Klasifikacija pomocu linearne funkcije korisnosti

Linearna funkcija korisnosti ima sledeći oblik:

$$u^c(x) = c^T x = \sum_{k=1}^n c_k x_k, c \in \mathbf{R}^n$$

Ako imamo linearnu kaznenu funkciju i linearnu funkciju korisnosti, tada se klasifikacija vrši rešavajuci sledeći problem:

$$\begin{aligned}
 & \min_{c, u, \mathbf{s}^-, \mathbf{s}^+} \sum_{i=1}^n (\mathbf{a}_i^- \mathbf{s}_i^- + \mathbf{a}_i^+ \mathbf{s}_i^+) \\
 & \sum_{k=1}^n c_k x_k^i - u_{j-1} + \mathbf{s}_i^- \geq \mathbf{s}, i \in I_j, j = 1, \dots, J \\
 & \sum_{k=1}^n c_k x_k^i - u_j - \mathbf{s}_i^+ \leq 0, i \in I_j, j = 1, \dots, J \\
 & u_j - u_{j-1} \leq s \\
 & \mathbf{s}_i^+ \geq 0, \mathbf{s}_i^- \geq 0, i = 1, \dots, m \\
 & c \in R^n
 \end{aligned}$$

3.2.5 Klasifikacija uz pomoc kvadratne funkcije korisnosti

Posmatramo sledecu kvadratnu funkciju korisnosti:

$$u^{D',c}(x) = x^T D' x + c^T x = \sum_{k=1}^n \sum_{l=1}^n d_{kl}^{\cdot} x_k x_l + \sum_{k=1}^n c_k x_k$$

Ova funkcija ima $n^2 + n$ elemenata, pošto matrica D' ima n^2 elemenata, a vektor c ima n koeficijenata. Ekvivalentna definicija kvadratne funkcije je sledeca:

$$u^{D',c}(x) = \sum_{k=1}^n d_{kk}^{\cdot} x_k^2 + \sum_{k=1}^{n-1} \sum_{l=k+1}^n (d_{kl}^{\cdot} + d_{lk}^{\cdot}) x_k x_l + \sum_{k=1}^n c_k x_k$$

Ako uvedemo nove oznake:

$$d_k = d_{kk}^{\cdot}$$

$$d_{kl} = d_{kl}^{\cdot} + d_{lk}^{\cdot}, k < l$$

kvadratna funkcija korisnosti može biti zapisana na sledeci nacin:

$$u^{D',c}(x) = \sum_{k=1}^{n-1} d_k x_k^2 + \sum_{k=1}^{n-1} \sum_{l=k+1}^n d_{kl} x_k x_l + \sum_{k=1}^n c_k x_k$$

Ako imamo linearu kaznenu funkciju i kvadratnu funkciju korisnosti rešavamo sledeci problem:

$$\begin{aligned}
 & \min_{D, c, u, \mathbf{s}^+, \mathbf{s}^-} F = \sum_{i=1}^m (\mathbf{a}_i^- \mathbf{s}_i^- + \mathbf{a}_i^+ \mathbf{s}_i^+) \\
 & \sum_{k=1}^n d_k x_k^i + \sum_{k=1}^{n-1} \sum_{l=k+1}^n d_{kl} x_k^i x_l^i + \sum_{k=1}^n c_k x_k^i - u_{j-1} + \mathbf{s}_i^- \geq \mathbf{d}, i \in I_j, j = 1, \dots, J \\
 & \sum_{k=1}^n d_k x_k^i + \sum_{k=1}^{n-1} \sum_{l=k+1}^n d_{kl} x_k^i x_l^i + \sum_{k=1}^n c_k x_k^i - u_j - \mathbf{s}_i^+ \leq 0, i \in I_j, j = 1, \dots, J \\
 & u_j - u_{j-1} \leq s \\
 & \mathbf{s}_i^+ \geq 0, \mathbf{s}_i^- \geq 0, i = 1, \dots, m, \quad d \in R^n, c \in R^n
 \end{aligned}$$

3.3 Logit analiza

Regresiona analiza se bavi proucavanjem odnosa izmedu zavisne i jedne ili više nezavisnih promenljivih. Zavisna promenljiva se uglavnom označava sa y , a nezavisne sa x_1, \dots, x_n . Jedan od mogucih ciljeva regresione analize je:

1. da oceni srednju odnosno ocekivanu vrednost zavisne promenljive, pri cemu su date vrednosti nezavisnih promenljivih
2. da testira hipoteze o prirodi zavisnosti izmedu promenljivih
3. da predvidi ocekivanu vrednost zavisne promenljive, pri cemu su date vrednosti nezavisnih promenljivih
4. kombinovano više prethodnih ciljeva

Linearna regresija je tip regresije koja prepostavlja da postoji linearna zavisnost izmedu zavisne i nezavisnih promenljivih. Jednodimenzionalna linearna regresija se koristi u slučaju kada imamo jednu zavisnu i jednu nezavisnu promenljivu. Ukoliko imamo više nezavisnih promenljivih koristimo višedimenzionalnu regresiju.

U slučaju **jednodimenzionalne regresije** posmatramo sledeću jednacinu:

$$y_i = \mathbf{a} + \mathbf{b}x_i + u_i, \quad i = 1, \dots, n$$

gde su α (*odsekak* ocekivana vrednost za y ukoliko je x jednako nula) i β (ocekvana vrednost za y ukoliko se x promeni za jednu jedinicu) *koeficijenti regresije* koje treba oceniti ako je dato n opservacija za y i x , a u je *greška regresije*. Da bi mogli da ocenimo parametre moramo da navedemo nekoliko prepostavki:

1. $E(u_i) = 0, \forall i$
2. $\text{var}(u_i) = s^2, \forall i$
3. u_i i u_j su nezavisne za svako $i \neq j$
4. u_i i x_j su nezavisne za svako i i j
5. u_i imaju normalnu raspodelu za svako i sa nula ocekivanjem i σ^2 varijansom (što sledi iz 1. i 2. prepostavke)
6. $\sum (x_i - \bar{x})^2 \neq 0$
7. $\sum (x_i - \bar{x})^2 < \infty$

Sada možemo da zapišemo kako izgleda *funkcija jednodimenzionalne regresije*:

$$\hat{y}_i = \hat{\mathbf{a}} + \hat{\mathbf{b}}x_i$$

gde je sa kapicom označena ocekivana vrednost, a znamo da je $E(u_i) = 0, \forall i$. Stohasticka verzija gornje jednacine je

$$y_i = \hat{\mathbf{a}} + \hat{\mathbf{b}}x_i + \hat{u}_i$$

gde \hat{u}_i predstavlja ocenjenu grešku koja se naziva *rezidual*. Iz ove jednacine se dobijaju ocene za parametre regresije α i β . Tri metode koje se najčešće koriste za ocenu parametara su:

1. metoda najmanjih kvadrata
2. metoda momenta
3. metoda maksimalne verovatnoće

Ovim metodama dobijamo nepristrasne parametre koji imaju najmanju varijansu u klasi linearnih nepristrasnih ocenjenih parametara. Mi u ovom radu necemo objašnjavati ove metode, sem treće koja se koristi pri oceni parametara u slučaju logit regresije.

U slučaju **višedimenzionalne regresije** se posmatra odnos zavisne y i više nezavisnih x_1, \dots, x_n promenljivih, pa imamo sledeći model:

$$y_i = \mathbf{a} + \mathbf{b}_1 x_{1i} + \mathbf{b}_2 x_{2i} + \dots + \mathbf{b}_k x_{ki} + u_i, i = 1, \dots, n$$

gde je α odsecak (ocekivana promena u y ako su svi $x_i, i = 1, \dots, n$ jednaki nula), a β_i parcijalni koeficijenti korelacije (ocekivana vrednost za y ukoliko dode do jedinicne promene u x_i pri cemu se pretpostavlja da su svi ostali x -evi nepromenjeni). Važe iste pretpostavke kao i u slučaju jednodimenzionalne regresije, samo što je dodata još jedna koja kaže da medu nezavisnim promenljivama nema linearne zavisnosti. Analogno imamo stohasticku jednacinu:

$$y_i = \hat{\mathbf{a}} + \hat{\mathbf{b}}_1 x_{1i} + \hat{\mathbf{b}}_2 x_{2i} + \dots + \hat{\mathbf{b}}_k x_{ki} + \hat{u}_i, i = 1, \dots, n$$

Metod najmanjih kvadrata daje ocene za $\alpha, \beta_1, \dots, \beta_k$ koje su nepristrasne i imaju najmanju varijansu u klasi svih linearnih nepristrasnih ocenjenih parametara.

Medutim, dešava se da je zavisna promenljiva binarna, tj. da uzima samo dve vrednosti, na primer uspeh i neuspeh. Postavlja se pitanje kako se modelira odnos nezavisnih promenljivih i binarne zavisne promenljive? Odgovor na to pitanje nam daje **logit regresija**.

Prepostavljemo da je Y Bernulijeva slučajna promenljiva koja uzima vrednosti 0 ili 1, u zavisnosti da li je ishod dobar ili loš. Verovatnoca da ce ishod biti loš u zavisnosti od datih nezavisnih promenljivih, tj. da ce Y biti jednako 1, se definiše kao $\pi = P(Y = 1 | X = x)$, a verovatnoca da ce ishod biti dobar sa $1 - \pi = P(Y = 0 | X = x)$. Posmatramo odnos (*odds*) ove dve verovatnoće:

$$odds(x) = \frac{P(Y = 1 | X = x)}{P(Y = 0 | X = x)} = \frac{p}{1-p}$$

Logaritam ovog odnosa, tzv. *logit* je linearna funkcija nezavisnih promenljivih $x_i, i = 0, 1, 2, \dots, n$:

$$\ln(odds(x)) = \mathbf{b}_0 + \mathbf{b}_1 x_1 + \dots + \mathbf{b}_n x_n$$

Dakle,

$$\ln\left(\frac{\mathbf{p}}{1-\mathbf{p}}\right) = \mathbf{b}_0 + \mathbf{b}_1 x_1 + \dots + \mathbf{b}_n x_n$$

Jednostavnim eksponencijalnim transformacijama dobijamo da je:

$$\mathbf{p} = \frac{\exp(\mathbf{b}_0 + \mathbf{b}_1 x_1 + \dots + \mathbf{b}_n x_n)}{1 + \exp(\mathbf{b}_0 + \mathbf{b}_1 x_1 + \dots + \mathbf{b}_n x_n)}$$

Ako je promenljiva $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})$, $i = 1, 2, \dots, N$, tada je verovatnoca pozitivnog ishoda jednaka:

$$\mathbf{p}(\mathbf{x}_i) = \frac{\exp\left(\sum_{k=1}^p \mathbf{b}_k x_{ik}\right)}{1 + \sum_{k=1}^p \mathbf{b}_k x_{ik}}$$

Prepostavke koje važe kod logit regresije su sledeće:

1. Y ima Bernulijevu raspodelu sa parametrom $\pi(x)$, tj.

$$Y : \begin{pmatrix} 0 & 1 \\ 1 - \mathbf{p}(x) & \mathbf{p}(x) \end{pmatrix}$$

2. Nijedna promenljiva od znacaja nije izostavljena i nijedna promenljiva koja nema znacaja nije ukljucena
3. Logaritam nezavisnih promenljivih i zavisna promenljiva su linearno nezavisne
4. Nema znacajne korelacije izmedu nezavisnih promenljivih

3.3.1 Selektioni kriterijumi

U situaciji kada imamo veliki broj nezavisnih promenljivih koje mogu, a i ne moraju biti relevantne za donošenje prepostavki o zavisnoj promenljivoj, korisno je imati mogucnost redukcije modela, tako da u njemu ostanu samo promenljive koje nam obezbeduju važne informacije o zavisnoj promenljivoj. Nije uvek trivijalno odluciti koju promenljivu treba ostaviti u modelu. **Maksimalni model** definišemo kao model koji sadrži sve nezavisne promenljive koje mogu biti prisutne u modelu. Neka k predstavlja maksimalni broj nezavisnih promenljivih, tada maksimalni model ima oblik:

$$Y_i = \mathbf{b}_0 + \mathbf{b}_1 x_{i1} + \dots + \mathbf{b}_k x_{ik} + \mathbf{e}_i$$

gde su x_{i1}, \dots, x_{ik} nezavisne promenljive, a \mathbf{e}_i su greške za koje važi da su nezavisne, da imaju normalnu raspodelu sa nula ocekivanjem i zajednickom varijansom.

Kada se definiše maksimalni model, važno je da on sadrži sve nezavisne promenljive koje mogu uticati na zavisnu promenljivu, ali se mora paziti da u model ne ude previše nezavisnih promenljivih koje nemaju znacaja. Ukoliko model sadrži previše nezavisnih promenljivih u poređenju sa brojem opservacija, standardne greške u ocenjenim parametrima regresije mogu da budu izuzetno velike, što dovodi do nepreciznih rezultata. Takođe, što je veci broj nezavisnih promenljivih to je veci rizik da dođe do medusobne korelacije izmedu promenljivih. U opštem slučaju, treba da se

uzme u obzir velicina uzorka, što je manji uzorak, to treba da bude manji i maksimalni model. Postoje razna pravila o tome kakav ovaj odnos treba da bude. Tako na primer jedno od takvih pravila kaže da bi trebalo da bude najmanje pet opservacija na jednu nezavisnu promenljivu, tj. $n \geq 5k$.

Kada je definisan maksimalan model, sledeći korak je uporediti dva modela i odrediti koji je od njih bolji. U tu svrhu koristimo **selekcione kriterijume**, ciji je zadatak da porede maksimalni model

$$Y = \mathbf{b}_0 + \mathbf{b}_1 x_1 + \dots + \mathbf{b}_m x_m + \mathbf{b}_{m+1} x_{m+1} + \dots + \mathbf{b}_k x_k + \mathbf{e}$$

sa redukovanim modelom

$$Y = \mathbf{b}_0 + \mathbf{b}_1 x_1 + \dots + \mathbf{b}_m x_m$$

koji se dobija od maksimalnog modela. Cilj je videti da li redukovani model odgovara podacima podjednako dobro kao maksimalni model, i u tom slučaju cemo se odluciti da koristimo redukovani model umesto maksimalnog. Sada cemo navesti par selekcionih kriterijuma.

- R_a^2 **kriterijum**. Koeficijent determinacije R^2 odreduje koliko dobro model odgovara podacima i on se racuna na sledeći nacin:

$$R^2 = \frac{S_{yy} - RSS}{S_{yy}}$$

gde je $S_{yy} = \sum_{i=1}^n (Y_i - \bar{Y})^2$, $RSS = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$ je suma kvadrata reziduala. Ako sa R_m^2 oznamo sumu kvadrata reziduala maksimalnog modela, a sa R_j^2 sumu kvadrata reziduala redukovaniog modela, tada je

$$R_j^2 = \frac{S_{yy} - RSS_j}{S_{yy}}$$

gde je

$$RSS_j = \sum_{i=1}^n (Y_i - \hat{\mathbf{b}}_{j0} - \hat{\mathbf{b}}_{j1}x_{i1} - \hat{\mathbf{b}}_{j2}x_{i2} - \dots - \hat{\mathbf{b}}_{jj}x_{ij})^2 \quad (*)$$

gde $\hat{\mathbf{b}}_{ji}$ označava ocenu za parametar \mathbf{b}_{ji} koja se dobija primenom metode najmanjih kvadrata. Što model bolje odgovara podacima, to je veci R^2 . Dakle, jedna od mogucnosti za poredenje dva modela je da se uporede odgovarajuci koeficijenti determinacije, pri cemu je model sa vecim koeficijentom determinacije precizniji. Međutim, kao problem sa javlja cinjenica da model sa vecim brojem nezavisnih promenljivih ima veci ovaj koeficijent, nezavisno od toga kakav uticaj te nezavisne promenljive imaju na zavisnu promenljivu. Da bi se ovaj problem izbegao uvodi se pojam *prilagodenog koeficijenta determinacije* R_a^2 koji se dobija na sledeći nacin:

$$R_a^2 = 1 - \frac{n-1}{n-k-1} (1 - R^2)$$

Kao što možemo da vidimo, sa porastom broja promenljivih, ne mora da znaci da će i R_a^2 porasti. Prema R_a^2 kriterijumu zaključujemo da *treba da se izabere model sa najvećim R_a^2 .*

- **F-test kriterijum.** Ideja ovog kriterijuma je testirati znacajnost $k-m$ nezavisnih promenljivih x_{m+1}, \dots, x_k u maksimalnom modelu, sa ciljem da se dobije redukovani model. Dakle, treba da testiramo nultu hipotezu:

$$H_0 : \mathbf{b}_{m+1} = \dots = \mathbf{b}_k = 0$$

F-test statistika je data sa:

$$F_m = \frac{(RSS_m - RSS_k)/(k-m)}{RSS_k/(n-k-1)}$$

gde su RSS_m i RSS_k definisane sa (*). Ako se hipoteza H_0 prihvati, redukovani model odgovara podacima podjednako dobro kao maksimalni, pa samim tim možemo da koristimo redukovani umesto maksimalnog modela. *F-test kriterijum za selekovanje promenljivih nalazi najmanji podskup nezavisnih promenljivih x_1, x_2, \dots, x_m takvih da test statistika F_m nije znacajna.*

Kada smo se upoznali sa selepcionim kriterijumima, možemo da navedemo i **selekciione metode** koje koriste ove se kriterijume da bi odredile da li se u modelu nalazi optimalan broj promenljivih.

Jedna ovakva metoda je tzv. **eliminacija unazad**. Ova metoda je u suštini niz testova znacajnosti nezavisnih promenljivih. Dakle, zaključujemo da ona koristi *F-test kriterijum*. Eliminacija unazad pocinje sa maksimalnim modelom i u svakom koraku eliminiše promenljivu sa najvećom *p-vrednosti*, pri cemu je unapred određen nivo znacajnosti. Ova metoda se završava onda kada nemamo mogucnost više nijednu promenljivu da isključimo iz modela.

Druga metoda je **metoda unapred**, koja na kraju treba da da iste rezultate kao i eliminacija unazad. Ova metoda pocinje sa praznim modelom i u svakom koraku dodaje promenljivu koja je najznacajnija i tako sve dok se dode do situacije da se više nema koristi od dodavanja novih promenljivih.

Kombinacija ove dve metode je predstavljena u **metodi po etapama**. U prvom koraku ove metode se dodaje promenljiva promenljiva sa najvećom vrednošću koeficijenta determinacije R^2 . Koeficijenti determinacije ostalih promenljivih se ispituju da se vidi da li one obezbeduju neke dodatne informacije za model. Kada je dodata nova promenljiva u model, pomocu *F-testa* se proverava znacajnost promenljivih koje se u tom trenutku nalaze u modelu. One promenljive koje ulaskom nove promenljive nisu više znacajne, izlaze iz modela. Može se desiti da promenljiva koja ušla u model u prethodnom koraku, u nekom od sledećih koraka izgubi znacaj, i u tom slučaju ona izlazi iz modela. Ovaj proces se nastavlja sve dok se ne dode do toga da nema više koristi dodavati nove promenljive.

3.3.2 Metod maksimalne verovatnoće

Kada imamo više od jednog opažanja za promenljvu \mathbf{x}_i , tada nam je bitno da znamo broj opažanja n_i i broj uspeha koje označavamo sa y_i . Tada su $\{Y_1, \dots, Y_N\}$ nezavisne binarne promenljive sa očekivanjem $E(Y_i) = n_i p(\mathbf{x}_i)$, gde je $n_1 + \dots + n_N = n$. Tada je njihova zajednicka funkcija verovatnoće jednaka:

$$L(\mathbf{b}) = \prod_{i=1}^N p(x_i)^{y_i} (1-p(x_i))^{n_i - y_i}$$

Primenom metode maksimalne verovatnoće dobijamo ocene za β . Kao što sama reč kaže, traži se maksimalna vrednost ove funkcije, koja se dobija kada prvi izvod izjednacimo sa nulom. Dobijamo sistem nelinearnih jednacina po β_k koji treba da rešimo:

$$\sum_{i=1}^N y_i x_{ik} - \sum_{i=1}^N n_i \hat{p}_i x_{ik} = 0, \quad k = 1, \dots, p \quad (*)$$

gde je

$$\hat{p}_i = \frac{\exp(\sum_{k=1}^p \hat{b}_k x_{ik})}{1 + \exp(\sum_{k=1}^p \hat{b}_k x_{ik})}$$

ocena za $\pi(\mathbf{x}_i)$. Za rešavanje ovog sistema koristi se neki iterativni postupak, najčešće Newton-Raphson algoritam.

Kada ovaj sistem rešimo po β_k , svako ovakvo rešenje, ako postoji, određuje jednu kriticnu tacku, ili minimum ili maksimum. Kriticna tacka će biti maksimum ako je matrica drugih parcijalnih izvoda negativno definitna, tj. ako je svaki element na dijagonali matrice manji od nule. Ova matrica će imati elemente sledećeg oblika:

$$-\sum_{i=1}^n x_{ik} x_{ia} n_i p_i (1 - p_i)$$

Matrica $-\frac{\partial^2 L(\mathbf{b})}{\partial \mathbf{b}_a \partial \mathbf{b}_b}$ se naziva *Fisher-ova matrica podataka*.

Matrica ocenjenih kovarijansi ima sledeći oblik:

$$\text{cov}(\hat{\mathbf{b}}) = (\mathbf{X}' \text{diag}(n_i \hat{p}_i (1 - \hat{p}_i)) \mathbf{X})^{-1}$$

gde $\text{diag}(n_i \hat{p}_i (1 - \hat{p}_i))$ predstavlja $N \times N$ dijagonalnu matricu sa elementima $(n_i \hat{p}_i (1 - \hat{p}_i))$ na glavnoj dijagonali. Kvadratni koren elemenata sa glavne dijagonale matrice kovarijansi su *ocenjene standardne greške koeficijenata* $\hat{\mathbf{b}}$.

Rezultati logit regresije mogu da se interpretiraju na više nacija. Kao prvo, ako posmatramo jednacinu u obliku

$$\ln\left(\frac{\hat{p}}{1 - \hat{p}}\right) = \mathbf{b}_0 + \mathbf{b}_1 x_1 + \dots + \mathbf{b}_n x_n$$

tada se sa jedinicnim porastom (smanjenjem) nezavisne promenljive x_i logaritam odnosa verovatnoca povecava (smanjuje) za β_i . Naravno, u slučaju kada su sve nezavisne promenljive jednake nuli, ovaj odnos ce biti jednak konstanti β_0 . Međutim, ova jednacina može da se transformiše eksponenciranjem:

$$\frac{\hat{p}}{1-\hat{p}} = \exp(\mathbf{b}_0 + \mathbf{b}_1 x_1 + \dots + \mathbf{b}_n x_n) = \exp(\mathbf{b}_0) * \exp(\mathbf{b}_1 x_1) * \dots * \exp(\mathbf{b}_n x_n)$$

u ovom slučaju ako x_i poraste za jednu jedinicu, ocjenjeni odnos ce porasti $\exp(\beta_i)$ puta. Ukoliko su sve nezavisne promenljive jednake nuli, tada ce ovaj odnos biti jednak $\exp(\beta_0)$. A ako posmatramo sledecu jednacinu:

$$\hat{p} = \frac{\exp\left(\sum_{k=1}^p \hat{\mathbf{b}}_k x_{ik}\right)}{1 + \exp\left(\sum_{k=1}^p \hat{\mathbf{b}}_k x_{ik}\right)}$$

nije lako protumaciti rezultate logit regresije.

3.3.3 Newton-Raphson algoritam (NRA)

Newton-Raphson algoritam je iterativni postupak za rešavanje nelinearnih jednacina. Sada cemo pokazati kako NRA odreduje vrednost $\hat{\mathbf{b}}$ kada se maksimizira funkcija $L(\mathbf{b})$. Neka je $\mathbf{u}' = (\partial L(\mathbf{b}) / \partial \mathbf{b}_1, \dots, \partial L(\mathbf{b}) / \partial \mathbf{b}_p)$, a sa \mathbf{H} označavamo Hesijan matricu koja ima sledeće elemente $h_{ab} = \partial^2 L(\mathbf{b}) / \partial \mathbf{b}_a \partial \mathbf{b}_b$. Neka je sa $\mathbf{u}^{(t)}$ i $\mathbf{H}^{(t)}$ označeno \mathbf{u} i \mathbf{H} u $\mathbf{b}^{(t)}$, što predstavlja t -ti pokušaj za $\hat{\mathbf{b}}$. Korak t iterativnog procesa ($t = 0, 1, 2, \dots$) koji aproksimira $L(\mathbf{b})$ u blizini $\mathbf{b}^{(t)}$ Tejlorovim polinomom drugog reda je:

$$L(\mathbf{b}) \approx L(\mathbf{b}^{(t)}) + \mathbf{u}^{(t)'}(\mathbf{b} - \mathbf{b}^{(t)}) + \frac{1}{2}(\mathbf{b} - \mathbf{b}^{(t)})' \mathbf{H}^{(t)} (\mathbf{b} - \mathbf{b}^{(t)})$$

Rešavajuci $\partial L(\mathbf{b}) / \partial \mathbf{b} \approx \mathbf{u}^{(t)} + \mathbf{H}^{(t)}(\mathbf{b} - \mathbf{b}^{(t)}) = 0$ po \mathbf{b} dolazimo do sledeće tacke, što može da se zapiše na sledeći nacin:

$$\mathbf{b}^{(t+1)} = \mathbf{b}^{(t)} - (\mathbf{H}^{(t)})^{-1} \mathbf{u}^{(t)} \quad (*)$$

i prepostavljamo da $\mathbf{H}^{(t)}$ nije singularna matrica.

U slučaju logaritamske regresije imamo sledeći slučaj:

$$u_j^{(t)} = \frac{\partial L(\mathbf{b})}{\partial \mathbf{b}_j} \Big|_{\mathbf{b}^{(t)}} = \sum_i (y_i - n_i \mathbf{p}_i^{(t)}) x_{ij}$$

$$h_{ab}^{(t)} = \frac{\partial^2 L(\mathbf{b})}{\partial \mathbf{b}_a \partial \mathbf{b}_b} \Big|_{\mathbf{b}^{(t)}} = - \sum_i x_{ia} x_{ib} n_i \mathbf{p}_i^{(t)} (1 - \mathbf{p}_i^{(t)})$$

gde je $\mathbf{p}^{(t)}$ t -ta aproksimacija za $\hat{\mathbf{p}}$ dobijena uz pomoc $\mathbf{b}^{(t)}$

$$\mathbf{p}_i^{(t)} = \frac{\exp(\sum_{k=1}^p \mathbf{b}_k^{(t)} x_{ik})}{1 + \exp(\sum_{k=1}^p \mathbf{b}_k^{(t)} x_{ik})}$$

Koristimo jednacinu (*) da bi dobili sledecu vrednost $\mathbf{b}^{(t+1)}$:

$$\mathbf{b}^{(t+1)} = \mathbf{b}^{(t)} + (X' \text{diag}(n_i \mathbf{p}_i^{(t)} (1 - \mathbf{p}_i^{(t)})) X)^{-1} X' (y - \mathbf{m}^{(t)})$$

gde je $\mathbf{m}^{(t)} = n_i \mathbf{p}_i^{(t)}$. Ova vrednost se dalje koristi da se dobije $\mathbf{p}^{(t+1)}$ i tako dalje. Vrednosti $\mathbf{p}^{(t)}$ i $\mathbf{b}^{(t)}$ konvergiraju ka ocenama koje se dobijaju uz pomoc funkcije maksimalne verovatnoće $\hat{\mathbf{p}}$ i $\hat{\mathbf{b}}$. A Hesijan $\mathbf{H}^{(t)}$ konvergira ka $\hat{H} = -X' \text{diag}(n_i \hat{\mathbf{p}}_i (1 - \hat{\mathbf{p}}_i)) X$.

3.3.4 Testiranje hipoteza

Za testiranje hipoteza koristi se **Waldov test** koji se najčešće koristi da pokaže da li efekat postoji ili ne, tj. on pokazuje da li nezavisna promenljiva ima statistički značajan odnos sa zavisnom promenljivom. Waldov test poređi ocene $\hat{\mathbf{b}}$ parametara od značaja β sa predloženim vrednostima u nultoj hipotezi β_0 , pod pretpostavkom da će se razlike ove dve vrednosti aproksimirati normalnom raspodelom, pa se kvadrat ove razlike aproksimira χ^2 -raspodelom. Dakle, Waldova statistika ima sledeći oblik:

$$\frac{(\hat{\mathbf{b}}_j - \mathbf{b}_{j0})^2}{\text{var}(\hat{\mathbf{b}})} \sim \mathbf{c}_1^2$$

$$\text{ili } \frac{\hat{\mathbf{b}}_j - \mathbf{b}_{j0}}{\hat{s}e(\hat{\mathbf{b}})} \sim N(0,1)$$

gde se $\hat{s}e(\hat{\mathbf{b}})$ dobija tako što se uzme inverzni element iz ocenjene Fisher-ove matrice podataka. Ukoliko testiramo hipotezu za više parametara u isto vreme, pod pretpostavkom da je $H_0: \beta = \beta_0$, Waldova statistika ima sledeći oblik:

$$W = (\hat{\mathbf{b}} - \mathbf{b}_0)' [\text{cov}(\hat{\mathbf{b}})]^{-1} (\hat{\mathbf{b}} - \mathbf{b}_0)$$

koja prati χ^2 raspodelu, a broj stepena slobode je jednak rangu matrice kovarijansi.

Drugi test koji se koristi za testiranje parametara je **test odnosa verovatnoca (likelihood-ratio test)**. Obično se obeležava sa grčkim slovom Λ (velikim lambda) i predstavlja odnos dve maksimalne vrednosti: prva je maksimum dobijen od parametara pod pretpostavkom nulte hipoteze koji se označava sa l_0 i opštег maksimuma, tj. maksimuma dobijenog pod pretpostavkom $H_0 \bigcup H_A$ i taj maksimum označavamo sa l_1 . Odnos $\Lambda = l_0 / l_1$ ne može da bude veci od jedan. Wilks je pokazao da $-2\log\Lambda$ prati χ^2 raspodelu gde je broj stepeni slobode jednak razlici dimenzija parametara pod pretpostavkama $H_0 \bigcup H_A$ i H_0 . Dakle,

$$-2\log\Lambda = -2\log(l_0 / l_1) = -2(L_0 - L_1)$$

Treci nacin za testiranje hipoteza je **skor test**. To je najmocniji test kada je prava vrednost parametra blizu ocenjene vrednosti. Neka je L funkcija verovatnoce koja zavisi od parametra β , tada je skor funkcija:

$$u(\mathbf{b}) = \frac{\partial L(\mathbf{b})}{\partial \mathbf{b}}$$

Fisher-ova matrica podataka je:

$$I(\mathbf{b}) = -\frac{\partial^2 L(\mathbf{b})}{\partial \mathbf{b}^2}$$

Tada je statistika koja testira hipotezu $H_0: \beta = \beta_0$ data sa:

$$S(\mathbf{b}) = \frac{u(\mathbf{b}_0)^2}{I(\mathbf{b}_0)} \sim \mathbf{c}_1^2$$

3.3.5 Konstruisanje intervala poverenja

Neka z_a predstavlja z -vrednost standardizovane normalne raspodele koja ima verovatnocu a , tj. to je $100(1-a)\%$ rapodele. Neka $\chi^2_{df}(a)$ predstavlja $100(1-a)\%$ χ^2 -raspodele sa df stepeni slobode. Tada je interval poverenja za Waldov test skup β_0 za koje je:

$$\frac{|\hat{\mathbf{b}} - \mathbf{b}_0|}{SE} < z_{a/2}$$

Što daje sledeci interval: $\hat{\mathbf{b}} \pm z_{a/2} SE$.

Za test odnosa verovatnoca interval poverenja je skup β_0 za koje je

$$-2[L(\mathbf{b}_0) - L(\hat{\mathbf{b}})] < \mathbf{c}_1^2(a)$$

Podsetimo se da je $\mathbf{c}_1^2(a) = z_{a/2}^2$.

Dok je za skor test interval poverenja skup β_0 za koje je

$$\left| \frac{u(\mathbf{b}_0)^2}{I(\mathbf{b}_0)} \right| < \mathbf{c}_1^2(a)$$

3.3.6 Primena softvera

Uz pomoc statistickog paket *Statgraphics* cemo prikazati jedan jednostavan primer logit regresije. Prepostavimo da imamo zavisnu promenljivu *Loš* koja uzima vrednost 1 ukoliko je klijent loš i vrednost 0 ako je klijent dobar, i pet nezavisnih promenljivih:

- ❖ *Starost*: koja pokazuje koliko klijent ima godina
- ❖ *Muško*: binarna promenljiva koja uzima vrednost 1 ukoliko je klijent muško
- ❖ *Prihodi*: koja predstavlja klijentove prihode
- ❖ *Oženjen*: binarna promenljiva koja uzima vrednost 1 ukoliko je klijent oženjen
- ❖ *Vlasnik stana*: binarna promenljiva koja uzima vrednost 1 ukoliko je klijent vlasnik stana u kome živi

	Loš	Starost	Musko	Prihodi	Oženjen	Vlasnik stana
1	1	25	1	500	0	1
2	0	40	1	1000	1	0
3	1	45	1	1500	1	0
4	0	21	0	300	0	0
5	0	25	1	300	0	0
6	1	35	0	550	1	1
7	0	30	0	600	0	0
8	1	40	1	1200	1	1
9	1	42	0	1200	1	0
10	1	37	1	700	0	1
11	0	51	0	1500	1	1
12	1	22	1	600	0	1
13	0	21	1	300	0	0
14	0	29	1	900	0	0
15	1	33	1	1000	1	1

Promenljive koje ulaze u model

Kada smo uneli ove podatke u radni list, možemo da primenimo logit regresiju i da vidimo kakve cemo rezultate dobiti. Prvo dobijamo prozor u kome se nalazi tabela sa ocenjenim parametrima regresije uz pomoc funkcije maksimalne verovatnoće, njihovim standardnim greškama i ocenjenim «*odds*» odnosom koji se dobija na sledeći nacin:

$$odds_ratio = \exp(\hat{b}_i)$$

i on predstavlja procenat u kome se povecava $odds = \frac{p}{1-p}$ ukoliko dode do jedinicnog poveca promenljive *X*. Takode dobijemo i tabelu sa rezultatima testa znacajnosti promenljivih. Kao što možemo da primetimo, ovaj paket koristi test odnosa verovatnoca da testira znacajnost promenljivih u modelu. Na kraju imamo kratak opis rezultata, kao i regresionu jednacinu.

Logistic Regression - Los

Dependent variable: Los

Factors:

Musko
Ozenjen
Prihodi
Starost
Vlasnik stana

Estimated Regression Model (Maximum Likelihood)

Parameter	Estimate	Standard Error	Estimated Odds Ratio
CONSTANT	0.913344	5.13631	
Musko	1.37705	1.86523	3.96318
Ozenjen	3.32295	2.91391	27.742
Prihodi	0.00275722	0.00472383	1.00276
Starost	-0.20772	0.254429	0.812435
Vlasnik stana	3.53999	1.94519	34.4667

Analysis of Deviance

Source	Deviance	Df	P-Value
Model	8.85067	5	0.1152
Residual	11.877	9	0.2203
Total (corr.)	20.7277	14	

Percentage of deviance explained by model = 42.6997

Adjusted percentage = 0.0

Likelihood Ratio Tests

Factor	Chi-Squared	Df	P-Value
Musko	0.614137	1	0.4332
Ozenjen	1.58047	1	0.2087
Prihodi	0.357043	1	0.5502
Starost	0.797817	1	0.3717
Vlasnik stana	5.79497	1	0.0161

The StatAdvisor

The output shows the results of fitting a logistic regression model to describe the relationship between Los and 5 independent variable(s). The equation of the fitted model is

$$\text{Los} = \exp(\eta) / (1 + \exp(\eta))$$

where

$$\eta = 0.913344 + 1.37705 * \text{Musko} + 3.32295 * \text{Ozenjen} + 0.00275722 * \text{Prihodi} - 0.20772 * \text{Starost} + 3.53999 * \text{Vlasnik stana}$$

Because the P-value for the model in the Analysis of Deviance table is greater or equal to 0.05, there is not a statistically significant relationship between the variables at the 95.0% or higher confidence level. In addition, the P-value for the residuals is greater than or equal to 0.05, indicating that the model is not significantly worse than the best possible model for this data at the 95.0% or higher confidence level.

The pane also shows that the percentage of deviance in Los explained by the model equals 42.6997%. This statistic is similar to the usual R-Squared statistic. The adjusted percentage, which is more suitable for comparing models with different numbers of independent variables, is 0.0%.

In determining whether the model can be simplified, notice that the highest P-value for the likelihood ratio tests is 0.5502, belonging to Prihodi. Because the P-value is greater or equal to 0.05, that term is not statistically significant at the 95.0% or higher confidence level. Consequently, you should consider removing Prihodi from the model.

Statgraphics nam daje opciju da vidimo i koliki su 95% intervali poverenja za ocenjene parametre.

95.0% confidence intervals for coefficient estimates

Parameter	Estimate	Standard Error	Lower Limit	Upper Limit
CONSTANT	0.913344	5.13631	-10.7058	12.5325
Musko	1.37705	1.86523	-2.8424	5.5965
Ozenjen	3.32295	2.91391	-3.2688	9.91469
Prihodi	0.00275722	0.00472383	-0.00792885	0.0134433
Starost	-0.20772	0.254429	-0.783279	0.36784
Vlasnik stana	3.53999	1.94519	-0.860334	7.94032

95.0% confidence intervals for odds ratios

Parameter	Estimate	Lower Limit	Upper Limit
Musko	3.96318	0.0582854	269.481
Ozenjen	27.742	0.0380521	20225.4
Prihodi	1.00276	0.992102	1.01353
Starost	0.812435	0.456905	1.44461
Vlasnik stana	34.4667	0.423021	2808.26

The StatAdvisor

This table shows 95.0% confidence intervals for the coefficients in the model. Confidence intervals show how precisely the coefficients can be estimated given the amount of available data and the noise which is present. Also shown are confidence intervals for the odds ratios. The odds ratio equals the inverse natural logarithm of the coefficient and shows the proportional change in the response variable as the independent variable is increased by 1 unit.

Da bi videli da li postoji odredena medusobna korelacija izmedu promenljivih možemo da koristimo opciju za dobijanje korelaceione matrice:

Correlation matrix for coefficient estimates					
	CONSTANT	Musko	Ozenjen	Prihodi	
CONSTANT	1.0000	-0.3334	0.4237	0.5799	
Musko	-0.3334	1.0000	0.3533	-0.1267	
Ozenjen	0.4237	0.3533	1.0000	0.0392	
Prihodi	0.5799	-0.1267	0.0392	1.0000	
Starost	-0.8735	0.0267	-0.5098	-0.8154	
Vlasnik stana	0.3512	0.2026	0.4478	0.4347	

	Starost	Vlasnik stana
CONSTANT	-0.8735	0.3512
Musko	0.0267	0.2026
Ozenjen	-0.5098	0.4478
Prihodi	-0.8154	0.4347
Starost	1.0000	-0.5831
Vlasnik stana	-0.5831	1.0000

The StatAdvisor

This table shows estimated correlations between the coefficients in the fitted model. These correlations can be used to detect the presence of serious multicollinearity, i.e., correlation amongst the predictor variables. In this case, there are 3 correlations with absolute values greater than 0.5.

Takode, imamo mogucnost primene ranije opisanih selekcionih metoda i to metode unapred i metode eliminacije unazad i za naš primer dobijamo sledece rezultate. Za metod unapred:

Logistic Regression - Los			
Estimated Regression Model (Maximum Likelihood)			
Parameter	Estimate	Standard Error	Estimated Odds Ratio
CONSTANT	-1.09861	0.816497	
Vlasnik stana	2.89037	1.35398	18.0

Analysis of Deviance			
Source	Deviance	Df	P-Value
Model	5.98871	1	0.0144
Residual	14.739	13	0.3239
Total (corr.)	20.7277	14	

Percentage of deviance explained by model = **28.8923**
Adjusted percentage = **9.59445**

Likelihood Ratio Tests

Factor	Chi-Squared	Df	P-Value
Vlasnik stana	5.98871	1	0.0144

Stepwise factor selection
Method: forward selection
P-to-enter: 0.05
P-to-remove: 0.05
Step 0:
0 factors in the model. 14 d.f. for error.
Percentage of deviance explained = 0.00% Adjusted percentage = 0.00%
Step 1:
Adding factor Vlasnik stana with P-to-enter = 0.0143954
1 factors in the model. 13 d.f. for error.
Percentage of deviance explained = 28.89% Adjusted percentage = 9.59%

Final model selected.

The StatAdvisor

The output shows the results of fitting a logistic regression model to describe the relationship between Los and 5 independent variable(s). The equation of the fitted model is

$$\text{Los} = \exp(\eta)/(1+\exp(\eta))$$

where

$$\eta = -1.09861 + 2.89037 * \text{Vlasnik stana}$$

Because the P-value for the model in the Analysis of Deviance table is less than 0.05, there is a statistically significant relationship between the variables at the 95.0% confidence level. In addition, the P-value for the residuals is greater than or equal to 0.05, indicating that the model is not significantly worse than the best possible model for this data at the 95.0% or higher confidence level.

The pane also shows that the percentage of deviance in Los explained by the model equals 28.8923%. This statistic is similar to the usual R-Squared statistic. The adjusted percentage, which is more suitable for comparing models with different numbers of independent variables, is 9.59445%.

In determining whether the model can be simplified, notice that the highest P-value for the likelihood ratio tests is 0.0144, belonging to Vlasnik stana. Because the P-value is less than 0.05, that term is statistically significant at the 95.0% confidence level. Consequently, you probably don't want to remove any variables from the model.

A za metod eliminacije unazad:

Logistic Regression - Los

Estimated Regression Model (Maximum Likelihood)

Parameter	Estimate	Standard	Estimated
		Error	Odds Ratio
CONSTANT	-1.09861	0.816497	
Vlasnik stana	2.89037	1.35398	18.0

Analysis of Deviance

Source	Deviance	Df	P-Value
Model	5.98871	1	0.0144
Residual	14.739	13	0.3239
Total (corr.)	20.7277	14	

Percentage of deviance explained by model = 28.8923

Adjusted percentage = 9.59445%

Likelihood Ratio Tests

Factor	Chi-Squared	Df	P-Value
Vlasnik stana	5.98871	1	0.0144

Stepwise factor selection

Method: backward selection

P-to-enter: 0.05

P-to-remove: 0.05

Step 0:

5 factors in the model. 9 d.f. for error.

Percentage of deviance explained = 42.70% Adjusted percentage = 0.00%

Step 1:

Removing factor Prihodi with P-to-remove = 0.550152

4 factors in the model. 10 d.f. for error.

Percentage of deviance explained = 40.98% Adjusted percentage = 0.00%

Step 2:

Removing factor Starost with P-to-remove = 0.496249

3 factors in the model. 11 d.f. for error.

Percentage of deviance explained = 38.74% Adjusted percentage = 0.15%

Step 3:

Removing factor Musko with P-to-remove = 0.355254

2 factors in the model. 12 d.f. for error.

Percentage of deviance explained = 34.62% Adjusted percentage = 5.67%

Step 4:

Removing factor Ozenjen with P-to-remove = 0.275858

1 factors in the model. 13 d.f. for error.

Percentage of deviance explained = 28.89% Adjusted percentage = 9.59%

Final model selected.

The StatAdvisor

The output shows the results of fitting a logistic regression model to describe the relationship between Los and 5 independent variable(s). The equation of the fitted model is

$$\text{Los} = \exp(\eta) / (1 + \exp(\eta))$$

where

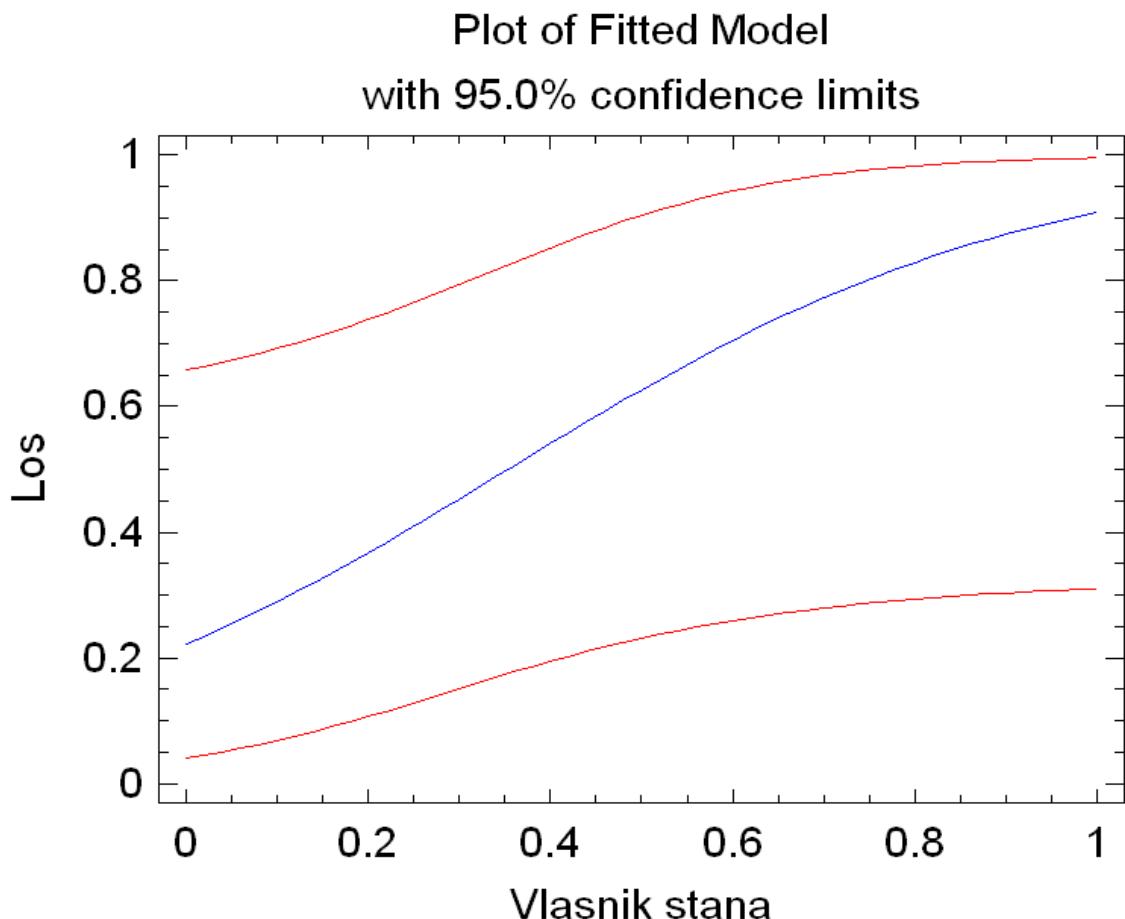
$$\eta = -1.09861 + 2.89037 * \text{Vlasnik stana}$$

Because the P-value for the model in the Analysis of Deviance table is less than 0.05, there is a statistically significant relationship between the variables at the 95.0% confidence level. In addition, the P-value for the residuals is greater than or equal to 0.05, indicating that the model is not significantly worse than the best possible model for this data at the 95.0% or higher confidence level.

The pane also shows that the percentage of deviance in Los explained by the model equals 28.8923%. This statistic is similar to the usual R-Squared statistic. The adjusted percentage, which is more suitable for comparing models with different numbers of independent variables, is 9.59445%.

In determining whether the model can be simplified, notice that the highest P-value for the likelihood ratio tests is 0.0144, belonging to Vlasnik stana. Because the P-value is less than 0.05, that term is statistically significant at the 95.0% confidence level. Consequently, you probably don't want to remove any variables from the model.

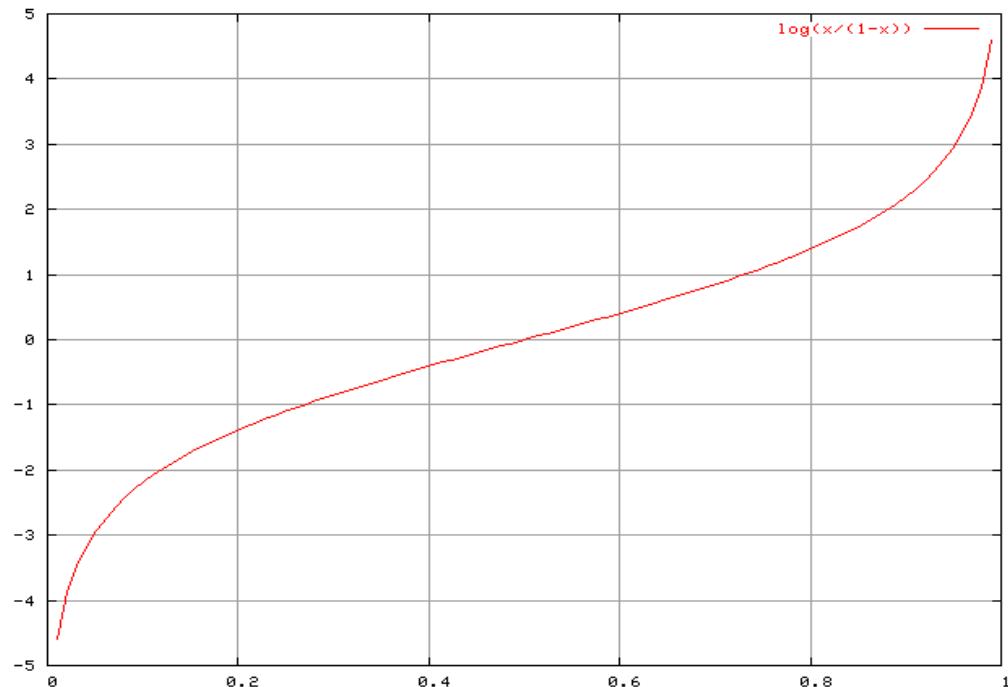
Statgraphics nam omogucava da graficki prikažemo ocenjenu verovatnocu da je klijent loš u odnosu na pojedinacnu nezavisnu promenljivu, dok ostale nezavine promenljive ostaju konstantne. Mi necemo dati grafike za sve nezavisne promenljive, nego samo za jednu, na primer *Vlasnik stana*:



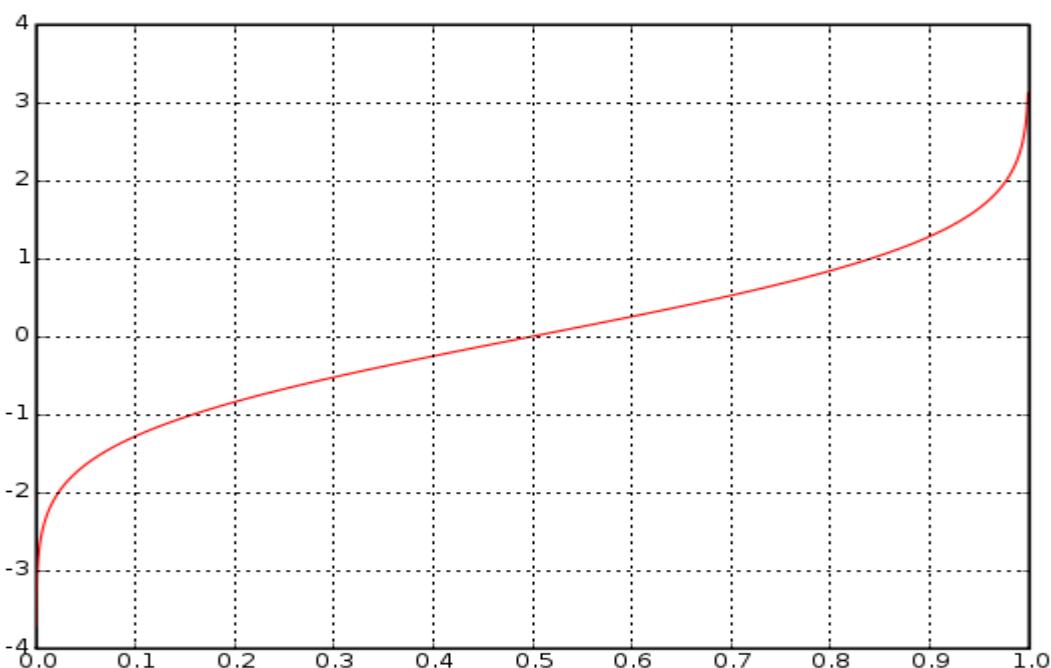
Paket *Statgraphics* je vrlo koristan softver koji nam pored ovih, gore navedenih, omogucava još mnoštvo drugih analiza na podacima koje nam mogu biti od koristi.

3.4 Probit analiza

Probit analiza je alternativa logit analizi. Logit i probit analizu su jako slicne, samo što logit analiza koristi logaritam proporcije verovatnoca, dok probit koristi kumulativnu normalnu raspodelu. Grafik kumulativne normalne raspodele ima oblik slova S i kreće se u granicama od 0 do 1, što je vrlo slicno grafiku logit funkcije. Dakle, zaključujemo da ove dve metode daju vrlo slicne rezultate.



Logit kriva



Probit kriva

Probit funkcija je inverzna funkcija raspodele koja je povezana sa standardnom normalnom raspodelom. **Probit model** je specijalni oblik generalizovanog linearног modela, i uglavnom se koristi u regresiji. Neka je Y zavisna promenljiva, a X vektor nezavisnih promenljivih, tada probit model kaže da je:

$$P(Y = 1 | X = x) = \Phi(x' \mathbf{b})$$

gde je Φ funkcija raspodele standardne normalne raspodele. Parametri β se ocenjuju pomocu funkcije maksimalne verovatnoće. *Višedimenzionalni probit model* je generalizacija probit modela koji se koristi kada treba zajedno da se oceni nekoliko povezanih binarnih promenljivih. Neka su Y_1 i Y_2 binarne zavisne promenljive tako da je

$$Y_1 = 1(Y_1^* > 0)$$

$$Y_2 = 1(Y_2^* > 0)$$

gde je

$$Y_1^* = X\mathbf{b}_1 + \mathbf{e}_1$$

$$Y_2^* = X\mathbf{b}_2 + \mathbf{e}_2$$

gde je X vektor nezavisnih promenljivih i prepostavlja se da one imaju normalnu raspodelu sa nula očekivanjem, jedinicnom varijansom i koeficijentom korelacije ρ .

Logaritamska funkcija verovatoće za probit model ima sledeći oblik:

$$\begin{aligned} \ln l = & \sum_{i=1}^N (1 - y_{1i}) \ln(1 - f(\mathbf{x}_{1i}\mathbf{a}_1)) + \sum_{i=1}^N y_{1i}(1 - y_{2i}) \ln(f(\mathbf{x}_{1i}\mathbf{a}_1) - f_2(\mathbf{x}_{1i}\mathbf{a}_1, \mathbf{x}_{2i}\mathbf{a}_2; \mathbf{r})) + \\ & \sum_{i=1}^N y_{1i}y_{2i} \ln f_2(\mathbf{x}_{1i}\mathbf{a}_1, \mathbf{x}_{2i}\mathbf{a}_2; \mathbf{r}) \end{aligned}$$

gde $\phi(\bullet)$ i $\phi(\bullet, \bullet; \rho)$ predstavljaju jednodimenzionalnu i binarnu standardnu normalnu funkciju raspodele. Ova funkcija nam dalje služi za ocenu parametara β_1 , β_2 , ρ uz pomoć funkcije maksimalne verovatoće.

3.4.1 Primena softvera

I u slučaju probit analize, izuzetno nam je koristan paket Statgraphics. Koristicemo se primerom koji smo dali u odeljku 3.3.6. Kada na date podatke primenimo probit analizu, dobijamo sledeće rezultate:

Probit Analysis - Los

Dependent variable: Los

Factors:

- Musko
- Ozenjen
- Prihodi
- Starost
- Vlasnik stana

Estimated Regression Model (Maximum Likelihood)

Parameter	Standard	
	Estimate	Error
CONSTANT	0.616715	2.90843
Musko	0.830676	1.06687
Ozenjen	2.10205	1.70154
Prihodi	0.00159893	0.00270175
Starost	-0.127185	0.144372
Vlasnik stana	2.12783	1.05592

Analysis of Deviance

Source	Deviance	Df	P-Value
Model	9.01866	5	0.1083
Residual	11.709	9	0.2302
Total (corr.)	20.7277	14	

Percentage of deviance explained by model = 43.5102

Adjusted percentage = 0.0

Likelihood Ratio Tests

Factor	Chi-Squared	Df	P-Value
Musko	0.706922	1	0.4005
Ozenjen	1.74547	1	0.1864
Prihodi	0.334239	1	0.5632
Starost	0.826978	1	0.3631
Vlasnik stana	5.95822	1	0.0146

The StatAdvisor

The output shows the results of fitting a probit regression model to describe the relationship between Los and 5 independent variable(s). The equation of the fitted model is

Los = normal(eta)

where

eta = 0.616715 + 0.830676*Musko + 2.10205*Ozenjen + 0.00159893*Prihodi - 0.127185*Starost + 2.12783*Vlasnik stana

Because the P-value for the model in the Analysis of Deviance table is greater or equal to 0.05, there is not a statistically significant relationship between the variables at the 95.0% or higher confidence level. In addition, the P-value for the residuals is greater than or equal to 0.05, indicating that the model is not significantly worse than the best possible model for this data at the 95.0% or higher confidence level.

The pane also shows that the percentage of deviance in Los explained by the model equals 43.5102%. This statistic is similar to the usual R-Squared statistic. The adjusted percentage, which is more suitable for comparing models with different numbers of independent variables, is 0.0%.

In determining whether the model can be simplified, notice that the highest P-value for the likelihood ratio tests is 0.5632, belonging to Prihodi. Because the P-value is greater or equal to 0.05, that term is not statistically significant at the 95.0% or higher confidence level. Consequently, you should consider removing Prihodi from the model.

Kada uporedimo ove rezultate sa rezultatima koje smo dobili putem logaritamske regresije, možemo da primetimo da su oni vrlo slični. I korelaciona matrica daje slike rezultate:

Correlation matrix for coefficient estimates

	CONSTANT	Musko	Ozenjen	Prihodi
CONSTANT	1.0000	-0.3495	0.4296	0.5856
Musko	-0.3495	1.0000	0.3466	-0.1661
Ozenjen	0.4296	0.3466	1.0000	0.0303
Prihodi	0.5856	-0.1661	0.0303	1.0000
Starost	-0.8801	0.0620	-0.5079	-0.8121
Vlasnik stana	0.3070	0.1652	0.3887	0.3873

	Starost	Vlasnik stana
CONSTANT	-0.8801	0.3070
Musko	0.0620	0.1652
Ozenjen	-0.5079	0.3887
Prihodi	-0.8121	0.3873
Starost	1.0000	-0.5239
Vlasnik stana	-0.5239	1.0000

The StatAdvisor

This table shows estimated correlations between the coefficients in the fitted model. These correlations can be used to detect the presence of serious multicollinearity, i.e., correlation amongst the predictor variables. In this case, there are 3 correlations with absolute values greater than 0.5.

I u slučaju probit analize možemo da primenimo selekciione metode. Pri tome, metoda unapred daje sledeće rezultate:

Probit Analysis - Los

Dependent variable: Los

Factors:

Musko
Ozenjen
Prihodi
Starost
Vlasnik stana

Estimated Regression Model (Maximum Likelihood)

		Standard
Parameter	Estimate	Error
CONSTANT	-0.674486	0.482718
Vlasnik stana	1.74205	0.760325

Analysis of Deviance

Source	Deviance	Df	P-Value
Model	5.98871	1	0.0144
Residual	14.739	13	0.3239
Total (corr.)	20.7277	14	

Percentage of deviance explained by model = 28.8923

Adjusted percentage = 9.59445

Likelihood Ratio Tests

Factor	Chi-Squared	Df	P-Value
Vlasnik stana	5.98871	1	0.0144

Stepwise factor selection

Method: forward selection

P-to-enter: 0.05

P-to-remove: 0.05

Step 0:

0 factors in the model. 14 d.f. for error.

Percentage of deviance explained = 0.00% Adjusted percentage = 0.00%

Step 1:

Adding factor Vlasnik stana with P-to-enter = 0.0143954

1 factors in the model. 13 d.f. for error.

Percentage of deviance explained = 28.89% Adjusted percentage = 9.59%

Final model selected.

The StatAdvisor

The output shows the results of fitting a probit regression model to describe the relationship between Los and 5 independent variable(s). The equation of the fitted model is

Los = normal(eta)

where

eta = -0.674486 + 1.74205*Vlasnik stana

Because the P-value for the model in the Analysis of Deviance table is less than 0.05, there is a statistically significant relationship between the variables at the 95.0% confidence level. In addition, the P-value for the residuals is greater than or equal to 0.05, indicating that the model is not significantly worse than the best possible model for this data at the 95.0% or higher confidence level.

The pane also shows that the percentage of deviance in Los explained by the model equals 28.8923%. This statistic is similar to the usual R-Squared statistic. The adjusted percentage, which is more suitable for comparing models with different numbers of independent variables, is 9.59445%.

In determining whether the model can be simplified, notice that the highest P-value for the likelihood ratio tests is 0.0144, belonging to Vlasnik stana. Because the P-value is less than 0.05, that term is statistically significant at the 95.0% confidence level. Consequently, you probably don't want to remove any variables from the model.

Dok metoda eliminacije unazad u slučaju probit analize izgleda ovako:

Probit Analysis - Los

Estimated Regression Model (Maximum Likelihood)

Parameter	Estimate	Standard	
		Error	
CONSTANT	-0.674486	0.482718	
Vlasnik stana	1.74205	0.760325	

Analysis of Deviance

Source	Deviance	Df	P-Value
Model	5.98871	1	0.0144
Residual	14.739	13	0.3239
Total (corr.)	20.7277	14	

Percentage of deviance explained by model = 28.8923

Adjusted percentage = 9.59445

Likelihood Ratio Tests

Factor	Chi-Squared	Df	P-Value
Vlasnik stana	5.98871	1	0.0144

Stepwise factor selection

Method: backward selection

P-to-enter: 0.05

P-to-remove: 0.05

Step 0:

5 factors in the model. 9 d.f. for error.

Percentage of deviance explained = 43.51% Adjusted percentage = 0.00%

Step 1:

Removing factor Prihodi with P-to-remove = 0.563173

4 factors in the model. 10 d.f. for error.

Percentage of deviance explained = 41.90% Adjusted percentage = 0.00%

Step 2:

Removing factor Starost with P-to-remove = 0.458766

3 factors in the model. 11 d.f. for error.

Percentage of deviance explained = 39.25% Adjusted percentage = 0.65%

Step 3:

Removing factor Musko with P-to-remove = 0.287249

2 factors in the model. 12 d.f. for error.

Percentage of deviance explained = 33.79% Adjusted percentage = 4.84%

Step 4:

Removing factor Ozenjen with P-to-remove = 0.313862

1 factors in the model. 13 d.f. for error.

Percentage of deviance explained = 28.89% Adjusted percentage = 9.59%

Final model selected.

The StatAdvisor

The output shows the results of fitting a probit regression model to describe the relationship between Los and 5 independent variable(s). The equation of the fitted model is

Los = normal(eta)

where

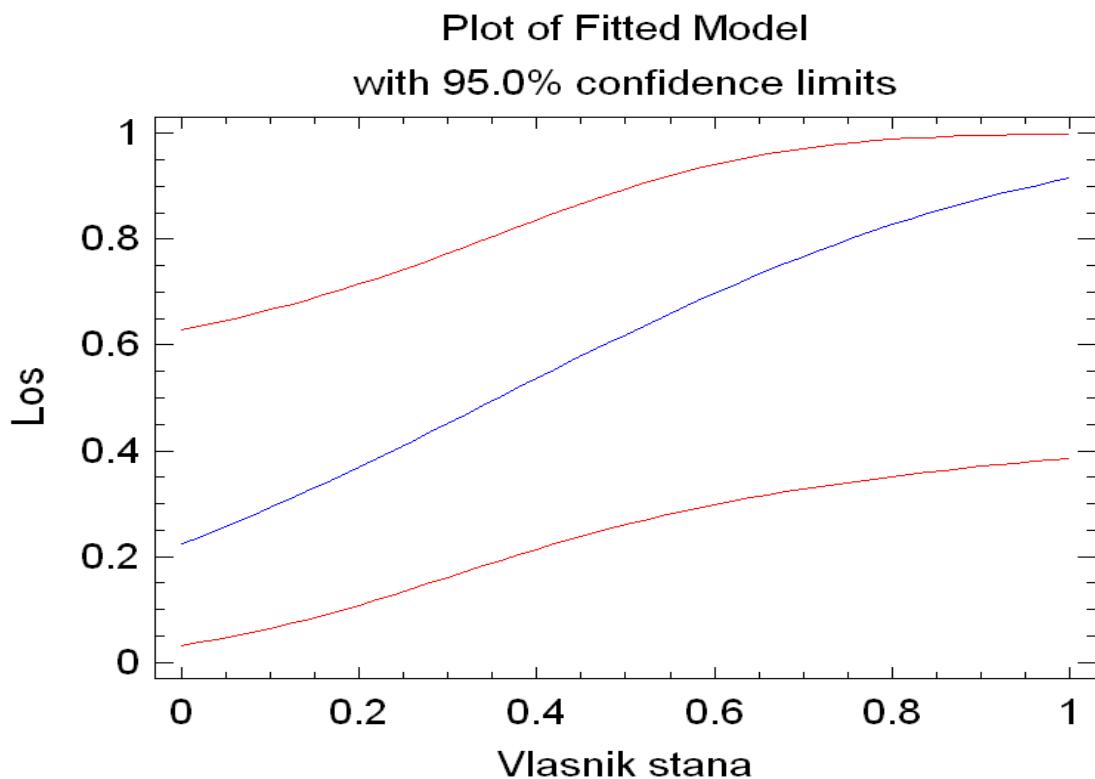
eta = -0.674486 + 1.74205*Vlasnik stana

Because the P-value for the model in the Analysis of Deviance table is less than 0.05, there is a statistically significant relationship between the variables at the 95.0% confidence level. In addition, the P-value for the residuals is greater than or equal to 0.05, indicating that the model is not significantly worse than the best possible model for this data at the 95.0% or higher confidence level.

The pane also shows that the percentage of deviance in Los explained by the model equals 28.8923%. This statistic is similar to the usual R-Squared statistic. The adjusted percentage, which is more suitable for comparing models with different numbers of independent variables, is 9.59445%.

In determining whether the model can be simplified, notice that the highest P-value for the likelihood ratio tests is 0.0144, belonging to Vlasnik stana. Because the P-value is less than 0.05, that term is statistically significant at the 95.0% confidence level. Consequently, you probably don't want to remove any variables from the model.

Kada graficki predstavimo, vidimo da su i grafici slicni logit regresiji, što smo vec napomenuli u teorijskom delu, a ovo je samo potvrda toj cinjenici:



3.5 Tobit model

Tobit model je tip regresije sa ogranicenjem koji odreduje odnos zavisne y_i i nezavisnih promenljivih x_i . Tobit model ima sledeci oblik:

$$y_i = \begin{cases} y_i^*, & y_i^* > 0 \\ 0, & y_i^* \leq 0 \end{cases}$$

gde je y_i^* neopažena promenljiva, tako da je:

$$y_i^* = \mathbf{b}x_i + u_i, u_i \sim N(0, \sigma^2)$$

Cesta modifikacija Tobit modela je da se odredi prag y_L koji je razlicit od nule, i tada posmatramo sledeci slucaj:

$$y_i = \begin{cases} y_i^*, & y_i^* > y_L \\ y_L, & y_i^* \leq y_L \end{cases}$$

Logaritamska funkcija verovatoce za tobit model ima sledeci oblik:

$$\ln L = \sum_{i=1}^N (d_i(-\ln s + \ln f(\frac{y_i - bx_i}{s})) + (1 - d_i)\ln(1 - f(\frac{bx_i}{s})))$$

gde je ϕ funkcija normalne raspodele, a d_i je indikatorska promenljiva:

$$d_i = \begin{cases} 0, & y_i^* \leq y_L \\ 1, & y_i^* > y_L \end{cases}$$

Uz pomoc funkcije maksimalne verovatnoce ocenjujemo parametre β .

4. Modeli kreditnog scoringa sistema

4.1 Primena problema separacije

Sada cemo na primeru videti kako se problem separacije primenjuje na kreditni scoring sistem. Koristimo se primerom datim u radu «*Credit Cards Scoring with Quadratic Utility Function*» [6]. Neka je zahtev za kreditnu karticu objekat, a na primer stavka «bracno stanje» koje se popunjava u zahtevu za kredit je indikatorska promenljiva. Prepostavimo da imamo skup od m objekata i neka je svaki objekat predstavljen vektorom $x^i = (x_1^i, \dots, x_n^i)$, $i = 1, \dots, m$ gde su x_1^i, \dots, x_n^i indikatorske promenljive. Neka je Ω_k skup svih mogucih vrednosti k -tog indikatora. U slucaju kada imamo diskretan indikator, prepostavljamo da klasifikacija pocinje sa nulom i prima uzastopne cele vrednosti. Tako na primer, «bracno stanje» možemo kodirati na sledeci nacin: 0 – razvedena, 1 – neodata, 2 – udata. U slucaju kada imamo neprekidnu promenljivu, nju možemo zameniti skupom diskretnih indikatora. Prepostavimo da želimo da diskretizujemo k -tu indikatorsku promenljivu za objekat $x^i = (x_1^i, \dots, x_n^i)$. Neka k -ti indikator prima vrednosti iz intervala $(0, t)$, gde je t ceo broj. Na primer, za neprekidni indikator kao što je «starost», možemo uvesti sledecu aproksimaciju:

godine	<18	18-20	21-24	25-29	30-55	56-59	60-75	76-84	>84
kod	0	1	2	3	4	5	6	7	8

Posmatrane funkcije korisnosti, pogotovo kvadratna funkcija korisnosti, mogu biti suviše fleksibilne (tj. mogu imati previše stepeni slobode) za skup podataka sa malim brojem tacaka. Uvodenje ogranicenja, kao što je na primer ogranicenje monotonosti na indikatorske promenljive može da smanji prekomernu fleksibilnost modela. Za linearu funkciju korisnosti, uslov monotonosti za indikatorske promenljive se definiše na sledeci nacin: linearna funkcija $u^c(x) = c^T x = \sum_{k=1}^n c_k x_k$ je rastuća po x_k ako i samo ako je $c_k > 0$. Na primer, za problem scoring sistema možemo nametnuti uslov monotonosti za promenljive koje odgovaraju indikatoru «bracno stanje»:

$$\text{razvedena} \setminus \text{neodata} \setminus \text{udata}$$

Za kvadratnu funkciju je mnogo teže nametnuti ogranicenje monotonosti. Posmatramo potklasu monotonih kvadratnih funkcija sa nenegativnim elementima matrice D i vektora c . Ovakva funkcija je monotona po svakoj promenljivoj na $\mathbf{R}_+^n = \{x = (x_1, \dots, x_n) / x_k \geq 0, \forall k \in \{1, \dots, n\}\}$. Da bismo mogli da koristimo takvu klasu funkcija, potrebno je da indikatorske promenljive budu nenegativne i rastuce po preferenciji.

Za neke indikatore je moguce nametnuti ogranicenje monotonosti, dok je za druge bolje koristiti neke druge osobine, kao što su konveksnost i konkavnost. Ako

posmatramo neprekidan indikator «starost», prepostavimo da su srednje godine bolje od mlađih i starijih i neka je preferencija data sa:

$$\begin{aligned} \{0-18\} &\setminus \{18-20\} \setminus \{21-24\} \setminus \{25-29\} \setminus \{85-\infty\} \\ &\setminus \{76-84\} \setminus \{60-75\} \setminus \{56-59\} \setminus \{30-55\} \end{aligned}$$

Preferencija za godine je odredena pomocu nemonotonog kriterijuma nad celim opsegom godina. Ovo možemo uvrstiti u model kreditnog bodovnog sistema na sledeći nacin:

- 1) Možemo da prepostavimo da je funkcija korisnosti konkavna po datoj promenljivoj. Ovo možemo obezbediti racunanjem drugog izvoda po datoj promenljivoj i uvodenjem sledeceg ogranicenja u model linearног programiranja:

$$\frac{\partial^2 u(x)}{\partial x_i^2} \leq 0$$

- 2) Kodiranje indikator promenljive «starost» možemo izvesti na sledeci nacin:

godine	<18	18-20	21-24	25-29	30-55	56-59	60-75	76-84	>84
kod	0	1	2	3	4	3	2	1	0

Na kraju želimo da odredimo pragove u_1 i u_2 i funkciju $u(x)$ koja klasificuje skup svih zahteva za kredit $X_i = \{x^i / i \in I\}$ na sledeci nacin:

- Zahtev i je odobren ako i samo ako je $u_2 \leq u(x^i)$
- Odluka o zahtevu i je odložena ako i samo ako je $u_1 \leq u(x^i) < u_2$
- Zahtev je odbijen ako i samo ako je $u(x^i) < u_1$

Razmatramo funkciju korisnosti koja je linearна по promenljivama koje su korišcene за odlucivanje, и linearна или kvadratna по indikatorskim promenljivama. Da bismo našli funkciju korisnosti koja klasificuje skup objekata sa minimalnom greškom rešavamo ranije objašnjene probleme linearног programiranja, u odeljku 3.2.4. i 3.2.5.

4.2 Dve faze kreditnog ocenjivanja u procesu odobravanja kredita

Pri podnošenju zahteva za kredit, od potencijalnog klijenta se zahteva da dostavi osnovne licne podatke, vlasništvo nad kucom ili nepokretnom imovinom, podatke o zaposlenju, kreditnu istoriju i stanje finansijske aktive i pasive. Pored ovih podataka dostavljenih od klijenta, banka traži izveštaj o klijentu i iz kreditnog biroa. Na osnovu svih ovih podataka, banka donosi odluku o kreditnoj sposobnosti klijenta. Ovaj postupak se najčešće zove *ekspertni sistem kreditne procene*. Ekspertni sistem je kasnije zamjenjen *kreditnim scoring sistemom*, koji je matematički model koji koristi karakteristike podnosioca zahteva da bi izracunao skor klijenta, koji se povezuje sa verovatnocom da klijent neće biti u mogućnosti da vraca kredit (*probability of default - PD*), ili da rangira klijente na osnovu njihovog *default rizika* (rizika da neće vratiti kredit). U ovom delu rada, koji se poziva na rad «*Two stages credit evaluation in bank loan appraisal*» [5], cemo videti da banka u procesu odlucivanja prolazi kroz dve faze, pri cemu cemo navesti i kriterijume koji se uzimaju u obzir da bi se ocenilo da li se isplati izvoditi drugu fazu odlucivanja, pošto se prva uvek izvodi.

Kao što smo vec napomenuli klijente klasifikujemo u dve grupe: dobre i loše (G i B) u zavisnosti od toga da li ce otplatiti kredit kada on dode na naplatu. Prepostavlja se da je procenat dobrih G klijenata koji su kreditno sposobni i koji ce skoro sigurno otplatiti kredit jednak α , gde je $0 < \alpha < 1$. Dok je procenat loših B klijenata koji ce kasniti u otplati kredita i verovatno ga neće ni otplatiti jednak $1-\alpha$. Međutim, kada banka odobrava kredit ona ne može sa sigurnošću da tvrdi kojoj grupi ce pripasti klijent. Nakon što banka primi kreditni zahtev sa zahtevanim dokumentima, ona ce doneti preliminarnu odluku na osnovu kreditnog scoring modela zasnovanog na standardizovanim podacima. Ova faza ocene se naziva *prva faza kreditne procene* (e_1). Banka zatim može da traži dodatne informacije o klijentu i tada se prelazi u *drugu fazu kreditne procene* (e_2).

Stopu preciznosti za prvu fazu kreditne procene cemo označavati sa q_1 , to je verovatnoca da ce klijent iz grupe G (odnosno B) biti tacno klasifikovan u grupu G (odnosno B). Ali takođe postoji $1-q_1$ šanse da banka klasificuje klijenta iz grupe G u grupu B i obrnuto.

Troškovi rezervisanja (zaštite) od gubitaka su jednaki procenatu k koji se primenjuje na iznos odobrenog kredita. Banke moraju za svakog klijenta da odvoje određeni iznos koji ce im barem deliminicno pokriti gubitke ukoliko klijent ne bude u mogućnosti da otplacuje kredit. Prepostavlja se da su troškovi prve faze zanemarljivi. Pošto banaka u drugoj fazi zahteva dodatne informacije o klijentu, prepostavlja se da ce stopa preciznosti druge faze q_2 biti veca od q_1 . Kao i u prvoj fazi, postoji $1-q_2$ šanse da ce banka pogrešno klasifikovati klijenta.

Odluka o tome da li ce se banka odluciti za drugu fazu je problem koristi i troškova. Troškovi procene A druge faze i stope preciznosti q_1 i q_2 igraju odlucujući ulogu. Ukoliko dolazi do velikog povecanja preciznosti u drugoj fazi u odnosu na prvu fazu, korisno je izvoditi drugu fazu. Međutim, ako je pomak mali, cak i ako su troškovi A relativno niski, banka ce odluciti da nije dovoljno ekonomicno da se izvodi druga faza.

Verovatnoca da li ce banka sprovesti drugu fazu procene zavisi od rezultata prve faze. Ako je klijent klasifikovan kao dobar u prvoj fazi, verovatnoca da ce banka sprovesti drugu fazu je jednaka a_G , a ukoliko je klasifikovan kao loš verovatnoca da ce biti sprovedena druga faza je jednaka a_B .

Model se sastoji iz sledecih koraka:

Korak 1: klijent podnosi zahtev za kredit sa zahtevanom dokumentacijom

Korak 2: banka sprovodi prvu fazu kreditne procene i deli klijente na dobre i loše

Korak 3: banka odlucuje da li da odobri kredit na osnovu rezultata iz prve faze ili da nastavi sa drugom fazom. Ako se sprovodi samo prva faze, kredit velicine D sa kamatnom stopom i ce biti odobren samo klijentima koji su svrstani u grupu G , dok ce klijenti koji su svrstani u grupu B odbijeni.

Korak 4: banka sprovodi drugu fazu kreditne procene.

Korak 5: kreditne odluke se donose na bazi rezultata iz druge faze procene, tj. klijenti koji su klasifikovani kao dobri, tj. u grupu G su dobili kredit, a oni koji su klasifikovani kao loši, tj. u grupu B su odbijeni.

Prepostavlja se da je kamatna stopa i veca od procenta troškova rezervisanja k . Stopa povracaja banke zavisi od visine kamatne stope (i), procenta rezervisanja za eventualne gubitke (k) i velicine kredita (D). Da bi se model pojednostavio prepostavlja se da je ocekivana stopa povracaja od davanja kredita dobrom klijentu ($E(r/G)$) pozitivna, a lošem klijentu ($E(r/B)$) negativna, tj.

$$E(r/G) = D(i-k) > 0$$

$$E(r/B) = -D(1+k) < 0$$

Lema 1 Bez upotrebe sistema procene kreditnog rizika, banka daje kredit sa kamatnom stopom $i \geq \frac{1+k-a}{a} \equiv i^*$, gde se i^* može posmatrati kao kamatna stopa u savršeno konkurentnom kreditnom tržištu koja osigurava da ce ocekivana stopa povracaja za banku biti nenegativna.

Dokaz Bez upotrebe sistema za ocenu kreditnog rizika, banka ce odobriti kredit ako i samo ako je ocekivana stopa povracaja $E(r)$ za banku nenegativna.

$$E(r) = aE(r|G) + (1-a)E(r|B) \geq 0 \Leftrightarrow aD(i-k) - (1-a)D(1+k) \geq 0$$

$$\Leftrightarrow D(ai - 1 - k + a) \geq 0, D \geq 0$$

$$\Leftrightarrow i \geq \frac{1+k-a}{a} \equiv i^*$$

U praksi, banka ce uvek sprovesti prvu fazu da bi procenila kreditnu sposobnost klijenta. Prepostavlja se da su troškovi sprovodenja prve faze zanemarljivi i da banka uvek izvodi ovu fazu. Lema 2 postavlja donju granicu na stopu preciznosti prve faze q_1 .

Lema 2 Prepostavljamo da su troškovi prve faze procenjivanja zanemarljivi, sledeći uslov za stopu preciznosti prve faze mora da bude zadovoljen da bi se sprovela prva faza procene:

$$q_1 \geq \frac{\mathbf{a}(i-k)}{\mathbf{a}(i-k) + (1-\mathbf{a})(1+k)} \equiv q_1^*$$

U uslovima savršene konkurencije na kreditnom tržištu, q_1^* je jednako $\frac{1}{2}$.

Dokaz Da bi banka sprovela prvu fazu procene, očekivani povracaj koji zavisi od procene prve faze $E(r|e_1)$ mora da bude veci ili jednak očekivanom povracaju $E(r)$ bez primene sistema za ocenjivanje.

$$\begin{aligned} E(r|e_1) \geq E(r) &\Leftrightarrow \mathbf{a}q_1D(i-k) - (1-\mathbf{a})(1-q_1)D(1+k) \geq \mathbf{a}D(i-k) - (1-\mathbf{a})D(1+k) \\ &\Leftrightarrow \mathbf{a}(i-k)q_1 + (1-\mathbf{a})(1+k)q_1 \geq \mathbf{a}(i-k) \\ &\Leftrightarrow q_1 \geq \frac{\mathbf{a}(i-k)}{\mathbf{a}(i-k) + (1-\mathbf{a})(1+k)} \equiv q_1^* \end{aligned}$$

U uslovima savršene konkurencije, kamatna stopa $i = \frac{1+k-\mathbf{a}}{\mathbf{a}}$. Kada to uvrstimo u q_1^* dobijamo da je $q_1^* = 1/2$.

Nakon što su ispunjeni uslovi iz Leme 1 i Leme 2, banka će odluciti da li joj se isplati da sprovodi drugu fazu, ili da donese odluku na osnovu rezultata iz prve faze. Teorema 1 nam daje uslov na osnovu koga se banka odluceće za drugu fazu u situaciji kada je klijent klasifikovan kao dobar u prvoj fazi. Dok nam Teorema 2 daje uslov na osnovu koga se banka odluceće za drugu fazu u situaciji kada je klijent klasifikovan kao loš u prvoj fazi.

Teorema 1 Neka je klijent klasifikovan kao dobar, tj. u grupu G , $q_1 \geq q_1^*$ i $i = i^*$, ako je $A < A_G^* \equiv \frac{(1-\mathbf{a})(q_2 - q_1)(1+k)D}{\mathbf{a}q_1 + (1-\mathbf{a})(1-q_1)}$ tada je verovatnoca da će se sprovesti druga faza procene jednaka $a_G^* = 1$, u suprotnom je $a_G^* = 0$.

Dokaz Ako je klijent u prvoj fazi klasifikovan kao dobar, što označavamo sa e_{IG} i ako je verovatnoca izvodenja druge faze procene definisana sa a_G , tada je očekivani povracaj za banku nakon što se izvede druga faza jednak:

$$\begin{aligned} E(r|e_{IG}, e_2) &= a_G \left(\frac{\mathbf{a}q_1}{\mathbf{a}q_1 + (1-\mathbf{a})(1-q_1)} (q_2(i-k) + (1-q_2) \cdot 0)D + \frac{(1-\mathbf{a})(1-q_1)}{\mathbf{a}q_1 + (1-\mathbf{a})(1-q_1)} (- (1-q_2)(1+k) \right. \\ &\quad \left. + q_2 \cdot 0)D - A \right) + (1-a_G) \left(\frac{\mathbf{a}q_1}{\mathbf{a}q_1 + (1-\mathbf{a})(1-q_1)} (i-k)D - \frac{(1-\mathbf{a})(1-q_1)}{\mathbf{a}q_1 + (1-\mathbf{a})(1-q_1)} (1+k)D \right) \end{aligned}$$

Banka će odluciti da izvede drugu fazu procene jedino ako povećanje verovatnoće sprovodenja druge faze procene ima za posledicu veće očekivane prinose, tj. $\frac{\partial E(r|e_{IG}, e_2)}{\partial a_G} > 0$, u suprotnom banka neće sprovesti drugu fazu, tj. $a_G^* = 0$.

$$\begin{aligned}
 \frac{\partial E(r | e_{1G}, e_2)}{\partial a_G} &= \left(\frac{D}{\mathbf{a}q_1 + (1-\mathbf{a})(1-q_1)} \right) (\mathbf{a}q_1 q_2(i-k) - \mathbf{a}q_1(i-k) - (1-\mathbf{a})(1-q_1)(1-q_2)(1+k) + \\
 &+ (1-\mathbf{a})(1-q_1)(1+k)) - A > 0 \\
 \Leftrightarrow A < ((1-\mathbf{a})(1-q_1)q_2(1+k) - \mathbf{a}q_1(1-q_2)(i-k)) \left(\frac{D}{\mathbf{a}q_1 + (1-\mathbf{a})(1-q_1)} \right) \\
 \Leftrightarrow A < \frac{(1-\mathbf{a})(q_2 - q_1)(1+k)D}{\mathbf{a}q_1 + (1-\mathbf{a})(1-q_1)} \equiv A_G^* \\
 \text{gde je } i = i^* = \frac{1+k-\mathbf{a}}{\mathbf{a}}.
 \end{aligned}$$

Teorema 2 Neka je klijent klasifikovan kao loš, tj. u grupu B , $q_1 \geq q_1^*$ i $i = i^*$, ako je $A < A_B^* \equiv \frac{(1-\mathbf{a})(q_2 - q_1)(1+k)D}{\mathbf{a}(1-q_1) + (1-\mathbf{a})q_1}$ tada je verovatnoca da će se sprovesti druga faza procene jednaka $a_B^* = 1$, u suprotnom je $a_B^* = 0$.

Dokaz Ako je klijent u prvoj fazi klasifikovan kao loš, što označavamo sa e_{1B} i ako je verovatnoca izvodenja druge faze procene definisana sa a_B , tada je očekivani povratak za banku nakon što se izvede druga faza jednak:

$$\begin{aligned}
 E(r | e_{1B}, e_2) &= a_B \left(\frac{\mathbf{a}(1-q_1)}{\mathbf{a}(1-q_1) + (1-\mathbf{a})q_1} (q_2(i-k) + (1-q_2) \cdot 0) D + \frac{(1-\mathbf{a})q_1}{\mathbf{a}(1-q_1) + (1-\mathbf{a})q_1} (-1-q_2)(1+k) \right. \\
 &\quad \left. + q_2 \cdot 0 \right) D - A + (1-a_B) \cdot 0
 \end{aligned}$$

Banka će odluciti da izvede drugu fazu procene jedino ako povećanje verovatnoće sprovodenja druge faze procene ima za posledicu veće očekivane prinose, tj.

$$\frac{\partial E(r | e_{1B}, e_2)}{\partial a_B} > 0, \text{ u suprotnom banka neće sprovesti drugu fazu, tj. } a_B^* = 0.$$

$$\begin{aligned}
 \frac{\partial E(r | e_{1B}, e_2)}{\partial a_B} &= \left(\frac{D}{\mathbf{a}(1-q_1) + (1-\mathbf{a})q_1} \right) (\mathbf{a}(1-q_1)q_2(i-k) - \mathbf{a}q_1(i-k) - (1-\mathbf{a})q_1(1-q_2)(1+k)) - A > 0 \\
 \Leftrightarrow A < & (\mathbf{a}(1-q_1)q_2(i-k) - (1-\mathbf{a})q_1(1-q_2)(1+k)) \left(\frac{D}{\mathbf{a}(1-q_1) + (1-\mathbf{a})q_1} \right) \\
 \Leftrightarrow A < & \frac{(1-\mathbf{a})(q_2 - q_1)(1+k)D}{\mathbf{a}(1-q_1) + (1-\mathbf{a})q_1} \equiv A_B^* \\
 \text{gde je } i = i^* &= \frac{1+k-\mathbf{a}}{\mathbf{a}}.
 \end{aligned}$$

Teoreme 1 i 2 daju granicne troškove druge faze ukoliko banka odluci da je sprovede. Ako su stvarni troškovi druge faze procene manji od granicnih troškova, isplati se sprovoditi drugu fazu. Brojilac granicnih troškova A_G^* i A_B^* je korist od izbegavanja gubitaka koji se stvaraju davanjem kredita lošim klijentima, a imenilac predstavlja procenat klijenata koji su klasifikovani kao dobri odnosno kao loši u prvoj fazi.

U ovom modelu, prepostavlja se da banka izjednacava kamatnu stopu sa kamatnom stopom u uslovima savršene konkurencije na tržištu kapitala, tj.
 $i = i^* = \frac{1+k-\alpha}{\alpha}$.

Lema 3 Ako je $\alpha < \frac{1}{2}$ i $q_1 \geq \frac{1}{2}$ tada je $A_G^* \geq A_B^*$. Ali ako je $\alpha > \frac{1}{2}$ i $q_1 \geq \frac{1}{2}$ tada je $A_G^* \leq A_B^*$.

Dokaz Ako je $A_G^* \geq A_B^*$, tada je

$$\begin{aligned} \frac{(1-\alpha)(q_2 - q_1)(1+k)}{\alpha q_1 + (1-\alpha)(1-q_1)} &\geq \frac{(1-\alpha)(q_2 - q_1)(1+k)D}{\alpha(1-q_1) + (1-\alpha)q_1} \\ \Leftrightarrow \alpha(1-q_1) + (1-\alpha)q_1 &\geq \alpha q_1 + (1-\alpha)(1-q_1) \\ \Leftrightarrow (1-2\alpha)q_1 &\geq (1-2\alpha)(1-q_1) \\ \Leftrightarrow q_1 &\geq 1/2 \text{ ako je } \alpha < 1/2 \end{aligned}$$

Dakle, ako je $q_1 \geq 1/2$ i $\alpha < 1/2$ tada je $A_G^* \geq A_B^*$.

Ako je $A_B^* \geq A_G^*$, tada je

$$\begin{aligned} \frac{(1-\alpha)(q_2 - q_1)(1+k)}{\alpha q_1 + (1-\alpha)(1-q_1)} &\leq \frac{(1-\alpha)(q_2 - q_1)(1+k)D}{\alpha(1-q_1) + (1-\alpha)q_1} \\ \Leftrightarrow \alpha(1-q_1) + (1-\alpha)q_1 &\leq \alpha q_1 + (1-\alpha)(1-q_1) \\ \Leftrightarrow (1-2\alpha)q_1 &\leq (1-2\alpha)(1-q_1) \\ \Leftrightarrow q_1 &\geq 1/2 \text{ ako je } \alpha > 1/2 \end{aligned}$$

Dakle, ako je $q_1 \geq 1/2$ i $\alpha > 1/2$ tada je $A_B^* \geq A_G^*$.

Lema 3 pokazuje da ako je procenat klijenata koji otplacuju kredit manji (veci) od onih koji kasne u otplati, granicni troškovi za dobre klijente su viši (niži) od tih troškova za loše klijente.

Teorema 3 Neka je $\alpha = \frac{1}{2}$, $i = i^*$ i $q_1 \geq q_1^* = \frac{1}{2}$,

(1) ako je $A < A_G^* = A_B^*$, tada je $a_G^* = a_B^* = 1$

(2) ako je $A \geq A_G^* = A_B^*$, tada je $a_G^* = a_B^* = 0$, gde je $A_G^* = \frac{(1-\alpha)(q_2 - q_1)(1+k)D}{\alpha q_1 + (1-\alpha)(1-q_1)}$,

$$a \quad A_B^* = \frac{(1-\alpha)(q_2 - q_1)(1+k)D}{\alpha(1-q_1) + (1-\alpha)q_1}.$$

Dokaz Ako je $\alpha = 1/2$ tada je

$$A_G^* = \frac{(1-\alpha)(q_2 - q_1)(1+k)D}{\alpha q_1 + (1-\alpha)(1-q_1)} = (q_2 - q_1)(1+k)D$$

$$A_B^* = \frac{(1-\mathbf{a})(q_2 - q_1)(1+k)D}{\mathbf{a}(1-q_1) + (1-\mathbf{a})q_1} = (q_2 - q_1)(1+k)D$$

Pošto je $A_G^* = A_B^*$,

(1) ako je $A < A_G^* = A_B^*$ onda je $a_G^* = a_B^* = 1$ na osnovu Teoreme 1 i Teoreme 2

(2) ako je $A \geq A_G^* = A_B^*$, tada je $a_G^* = a_B^* = 0$ na osnovu Teoreme 1 i Teoreme 2.

Teorema 3 pokazuje da kada je procenat dobrih i loših klijenata jednak, tada su i granicni troškovi jednakci.

Teorema 4 Neka je $\alpha < \frac{1}{2}$, $i = i^*$ i $q_1 \geq q_1^* = \frac{1}{2}$,

(1) ako je $A < A_B^*$, tada je $a_G^* = a_B^* = 1$;

(2) ako je $A_B^* \leq A < A_G^*$, tada je $a_G^* = 1$ i $a_B^* = 0$;

(3) ako je $A \geq A_G^*$, tada je $a_G^* = a_B^* = 0$, gde je $A_G^* = \frac{(1-\mathbf{a})(q_2 - q_1)(1+k)D}{\mathbf{a}q_1 + (1-\mathbf{a})(1-q_1)}$ i

$$A_B^* = \frac{(1-\mathbf{a})(q_2 - q_1)(1+k)D}{\mathbf{a}(1-q_1) + (1-\mathbf{a})q_1}.$$

Dokaz Na osnovu Leme 3, ako je $\mathbf{a} < 1/2$ i $q_1 \geq 1/2$ tada je $A_G^* \geq A_B^*$. Zatim,

(1) ako je $A < A_B^*$, tada je $a_G^* = a_B^* = 1$ na osnovu Teoreme 1 i Teoreme 2

(2) ako je $A_B^* \leq A < A_G^*$, tada je $a_G^* = 1$ i $a_B^* = 0$ na osnovu Teoreme 1 i Teoreme 2

(3) ako je $A \geq A_G^*$, tada je $a_G^* = a_B^* = 0$ na osnovu Teoreme 1 i Teoreme 2

Teorema 5 Neka je $\alpha > \frac{1}{2}$, $i = i^*$ i $q_1 \geq q_1^* = \frac{1}{2}$,

(1) ako je $A < A_G^*$, tada je $a_G^* = a_B^* = 1$;

(2) ako je $A_G^* \leq A < A_B^*$, tada je $a_G^* = 0$ i $a_B^* = 1$;

(3) ako je $A \geq A_B^*$, tada je $a_G^* = a_B^* = 0$, gde je $A_G^* = \frac{(1-\mathbf{a})(q_2 - q_1)(1+k)D}{\mathbf{a}q_1 + (1-\mathbf{a})(1-q_1)}$ i

$$A_B^* = \frac{(1-\mathbf{a})(q_2 - q_1)(1+k)D}{\mathbf{a}(1-q_1) + (1-\mathbf{a})q_1}.$$

Dokaz Na osnovu Leme 3, ako je $\mathbf{a} > 1/2$ i $q_1 \geq 1/2$ tada je $A_G^* \leq A_B^*$. Zatim,

(1) ako je $A < A_G^*$, tada je $a_G^* = a_B^* = 1$ na osnovu Teoreme 1 i Teoreme 2

(2) ako je $A_G^* \leq A < A_B^*$, tada je $a_G^* = 0$ i $a_B^* = 1$ na osnovu Teoreme 1 i Teoreme 2

(3) ako je $A \geq A_B^*$, tada je $a_G^* = a_B^* = 0$ na osnovu Teoreme 1 i Teoreme 2

Teoreme 4 i 5 demonstriraju, da u zavisnosti od odnosa dobrih i loših klijenata, banka sprovodi optimalnu strategiju da li da izvede drugu fazu na (1) obe grupe, (2) samo jednoj grupi ili (3) ni na jednoj grupi. Ako je procenat klijenata koji otplacuju kredit manji od procenta klijenata koji kasne u otplati, granicni troškovi za dobre klijente su viši od onih za loše klijente. U ovom slučaju postoji mogucnost (ako je $A_B^* \leq A < A_G^*$) da banka odluci da sprovede drugu fazu samo na dobre klijente. Tada banke teže da sprovedu drugu fazu procene na klijente koji su klasifikovani kao dobri, u suprotnom, namera da se dalje procenjuje je mnogo slabija.

Ukoliko je procenat dobrih klijenata viši od procenta loših klijenata, granicni troškovi ce biti niži od onih za loše klijente. U tom slučaju postoji mogucnost (ako je $A_G^* \leq A < A_B^*$) da banka odluci da sprovede drugu fazu samo na loše klijente. Tada banke teže da sprovedu drugu fazu procene na klijente koji su klasifikovani kao loši, u suprotnom, namera da se dalje procenjuje je mnogo slabija.

Sada cemo na kratkom numerickom primeru videti kakva je osetljivost rezultata u zavisnosti od vrednosti parametara (k, D, q_1, q_2)

Pretpostavke: $k = 5\%$, $D = \$100$, $q_1 = 0.6$, $q_2 = 0.7$		
Rezultati	$a = 0.95$	$a = 0.40$
i^*	10.53%	162.5%
q_1^*	0.50	0.50
A_G^*	\$ 0.89	\$ 13.13
A_B^*	\$ 1.28	\$ 12.12
Pretpostavke: $k = 5\%$, $D = \$100$, $q_1 = 0.6$, $q_2 = 0.9$		
Rezultati	$a = 0.95$	$a = 0.40$
i^*	10.53%	162.50%
q_1^*	0.50	0.50
A_G^*	\$ 2.67	\$ 39.38
A_B^*	\$ 3.84	\$ 36.35

Kao što vidimo iz tabele, ako se pretpostavi da je $\alpha = 0.95$ i da je mala promena u stopi preciznosti, A_G^* i A_B^* su veoma niski pri cemu je A_G^* niža od A_B^* , što je u saglasnosti sa Teoremom 5. Ako procenat dobrih klijenata padne ispod 0.50 ($\alpha = 0.40$), zajedno sa A_G^* i A_B^* se drastично povecava i u tom slučaju je A_G^* veće od A_B^* , što je u saglasnosti sa Teoremom 4. Medutim, u tom slučaju, kamatna stopa mora da bude izuzetno visoka (162.5 %). U drugom primeru vidimo da znacajan napredak u stopi preciznosti rezultira u višim granicnim troškovima A_G^* i A_B^* , a samim tim je veca težnja da se sprovede duga faza procenjivanja.

4.3 Primena probit i tobit modela

U praksi, vecina kreditnih scoring sistema pati od pristrasnosti pri odabiru uzorka zbog toga što se modeli ocenjuju na uzorku koji se sastoji od odobrenih kredita, pri cemu se ne uzima u obzir kriterijum po kojem se klijenti odbijaju. Glavni razlog tome je nedostatak podataka o odbijenim klijentima koji su dostupni javnosti. Boyes je izbegao ovu pristrasnost tako što je formirao binarni probit model sa dva uzastopna dogadaja kao zavisne promenljive: odluka da li odobrati kredit ili ne, i zavisno od toga da li je kredit odobren, klijentova sposobnost da ga otplati. On je, koristeci svoj model, došao da zaključka da odobravanje kredita nije u saglasnosti sa minimizacijom rizika da ce klijent biti loš.

Za pocetak treba oceniti nepristrasan kreditni scoring sistem, tj. binarni probit model. U ovom primeru korišcena je velika baza podataka koja sadrži veliki skup finansijskih i licnih informacija kako o odobrenim tako i o odbijenim klijentima.

Skup podataka se sastoji od 13,338 zahteva za kredit u jednoj Švedskoj banci u periodu izmedu septembra 1994. i avgusta 1995. godine (pozivamo se na rad «*Bank lending policy, credit scoring and value-at-risk*» [1]). Krediti su davani u obliku kreditnih kartica. To su revolving krediti tako da nemaju odredeno vreme do kada treba da se otplate. U ovom primeru, kredit se klasificuje kao *loš* ukoliko je u kašnjenju dužem od devedeset dana i samim tim je završio u agenciji za naplatu kredita koji su u problemu. Od pocetnih 57 promenljivih, korišceno je samo 17 za konacnu ocenu modelu. Mnoge promenljive su iskljucene iz modela zbog odsustva odnosa sa promenljivama od znacaja - odluka da li odobrati kredit i ponašanje pri otplati kredita. Drugi razlog je izuzetno jaka korelacija sa nekom drugom promenljivom koja se odnosi na istu stvar (kao što su na primer bruto i neto prihodi), a ova druga ima vecu objašnjavajucu moc. Promenljive koje se nalaze u konacnom modelu i njihove definicije su sledeće:

1. *Starost* – koliko godina ima klijent
2. *Muško* – binarna promenljiva, uzima vrednost 1 ako je klijent muško
3. *Razveden* - binarna promenljiva, uzima vrednost 1 ako je klijent razveden
4. *Kuca*- binarna promenljiva, uzima vrednost 1 ako je klijent vlasnik kuce
5. *Veliki grad* - binarna promenljiva, uzima vrednost 1 ako klijent živi u jednom od tri najveca grada
6. *Broj zahteva* – broj zahteva za informacije o klijentu koje je kreditni biro primio u poslednjih 36 meseci
7. *Preduzece* - binarna promenljiva, uzima vrednost 1 ako klijent ima oporezujuci prihod od registrovanog posla
8. *Prihod* – prijavljeni godišnji prihod od plata
9. *Razlika prihoda* – razlika izmedu godišnjih prihoda tekuce i prethodne godine
10. *Kapital* - binarna promenljiva, uzima vrednost 1 ako klijent ima oporezujuci prihod od kapitala

11. *Balans* – odnos prihoda i ukupnih kreditnih kapaciteta, izraženo u procentima
12. *Nula limit* - binarna promenljiva, uzima vrednost 1 ako klijent nema neizmirenih obaveza po kreditima
13. *Limit* – ukupan odobreni iznos kredita
14. *Broj kredita* – broj kredita
15. *Upotrebljen limit* - procenat *Limita* koji je iskorišćen
16. *Iznos kredita* – iznos odobrenog kredita
17. *Jemac* - binarna promenljiva, uzima vrednost 1 ako klijent ima jemca

Od svih 13,338 klijenata, 6899 ili 51.7% je odbijeno za kredit. Dakle, 6439 klijenata je dobilo kredit. U *Tabeli 1* su date deskriptivne statistike za gore navedene promenljive.

Tabela 1: Deskriptivne statistike za promenljive korišcene u modelu

Promenljive	Odobreni klijenti (N = 6899)				Odbijeni klijenti (N = 6439)			
	Srednja vrednost	Standardna devijacija	Minimalna vrednost	Maksimalna vrednost	Srednja vrednost	Standardna devijacija	Minimalna vrednost	Maksimalna vrednost
<i>Starost</i>	38.65	12.76	18	84	41.02	12.08	20	83
<i>Muško</i>	0.62	0.48	0	1	0.65	0.48	0	1
<i>Razveden</i>	0.13	0.34	0	1	0.14	0.35	0	1
<i>Kuca</i>	0.34	0.47	0	1	0.47	0.50	0	1
<i>Veliki grad</i>	0.41	0.49	0	1	0.37	0.48	0	1
<i>Broj zahteva</i>	4.69	2.60	1	10	4.81	2.68	1	19
<i>Preduzece</i>	0.04	0.21	0	1	0.02	0.16	0	1
<i>Prihod</i>	129.93	70.38	0	737.9	189.47	75.70	0	1093.0
<i>Razlika prihoda</i>	5.37	34.06	-438.5	252.6	9.03	34.63	-6226	5006.0
<i>Kapital</i>	0.12	0.32	0	1	0.07	0.25	0	1
<i>Balans</i>	91.04	894.53	0	41533	31.01	386.15	0	22387
<i>Nula limit</i>	0.15	0.36	0	1	<0.01	0.05	0	1
<i>Limit</i>	79.89	93.69	0	1703.0	50.47	51.07	0	949.2
<i>Broj kredita</i>	2.99	2.42	0	18	3.65	2.04	0	16
<i>Upotrebljen limit</i>	64.34	38.88	0	278.0	53.22	33.94	0	124.0
<i>Jemac</i>	0.16	0.36	0	1	0.14	0.35	0	1

Primetimo da je u ovoj tabeli izostavljena promenljiva «*Iznos kredita*», razlog je taj što ona nema vrednost za odbijene klijente, ali se nalazi u sledecoj tabeli u kojoj su predstavljene vrednosti za odobrene klijente.

U oktobru 1996. je uraden monitoring klijenata kojima je odobren kredit i došlo se do zakljucka da je 6% ili 388 klijenata zapalo u problem nemogucnosti otplacivanja

kredita, tj. bili su prosledeni u agenciju za naplatu duga, dok su ostali klijenti redovno otpacivali kredit. U Tabeli 2 predstavljene su deskriptivne statistike za ove klijente:

Tabela 2: Deskriptivne statistike za odobrene kredite

Promenljive	Loši klijenti ($N = 388$)				Dobili klijenti ($N = 6051$)			
	Srednja vrednost	Standardna devijacija	Minimalna vrednost	Maksimalna vrednost	Srednja vrednost	Standardna devijacija	Minimalna vrednost	Maksimalna vrednost
Starost	36.11	11.03	21	75	41.33	12.07	20	83
Muško	0.67	0.47	0	1	0.65	0.48	0	1
Razveden	0.20	0.40	0	1	0.14	0.35	0	1
Kuca	0.28	0.45	0	1	0.48	0.50	0	1
Veliki grad	0.41	0.49	0	1	0.36	0.48	0	1
Broj zahteva	6.15	2.85	1	14	4.72	2.64	1	19
Preduzece	0.02	0.13	0	1	0.03	0.16	0	1
Prihod	165.36	82.35	0	1093.0	191.01	75.00	0	1031.7
Razlika prihoda	3.52	39.01	-135.0	439.7	9.38	34.30	-622.6	500.6
Kapital	0.04	0.20	0	1	0.07	0.26	0	1
Balans	39.92	313.51	0	6041	30.44	390.36	0	22387
Nula limit	0.04	0.20	0	1	<0.01	0.02	0	1
Limit	41.44	57.98	0	511.5	51.05	50.54	0	949.2
Broj kredita	2.34	1.64	0	11	3.74	2.04	0	16
Upotrebljen limit	75.69	33.37	0	124.0	51.78	33.47	0	112.0
Iznos kredita	7.08	3.95	3.0	24.5	7.12	3.83	3.0	30.0
Jemac	0.07	0.26	0	1	0.14	0.35	0	1

U ovom primeru koristi se binarni probit model koji se sastoji od dve simultane jednacine. Prva y_{1i} za odluku da li odobriti kredit ili ne, a druga y_{2i} za binarni rezultat, «dobar» ili «loš» klijent. Neka je sa $*$ označena neopažena promenljiva i pretpostavljamo da je:

$$y_{1i}^* = \mathbf{x}_{1i}\mathbf{a}_1 + \mathbf{e}_{1i}$$

$$y_{2i}^* = \mathbf{x}_{2i}\mathbf{a}_2 + \mathbf{e}_{2i} \quad \text{za } i = 1, 2, \dots, N$$

gde su $\mathbf{x}_{ji}, j = 1, 2, 1 \times k_j$ vektori nezavisnih promenljivih i pretpostavlja se da one imaju normalnu raspodelu sa nula očekivanjem, jedinicnom varijansom i koeficijentom korelacije ρ .

Binarna promenljiva y_{1i} uzima vrednost 1 ako je kredit odobren, a ukoliko je klijent odbijen uzima vrednost 0:

$$y_{1i} = \begin{cases} 0, & \text{odbijen}(y_{1i}^* < 0) \\ 1, & \text{odobren}(y_{1i}^* \geq 0) \end{cases}$$

Binarna promenljiva y_{2i} uzima vrednost 1 ako je klijent dobar, a 0 ukoliko je klijent loš:

$$y_{2i} = \begin{cases} 0, loš(y_{2i}^* < 0) \\ 1, dobar(y_{2i}^* \geq 0) \end{cases}$$

Ocenjeni parametri i njihove standardne greške su predstavljene u *Tabeli 3*.

Tabela 3: ocenjeni parametri binarnog probit modela

Promenljive	P (da je kredit odobren)		P (da je klijent dobar)	
	$\hat{\alpha}_1$	t - statistika	$\hat{\alpha}_2$	t - statistika
Konstatnta	-0.2374	-3.57***	2.2900	15.65***
Starost	-0.004303	-3.69***	0.006892	2.63***
Muško	-0.2003	-7.10***	-0.02456	-0.43
Razveden	-0.02588	-0.70	-0.2380	-3.34***
Kuca	0.06391	2.32**	-0.02019	0.35
Veliki grad	-0.2382	-8.96***	-0.03724	-0.69
Broj zahteva	-0.008123	-1.58	-0.1000	-9.84***
Preduzece	0.5223	8.30***	0.2065	1.28
Prihod	0.008929	49.17***	-0.002392	-4.81***
Razlika prihoda	-0.002336	-6.78***	0.002233	3.04***
Kapital	-0.2776	-5.48***	0.1189	0.97
Balans	0.00006548	3.48***	-0.00009135	-1.54
Nula limit	-2.244	-54.94***	0.005064	-2.23**
Limit	-0.008381	-54.94***	0.005064	10.28***
Broj kredita	0.08420	12.23***	0.2698	14.41***
Upotrebljen limit	-0.007746	-17.72***	-0.01197	-12.91***
Iznos kredita	-	-	-0.006637	-0.98
Jemac	0.1300	3.83***	0.4374	4.50***
r	-	-	-0.9234	-17.34***

gde *, **, *** predstavljaju 10%, 5% i 1% nivo znacajnosti.

Primetimo da «Iznos kredita» ne može da se koristi kao nezavisna promenljiva u prvoj jednacini, pošto ne postoje podaci za ovu promenljivu u slučaju odbijenih klijenata. Uticaj vecine promenljivih na verovatnoca da će klijent biti loš je u saglasnosti sa politikom banke. Tako na primer «Prihodi», «Kuca», «Preduzece», «Broj kredita» i «Jemac» su vrlo bitni faktori koji pozitivno uticu, dok «Nula limit», «Limit» i «Upotrebljen limit» imaju negativan uticaj kada banka donosi odluku.

Promenljiva «Iznos kredita» nema znacajan uticaj na rizik da će klijent biti loš, dok «Limit» ima mnogo veći uticaj. Zato se mora biti obazriv u generalizaciji ovih

rezultata u slučaju kada je odnos «*Iznosa kredita*» i «*Limita*» blizu jedan. Ali ipak, rezultati pokazuju da kada se posmatraju manji krediti, iznos kredita ne utice na rizik koji nosi taj plasman.

Koefficijent korelacije -0.9234 pokazuje skoro savršenu korelaciju između odluke da se odobri kredit i rizika da će klijent biti loš.

Kada se posmatraju rezultati u *Tabeli 3*, dolazi se do sledećih zaključaka: prvo, većina promenljivih koje pozitivno uticu na odobravanje kredita nisu među onima koje smanjuju rizik da će klijent biti loš, dakle, banka ne teži ka tome da minimizira rizik; takođe, rezultati pokazuju da iznos kredita ne utice na rizik. Dakle, banke teže da odobravaju veće kredite cak i ako su oni rizičniji da bi maksimizirale profit, pošto se za rizičnije klijente obračunava viša kamatna stopa, pa je samim tim i viša očekivana stopa prinosa.

Pored toga da li će klijent biti loš, za banke je vrlo bitno da znaju i kada će se desiti da klijent upadne u problem i prestane da otplacuje kredit, tj. da ocene vreme «**preživljavanja kredita**». Binarni Tobit model je vrlo efikasan u razdvajanju klijenata koji brzo postanu loši i onih koji imaju duže vreme preživljavanja ili nikada ne postanu loši. Uzeti su rezultati iz rada «*Bank lending policy, credit scoring and the survival of loans*» [2].

Koristimo se vec datim primerom, jedino što je u modelu umesto promenljive «*Starost*» uključena binarna promenljiva «*Oženjen*» koja uzima vrednost jedan ako je klijent razveden, dok u ovom slučaju promenljiva «*Razveden*» znači da je klijent udovac ili da je razdvojen. U ovom slučaju, klijent je loš ako kasni u otplati kredita duže od devedeset dana, a dobar ako i dalje izmiruje svoje obaveze, tj. i dalje je aktivran. Za loše klijente, vreme preživljavanja predstavlja broj dana između odobravanja kredita i dana kada je postao loš, dok je za dobre klijente vreme preživljavanja razlika između dana kada je odobren kredit i dana kada je raden monitoring, u ovom slučaju 9. oktobra 1996.

Model se sastoji od dve simultane jednacine, prva za odluku da li odobriti kredit ili ga odbiti, y_i , a druga je prirodni logaritam od vremena preživljavanja kredita (izraženo u danima), t_i . Neka je sa * označena neopažena promenljiva i prepostavljamo da je:

$$\begin{aligned} y_i^* &= \mathbf{x}_{1i} \mathbf{b}_1 + \mathbf{e}_{1i} \\ t_i^* &= \mathbf{x}_{2i} \mathbf{b}_2 + \mathbf{e}_{2i} \quad \text{za } i = 1, 2, \dots, N \end{aligned} \tag{*}$$

i prepostavimo da su raspodele dvodimenzionalne normalne raspodele:

$$\begin{pmatrix} \mathbf{e}_{1i} \\ \mathbf{e}_{2i} \end{pmatrix} \sim N \begin{pmatrix} 0 & 1 & \mathbf{s}_{12} \\ 0 & \mathbf{s}_{12} & \mathbf{s}_2^2 \end{pmatrix}$$

Binarna promenljiva y_i uzima vrednost 1 ako je kredit odobren, a ukoliko je klijent odbijen uzima vrednost 0:

$$y_i = \begin{cases} 0, & \text{odbijen}(y_i^* < 0) \\ 1, & \text{odobren}(y_i^* \geq 0) \end{cases}$$

Za loše kredite, može da se uoci tacno vreme preživaljavanja. Dok za kredite koji se i dalje redovno otplacuju u vreme monitoringa, preživljavanje je ocenjeno pošto ne znamo da li ce i kada ce klijent postati loš. Tako na primer, kreditu koji je odobren 1.septembra 1994. godine prag za ocenjivanje ce biti 768 dana, a za kredit koji je odobren 31.avgusta 1995. godine on ce biti 434 dana. Ova vrednost praga za ocenjivanje ce biti oznacena sa \bar{t}_i . Dakle, pravilo za ocenjivanje je sledece:

$$t_i^* = \begin{cases} t_i^*, \text{ako } t_i^* < \bar{t}_i \\ \bar{t}_i, \text{ako } t_i^* \geq \bar{t}_i \end{cases}$$

Binarna promenljiva d_i razdvaja skup odobrenih kredita na dobre i loše i to na sledeci nacin, ukoliko je vreme preživljavanja kredita manje od praga \bar{t}_i tada je klijent loš, u suprotnom je dobar, tj.

$$d_i = \begin{cases} 0, \text{ako } t_i^* \leq \bar{t}_i \\ 1, \text{ako } t_i^* > \bar{t}_i \end{cases}$$

Za ocenjivanje parametara koristimo sledecu jednacinu:

$$\ln l = \sum_{i=1}^N (1-y_i) \ln(1-\Phi(x_{1i}\mathbf{b}_1)) + \sum_{i=1}^N y_i (1-d_i) \left\{ \ln \Phi\left(\frac{x_{1i}\mathbf{b}_1 - \frac{\mathbf{s}_{12}}{\mathbf{s}_2} (t_i - x_{2i}\mathbf{b}_2)}{\sqrt{(1-\mathbf{r}^2)}}\right) - \frac{1}{2} \ln 2\mathbf{p} + \ln\left(\frac{1}{\mathbf{s}_2}\right) - \frac{1}{2} \left(\frac{t_i - x_{2i}\mathbf{b}_2}{\mathbf{s}_2}\right)^2 \right\} + \sum_{i=1}^N y_i d_i \ln \Phi_2(x_{1i}\mathbf{b}_1, \frac{x_{2i}\mathbf{b}_2 - \bar{t}_i}{\mathbf{s}_2}; \mathbf{r})$$

gde $\Phi(\bullet)$ i $\Phi(\bullet, \bullet; \rho)$ predstavljaju jednodimenzionalnu i binarnu standardnu normalnu funkciju raspodele, poslednja sa koeficijentom korelacije ρ .

U Tabeli 4 su prikazani jednodimenzionalne i binarne ocene za parametre β_1 i β_2 koje su dobijene korišcenjem funkcije maksimalne verovatnoće. Jednodimenzionalne ocene za β_1 i β_2 se dobijaju kada se pojedinačno ocenjuju prva i druga jednacina u (*), pri cemu se prepostavlja da je $\rho = 0$. Dvodimenzionalne ocene su dobijene kada su jednacine u (*) rešavane zajedno, tako da su uzeti u obzir efekti odabira uzorka, takode neophodna je i ocena koeficijenta korelacije ρ . Binarna probit ocena α_2 se dobija kada se u (*) umesto druge jednacine koristi jednacina koja meri default rizik, tj. rizik da ce klijent biti loš.

Razliciti pristupi kreditnom scoring sistemu

Tabela 4: Jednodimenzionalne i binarne ocene za parametre; standardne greške su date u zagradama, a promenljive koje su znacajne na nivou znacajnosti 10% su podebljane

Promenljive	Jednacine				
	Odobravanje kredita		Preživljavanje kredita		Nije loš
	b ₁	b ₂	a ₂	Binarna	Binarna
Jednodimenzionalna	Binarna	Jednodimenzionalna	Binarna	Binarna	Binarna
Konstatnta	-0.336 (0.051)	-0.328 (0.051)	8.246 (0.156)	9.065 (0.193)	2.460 (0.113)
Veliki grad	-0.232 (0.027)	-0.222 (0.027)	-0.128 (0.058)	0.233 (0.053)	-0.042 (0.055)
Muško	-0.207 (0.028)	-0.196 (0.028)	-0.106 (0.061)	0.024 (0.058)	0.024 (0.059)
Razveden	-0.186 (0.040)	-0.179 (0.039)	-0.124 (0.078)	-0.009 (0.069)	-0.071 (0.074)
Kuca	0.110 (0.028)	0.102 (0.028)	0.061 (0.061)	-0.023 (0.062)	-0.011 (0.060)
Oženjen	-0.242 (0.030)	-0.233 (0.030)	0.187 (0.068)	0.344 (0.066)	0.253 (0.066)
Broj zahteva	-0.007 (0.005)	-0.004 (0.005)	-0.116 (0.012)	-0.097 (0.011)	-0.104 (0.010)
Preduzece	0.570 (0.064)	0.570 (0.064)	0.135 (0.181)	0.162 (0.160)	0.147 (0.157)
Prihod	0.901 (0.018)	0.886 (0.018)	0.038 (0.042)	-0.286 (0.050)	-0.226 (0.050)
Razlika prihoda	-0.243 (0.035)	-0.237 (0.035)	0.147 (0.075)	0.185 (0.079)	0.206 (0.073)
Kapital	-0.284 (0.051)	-0.272 (0.050)	-0.057 (0.123)	0.195 (0.100)	0.142 (0.127)
Nula limit	-2.253 (0.106)	-2.218 (0.114)	-2.244 (0.401)	-0.328 (0.140)	-0.703 (0.306)
Limit	-0.861 (0.019)	-0.848 (0.021)	0.006 (0.059)	0.561 (0.050)	0.486 (0.056)
Broj kredita	0.086 (0.007)	0.086 (0.007)	0.329 (0.026)	0.259 (0.023)	0.270 (0.020)
Upotrebljen limit	-0.747 (0.045)	-0.759 (0.045)	-1.295 (0.120)	-1.223 (0.116)	-1.210 (0.093)
Iznos kredita	-	-	-0.069 (0.073)	-0.070 (0.070)	-0.067 (0.069)
s^2	-	-	0.919 (0.045)	0.196 (0.057)	-
ρ	-	-	-	-0.986 (0.021)	-0.911 (0.056)

U Tabeli 4 vidimo da osobe sa vecim prihodom imaju veci rizik da ce postati loši od osoba sa manjim prihodom. Ostale promenljive koje imaju koeficijente od znacaja u ovom primeru su «*Broj zahteva*», «*Nula limit*», «*Broj kredita*», «*Limit*» i «*Upotrebljeni limit*». Broj zahteva za informacije o klijentu koje kreditni biro primi je samo znak koliko je klijentu potreban novi kredit, a to samim tim negativno utice na preživljavanje tog kredita, jer je to rani znak da je klijent u problemu. Ukoliko klijent do sada nema kredita, što je predstavljeno promenljivom «*Nula limit*», je znak da je klijent neiskusan u otplati kredita, što može dovesti do toga da brzo ode u default, pa samim tim negativno utice na vreme preživljavanja. Za klijente koji vec imaju kredit se prepostavlja da vec imaju dovoljno iskustva u servisiranju duga, pa samim tim «*Broj kredita*» i «*Limit*» pozitivno uticu na preživljavanje kredita. Promenljiva «*Upotrebljeni limit*» označava u kom opsegu klijent koristi odobreni kredit, a samim tim i utice na smanjenje preživljavanja.

Ako posmatramo poslednju kolonu *Tabele 4* vidimo da probit parametri koji odreduju verovatnocu da klijent biti loš (α_2) imaju isti znak, sem u slučaju promenljive «*Kapital*», kao parametri β_2 u tobit modelu koji ocenjuje vreme preživljavanja. Ipak, promenljive kao što su «*Muško*», «*Razveden*», «*Kuca*», «*Veliki grad*» i «*Preduzece*», koje su znacajne u donošenju odluke nemaju uticaja pri analizi rizika da ce klijent biti loš niti pri analizi vremena preživljavanja, dok promenljiva «*Broj kredita*» ima uticaja na preživljavanje kredita, a ne igra nikakvu ulogu pri odlucivanju da li treba odobriti kredit. Pošto svi sem jednog parametra imaju isti znak u jednacini preživljavanja kao odgovarajuci parametri u binarnom probit modelu, ove promenljive u isto vreme smanjuju, odnosno povecavaju ocekivano vreme preživljavanja i stopu prinosa od kredita. Drugim recima, ako banka ne teži da minimizira rizik da ce klijent biti loš, razlog tome je taj što krediti sa vecim rizikom imaju duže vreme preživljavanja, a samim tim se ocekaju i veci prinosi.

Ovi rezultati su potvrdili da banka ne teži da minimizira rizik. Neke promenljive koje povecavaju (smanjuju) verovatnocu da klijentu bude odobren kredit u isto vreme smanjuju (povecavaju) ocekivano vreme preživljavanja (a samim tim i prinos) i povecavaju (smanjuju) verovatnocu da ce klijent biti loš. Takođe, pokazalo se da na odluku banke ne utice iznos kredita. Ovo nas vodi do zaključka da banka ne teži ni ka tome da maksimizira profit, a samim tim ni vreme preživljavanja kredita. Ali sigurno je da banka u svojoj politici ima definisan cilj kome teži, ka maksimizaciji ili minimizaciji nekih drugih parametara, kao što su na primer broj klijenata, profit od odredene vrste proizvoda, ili sличno.

4.4 Primena logit modela

Kako se bankarsko tržište sve više širi, ne postoji konkurenčija samo među domaćim bankama, nego je ona sve jaca i između stranih banka. Sve je veća tražnja za potrošackim kreditima, pa samim tim je neophodno usmeriti pažnju na upravljenje rizikom, što nas vodi do razvijanja dobrog kreditnog skoring modela kojim se ocenjuje kreditna sposobnost klijenta, kolika je verovatnoca da će isti redovno otplatiti.

Kreditni skoring sistem koristi istorijske podatke o kreditima i podatke o klijentu da odredi koje klijentove karakteristike najbolje razdvajaju dobre od loših klijenata. Kada se napravi jedan ovakav model, on se zatim može primeniti na nove klijente za koje nije poznato sa kojom verovatnocom će postati loši *PD* (*probability of default*). Na osnovu ocenjenog kreditnog skoring modela mogu da se izracunaju skorovi za svakog novog klijenta, pri cemu viši skor označava manji PD. Ovaj skor zatim treba da se uporedi sa granicnom vrednošću kreditnog skoring sistema koja određuje koji zahtev za kredit će biti odobren, a koji odbijen, a koju određuje sama banka u zavisnosti od toga koliki je rizik spremna da prihvati.

Kao što smo vec videli postoje razni metodi koji mogu da se koriste za razvijanje kreditnog skoring modela, kao što su na primer probit analiza, diskriminantna analiza, linearna regresija i sличno, međutim ovde se koristi logaritamska regresija, tj. *logit model* primenjen na uzorak klijenata koji cine Vijetnamsko finansijsko tržište («*A credit scoring model for Vietnam's retail banking market*» [5]).

Da bi mogao da se primeni ovaj model prvo treba definisati koje će karakteristike da uđu u model. Ne postoji nikakvo pisano pravilo o broju i tipu promenljivih koje će ući u model. Uglavnom se tu nalaze karakteristike koje pokazuju finansijske mogućnosti klijenta kao što se prihodi i vrednost kuće koja je u njegovom vlasništvu, a takođe i promenljive koje to indirektno pokazuju, kao što su obrazovanje, zatim broj nekretnina, broj godina na istom poslu i sличno. Treba obratiti pažnju i na zakon koji u nekim zemljama zabranjuje da se neke karakteristike, kao što su pol, religija i sличno, uključe u model zbog mogućnosti diskriminacije.

Nakon što se ocene koeficijenti β_i uz pomoć metoda ocene maksimalne verovatnoće treba testirati model i odrediti sa kolikom preciznošću je ocenjen model. Ovaj proces se naziva *validacija*. Uglavnom se uzima uzorak koji nije korišćen za izvođenje modela, *out-of-sample*. Prvo se ocenjuju verovatnoće da će klijent biti loš za ovaj uzorak. Ovi PD-evi se upoređuju sa granicnom vrednošću da bi odredili da li će podnositelj zahteva biti dobar ili loš klijent. Granicnu vrednost određuje sama banka, na osnovu toga koliki rizik je spremna da prihvati. Ako na primer izaberemo da je granicna vrednost 50%, tada će klijent ciji je ocenjeni PD veci (manji) od 50% biti klasifikovan kao loš (dobar). U ovoj fazi se koristi *matrica klasifikacije*, gde G_g predstavlja broj tacno klasifikovanih dobrih kredita, a G_b broj dobrih kredita koji su pogrešno klasifikovani. Slicno, B_b predstavlja broj dobro klasifikovanih loših kredita, a B_g broj pogrešno klasifikovanih loših kredita. Procenat tacno klasifikovanih kredita (PCC) služi kao mera preciznosti. Procenat tacno klasifikovanih dobrih kredita

(PCC_{good}) je definisan kao odnos tacno klasifikovanih dobrih kredita i ukupnog broja posmatranih kredita. Slicno, procenat tacno klasifikovanih loših kredita (PCC_{bad}) je definisan kao odnos tacno klasifikovanih loših kredita i ukupnog broja posmatranih kredita. I na kraju, procenat tacno klasifikovanih kredita (PCC_{total}) je definisan kao odnos tacno klasifikovanih kredita i ukupnog broja kredita.

Matrica klasifikacije

Stvarna opažanja	Ocekivana opažanja		PCC
	dobili	loši	
dobri	G_g	G_b	$PCC_{good}=G_g/(G_g+G_b)$
loši	B_g	B_b	$PCC_{bad}=B_b/(B_b+B_g)$
ukupno			$PCC_{total}=(G_g+B_g)/(G_g+G_b+B_b+B_g)$

Banka možda želi da minimizira obe greške, B_g i G_b . Medutim, smanjivanje B_g dovodi do povecanja G_b i suprotno. Iz tog razloga treba razmotriti kakav je odnos u troškovima (gubicima) koji nastaju pri pogrešnoj klasifikaciji i na osnovu tog kriterijuma odluciti koju grešku treba minimizirati.

Sada cemo se upoznati sa primerom vijetnamskog bankarskog tržišta koje je i dalje dosta nerazvijeno i još uvek u mnogome zaostaje za razvijenim industrijskim zemljama. Bankarski sistem i dalje pati od nedostatka kapitala, adekvatne zaštite od mogucih gubitaka, niske profitabilnosti. Ali on se konstantno razvija.

Da bi razvili kreditni skoring sistem uzeti su podaci o kreditima iz jedne Vijetnamske komercijalne banke u periodu izmedu 1992. i 2005. godine. Dobijamo uzorak od 56,037 kredita. Banka klasificuje klijenta kao lošeg ukoliko je u kašnjenju dužem od 90 dana. Naš uzorak sadrži 3.3% loših klijenata. Pošto uzorak sadrži samo informacije o klijentima kojima je kredit odobren, on nije dovoljno reprezentativan za buduce klijente. Takođe, ne možemo da predvidimo kako bi se odbijeni klijenti ponašali da su odobreni. Jedno od rešenja je out-of-sample kalibracija, pri cemu uzorak od 56,037 kradita delimo na dva poduzorka. Pocetni uzorak sadrži 30,994 kredita od cega su 1026 (3.3%) loši krediti i *hold-out* uzorak koji sadrži 25,043 kredita od cega su 798 (3.2%) loši krediti.

Da bi se ocenio kreditni skoring sistem, koristi se *stepenasta metoda unapred*. Ova metoda pocinje sa modelom koji ne sadrži nijednu nezavisnu promenljivu i postepeno dodaje promenljive. U svakom koraku, dodaje po jedna promenljiva koja ima najvecu moc u poboljšanju tacnosti ocenjivanja, i pri tome se posmatra nivo znacajnosti koji je manji od 5%. Proces se zaustavlja kada se više ne može dodati nijedna promenljiva sa nivoom znacajnosti manjim od 5%. Na ovaj nacin smo od pocetne 22 promenljive, skup sveli na 16 nezavisnih promenljivih. Da bi se osigurali da su ovo stvarno najznacajnije promenljive, koristi se stepenasti metod unazad. Ovaj model u pocetku sadrži 22 promenljive, i u svakom koraku izbacuje po jednu koja ima najslabiji uticaj i ovako dobijamo konacnih 16 promenljivih kao i stepenastom metodom unapred.

Promenljive koje su korišcene u modelu su sledeće:

- ◆ *Obrazovanje* – možemo da ocekujemo da su obrazovaniji ljudi mnogo stabilniji, imaju veće prihode i samim tim niži PD.
- ◆ *Pol* – u nekim zemljama ova promenljiva mora da se iskljuci iz modela, zbog mogucnosti diskriminacije.
- ◆ *Region* – predstavlja deo zemlje u kojoj klijent živi. Smatra se da ljudi sa sличnim bogatstvom žive na istoj lokaciji. Dakle, geografski kriterijum može da pokazuje nivo finansijskog bogatstva klijenta
- ◆ *Vreme na trenutnoj adresi* – predstavlja broj godina koliko klijent živi na istoj adresi. Ova promenljiva može da označava klijentovu zrelost, stabilnost i izbegavanje rizika. Međutim, u Vijetnamu PD se povećava sa brojem godina na trenutnoj adresi, pošto ljudi što više zaraduju i napreduju, teže da žive u što boljim uslovima, pa samim tim i cesto menjaju adresu stanovanja.
- ◆ *Stambeno pitanje* – pokazuje da li je klijent vlasnik svog doma, iznajmljuje stan, ili živi sa roditeljima.
- ◆ *Bracno stanje* – u našem uzorku verovatnoca da će klijent biti loš je veća kod ljudi koji su u braku nego za samce.
- ◆ *Broj ljudi koje klijent izdržava* – što je veći broj ovih zavisnika, to se i PD povećava, jer je veće opterenje na klijentove prihode.
- ◆ *Telefon* – meri da li klijent ima kućni telefon, a time i koliko lako banka može da održava kontakt sa klijentom. Ukoliko klijent nema telefon, to mu je viši PD.
- ◆ *Svrha kredita* – opisuje kako će odobrena sredstva biti iskorišćena
- ◆ *Tip kolateral-a* – opisuje koji tip obezbedenja pokriva kredit.
- ◆ *Vrednost kolateral-a*
- ◆ *Trajanje kredita* – u Vijetnamu trajanje kredita predlaže klijent, pa samim tim i ova promenljiva pokazuje koliko je klijent spreman da prihvati rizik, kao i samoocenjivanje sposobnosti otplate kredita.
- ◆ *Vreme u banci* – meri koliko godina klijent posluje preko te banke.
- ◆ *Broj kredita* – racuna broj kredita koliko je klijent primio od kada je u banci
- ◆ *Tekuci racun* – je binarna promenljiva koja predstavlja da li klijent održava tekuci racun
- ◆ *Štedni racun* – je binarna promenljiva koja predstavlja da li klijent održava štedni racun

Tabela 1: ocjenjeni kreditni scoring model

Promenljive	Ocenjeni koeficijenti	Standardne greške	Nivo znacajnosti
Vreme u banci	-1.774	0.121	0.0%
Pol	-1.557	0.222	1.0%
Broj kredita	-0.938	0.051	1.4%
Trajanje kredita	-0.845	0.080	3.7%
Štedni racun	-0.750	0.104	3.1%
Region	-0.652	0.030	13.6%
Stambeno pitanje	-0.551	0.278	44.6%
Tekuci racun	-0.492	0.208	10.4%
Vrednost kolateralna	-0.402	0.096	9.8%
Broj zavisnika	-0.356	0.096	9.9%
Vreme na trenutnoj adresi	-0.285	0.054	2.5%
Bracno stanje	-0.233	0.101	68.1%
Tip kolateralna	-0.190	0.057	53.0%
Telefon	-0.181	0.047	3.4%
Obrazovanje	-0.156	0.067	60.3%
Svrha kredita	-0.125	0.054	3.3%
Konstanta	-3.176	0.058	4.6%

Kao što vidimo u ovoj tabeli, od datih 16 promenljivih, vreme u banchi je najznačajniji pokazatelj, a zatim slede pol, broj kredita i trajanje kredita. Negativni koeficijent koji stoji uz promenljivu pol znači da žene imaju manju verovatnocu da će kasniti sa placanjem od muškaraca. Trajanje kredita kao mera klijentove spremnosti prihvatanja rizika i samoocenjivanja je jedinstveno u Vijetnamu. Ovo pokazuje da banke ne mogu da se pouzdaju samo u sopstvene procene nego se oslanjaju i na klijentovu poštenu ocenu trenutnog stanja.

Da bi pravilno ocenili preciznost modela koristimo out-of-sample. Tada je srednja vrednost predvidjene verovatnoće da će klijent biti loš jednaka 1.73% za dobre klijente za razliku od 49.05% za loše klijente. Za loše klijente PD se kreće u opsegu od 0.01% do 96.81%, a za dobre od 0.00% do 73.54%. U sledećoj tabeli su date vrednosti koje ukazuju na preciznost predviđanja modela.

Tabela 2: preciznost predviđanja sa granicnom vrednošću 0.50

Uoceno	Predvidene vrednosti			
	Dobili	Loši	PCC	
Dobili	24,136	109	99.55%	= PCCgood
Loši	397	401	50.25%	= PCCbad
Ukupno			97.98%	= PCCtotal

4.5 Primer genetskog programiranja

Genetsko programiranje je proširenje tehnike genetskih algoritama. Mi necemo dublje ulaziti u ovaj problem, jer to prevazilazi okvire ovog rada, nego cemo samo navesti primer (uzet iz rada «*Genetic programming for credit scoring: the case of Egyptian public sector banks*» [7]) u kome cemo uporediti rezultate koji se dobijaju genetskim programiranjem i probit analizom. Koristimo dva tipa modela genetskog programiranja: *prost model* i *tim model* koji je kombinacija više prostih modela ciji je cilj da se postignu bolji rezultati nego kad se primenjuje pojedinacni prost model.

Za potrebe pravljenja modela korišcena je baza podataka o potrošackim kreditima dostavljena od strane Egipatskog bankarskog sektora. Baza podataka se sastoji od 1262 kredita, od cega je 851 kredit (67.43%) dobar, a 411 (32.57%) loš.

Prvo se pravi model na celom uzorku, a zatim se za potrebe testiranja tacnosti ocenjivanja skoring sistema koristi *training uzorak* koji sadrži 846 (67%) slučajeva, a testiranje modela se vrši na *uzorku za testiranje* koji sadrži 416 (33%) slučaj.

Za uporedivanje rezultata uzeta su u obzir cetiri kriterijuma: procenat tacne klasifikacije (*PCC*), kriterijum ocenjenih troškova koji nastaju pogrešnom klasifikacijom (*OT*) sa koeficijentom troškova pogrešne klasifikacije 5:1, kriterijum *OT* sa koeficijentom troškova pogrešne klasifikacije 7:1, kriterijum *OT* sa koeficijentom troškova pogrešne klasifikacije 10:1.

Sa pojmom procenta pogrešne klasifikacije i matricom klasifikacije smo se upoznali u delu 4.3.2, a sada cemo dati osnovu ideju kriterijuma ocenjenih troškova koji nastaju pogrešnom klasifikacijom. Ovaj kriterijum daje ocenu efikasnosti skoring sistema. Za izracunavanje *OT*-a koristi se sledeća jednacina:

$$OT = C(I)x(G_b/TG)x(TG/TN) + C(II)x(B_g/TB)x(TB/TN)$$

gde je *C(I)* trošak nastao pogrešnom klasifikacijom koja je povezana sa greškom prvog reda; (*G_b/TG*) je verovatnoca da ce se napraviti greška prvog reda, predstavljena kao odnos broja dobrih klijenata koji su ocenjeni kao loši (*G_b*) i ukupnog broja dobrih klijenata; *TG/TN* je verovatnoca da je klijent dobar, tj. odnos broja dobrih *TG* i ukupnog broja klijenata *TN*; *C(II)* je trošak nastao pogrešnom klasifikacijom koja je povezana sa greškom drugog reda; (*B_g/TB*) je verovatnoca da ce se napraviti greška drugog reda, predstavljena kao odnos broja loših klijenata koji su ocenjeni kao dobri (*B_g*) i ukupnog broja loših klijenata; *TB/TN* je verovatnoca da je klijent loš, tj. odnos broja loših *TB* i ukupnog broja klijenata *TN*.

Gornja jednacina može da se napiše i na sledeći nacin:

$$OT = C(I)x(G_b/TN) + C(II)x(B_g/TN)$$

4.5.1 Upoređivanje genetskog programiranja i probit analize na celom uzorku

Tabela 1 sumira stope tacne klasifikacije i ocenjene troškove koji nastaju pogrešnom klasifikacijom, i to primenom genetskog programiranja i probit analize. Vidimo da GP_t (tim model genetskog programiranja) ima najveću stopu tacne klasifikacije. GP_p (prost model genetskog programiranja) je model koji najbolje klasificiše dobre klijente (91.89%), a GP_t je model koji najbolje klasificiše loše klijente (74.94%). GP_t je izabrani model pošto ima najnižu ocenu troškova pogrešne klasifikacije, a takođe ima i najbolju stopu tacne klasifikacije.

Tabela 1: uporedivanja rezultata klasifikacije, grešaka i ocenjenih troškova na celom uzorku

Model	Tacno klasifikovani rezultati			Rezultati o greškama		OT	OT	OT
	G%	B%	PCC	I tipa	II tipa	(5:1)	(7:1)	(10:1)
Ceo uzorak								
PA	88.95	67.40	81.93	0.1105	0.3260	0.6054	0.8178	1.1363
GP_p	91.89	65.45	83.28	0.0811	0.3455	0.6173	0.8424	1.1800
GP_t	91.07	74.94	85.82	0.0893	0.2506	0.4683	0.6316	0.8764

* G – dobri klijenti; B – loši klijenti

4.5.2 Upoređivanje genetskog programiranja i probit analize na poduzorku

Koristi se uzorak za testiranje da bi se testirala moć ocenjivanja razvijenog skoring modela. Vidimo da GP_t ima najvišu stopu tacne klasifikacije. Takođe, GP_t model pravi najveću grešku pri klasifikaciji dobrih klijenata, a GP_p pri klasifikaciji loših klijenata. U ovom slučaju se nije lako odluciti koji model je najbolji, pošto na primer GP_t ima najvišu stopu tacne klasifikacije ali u isto vreme i dosta visoke troškove pogrešne klasifikacije. U obzir bi mogao doci GP_p model pošto ima podnošljivo male troškove pogrešne klasifikacije i dosta visoku stopu tacne klasifikacije.

Tabela 2 : uporedivanja rezultata klasifikacije, grešaka i ocenjenih troškova na poduzorku

Model	Tacno klasifikovani rezultati			Rezultati o greškama		OT	OT	OT
	G%	B%	PCC%	I tipa	II tipa	(5:1)	(7:1)	(10:1)
Poduzorak								
PA	90.00	65.87	82.69	0.1000	0.3413	0.6232	0.8456	1.1790
GP_p	88.62	69.84	82.93	0.1138	0.3016	0.5679	0.7644	1.0590
GP_t	94.48	59.52	83.89	0.0552	0.4048	0.6964	0.9601	1.3557

* G – dobri klijenti; B – loši klijenti

4.6 Primena analize obavijanja podataka (DEA)

U dosadašnjim primerima smo videli kako se kreditni scoring model primjenjuje kada se analizira kreditna sposobnost pojedinacnog klijenta. U ovom poglavlju cemo se malo upoznati i sa primerom kreditnog scoring sistema u situaciji kada firma podnosi zahtev za kredit. Kao što smo videli, u slučaju individue posmatraju se karakteristike kao što su starost klijenta, njegovo finansijsko stanje i obaveze i slično. U slučaju firmi, bitni su nam podaci kao što o njihovoj aktivi, pasivi, rashodima i prihodima, i na osnovu tih podataka dolazimo do zaključka o likvidnosti te firme, njihovoj profitabilnosti, produktivnosti, kao i o strukturi troškova. U bankama još uvek nije zastupljen ovaj tip kreditnog scoringa, pošto je on i dalje u fazi razvoja, ali samo je pitanje trenutka kada će kreditni scoring poceti da se primjenjuje i za analizu kreditne sposobnosti firmi. U ovom radu cemo predstaviti kako se pitem analize obavijanja podataka dobijaju skorovi koji nam kasnije služe za rangiranje firmi.

DEA (Data Envelopment Analysis), tj. analiza obavijanja podataka je metoda linearog programiranja za ocenu relativne efikasnosti organizacionih jedinica koje koriste više razlicitih inputa za stvaranje više razlicitih outputa i ona za razliku od predašnjih modela kojima trebaju istorijski podaci da bi napravili model, zahteva samo *observed* (uoceni) skup ulaznih i izlaznih podataka da bi izracunala kreditni scor. DEA analiza je neparametarski metod operacionih istraživanja koji se koristi u ekonomiji za ocenjivanje granične proizvodnje.

Proces se sastoji od šest koraka, od čega se prva tri bave odabiru firmi koje će ući u model i indikatora koji se koriste za ocenu finansijske situacije u toj firmi, četvrti korak koristi DEA da bi se izracunali kreditni skorovi firmi, peti korak proverava valjanost ovako dobijenih skorova tako što ih poređi sa onima koji su dobijeni putem regresije i diskriminantne analize, i konacno šesti korak predlaže metod za kreditni rejting uz pomoć raspodele odnosa dobrih i loših klijenata.

Uzimaju se u obzir firme koje podnose zahtev za novi kredit ili one koje hoće da promene vec postojeći kreditni limit. Na početku istraživanja posmatrano je 1400 firmi, ali *outlajeri*, tj. firme koje su imale koeficijente koji su znacajno odstupali (više od dve standardne devijacije) od odgovarajućih srednjih vrednosti su isključene iz modela i ostala je 1061 firma.

4.6.1 Odabir finansijskih koeficijenata

Za definisanje finansijskih koeficijenata uglavnom se koriste opšte prihvacene finansijske dimenzije kao što su rast prihoda, likvidnost, profitabilnost, produktivnost i struktura troškova. Tako da se razliciti finansijski koeficijenti grupišu po ovim dimenzijama. U ovom primeru («*A practical approach to credit scoring*» [3]) je odabrano šest koeficijenata koji su klasifikovani kao ulazne i izlazne promenljive za DEA.

Ulazni podaci, tj. koeficijenti koji treba da se minimiziraju su *FE*-finansijski troškovi u odnosu na prodaju, *CL*-koeficijent obaveza (trenutne obaveze u odnosu na kapital) i *TB*-ukupno zaduženje u odnosu na ukupnu aktivu. Koeficijent finansijskih troškova u odnosu na prodaju (*FE*) pokazuje sposobnost firme da plati svoje troškove,

što pokazuje kreditnu sposobnost firme. Koeficijent obaveza (CL) je indikator stabilnosti strukture kapitala. Povecanje ovog koeficijenta pokazuje nestabilnost strukture kapitala i likvidnosti firmi. Ukoliko se koeficijent ukupne zaduženosti (TB) povecava to pokazuje da se profitabilnost i stabilnost firme smanjuju.

Izlazni podaci koji treba da se maksimiziraju su koeficijent adekvatnosti kapitala (CA) koji predstavlja odnos kapitala firme i ukupne aktive, zatim pokazatelj trenutne likvidnosti (CR) i koeficijent pokrivenosti kamata (IC). Što više firme pokriva obaveze iz sopstvenih sredstava, tj. što je viši CA, to je ona ocenjena kao manje rizicna. Takođe, što je veci CR to je firma više likvidna i manja je verovatnoca da će upasti u probleme. Koeficijent pokrivenosti kamata pokazuje sposobnost firme da placa troškove kamata iz operativnog prihoda, a samim tim što je ovaj koeficijent veci to je veca profitabilnost firme.

4.6.2 Racunanje skorova uz pomoc DEA

Kada smo definisali ulazne i izlazne promenljive, možemo da primenimo *CCR model* da bi izracunali DEA skorove:

$$\begin{aligned} \min \mathbf{q}_k - \mathbf{e} \sum_{r=1}^s s_r^+ - \mathbf{e} \sum_{i=1}^m s_i^- \\ \text{uz uslove } \sum_{j=1}^N I_j x_{ij} = \mathbf{q}_k x_{ik} - s_i^-; i = 1, \dots, m \\ \sum_{j=1}^N I_j y_{ij} = y_{rk} + s_r^+; r = 1, \dots, s \\ I_j, s_r^+, s_i^- \geq 0; \forall j, r, i, \mathbf{q}_k \end{aligned}$$

gde su

θ_k – skor koji pokazuje kreditne sposobnost za firmu k

N – broj firmi u uzorku

λ_j – težišna vrednost firme j

y_{rj} – r -ti izlazni koeficijent za firmu j

y_{rk} – r -ti izlazni koeficijent za firmu k

x_{ij} – i -ti ulazni koeficijent za firmu j

x_{ik} – i -ti ulazni koeficijent za firmu k

s_r^+ , s_i^- – dopunske promenljive za r -to ogranicenje i i -to ogranicenje

DEA skorovi su dati u obliku procenta, pri cemu firme sa DEA skorom od 100 predstavljaju najbolje firme i one pripadaju «DEA efikasnoj granici».

4.6.3 Provera valjanosti DEA skorova

Cilj ovog koraka je da oceni u kom procentu se DEA skorovi poklapaju sa skorovima koji se dobijaju putem regresione ili diskriminantne analize.

Linearna regresija se koristi kao kriterijum za testiranje objašnjavajuce moci indikator promenljivih u DEA. Za ovu svrhu, DEA skorovi se uzimaju kao zavisne, a DEA koeficijenti kao nezavisne promenljive.

Tabela 1: rezultati regresione analize

$R^2 = 0.741; F = 491.803$				
	Koeficijenti	Standardne greške	t-vrednost	p-vrednost
Konstanta	72.94035	3.772205	19.33626	0.0000
FE	-109.653	21.6798	-5.05785	0.0000
CL	-50.8614	3.123891	-16.2814	0.0000
TB	-56.3548	3.039185	-18.5427	0.0000
CA	47.45822	4.30213	11.03133	0.0000
CR	17.36333	1.085729	15.99233	0.0000
IC	1.146138	0.184531	6.211091	0.0000

Kao što se vidi u Tabeli 1 znaci koeficijenata koji stoje uz promenljive su očekivani (tj. minus uz promenljive koje negativno uticu i plus uz one koje pozitivno uticu) i sve su statisticki znacajne (p – vrednost je 0.0000), što pokazuje da je DEA algoritam uspešno izracunao vrednosti za ovih šest koeficijenata. Kada vrednosti koje smo dobili putem regresije uvrstimo u jednacinu, imamo:

$$\text{DEA} = 73 - 109.7 \text{ FE} - 50.9 \text{ CL} - 56.4 \text{ TB} + 47.5 \text{ CA} + 17.4 \text{ CR} + 1.2 \text{ IC}$$

Koristeci ovu jednacnu može da se izracuna linearna aproksimacija DEA skora, tako da kada se pojavi novi klijent ne mora da se ponavlja ceo DEA algoritam.

Diskriminatna analiza se koristi da utvrdi koliko dobro DEA skorovi klasifikuju firme u dve grupe: dobre i loše. Firme su podeljene u dve grupe na osnovu DEA skorova. Granicna vrednost izmedu ovih grupa se odreduje uzimajuci u obzir raspodelu DEA skorova. U ovom primeru za granicnu vrednost je uzeta srednja vrednost, dakle, 518 firme su klasifikovane kao dobre, a ostatak kao loše. Diskriminantna analiza izvodi funkciju diskriminacije koja ukuljuje pet od šest koeficijenata kao nezavisne promenljive (koeficijent IC je iskljucen).

Table 2: rezultati diskriminantne analize

Prognozirana grupa	Odabrana grupa		
	dobili	loši	ukupno
dobili	400 (84.9%)	57 (11.0%)	497
loši	78 (15.1%)	461 (89.0%)	539
ukupno	518	518	1036

Tabela 2 pokazuje da je ukupno dobro klasifikovano 87.0% populacije $((440+461)/(518+518) = 0.87)$. Dakle, greška klasifikacije je 13%.

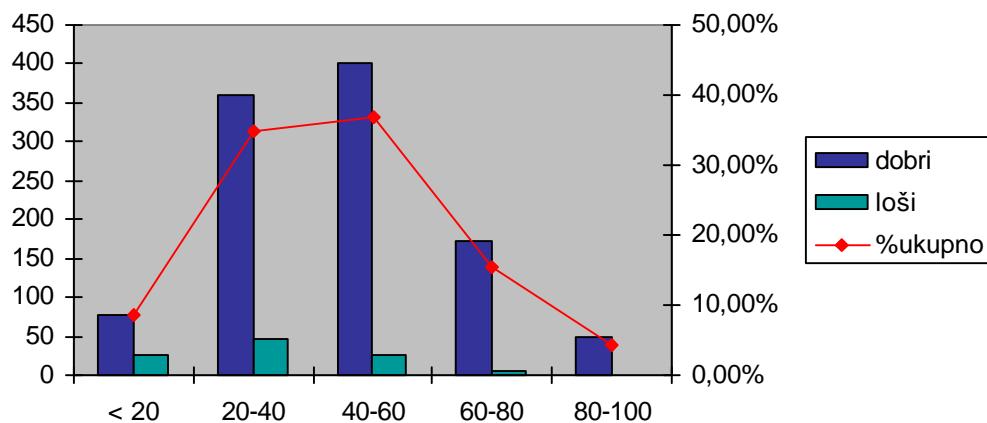
4.6.4 Metod kreditnog rejtinga

DEA skorovi klijenta bi trebalo da budu klasifikovani u klase kao što su A, B, C i tako dalje. Pošto se DEA skorovi kreću u intervalu od 0 do 100, može se napraviti podela da skorovi od 0 do 20 pripadaju jednoj klasi, od 20 do 40 drugoj klasi i tako dalje.

Tabela 3: raspodela ucestalosti dobrih i loših firmi

Klase	dobri	loši	Ukupno	% loših	% ukupno
Ispod 20	76	25	101	24.8%	8.7%
20-40	360	46	406	11.3%	34.9%
40-60	402	27	429	6.3%	36.9%
60-80	173	5	178	2.8%	15.3%
80-100	50	0	50	0.0%	4.3%
Ukupno	1061	103	1164	8.8%	

Rezultati se mogu prikazati i graficki, tako što je stubicima prikazan broj dobrih i loših klijenata, a kriva prikazuje raspodelu DEA skorova.

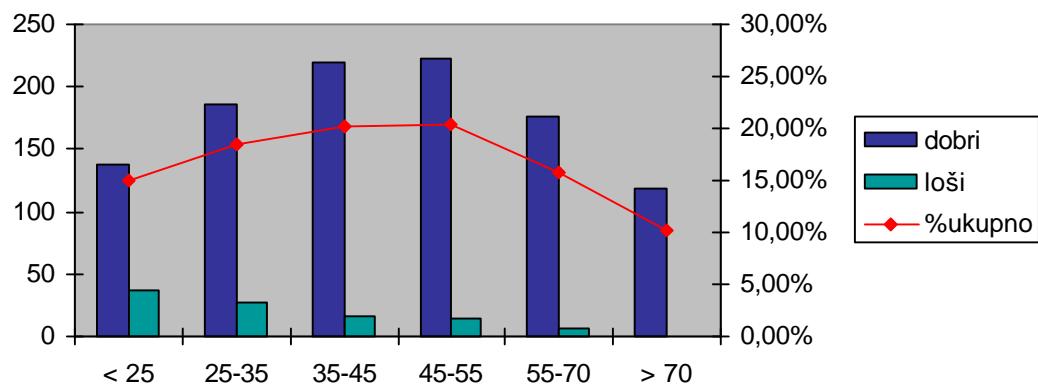


Kao što se vidi u ovoj tabeli, ucestalost loših klijenata (%loših) opada kako DEA skorovi rastu. To znači da ovi skorovi služe kao veoma korisni pokazatelji za rangiranje klijenata po njihovoј kreditnoј sposobnosti. Tako na primer, klasa 80-100 može da bude označena kao «A», klasa 60-80 kao «B» i tako dalje. Tabela 3 (kao i grafik) pokazuje da je raspodela DEA skorova (%ukupno) iskošena na desno, što ukazuje na to da mnogo firmi ima niske skorove. A većina komercijalnih banaka teže ka tome da rangiraju klijente tako da dobiju normalnu raspodelu. Ovaj cilj se postiže modifikacijom klasa DEA skorova. U Tabeli 4 su prikazane modifikovane klase kojima se postiže normalna raspodela.

Tabela 4: modifikovana raspodela ucestalosti dobrih i loših klijenata

Klasa	dobri	loši	Ukupno	% loših	% ukupno
Ispod 25	138	37	175	21.1%	15.0%
25-35	186	28	214	13.1%	18.4%
35-45	219	16	235	6.8%	20.2%
45-55	222	15	237	6.3%	20.4%
55-70	177	7	184	3.8%	15.8%
Iznad 70	119	0	119	0.0%	10.2%
Ukupno	1061	103	1164	8.8%	

Kada rezultate predstavimo graficki vidimo da je raspodela DEA skorova stvarno normalna.



5. Zaključak

Kreditni skoring sistem ne daje objašnjenje zašto je neki klijent odbijen, nego samo daje ocenu koliko je neki klijent rizican i upozorava da bi tog klijenta trebalo odbiti. Međutim, ko nereskira taj neprofitira. Zato je na banci da odluci koliki rizik želi da prihvati, a rezultati kreditnog skoringa služe kao odredene smernice.

Kao što smo videli kroz rad, postoje dva osnovna pristupa kredit skoringu. Prvi je pomocu separacije, a drugi je putem regresije. Cilj nam je da što tacnije klasifikujemo klijente na dobre i loše, jer nam pogrešna klasifikacija donosi i veće troškove koji tom prilikom nastaju. Ukoliko nekog klijenta koji je dobar klasifikujemo kao loš, samim tim ce taj klijent biti odbijen i time gubimo profit koji smo mogli da ostvarimo da smo mu odobrili kredit. Međutim, mnogo je veca greška da klijenta koji je loš klasifikujemo kao dobrog i odobrimo mu kredit, jer time nastaju odredeni gubici kada klijent više ne bude u mogućnosti da otplacuje kredit. Zato nam je vrlo bitno da ove greške svedemo na minimum. Međutim, ni to nije tako jednostavno. Smanjenje greške koja nastaje pogrešnom klasifikacijom loših klijenata dovodi do povecanja greške pogrešne klasifikacije dobroih klijenata i obrnuto. Iz tog razloga treba razmotriti kakav je odnos u gubicima koji nastaju pri pogrešnoj klasifikaciji i na osnovu tog kriterijuma odluciti koju grešku treba minimizirati.

U procesu minimizacije troškova može da nam pomogne metod koji se sastoji iz dve faze odlucivanja. Kada banka dobije zahtev za kredit, ona donosi preliminarnu odluku o tome da li će klijenta svrstati u grupu dobroih ili grupu loših klijenata. Zatim, na osnovu procene koristi i troškova donosi odluku da li će tražiti dodatne informacije o klijentu i na osnovu njih doneti konačnu odluku. Na osnovu toga koliko se stopa preciznosti povecala u drugoj fazi i koliki su troškovi, banka odlučuje da li joj se isplati da izvede drugu fazu ili da doneše odluku na osnovu prve faze ciji su troškovi zanemarljivi i koju banka uvek izvodi da bi ocenila kreditnu sposobnost klijenta.

U današnje vreme, u vecini banaka, sve je veca upotreba modela baziranih na regresiji, i to su najčešće u primeni logit modeli. Iako deluju vrlo jednostavno i razumljivo, da bi se napravio dobar i precizan model mora mu se posvetiti puno vremena i pažnje. Vrlo je bitno dobro izabrati promenljive koje ce ući u model. Polazi se od velikog skupa promenljivih, međutim taj broj se na kraju drastično smanjuje, najčešće zbog jake korelacije medu promenljivama, zatim zbog toga što neke promenljive predstavljaju istu ili sličnu stvar ali ipak jedna od njih ima veci znacaj. Dati primeri nam daju vrlo realne i logicne rezultate. Međutim, ima nekih stvari koje nas navode na razmišljanje. Tako na primer, posmatrajući Tabelu 3 u odeljku 4.3 vidimo da vecina promenljivih koje pozitivno uticu na odobravanje kredita nisu medu onima koje smanjuju rizik da ce klijent biti loš, što nas navodi na zaključak da banka ne teži ka tome da minimizira rizik. Kada uzmemos u razmatranje i vreme preživljavanja kredita, vidimo da neke promenljive koje povecavaju verovatnoca da klijentu bude odobren kredit, u isto vreme smanjuju očekivano vreme preživljavanja kredita i povecavaju verovatnoca da ce klijent biti loš. Tako na primer, osobe sa vecim prihodom imaju veci rizik da ce postati loše od osoba sa manjim prihodom, dok je verovatnoca da im bude odobren kredit veca. Treba biti vrlo obazriv kod ovakvih

situacija i ukoliko je moguce, izbaciti ih iz modela. Takođe, pokazalo se da na odluku banke ne utice iznos kredita, što nas vodi do zaključka da banka ne teži ka tome da maksimizira profit, a samim tim ni vreme preživljavanja kredita. Ali sigurno je da banka ima neku strategiju kojom se vodi pri odobravanju kredita, bila to minimizacija ili maksimizacija nekih drugih parametara, kao što su na primer broj klijenata, profit od odredene vrste proizvoda ili nešto sличno.

U primeru iz Vijetnama, u koje je kreditno tržište još uvek u fazi razvoja, dolazimo do vrlo zanimljivih rezultata. Na primer, promenljiva «Vreme na trenutnoj adresi» koja označava broj godina koje je klijent proveo na istoj adresi, u vecini zemalja označava klijentovu stabilnost i zrelost, tako da se teži da ova promenljiva bude što veća. Međutim, u Vijetnamu je malo drugacijia situacija, tamo se verovatnoca da će klijent biti loš povećava sa brojem godina koje je proveo na istoj adresi, pošto ljudi što više zaraduju i napreduju teže da žive u što boljim uslovima, pa samim tim i cesto menjaju adresu stanovanja. U Vijetnamu, klijent je taj koji predlaže trajanje kredita, on sam procenjuje u kom roku je sposoban da otplati kredit i koliki rizik je spreman da prihvati. Ovo pokazuje da banke ne mogu da se pouzdaju u sopstvenu procenu nego se oslanjaju i na klijentovu poštenu ocenu trenutnog stanja.

Uporedivanjem genetskog programiranja i probit analize videli smo da nam model genetskog programiranja daje preciznije rezultate. To je samo znak da se razvojem kreditnog tržišta razvijaju i precizniji i pouzdaniji modeli kreditnog scoringa.

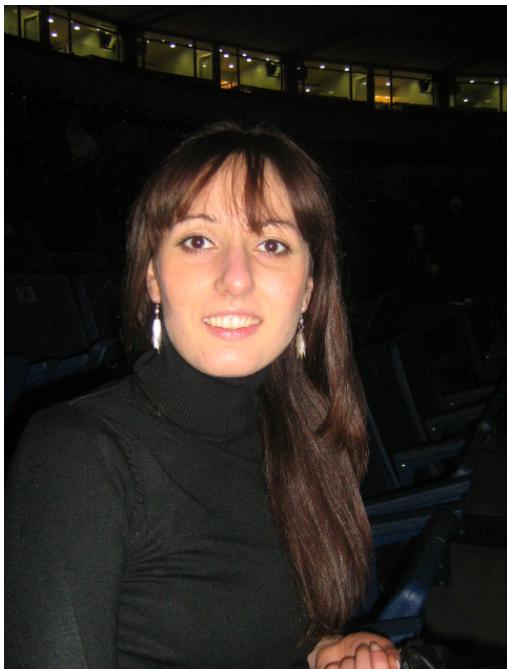
Iako se za sada kreditni scoring najčešće primenjuje u slučaju kada treba da se odredi kreditna sposobnost pojedinacnog klijenta, sve je veća primena scoringa i pri oceni rizika koje nose firme. U tom slučaju možemo da primenimo analizu obavijanja podataka (DEA), koja nam daje dosta tacne rezultate, što smo videli tako što smo rezultate uvrstili u linearnu regresiju i diskriminantnu analizu. Takođe smo videli da nam DEA skorovi služe kao veoma dobri pokazatelji pri rangiranju klijenata na osnovu njihove kreditne sposobnosti.

Opšte pravilo ne postoji, modeli se razlikuju od banke do banke, od proizvoda do proizvoda. Pošto se uzimaju istorijski podaci, a model se periodično ažurira, dešava se da je neka promenljiva u jednom periodu bila značajna, dok je u sledecem bila isključena iz modela.

6. Literatura

1. Tor Jacobson, Kasper Roszbach, “*Bank lending policy, credit scoring and value-at-risk*”, Journal of Banking & Finance 27 (2003) 615-633
2. Kasper Roszbach, “*Bank lending policy, credit scoring and the survival of loans*”, Sveriges Riksbank Working Paper Series No.154 (2003)
3. Jae H. Min, Young-Chan Lee, “*A practical approach to credit scoring*”, Expert System with Applications 35 (2008) 1762-1770
4. Yenpao Chen, Ruey-Ji Guo, Rao-Li Huang, “*Two stages credit evaluation in bank loan appraisal*”, Economic Modelling 26 (2009) 63-70
5. Thi Huyen Thanh Dinh, Stefanie Kleimeier, “*A credit scoring model for Vietnam’s retail banking market*”, International Review of Financial Analisys 16 (2007) 471-495
6. Vladimir Bugera, Hiroshi Konno, Stanislav Uryasev, “*Credit Cards Scoring with Quadratic Utility Function*”, Journal of Multi-criteria Decision Ananlysis 11 (2002) 197-211
7. Hussein A.Abdou, “*Genetic programming for credit scoring: the case of Egyptian public sector banks*”, Expert System with Applications (2009)

Biografija



Jelena Burgijašev je rođena 09.10.1984. godine u Novom Sadu od oca Mladena i majke Branislave. Završila je osnovnu školu «Jovan Popović» u Novom Sadu. Gimnaziju «Jovan Jovanović Zmaj», prirodno-matematički smer je završila 2003. godine sa odlicnim prosekom.

Studije na Prirodno-matematičkom fakultetu, Univerzitet u Novom Sadu, smer matematika-finansija je upisala 2003. godine. Diplomirala je 2007. godine sa prosekom 9,17. U školskoj 2005/2006. godini je bila stipendista Ministarstva prosvete i sporta, a u školskoj 2006/2007. godini primala stipendiju iz Fonda za mlade talente. Master studije, smer diplomirani matematičar-master matematike finansija, upisala 2007. godine.

Posle diplomiranja, 2007. godine se zaposlila u Ministarstvu ekonomije i regionalnog razvoja, kao saradnik u Sektoru za strana ulaganja i koncesije. Od avgusta 2008. godine radi u Raiffeisen banci, kao mladi saradnik u Sektoru za Upravljanje rizicima u poslovanju sa malim i srednjim preduzećima.

UNIVERZITET U NOVOM SADU
PRIRODNO-MATEMATICKI FAKULTET
KLJUCNA DOKUMENTACIJSKA INFORMACIJA

Redni broj:

RBR

Identifikacioni broj:

IBR

Tip dokumentacije: Monografska dokumentacija

TZ

Vrsta rada: Master rad

VR

Autor: Jelena Burgijašev

AU

Mentor: Prof. Dr. Zorana Lužanin

MN

Naslov rada: Razliciti pristupi kreditnom skoring sistemu

NR

Jezik publikacije: srpski

JI

Zemlja publikacije: Srbija

ZP

Uže geografsko područje: Vojvodina

UGP

Godine: 2009.

GO

Izdavac: Autorski reprint

IZ

Mesto izdavanja: Novi Sad, Departman za matematiku i informatiku, PMF, Trg Dositeja Obradovica 4

MA

Fizicki opis rada: (6, 66, 0, 0, 34, 0)

FO

Naucna oblast: Matematika

NO

Naucna disciplina: Ekonometrija

ND

Predmetna odrednica: Kreditni skoring sistem, Logit regresija, Probit i Tobit model

PO

UKD:

Cuva se: U biblioteci Departmana za matematiku i informatiku

CU

Važna napomena:

VN

Izvod: U ovom radu su prezentovani razliciti pristupi kreditnom skoring sistemu, primenom logit regresije, probit i tobit analize, problema separacije. Ovaj model se najcešće sreće u bankama kada odobravaju kredit, pri proceni kreditne sposobnosti klijenta i rizika koji taj klijent nosi

IZ

Datum prihvatanja teme od stane NN veca:

DP

Datum odbrane:

DO

Clanovi komisije:

KO

Dr. Danijela Rajter-Ciric, *redovni profesor*

Dr. Zorana Lužanin, *redovni profesor*

Dr. Andreja Tepavcevic, *redovni profesor*

UNIVERSITY OF NOVI SAD
FACULTY OF NATURAL SCIENCES AND MATHEMATICS
KEY WORD DOCUMENTATION

Accession number:

ANO

Identification number:

IDN

Type of record: Monograph type

DT

Contents code: Master thesis

CC

Author: Jelena Burgijašev

AU

Mentor: Prof. Dr. Zorana Lužanin

MN

Title: Different approaches to credit scoring

TI

Language of text: Serbian

LT

Language of abstract: English

LA

Country of publication: Serbia

CP

Locality of publication: Vojvodina

LP

Publication year: 2009.

PY

Publisher: Author's reprint

PU

Publ. place: Novi Sad, Department of Mathematics and Informatics, Faculty of Sciences and Mathematics, Trg Dositeja Obradovica 4

PP

Physical description: (6, 66, 0, 0, 34, 0)

PD

Scientific field: Mathematics

SF

Scientific discipline: Econometrics

SD

Subject – Key words: Credit scoring system, Logit regression, Probit i Tobit model

SKW

UC:

Holding data: Library of Department of Mathematics and Informatics

HD

Note:

N

Abstract: Different approaches to credit scoring are presented in this thesis, using logit regression, probit and tobit model and separation. This model is used in banks for ranking clients who applying for credit. Their aim is to determine whether an applicant has the capacity to repay his debt.

AB

Accepted by the Scientific Board on:

ASB

Defended:

DE

Thesis defened board:

DB

Dr. Danijela Rajter-Ciric, *full professor*

Dr. Zorana Lužanin, *full professor*

Dr. Andreja Tepavcevic, *full professor*