



Dragan, Strnad, Rosi

Fakulteta za logistiko
Mariborska cesta 7
3000 Celje, Slovenija

Raziskava: Statistična analiza dokumentov Incoterm

UNIVERZA V MARIBORU
FAKULTETA ZA LOGISTIKO

**Statistična analiza dokumentov Incoterm s pomočjo
statističnega modeliranja tipa regresijski LOGIT in
ugotavljanje ključnih preferenc**

Delno poročilo študije

izr. prof. dr. Dejan Dragan, doc. dr. Marjan Strnad, prof. dr. Bojan Rosi (dekan FL UM)

Fakulteta za logistiko UM

Celje, 06.02.2019



Univerza v Mariboru

Dragan, Strnad, Rosi

Fakulteta za logistiko
Mariborska cesta 7
3000 Celje, Slovenija

Raziskava: Statistična analiza dokumentov Incoterm

Dokument: Statistična analiza dokumentov Incoterm s pomočjo statističnega modeliranja tipa regresijski LOGIT in ugotavljanje ključnih preferenc

Avtorji: Dejan Dragan, Marjan Strnad, Bojan Rosi

Izdelano: Leto 2019



KAZALO

1. UVOD IN OZADJE RAZISKAVE	1
2. PREDSTAVITEV PODATKOV IN DEFINICIJA SPREMENLJIVK	2
2.1 Predstavitev podatkov za uvoz 2016.....	3
2.2 Predstavitev podatkov za izvoz 2016	8
3. MONTE CARLO SIMULACIJA IN KRITERIJI	8
4. LOGISTIČNA REGRESIJA IN LOGIT MODEL	9
4.1 Logit model.....	9
4.2. Validacija modela	15
4.2.1. Devianca in Likelihood ratio R2L	15
4.2.2. Waldov test za testiranje posameznih regresorjev	16
4.2.3. Residualni testi	17
5. IZRAČUNAN REGRESIJSKI LOGIT MODEL.....	19
5.1. Logit model za uvoz 2016.....	19
5.2. Logit model za izvoz 2016	22
6. KVALITETETA PRILEGANJA MODELA REALNIM PODATKOM.....	27
6.1. Kvaliteteta Prileganja za Logit model za uvoz 2016	27
6.2. Kvaliteteta Prileganja za Logit model za izvoz 2016.....	28
7. VIRI IN LITERATURA	31

1. UVOD IN OZADJE RAZISKAVE

Raziskava obravnava Klavzule Incoterms, ki natančno definirajo odgovornosti in dolžnosti kupcev in prodajalcev v mednarodni trgovini pri prevozu blaga na podlagi prodajne pogodbe. Prvič so se pojavile že leta 1936, njihovo tolmačenje pa je uvedla Mednarodna trgovinska zbornica v Parizu.

Klavzule vsebujejo več neodvisnih spremenljivk (regresorjev), kot npr. 'razdalja', 'Zabojnik', 'Šifra_lege_kraja', 'Vrsta posla', 'promet na meji', 'promet v notranjosti', 'Carinski postopek', 'število EUL', 'število postavk', 'Statistična vrednost blaga', 'Neto masa', ter 'vrednost na kg'. Ključni podatek Incoterma pa je njegova koda (odvisna spremenljivka y).

Glavni namen raziskave je bil statistično ugotoviti, katere neodvisne spremenljivke signifikantno vplivajo na vrednost binomsko porazdeljene naključne spremenljivke y . Analizirani so bili podatki **uvoza in izvoza za Slovenijo za leto 2016**. Kot statistični model je bil izbran **logistični regresijski model logit**.

Ker je bil v prvotni obliki y porazdeljen po Bernoullijevi porazdelitvi (0 ali 1 – pripadnost enem ali drugemu razredu), smo morali pri izbrani kombinaciji neodvisnih spremenljivk prešteti število »ugodnih« in »neugodnih« izidov spremenljivke y . Pri tem smo tudi neodvisne spremenljivke 'Statistična vrednost blaga', 'Neto masa', ter 'vrednost na kg' razvrstili v razrede, da smo sploh lahko prešteli število »ugodnih« in »neugodnih« izidov spremenljivke y .

Najzahtevnejši del raziskave je predstavljalo iskanje optimalne kombinacije regresorjev, ki statistično signifikantno vpliva na vrednost y , pri čemer morajo biti hkrati izpolnjeni vsi statistični testi in pogoji glede kvalitete prileganja izhoda modela realnim podatkom, kot npr. Waldova signifikantnost regresorjev, Devianca manjša od kritične vrednosti, McFaddenov R^2 in ostali pseudo- R^2 morajo biti dovolj veliki, ocenjeni logit se mora dobro prilegati dejanskemu (iz podatkov), ocenjene verjetnosti, da y zavzame določeno vrednost, pa dejanskim (na osnovi podatkov).

Za potrebe iskanja optimalne kombinacije regresorjev je bila uvedena Monte Carlo (MC) simulacija, ki je za vsako izbrano število regresorjev (med tri in šest – pri večjem številu so se pojavile težave s slabo pogojenostjo matrik za ocenjevanje regresorjev), preigrala 100.000 možnih scenarijev, tako za podatke uvoza kot tudi izvoza. V MC postopku so se poleg omenjenih statističnih testov izračunavali tudi številni drugi testi, kot npr. vrednost

normaliziranega korelacijskega koeficienta med ocenjenimi in dejanskimi vrednostmi logita, oz. ocenjenimi in dejanskimi verjetnostmi, mere povezane s Pearsonovim residualom, itn.

Po koncu postopka in celotne raziskave se je izkazalo, da **pri uvozu** na vrednosti y najbolj (signifikantno!) vplivajo naslednji regresorji: '**konstanta**', '**Zabojnik**', '**Šifra_lege_kraja**', '**Statistična vrednost blaga (v razredih)**', ter '**Neto masa (v razredih)**'. **Pri izvozu** pa najbolj vplivajo: '**konstanta**', '**Zabojnik**', '**Šifra_lege_kraja**', ter '**Statistična vrednost blaga (v razredih)**'.

Pomembno je še omeniti, da smo zaradi preprečevanja nepotrebne redundantnosti izločili dogodke, ki so se zgodili le enkrat. Kot se izkaže, se tako pri uvozu kot tudi izvozu, poleg vseh izpolnjenih statističnih testov, model dobro prilega danim podatkom (logit, verjetnosti, dovolj velik normaliziran korelacijski koeficient, »število zadetkov« modela).

Glavna uporabna vrednost te raziskave je, da bi na osnovi statističnega modeliranja bolje razumeli subjektivne preference, ki so v letu 2016 vodile odločevalce k izbiri določenega Incoterma, ter preučili, kateri regresorji pri tovrstnih odločitvah sploh igrajo pomembno vlogo. Gre za eno prvih tovrstnih raziskav v svetu, kar je gotovo doprinos raziskave.

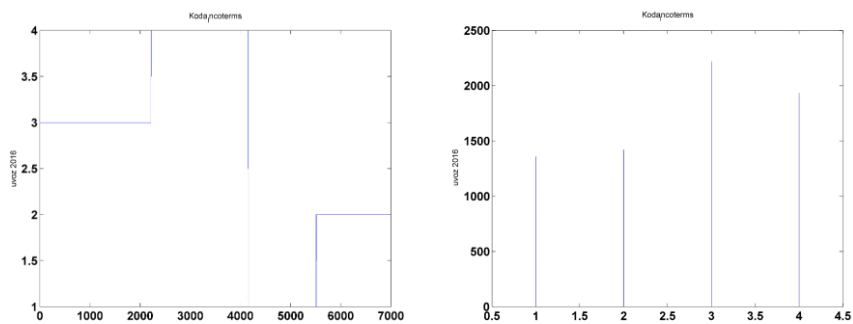
Proces modeliranja je bil implementiran v programskem orodju Matlab, pri čemer so bili uporabljeni tudi Statistics and Machine Learning toolbox, ter Econometrics Toolbox. Koncept in rezultati modeliranja, ki so bili pridobljeni s pomočjo razvite programske opreme, bodo predstavljene tudi v okviru končnega poročila študije. Poleg bodo rezultati, pridobljeni s pomočjo razvite programske opreme, predstavljeni tudi v izvirnem znanstvenem JCR ali SNIP članku.

2. PREDSTAVITEV PODATKOV IN DEFINICIJA SPREMENLJIVK

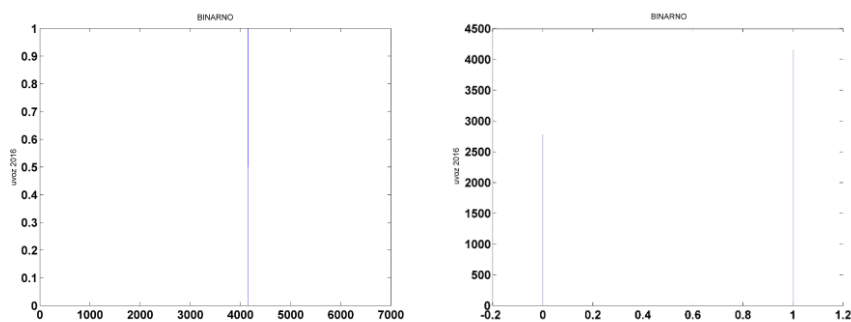
V tem poglavju bomo na kratko predstavili grafični prikaz danih podatkov Incotermov, njihove verjetnostne porazdelitve, ter ustrezno priredili simbole posameznim spremenljivkam.

2.1 Predstavitev podatkov za uvoz 2016

Izhodna spremenljivka y v originalni obliki (delitev za 4 razrede) je prikazana na sliki 1, njena binarna inačica pa na sliki 2. Na desni strani obeh slik je prikazana tudi ustrezna porazdelitev.

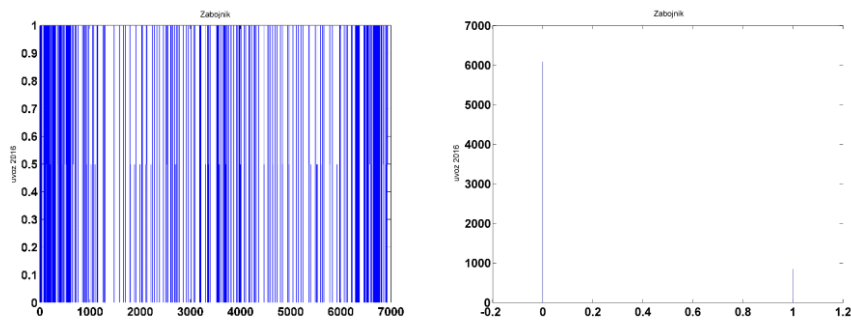


Slika 1: Izhodna spremenljivka y v originalni obliki (delitev za 4 razrede) in njena porazdelitev.



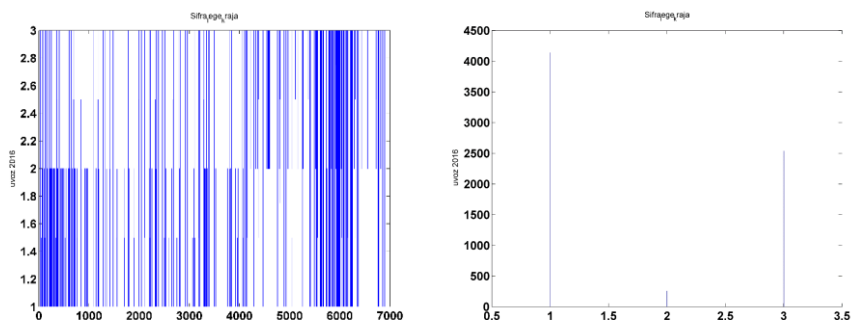
Slika 2: Izhodna spremenljivka y v binarni obliki in njena porazdelitev.

V nadaljevanju predstavimo vhodne spremenljivke. Zaradi večje preglednosti je zgornja meja oordinatne osi povsod omejena (na prikazih slik). Vhodna spremenljivka $X_1 = X_{ZABOJNIK}$ in njena porazdelitev je predstavljena na sliki 3.



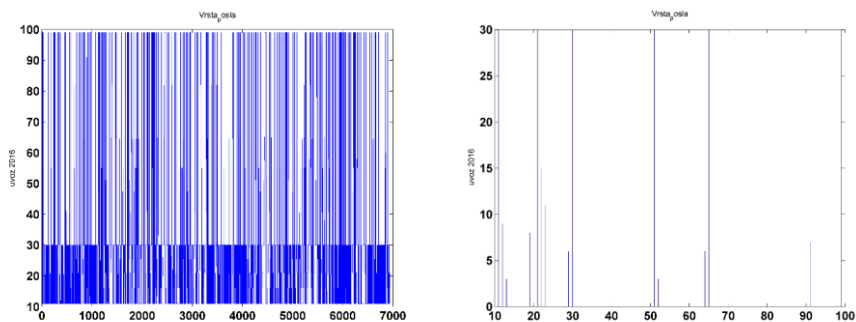
Slika 3: Vhodna spremenljivka $X_1 = X_{ZABOJNIK}$ in njena porazdelitev.

Vhodna spremenljivka $X_2 = X_{SIFRA_LEGE_KRAJA}$ in njena porazdelitev je predstavljena na sliki 4.



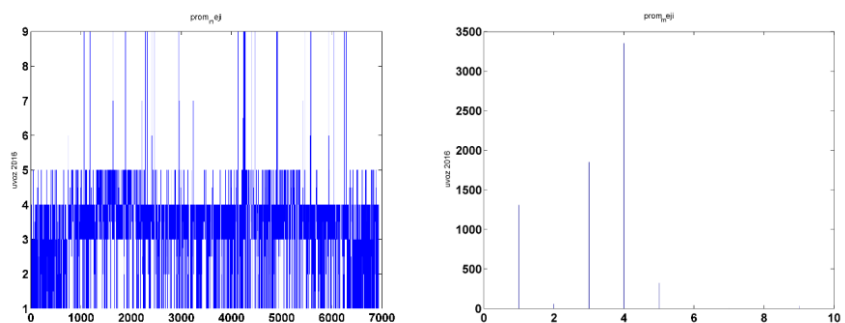
Slika 4: Vhodna spremenljivka $X_2 = X_{SIFRA_LEGE_KRAJA}$ in njena porazdelitev.

Vhodna spremenljivka $X_3 = X_{VRSTA_POSILA}$ in njena porazdelitev je predstavljena na sliki 5.



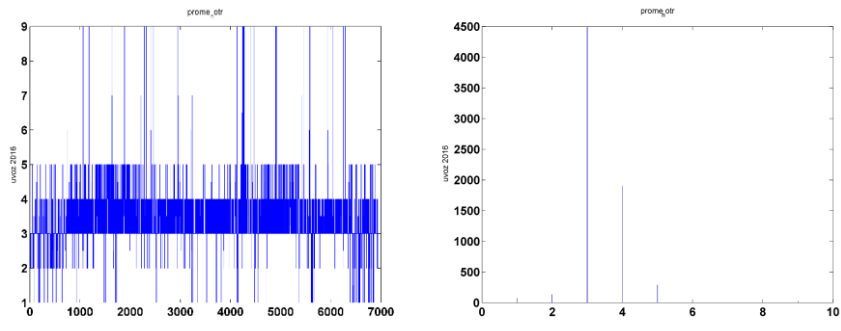
Slika 5: Vhodna spremenljivka $X_3 = X_{VRSTA_POSILA}$ in njena porazdelitev.

Vhodna spremenljivka $X_4 = X_{PROM_NA_MEJI}$ in njena porazdelitev je predstavljena na sliki 6.



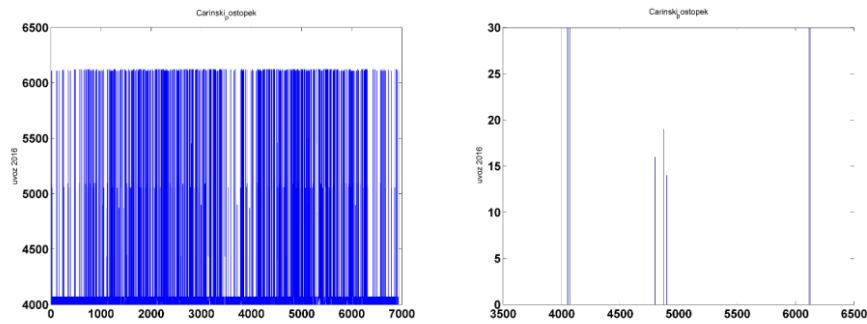
Slika 6: Vhodna spremenljivka $X_4 = X_{PROM_NA_MEJI}$ in njena porazdelitev.

Vhodna spremenljivka $X_5 = X_{PROM_V_NOTR}$ in njena porazdelitev je predstavljena na sliki 7.



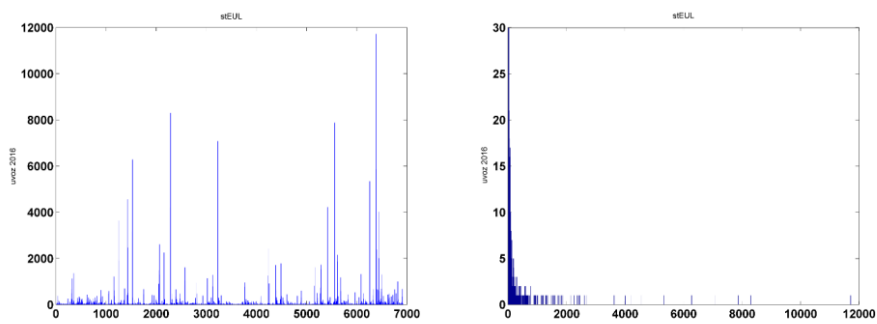
Slika 7: Vhodna spremenljivka $X_5 = X_{PROM_V_NOTR}$ in njena porazdelitev.

Vhodna spremenljivka $X_6 = X_{CARIN_POST}$ in njena porazdelitev je predstavljena na sliki 8.



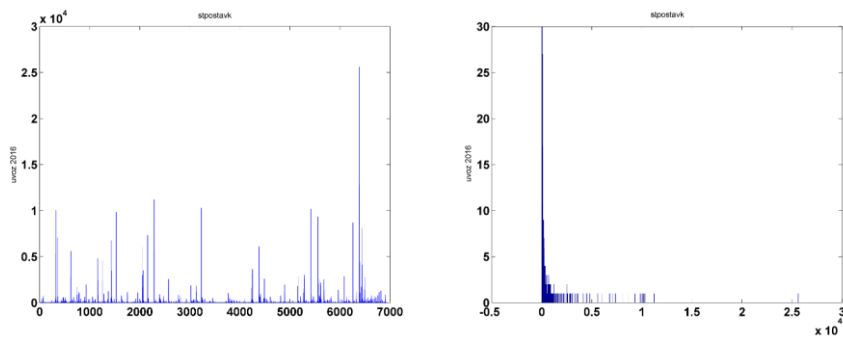
Slika 8: Vhodna spremenljivka $X_6 = X_{CARIN_POST}$ in njena porazdelitev.

Vhodna spremenljivka $X_7 = X_{ST_EUL}$ in njena porazdelitev je predstavljena na sliki 9.



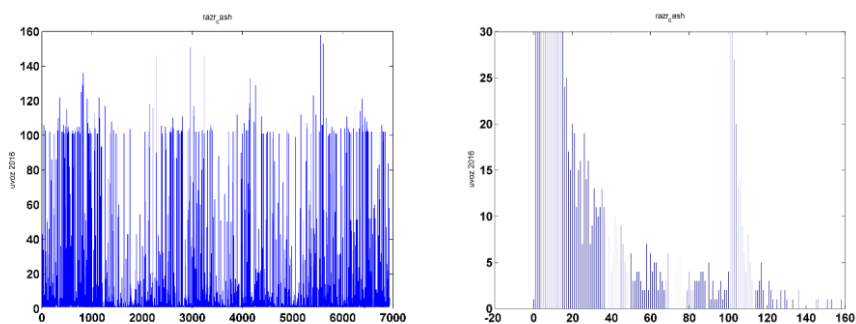
Slika 9: Vhodna spremenljivka $X_7 = X_{ST_EUL}$ in njena porazdelitev.

Vhodna spremenljivka $X_8 = X_{ST_POSTAVK}$ in njena porazdelitev je predstavljena na sliki 10.



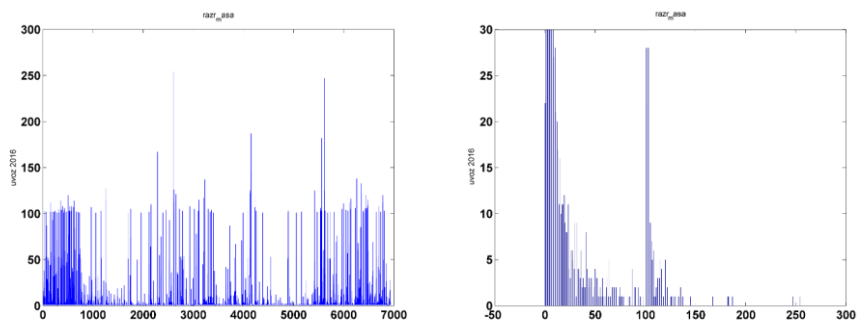
Slika 10: Vhodna spremenljivka $X_8 = X_{ST_POSTAVK}$ in njena porazdelitev.

Vhodna spremenljivka $X_9 = X_{CASH_RAZR}$ - Statistična vrednost blaga (v razredih) in njena porazdelitev je predstavljena na sliki 11. Na x osi leve slike 11 je število vzorcev, na y osi pa številka razreda denarne vrednosti. Višje denarne vrednosti imajo seveda višjo številko razreda. Interpretacija desne slike 11 pa je: na x osi je velikostni razred denarja, na ordinati pa frekvenca (število nastopov) vsakega posameznega velikostnega razreda.



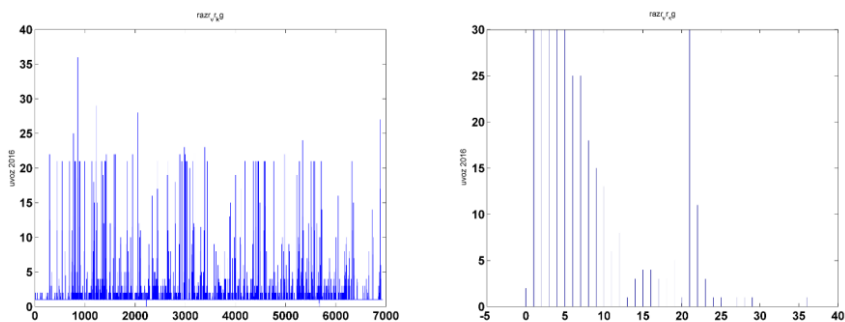
Slika 11: Vhodna spremenljivka $X_9 = X_{CASH_RAZR}$ in njena porazdelitev.

Vhodna spremenljivka $X_{10} = X_{MASA_RAZR}$ - Masa blaga (v razredih) in njena porazdelitev je predstavljena na sliki 12.



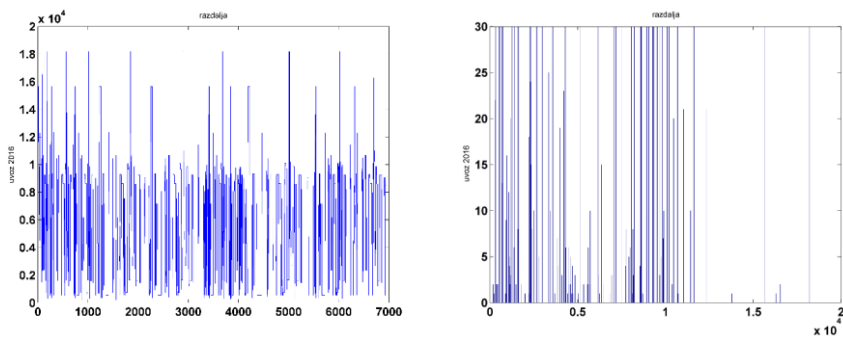
Slika 12: Vhodna spremenljivka $X_{10} = X_{MASA_RAZR}$ in njena porazdelitev.

Vhodna spremenljivka $X_{11} = X_{VRED_NA_KG_RAZR}$ - Vrednost/kg blaga (v razredih) in njena porazdelitev je predstavljena na sliki 13.



Slika 13: Vhodna spremenljivka $X_{11} = X_{VRED_NA_KG_RAZR}$ in njena porazdelitev.

Kako so vrednosti spremenljivk X_9, X_{10}, X_{11} točno porazdeljene na razrede, je razvidno iz poročila o programski opremi za dotično raziskavo, ter seveda samega programa. Slika 14 pa prikazuje še neodvisno spremenljivko razdalja in njeno porazdelitev, za katero pa se je izkazalo, da ni relevantna pri izgradnji modela.



Slika 14: Vhodna spremenljivka razdalja in njena porazdelitev.

2.2 Predstavitev podatkov za izvoz 2016

Podatkov za izvoz 2016 ne bomo predstavljali, saj so podobne oblike in značilnosti. Seveda so določene razlike, tako v velikosti vzorca, kot tudi amplitudah in frekvencah spremenljivk. Opravka pa imamo seveda s popolnoma enakimi tipi spremenljivk, kot so bile predstavljene v prejšnjem poglavju.

3. MONTE CARLO SIMULACIJA IN KRITERIJI

V splošnem lahko zapišemo naslednjo relacijo, ki povezuje binomski izhod (tip incoterma – 0 ali 1) z vhodnimi neodvisnimi spremenljivkami $X_p, p = 1, 2, \dots, P = 11$:

$$Y_i \in B(m_i, p_i(X_p)) \text{ za } i = 1, \dots, N; p = 1, \dots, P = 11 \quad (1)$$

torej se njegova vrednost porazdeljuje po binomski porazdelitvi B, kjer je N velikost vzorca. Pri tem lahko tvorimo polno matriko vseh vhodov:

$$\mathbf{X}_{\text{polni}} = [\mathbf{X}_1, \dots, \mathbf{X}_{11}] = \begin{bmatrix} X_1(1) & X_2(1) & \dots & X_{11}(1) \\ X_1(2) & X_2(2) & \dots & X_{11}(2) \\ \dots & \dots & \dots & \dots \\ X_1(N) & X_2(N) & \dots & X_{11}(N) \end{bmatrix} \quad (2)$$

Seveda se v praksi izkaže, da izhod ni statistično signifikantno odvisen od vseh vhodov, pač pa le nekaterih. Torej je potrebno poiskati takšno kombinacijo vhodov $X_l, l < 6 < p = 11$, pri kateri se bo uporabljeni statistični model, ki povezuje vhode z izhodom, najbolje obnesel. Kot se izkaže, ima lahko model največ 6 regresorjev (vhodov), sicer pride so slabe pogojenosti sistema, postopek iskanja parametrov modela pa divergira.

Za potrebe iskanja optimalne kombinacije regresorjev je bila uvedena Monte Carlo (MC) simulacija, ki je za vsako izbrano število regresorjev (med tri in šest), preigrala 100.000 možnih scenarijev, tako za podatke uvoza kot tudi izvoza. Z drugimi besedami to pomeni, da je bila v vsakem scenariju naključno izbrana drugačna sekvenca regresorjev, v paketih po tri, štiri, pet, ali šest regresorjev. Pri tem se je pri vsakem scenariju testiralo, če so hkrati izpolnjeni vsi statistični testi in pogoji glede kvalitete prilaganja izhoda modela realnim podatkom, kot npr.:

- Waldova signifikantnost regresorjev,

- Devianca manjša od kritične vrednosti,
- McFaddenov R2 in ostali pseudo-R2 morajo biti dovolj veliki,
- Ocenjeni logit se mora dobro prilegati dejanskemu (iz podatkov),
- Ocenjene verjetnosti, da y zavzame določeno vrednost, se morajo dobro prilegati dejanskim (na osnovi podatkov),
- Vrednost normaliziranega korelacijskega koeficienta med ocenjenimi in dejanskimi vrednostmi logita mora biti dovolj velika,
- Vrednost normaliziranega korelacijskega koeficienta med ocenjenimi in dejanskimi verjetnostmi mora biti dovolj velika,
- Mere povezane s Pearsonovim residualom morajo biti ustrezne,
- Ostale residualne metrike morajo biti ustrezne v smislu čim manjšega residuala (pogreška modela).

Ko je bil postopek preigravanja scenarijev $j = 1, 2, 3, \dots, 100000$ končan, se je tako za uvoz kot tudi izvoz poiskala tista optimalna reducirana matrika vhodov $\mathbf{X} = \mathbf{X}^* = \mathbf{X}(j^*)_{N \times l}$, $j^* = j_{OPTIM}$, ki je vsebovala takšno kombinacijo l regresorjev, pri kateri je statistični model imel najboljše lastnosti.

4. LOGISTIČNA REGRESIJA IN LOGIT MODEL

V tem poglavju si bomo pogledali nekaj kratkih teoretičnih osnov za logit model. Več podrobnosti pa lahko bralec zasledi v viru: Dejan Dragan, 2016: *Logistična regresija s programskim orodjem Matlab*.

4.1 Logit model

Za matriko v izrazu (2) lahko za vzorec velikosti N podamo vrednosti statističnih enot:

$$\mathbf{x}_i^T = \begin{bmatrix} x_1(i) \\ \vdots \\ x_p(i) \end{bmatrix}^T, \quad i = 1, \dots, N, \quad \mathbf{x}_i^T = [x_1(i), \dots, x_p(i)] \quad (3)$$

oziroma:

$$\begin{aligned}\mathbf{x}_1^T &= [x_1(1), \dots, x_p(1)] \\ \mathbf{x}_2^T &= [x_1(2), \dots, x_p(2)] \\ &\vdots \\ \mathbf{x}_N^T &= [x_1(N), \dots, x_p(N)]\end{aligned}\tag{4}$$

Za binarno odvisno spremenljivko Y v splošnem velja:

$$\begin{aligned}Y_1 = y_1, Y_2 = y_2, \dots, Y_N = y_N \\ y_i = \begin{cases} 1 \\ 0 \end{cases} \text{ za vsak } i = 1, \dots, N\end{aligned}\tag{5}$$

Ker pa smo pri izbranem vektorju neodvisnih spremenljivk \mathbf{x}_i^T naredili m_i realizacij, pa y_i pomeni število nastopov vrednosti 1. Če predpostavimo, da so naključne spremenljivke Y_i med seboj neodvisne binomske naključne spremenljivke, lahko zapišemo:

$$Y_i \in B(m_i, p_i) \text{ za } i = 1, \dots, N\tag{6}$$

kjer za prave verjetnosti na osnovi podatkov velja:

$$p_i = p(x_i)_{i=1, \dots, N} = \frac{e^{b_0 + b_1 x_1(i) + b_2 x_2(i) + \dots + b_p x_p(i)}}{1 + e^{b_0 + b_1 x_1(i) + b_2 x_2(i) + \dots + b_p x_p(i)}} = \frac{e^{\mathbf{x}_{ia}^T \cdot \mathbf{B}}}{1 + e^{\mathbf{x}_{ia}^T \cdot \mathbf{B}}}; \quad \mathbf{B} = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_p \end{bmatrix}\tag{7}$$

kjer:

$$\mathbf{x}_{ia}^T = \begin{bmatrix} 1 & \mathbf{x}_i^T \end{bmatrix}$$

Za logit funkcijo na osnovi podatkov velja izraz:

(8)

$$\log it(p_i) = \ln \frac{p_i}{1-p_i} = \ln \frac{e^{\mathbf{x}_{ia}^T \cdot \mathbf{B}}}{1 + e^{\mathbf{x}_{ia}^T \cdot \mathbf{B}}} = \ln \frac{e^{\mathbf{x}_{ia}^T \cdot \mathbf{B}}}{1} = \mathbf{x}_{ia}^T \cdot \mathbf{B} = b_0 + b_1 x_1(i) + \dots + b_p x_p(i)$$

Vektor regresijskih koeficientov, ki ga je treba oceniti, je:

$$\mathbf{B} = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_p \end{bmatrix} \rightarrow \text{Vektor regresijskih koeficientov} \quad (9)$$

Za ocenjene verjetnosti velja:

$$\hat{p}_i = \frac{e^{\mathbf{x}_{ia}^T \cdot \hat{\mathbf{B}}}}{1 + e^{\mathbf{x}_{ia}^T \cdot \hat{\mathbf{B}}}}; \quad i = 1, 2, \dots, N \quad (10)$$

Če tvorimo funkcijo največjega verjetja, po izpeljavi dobimo:

$$\ln L(\hat{\mathbf{B}}) = C + \sum_{i=1}^N [y_i \cdot \mathbf{x}_{ia}^T \cdot \hat{\mathbf{B}} + m_i \ln(1 - \hat{p}_i)] = \tilde{L}(\hat{\mathbf{B}}) \quad (11)$$

Funkcijo največjega verjetja moramo odvajati po parametrih $\hat{b}_0, \hat{b}_1, \dots, \hat{b}_p$. Izvesti moramo torej naslednjo operacijo:

$$\mathbf{S}(\hat{\mathbf{B}}) = \begin{bmatrix} S_0(\hat{\mathbf{B}}) \\ \vdots \\ S_p(\hat{\mathbf{B}}) \end{bmatrix} = \begin{bmatrix} \frac{\partial \tilde{L}(\hat{\mathbf{B}})}{\partial \hat{b}_0} \\ \frac{\partial \tilde{L}(\hat{\mathbf{B}})}{\partial \hat{b}_1} \\ \vdots \\ \frac{\partial \tilde{L}(\hat{\mathbf{B}})}{\partial \hat{b}_p} \end{bmatrix} = \frac{\partial}{\partial \hat{\mathbf{B}}} \tilde{L}(\hat{\mathbf{B}}) = \mathbf{0} \quad (12)$$

Po daljši izpeljavi dobimo za vektor \mathbf{S} :

$$\mathbf{S}(\hat{\mathbf{B}}) = \frac{\partial}{\partial \hat{\mathbf{B}}} \tilde{L}(\hat{\mathbf{B}}) = \sum_{i=1}^N (y_i \cdot \mathbf{x}_{ia} - m_i \cdot \mathbf{x}_{ia} \cdot \hat{p}_i) = \sum_{i=1}^N \mathbf{x}_{ia} \cdot (y_i - m_i \cdot \hat{p}_i) \quad (13)$$

Sistem enačb $\mathbf{S}(\hat{\mathbf{B}}) = \mathbf{0}$ ponavadi ni rešljiv z analitičnimi postopki. Zato moramo uporabiti postopke numerične matematike za rešitev sistema enačb (13). V ta namen bomo uporabili informacijsko matriko razsežnosti $(p+1) \cdot (p+1)$, ki ima obliko:

$$\mathbf{I}(\hat{\mathbf{B}}) = -\frac{\partial}{\partial \hat{\mathbf{B}}^T} \mathbf{S}(\hat{\mathbf{B}}) = - \begin{bmatrix} \frac{\partial S_0(\hat{\mathbf{B}})}{\partial \hat{b}_0} & \dots & \frac{\partial S_0(\hat{\mathbf{B}})}{\partial \hat{b}_p} \\ \vdots & \dots & \vdots \\ \frac{\partial S_p(\hat{\mathbf{B}})}{\partial \hat{b}_0} & \dots & \frac{\partial S_p(\hat{\mathbf{B}})}{\partial \hat{b}_p} \end{bmatrix} \quad (14)$$

Dokazati se da, da je njena vrednost enaka:

$$\mathbf{I}(\hat{\mathbf{B}}) = \underbrace{\begin{bmatrix} 1 & \dots & 1 \\ x_1(1) & \dots & x_1(N) \\ \vdots & \dots & \vdots \\ x_p(1) & \dots & x_p(N) \end{bmatrix}}_{\mathbf{X}^T} \cdot \underbrace{\begin{bmatrix} m_1 \cdot \hat{p}_1 \cdot (1 - \hat{p}_1) & \emptyset \\ \emptyset & \dots \\ \dots & m_N \cdot \hat{p}_N \cdot (1 - \hat{p}_N) \end{bmatrix}}_{\mathbf{W}} \cdot \underbrace{\begin{bmatrix} 1 & x_1(1) & \dots & x_p(1) \\ \dots & \vdots & \ddots & \vdots \\ 1 & x_1(N) & \dots & x_p(N) \end{bmatrix}}_{\mathbf{X}} \quad (15)$$

$$\mathbf{I}(\hat{\mathbf{B}}) = \mathbf{X}^T \cdot \mathbf{W} \cdot \mathbf{X}$$

kjer velja relacija za matriko \mathbf{X} :

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_{1a}^T \\ \dots \\ \mathbf{X}_{Na}^T \end{bmatrix} \quad (16)$$

Tudi vektor \mathbf{S} se da še nekoliko preoblikovati, pri čemer po krajši izpeljavi dobimo:

$$\mathbf{S}(\hat{\mathbf{B}}) = \mathbf{X}^T \cdot (\mathbf{Y} - \mathbf{E}) \quad (17)$$

kjer je \mathbf{Y} vektor izmerjenih vrednosti odvisne spremenljivke (izmerjenih frekvenc), \mathbf{E} je pa vektor ocenjenih vrednosti odvisne spremenljivke (teoretičnih frekvenc) [Jesenko].

Največjo vrednost funkcije največjega verjetja bomo iskali z Newton-Raphsonovo iteracijsko metodo, pri čemer bomo poskusili numerično rešiti sistem enačb $\mathbf{S}(\hat{\mathbf{B}}) = \mathbf{X}^T \cdot (\mathbf{Y} - \mathbf{E}) = \mathbf{0}$. Za začetne vrednosti ponavadi izberemo vektor koeficientov 0, to je:

$$\hat{\mathbf{B}}(1) = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \quad (90)$$

Iteracijski postopek pa temelji na naslednjem zaporedju približkov rešitev za vektor ocenjenih koeficientov [Jesenko]:

$$\hat{\mathbf{B}}(k+1) = \hat{\mathbf{B}}(k) + \left(\mathbf{I} \left[\hat{\mathbf{B}}(k) \right] \right)^{-1} \cdot \mathbf{S}(\hat{\mathbf{B}}(k)), \quad k = 1, 2, \dots \quad (18)$$

Iteracijski postopek izvajamo toliko časa, dokler ni dosežena prepisana toleranca ε in se ocenjeni parametri \hat{b}_p v dveh zaporednih iteracijah skorajda ne razlikujejo več med seboj:

$$\left\| \hat{\mathbf{B}}(k_{končni}) - \hat{\mathbf{B}}(k_{končni-1}) \right\| \leq \varepsilon \quad (19)$$

Končni rezultat je enak:

$$\hat{\mathbf{B}}_{OPT} = \hat{\mathbf{B}}^* = \hat{\mathbf{B}}(k_{končni}) \quad (20)$$

Ocena kovariančne matrike optimalnih ocenjenih koeficientov $\hat{\mathbf{b}}_0^*, \dots, \hat{\mathbf{b}}_p^*$ je enaka:

$$\hat{\Sigma}^* = \mathbf{I}^{-1}(\hat{\mathbf{B}}^*) = \begin{bmatrix} S_{11}^* & \cdots & S_{1,p+1}^* \\ \vdots & \ddots & \vdots \\ S_{p+1,1}^* & \cdots & S_{p+1,p+1}^* \end{bmatrix} \quad (21)$$

Kvadratni koreni diagonalnih elementov te matrike pa so standardne napake optimalnih ocen koeficientov:

$$\sigma(\hat{b}_0^*) = \sqrt{S_{11}^*}, \quad \sigma(\hat{b}_1^*) = \sqrt{S_{22}^*}, \quad \dots, \quad \sigma(\hat{b}_p^*) = \sqrt{S_{p+1,p+1}^*} \quad (22)$$

Ocenjene verjetnosti \hat{p}_i^* (točkaste ocene) in ocenjene logit vrednosti so enake:

$$\hat{p}_i^* = \frac{e^{\mathbf{x}_{ia}^T \cdot \hat{\mathbf{B}}^*}}{1 + e^{\mathbf{x}_{ia}^T \cdot \hat{\mathbf{B}}^*}} = \frac{e^{b_0^* + \hat{b}_1^* x_1(i) + \dots + \hat{b}_p^* x_p(i)}}{1 + e^{b_0^* + \hat{b}_1^* x_1(i) + \dots + \hat{b}_p^* x_p(i)}}, \quad i = 1, \dots, N$$

$$\log it(\hat{p}_i^*) = \ln \frac{\hat{p}_i^*}{1 - \hat{p}_i^*} = \ln \left(\frac{\frac{e^{\mathbf{x}_{ia}^T \cdot \hat{\mathbf{B}}^*}}{1 + e^{\mathbf{x}_{ia}^T \cdot \hat{\mathbf{B}}^*}}}{1 - \frac{e^{\mathbf{x}_{ia}^T \cdot \hat{\mathbf{B}}^*}}{1 + e^{\mathbf{x}_{ia}^T \cdot \hat{\mathbf{B}}^*}}} \right) = \ln e^{\mathbf{x}_{ia}^T \cdot \hat{\mathbf{B}}^*} = \mathbf{x}_{ia}^T \cdot \hat{\mathbf{B}}^* \quad (23)$$

Ocenjene vrednosti odvisne spremenljivke pri njih pa so enake:

$$\hat{y}_i^* = \hat{p}_i^* \cdot m_i, \quad i = 1, \dots, N \quad (24)$$

Pri logit modelu so pomembna tudi tki. ODDS razmerja za ocenjene parametre, ki izražajo elastičnost odvisne spremenljivke glede na spremembo neodvisne spremenljivke. Razmerje ODDS ratio potemtakem pomeni, za koliko se procentualno poveča logit funkcije, če katerega izmed regresorjev povečamo za eno enoto.

Denimo neko vhodno spremenljivko $x_p(i)$ povečamo na $x_p(i) + 1$. Potem bi za parameter \hat{b}_p^* veljal naslednji izraz:

$$OR_p^* = e^{\hat{b}_p^*} \quad (25)$$

Sklep bi torej bil: Če torej regresor x_p povečamo za eno enoto, se ODDS RATIO poveča za $e^{\hat{b}_p^*}$ enot, pri čemer gre $p = 0, \dots, P$. Na primer, če bi prišel $e^{\hat{b}_p^*} = 1,175$, bi to pomenilo, da se je OR_p^* povečal za 17,5%

4.2. Validacija modela

4.2.1. Devianca in Likelihood ratio R2L

Devianca polnega modela logistične regresije je enaka:

$$\begin{aligned}
 D &= -2 \left[\sum_{i=1}^N \left[y_i \cdot \ln \frac{\hat{y}_i^*}{y_i} + (m_i - y_i) \ln \frac{(m_i - \hat{y}_i^*)}{m_i - y_i} \right] \right] = \\
 &= 2 \left[\sum_{i=1}^N \left\{ y_i \cdot \ln \frac{y_i}{\hat{y}_i^*} + (m_i - y_i) \ln \frac{m_i - y_i}{m_i - \hat{y}_i^*} \right\} \right] = D_{POLNI}
 \end{aligned}
 \tag{26}$$

kjer je $\hat{y}_i^* = m_i \cdot \hat{p}_i^*$,

$$\hat{p}_i^* = \frac{e^{\mathbf{x}_{ia}^T \cdot \hat{\mathbf{B}}^*}}{1 + e^{\mathbf{x}_{ia}^T \cdot \hat{\mathbf{B}}^*}},$$

$i = 1, \dots, N$

Devianca je vedno nenegativna in je zato lahko tudi neka mera za prilagajanje modela logistične regresije podatkom. Za velike m_i in pri dobrem prilagajanju modela podatkom (ko je majhna devianca) je D približno χ^2 naključna spremenljivka z $N - p - 1$ prostostnimi stopnjami. Iz tega izhaja, da se za velike vrednosti deviance model slabo prilagodi podatkom. Devianco pa pogostokrat uporabljamo tudi za primerjavo različnih modelov logistične regresije z različnim številom parametrov.

Kot se izkaže, velja za velike m_i pri dobrem prilagajanju modela podatkom:

$$D \approx \chi^2(N - p - 1)
 \tag{27}$$

Celovit test kakovosti polnega modela pa izvedemo z naslednjo primerjavo:

H_0 : Ocenjeni model se dobro prilagodi podatkom (28)

H_1 : Ocenjeni model se slabo prilagodi podatkom

$$D = D_p > \chi_{krit}^2(\alpha, N - p - 1)$$

Če je $D_p < D_{pkrit} = \chi_{krit}^2(\alpha, N-p-1)$, ne zavrnemo ničelne hipoteze, torej se ocenjeni model dobro prilagodi podatkom. Likelihood ratio R2L dobimo, če izračunamo tudi devianco reduciranega (ničelnega) sistema D_{null} :

$$R^2_L = 1 - \frac{D_p}{D_{null}} \quad (29)$$

Na podoben način dobimo tudi ostale pseudo-R2 metrike, kot npr. McFaddenov R2, itn.

4.2.2. Waldov test za testiranje posameznih regresorjev

Za test enega samega parametra oz. testiranje parametrov posameznih regresorjev običajno uporabimo takoimenovano **Waldovo statistiko**. Za velik N in pri pravilni ničelni hipotezi $H_0 : b_p = 0$ je Waldova statistika približno standardizirana normalna naključna spremenljivka. Pri izbrani stopnji pomembnosti α in nasprotni hipotezi $H_1 : b_p \neq 0$ sprejmemo ničelno hipotezo, če velja $w_p^* \in \left(-\frac{z_{\alpha}}{2}, \frac{z_{\alpha}}{2} \right)$. Test hipoteze lahko zgradimo tudi s statistiko w_p^{*2} , ki je χ^2 naključna spremenljivka s stopnjo prostosti 1 [Jesenko].

Dokazati se da, da pri dovolj velikem N približno velja:

$$\frac{\hat{b}^*(p) - b_p}{\sigma(\hat{b}^*(p))} = \frac{\hat{b}^*(p) - b_p}{\sqrt{S_{p+1,p+1}^*}} \in Z_p^* = N(0,1), \quad (30)$$

$p=0, \dots, P$

kjer so $\sigma(\hat{b}^*(p))$ standardne napake optimalnih ocen koeficientov iz izraza (22). Postavimo naslednji hipotezi:

$$\begin{aligned} H_0 : b_p &= 0 \\ H_1 : b_p &\neq 0, \\ p &= 0, \dots, P \end{aligned} \quad (31)$$

Če ničelna hipoteza velja, sledi:

$$\frac{\hat{b}^*(p) - 0}{\sqrt{S_{p+1,p+1}^*}} = \frac{\hat{b}^*(p)}{\sqrt{S_{p+1,p+1}^*}} = Z_p^* = w_p^* \quad (32)$$

kjer je w_p^* takoimenovana Waldova statistika. Za ta test torej velja:

$$\begin{aligned} H_0 \text{ sprejmemo, \u0107e velja: } w_p^* \in \left(-\frac{z_\alpha}{2}, \frac{z_\alpha}{2} \right) &\Rightarrow \text{Potem } \hat{b}_p^* \text{ ni SIGNIFIKANTEN} \\ H_0 \text{ zavrnemo, \u0107e velja: } w_p^* \notin \left(-\frac{z_\alpha}{2}, \frac{z_\alpha}{2} \right) &\Rightarrow \hat{b}_p^* \text{ je SIGNIFIKANTEN} \end{aligned} \quad (33)$$

V praksi to pomeni, da je nek regresor statisti\u0107no signifikanten, \u0107e velja:

$$\left| \frac{\hat{b}^*(p)}{\sqrt{S_{p+1,p+1}^*}} \right| = |Z_p^*| = |w_p^*| > w_p^* \text{ krit} = \frac{z_\alpha}{2} \quad (34)$$

4.2.3. Residualni testi

Za testiranje kvalitete modela so pomembni tudi razli\u010dni tki. residualni testi, povezani z napako (pogre\u0161kom) prileganja izpeljanega modela danim podatkom. Analiziramo lahko razli\u010dne tipe pogre\u0161kov (residualov) modela in odtod izpeljemo razne metrike za merjenje kvalitete modela. V izrazih (35) so prikazane nekatere najbolj tipi\u010dne residualne metrike, ki smo jih dobili za najbol\u0161i model (pri najbol\u0161em scenariju j^* po kon\u010danem Monte Carlo postopku), kjer znak ‘*’ pomeni optimalne vrednosti, $n(j^*) \ll N$ je velikost binomskega vzorca pri zdru\u017eevanju ugodnih Bernoulijevih binarnih izidov in izklju\u010dtvi enkratnih dogodkov, $N_{SIGN_REGR}(j^*)$ pa je \u0161tevilo signifikantnih regresorjev. Poleg tega se v (35) pojavijo tudi normaliziran korelacijski koeficient CORR, kovarianca COV, ter varianca VAR. V izrazu (35) se pojavi tudi distan\u010dna metrika za simetri\u010dno Kullback-Leibler mero, ki se obi\u010dajno uporablja za merjenje oddaljenosti dveh signalov v postopku Dynamic Time Warping pri signalnem procesiranju (Matlab, Signal Processing Toolbox, 2018).

$$e(i, j^*) = \begin{cases} p(i) - \hat{p}(i, j^*) = e_1 \\ \log it(i) - \log it(i, j^*) = e_2 \\ e_{pearson}(i, j^*) = e_3 = r_{pi} = \frac{\tilde{p}_i - \hat{p}_i^*}{\sqrt{\hat{p}_i^* \cdot (1 - \hat{p}_i^*)}} = \frac{\frac{y_i}{m_i} - \hat{p}_i^*}{\sqrt{\hat{p}_i^* \cdot (1 - \hat{p}_i^*)}}, \quad i = 1, \dots, N \end{cases}$$

$$MSE(j^*) = \frac{1}{n} \cdot \sum_{i=1}^n e^2(i, j^*), \quad RMSE(j^*) = \sqrt{MSE(j^*)}, \quad MAE(j^*) = \frac{1}{n} \cdot \sum_{i=1}^n |e(i, j^*)|$$

$$MAPE_g(j^*) = \frac{1}{n} \cdot \sum_{i=1}^n \left| \frac{e(i, j^*)}{g(i)} \right| \cdot 100; \quad g(i) = \begin{cases} p(i) | e_1 \\ \log it(i) | e_2 \end{cases}, \quad \max_err(j^*) = \max |e(i, j^*)|$$

$$KL_{\log it}(j^*) = \sum_{i=1}^n \left[\log it(i) - \log it(i, j^*) \right] \cdot \left[\log[\log it(i)] - \log[\log it(i, j^*)] \right]$$

$$KL_p(j^*) = \sum_{i=1}^n \left[p(i) - \hat{p}(i, j^*) \right] \cdot \left[\log[p(i)] - \log[\hat{p}(i, j^*)] \right] \Rightarrow \text{Kullback - Leibler}$$

$$S_1(j^*) = \sum_{i=1}^n \left[|p(i) - \hat{p}(i, j^*)| > 0.1 \right]; \quad S_2(j^*) = \sum_{i=1}^n \left[|p(i) - \hat{p}(i, j^*)| < 0.1 \right]$$

$$RAZ_{DOBRI/SLABI}(j^*) = R_L^2(j^*) \cdot N_{SIGN_REGR}(j^*) \cdot n(j^*) \cdot \frac{S_2(j^*)}{S_1(j^*)}$$

$$CORR(p(i), \hat{p}(i, j^*)) = \frac{COV(p(i), \hat{p}(i, j^*))}{\sqrt{VAR(p(i))} \cdot \sqrt{VAR(\hat{p}(i, j^*))}} \tag{35}$$

$$CORR\left(\log it(i), \log it(i, j^*)\right) = \frac{COV\left(\log it(i), \log it(i, j^*)\right)}{\sqrt{VAR(\log it(i))} \cdot \sqrt{VAR\left(\log it(i, j^*)\right)}}$$

Poleg vseh metrik, podanih v izrazih (35), je izvedeno tudi štetje ugodnih zadetkov. Matlabova koda (K1) je naslednja:

% Matlab koda K1:

Ctrue=0;

Cfalse=0;

Cuncert=0;

for i = 1:length(pi)

if (pi(i) < 0.5) && (p_oc(i) < 0.5)

```

    Ctrue=Ctrue+1; % zadetki modela
elseif (pi(i) > 0.5)&&(p_oc(i) > 0.5)
    Ctrue=Ctrue+1; % zadetki modela
elseif (pi(i) == 0.5)
    Cuncert = Cuncert + 1; % uncertain
else
    Cfalse = Cfalse + 1; % zgrešitve modela
end
end
Chits = 1 - Cfalse/Ctrue; % procent uspešnosti zadetkov modela

```

5. IZRAČUNAN REGRESIJSKI LOGIT MODEL

5.1. Logit model za uvoz 2016

V raziskavi smo izvedli tudi **Belsleyev test multikolinearnosti**, ki je za dane regresorje bil ustrezen. Dekompozicija razmerij varianc je pokazala, da je maksimalni pogojni indeks (Conditional index) le za malenkost večji od od 30 ($CI_{\max} = 36 > 30$), vendar so razmerja varianc za vse pogojne indekse CI_1, \dots, CI_{\max} bila ustrezna. Tudi faktorji inflacije variance VIF so bili ustrezni ($VIF(p) \leq 7.8 < 10$, za $\forall p$). Po koncu Monte Carlo postopka in celotne raziskave se je izkazalo, da **pri uvozu** na vrednosti y najbolj (signifikantno!) vplivajo naslednji regresorji: '**konstanta**', '**Zabojn**ik', '**Šifra_lege_kraja**', '**Statistična vrednost blaga (v razredih)**', ter '**Neto masa (v razredih)**'. To pomeni, da imajo signifikanten vpliv na tip incoterma, to je binarne spremenljivke, gledane skozi prizmo šestih zaporedij ugodnih Bernouljevih poskusov kot izhodne binomske spremenljivke, naslednje vhodne spremenljivke:

- Vhodna spremenljivka $X_1 = X_{ZABOJNIK}$ (glej sliko 3),
- Vhodna spremenljivka $X_2 = X_{SIFRA_LEGE_KRAJA}$ (glej sliko 4),
- Vhodna spremenljivka $X_9 = X_{CASH_RAZR}$ - Statistična vrednost blaga (v razredih) (glej sliko 11), ter

- Vhodna spremenljivka $X_{10} = X_{MASA_RAZR}$ - Masa blaga (v razredih) (glej sliko 12).

Simbolično bi to lahko označili na naslednji način (glej tudi izraza (1), (2)):

$$y = f(X_1, X_2, X_9, X_{10}, link = logit)$$

$$\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_9, \mathbf{X}_{10}] = \begin{bmatrix} X_1(1) & X_2(1) & X_9(1) & X_{10}(1) \\ X_1(2) & X_2(2) & X_9(2) & X_{10}(2) \\ \dots & \dots & \dots & \dots \\ X_1(n(j^*)) & X_2(n(j^*)) & X_9(n(j^*)) & X_{10}(n(j^*)) \end{bmatrix} \quad (36)$$

Na osnovi izraza (23) lahko za ocenjene verjetnosti in ocenjeni logit zapišemo:

$$\hat{p}_i^* = \frac{e^{\mathbf{x}_{ia}^T \hat{\mathbf{B}}^*}}{1 + e^{\mathbf{x}_{ia}^T \hat{\mathbf{B}}^*}} = \frac{e^{\hat{b}_0^* + \hat{b}_1^* x_1(i) + \hat{b}_2^* x_2(i) + \hat{b}_3^* x_9(i) + \hat{b}_4^* x_{10}(i)}}{1 + e^{\hat{b}_0^* + \hat{b}_1^* x_1(i) + \hat{b}_2^* x_2(i) + \hat{b}_3^* x_9(i) + \hat{b}_4^* x_{10}(i)}}, \quad i = 1, \dots, n(j^*) = 134 \quad (37)$$

$$\log it(\hat{p}_i^*) = \hat{b}_0^* + \hat{b}_1^* x_1(i) + \hat{b}_2^* x_2(i) + \hat{b}_3^* x_9(i) + \hat{b}_4^* x_{10}(i)$$

Ocenjeni parametri najboljšega Monte Carlo scenarija pridejo:

$$\begin{aligned} \tilde{b}_0^* = \hat{b}_0^* &= \mathbf{5.9497}, \quad \tilde{b}_1^* = \hat{b}_1^* = \mathbf{-1.1012}, \quad \tilde{b}_2^* = \hat{b}_2^* = \mathbf{-2.0852}, \\ \tilde{b}_3^* = \hat{b}_3^* &= \mathbf{-0.25753}, \quad \tilde{b}_4^* = \hat{b}_4^* = \mathbf{-1.0053} \end{aligned} \quad (38)$$

ODDS razmerja za ocenjene parametre pridejo (glej izraz (25)):

$$\begin{aligned} OR_1^* &= e^{\hat{b}_1^*} = \mathbf{0.33246} \\ OR_2^* &= e^{\hat{b}_2^*} = \mathbf{0.12428} \\ OR_3^* &= e^{\hat{b}_3^*} = \mathbf{0.77296} \\ OR_4^* &= e^{\hat{b}_4^*} = \mathbf{0.36594} \end{aligned} \quad (39)$$

Devianca polnega modela logistične regresije pride enaka (glej (26)):

$$D_p = 2 \left[\sum_{i=1}^{134} \left\{ y_i \cdot \ln \frac{y_i}{\hat{y}_i^*} + (m_i - y_i) \ln \frac{m_i - y_i}{m_i - \hat{y}_i^*} \right\} \right] = \mathbf{125.47}$$

kjer je $\hat{y}_i^* = m_i \cdot \hat{p}_i^*$,

$$\hat{p}_i^* = \frac{e^{\hat{b}_0^* + \hat{b}_1^* x_1(i) + \hat{b}_2^* x_2(i) + \hat{b}_3^* x_9(i) + \hat{b}_4^* x_{10}(i)}}{1 + e^{\hat{b}_0^* + \hat{b}_1^* x_1(i) + \hat{b}_2^* x_2(i) + \hat{b}_3^* x_9(i) + \hat{b}_4^* x_{10}(i)}}, \quad i = 1, \dots, 134 \quad (40)$$

Ker velja: $D_P < D_{P_{krit}} = \chi_{krit}^2(\alpha, N - p - 1) = \mathbf{171.8}$ (glej izraz (28)), se očitno dobljeni logit model dobro prilagodi podatkom ($\alpha = 0.05$). Likelihood ratio R2L dobimo, če izračunamo tudi devianco reduciranega (ničelnega) sistema D_{null} (glej izraz (29)) :

$$R_L^2 = 1 - \frac{D_P}{D_{null}} = 1 - \frac{125.47}{1265.7} = 0.90087 \quad (41)$$

Torej tudi dokaj visoka vrednost Likelihood ratio R2L potrdi, da se dobljeni logit model dobro prilagodi podatkom. Za Waldovo signifikantnost regresorjev velja (glej (34)):

$$w_{p_{krit}}^* = z_{\frac{\alpha}{2}} = \mathbf{1.96}$$

$$\left| \frac{\hat{b}^*(0)}{\sqrt{S_{1,1}^*}} \right| = \left| \frac{\hat{b}_0^*}{\sqrt{S_{1,1}^*}} \right| = |Z_0^*| = |w_0^*| = \mathbf{7.4251} > \mathbf{1.96}$$

$$\left| \frac{\hat{b}^*(1)}{\sqrt{S_{2,2}^*}} \right| = \left| \frac{\hat{b}_1^*}{\sqrt{S_{2,2}^*}} \right| = |Z_1^*| = |w_1^*| = \mathbf{-2.7177} > \mathbf{1.96}$$

$$\left| \frac{\hat{b}^*(2)}{\sqrt{S_{3,3}^*}} \right| = \left| \frac{\hat{b}_2^*}{\sqrt{S_{3,3}^*}} \right| = |Z_2^*| = |w_2^*| = \mathbf{-26.596} > \mathbf{1.96}$$

$$\left| \frac{\hat{b}^*(3)}{\sqrt{S_{4,4}^*}} \right| = \left| \frac{\hat{b}_3^*}{\sqrt{S_{4,4}^*}} \right| = |Z_3^*| = |w_3^*| = \mathbf{-2.6969} > \mathbf{1.96}$$

$$\left| \frac{\hat{b}^*(4)}{\sqrt{S_{5,5}^*}} \right| = \left| \frac{\hat{b}_4^*}{\sqrt{S_{5,5}^*}} \right| = |Z_4^*| = |w_4^*| = \mathbf{-2.1711} > \mathbf{1.96} \quad (42)$$

Tudi Waldovo testiranje signifikantnosti regresorjev potrdi ustrezno veljavnost dobljenega logit modela.

Nazadnje si pogledjmo še, kako so se izračunali nekateri najbolj tipični residualni testi (glej izraz (35)):

$$e(i, j^*) = \begin{cases} p(i) - \hat{p}(i, j^*) = e_1 \\ \log it(i) - \log it(i, j^*) = e_2 \\ e_{pearson}(i, j^*) = e_3 = r_{pi} = \frac{\tilde{p}_i - \hat{p}_i^*}{\sqrt{\hat{p}_i^* \cdot (1 - \hat{p}_i^*)}} = \frac{\frac{y_i}{m_i} - \hat{p}_i^*}{\sqrt{\hat{p}_i^* \cdot (1 - \hat{p}_i^*)}}, \quad i = 1, \dots, 134 \end{cases}$$

$$KL_p(j^*) = \sum_{i=1}^{134} \left[p(i) - \hat{p}(i, j^*) \right] \cdot \left[\log[p(i)] - \log[\hat{p}(i, j^*)] \right] = \mathbf{0.24017}$$

$$S_1(j^*) = \sum_{i=1}^{134} \left[|p(i) - \hat{p}(i, j^*)| > 0.1 \right] = \mathbf{0.4403}; \quad S_2(j^*) = \sum_{i=1}^n \left[|p(i) - \hat{p}(i, j^*)| < 0.1 \right] = \mathbf{0.5597}$$

$$RAZ_{DOBRI/SLABI}(j^*) = R_L^2(j^*) \cdot N_{SIGN_REGR}(j^*) \cdot n(j^*) \cdot \frac{S_2(j^*)}{S_1(j^*)} = \mathbf{0.90} \times \mathbf{5} \times \mathbf{134} \times \frac{\mathbf{0.5597}}{\mathbf{0.4403}} = \mathbf{766.5}$$

$$CORR(p(i), \hat{p}(i, j^*)) = \frac{COV(p(i), \hat{p}(i, j^*))}{\sqrt{VAR(p(i))} \cdot \sqrt{VAR(\hat{p}(i, j^*))}} = \mathbf{0.86306}$$

$$CORR\left(\log it(i), \log it(i, j^*)\right) = \frac{COV\left(\log it(i), \log it(i, j^*)\right)}{\sqrt{VAR(\log it(i))} \cdot \sqrt{VAR\left(\log it(i, j^*)\right)}} = \mathbf{0.81326} \quad (43)$$

$$C_{true} = \mathbf{103}, \quad C_{false} = \mathbf{4}, \quad C_{uncert} = \mathbf{27}, \quad C_{hits} = \mathbf{0.96117}$$

Tudi prikazani residualni testi potrdijo dobro kvaliteto modela, kar se tiče residualov oz. pogreškov prileganja realnim podatkom. Enako velja za tiste rezultate testov, ki niso prikazani. KL metrika je zelo majhna (0.24017), normirano razmerje med dobrimi in slabimi izidi (pravilen ali napačen razred – 0 ali 1) napovedi modela je visoko (766.5), oba korelacijska koeficienta zadovoljivo visoka (0.86306 oz. 0.81326), model pa pravilno zadane napoved kar 103-krat, štirikrat zgreši, 27 je pa nedoločenih situacij (ko so realne verjetnosti točno 0.5). Tako procentualno zagotovi kar 96.117% zadetkov izida razreda (0 ali 1).

5.2. Logit model za izvoz 2016

V raziskavi smo ponovno izvedli tudi **Belsleyev test multikolinearnosti**, ki je za dane regresorje tudi tu bil ustrezen. Dekompozicija razmerij varianc je pokazala, da je maksimalni

pogojni indeks (Conditional index) manjši od 30 ($CI_{\max} = 24 < 30$), razmerja varianc za vse pogojne indekse CI_1, \dots, CI_{\max} pa so bila ustrezna. Tudi faktorji inflacije variance VIF so bili ustrezni ($VIF(p) \leq 8.7 < 10$, za $\forall p$). Po koncu Monte Carlo postopka in celotne raziskave se je izkazalo, da **pri izvozu** na vrednosti y najbolj (signifikantno!) najbolj vplivajo naslednji regresorji: 'konstanta', 'Zabojnik', 'Šifra_lege_kraja', ter 'Statistična vrednost blaga (v razredih)'. To pomeni, da imajo signifikanten vpliv na tip incoterma, to je binarne spremenljivke, gledane skozi prizmo šestih zaporedij ugodnih Bernoulijevih poskusov kot izhodne binomske spremenljivke, naslednje vhodne spremenljivke:

- Vhodna spremenljivka $X_1 = X_{ZABOJNIK}$ (glej sliko 3),
- Vhodna spremenljivka $X_2 = X_{SIFRA_LEGE_KRAJA}$ (glej sliko 4),
- Vhodna spremenljivka $X_9 = X_{CASH_RAZR}$ - Statistična vrednost blaga (v razredih) (glej sliko 11).

Kot vidimo, dobimo dokaj podobno situacijo kot pri uvozu, le da smo tam imeli še regresor za maso.

Simbolično bi funkcijsko odvisnost izhoda od vhodov lahko označili na naslednji način (glej tudi izraza (1), (2)):

$$y = f(X_1, X_2, X_9, link = logit)$$

$$\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_9] = \begin{bmatrix} X_1(1) & X_2(1) & X_9(1) \\ X_1(2) & X_2(2) & X_9(2) \\ \dots & \dots & \dots \\ X_1(n(j^*)) & X_2(n(j^*)) & X_9(n(j^*)) \end{bmatrix} \quad (44)$$

Na osnovi izraza (23) lahko za ocenjene verjetnosti in ocenjeni logit zapišemo:

$$\hat{p}_i^* = \frac{e^{\mathbf{x}_{ia}^T \cdot \hat{\mathbf{B}}^*}}{1 + e^{\mathbf{x}_{ia}^T \cdot \hat{\mathbf{B}}^*}} = \frac{e^{\hat{b}_0^* + \hat{b}_1^* x_1(i) + \hat{b}_2^* x_2(i) + \hat{b}_3^* x_9(i)}}{1 + e^{\hat{b}_0^* + \hat{b}_1^* x_1(i) + \hat{b}_2^* x_2(i) + \hat{b}_3^* x_9(i)}}, \quad i = 1, \dots, n(j^*) = 237 \quad (45)$$

$$\log it(\hat{p}_i^*) = \hat{b}_0^* + \hat{b}_1^* x_1(i) + \hat{b}_2^* x_2(i) + \hat{b}_3^* x_9(i)$$

Ocenjeni parametri najboljšega Monte Carlo scenarija pridejo:

$$\begin{aligned}\tilde{b}_0^* = \hat{b}_0^* &= \mathbf{-2.3366}, \quad \tilde{b}_1^* = \hat{b}_1^* = \mathbf{-0.57533}, \quad \tilde{b}_2^* = \hat{b}_2^* = \mathbf{1.3162}, \\ \tilde{b}_9^* = \hat{b}_3^* &= \mathbf{-0.10589}\end{aligned}\tag{46}$$

ODDS razmerja za ocenjene parametre pridejo (glej izraz (25)):

$$\begin{aligned}OR_1^* &= e^{\hat{b}_1^*} = \mathbf{0.56252} \\ OR_2^* &= e^{\hat{b}_2^*} = \mathbf{3.7294} \\ OR_3^* &= e^{\hat{b}_3^*} = \mathbf{0.89953}\end{aligned}\tag{47}$$

Devianca polnega modela logistične regresije pride enaka (glej (26)):

$$D_p = 2 \left[\sum_{i=1}^{237} \left\{ y_i \cdot \ln \frac{y_i}{\hat{y}_i^*} + (m_i - y_i) \ln \frac{m_i - y_i}{m_i - \hat{y}_i^*} \right\} \right] = \mathbf{167.54}$$

kjer je $\hat{y}_i^* = m_i \cdot \hat{p}_i^*$,

$$\hat{p}_i^* = \frac{e^{b_0^* + \hat{b}_1^* x_1(i) + \hat{b}_2^* x_2(i) + \hat{b}_9^* x_9(i)}}{1 + e^{b_0^* + \hat{b}_1^* x_1(i) + \hat{b}_2^* x_2(i) + \hat{b}_9^* x_9(i)}},$$

$i = 1, \dots, 237$

Ker velja: $D_p < D_{Pkrit} = \chi_{krit}^2(\alpha, N - p - 1) = \mathbf{288.99}$ (glej izraz (28)), se očitno dobljeni logit model tudi pri izvozu dobro prilagodi podatkom ($\alpha = 0.05$). Likelihood ratio R2L dobimo, če izračunamo tudi devianco reduciranega (ničelnega) sistema D_{null} (glej izraz (29)) :

$$R_L^2 = 1 - \frac{D_p}{D_{null}} = 1 - \frac{167.54}{746.06} = \mathbf{0.7754}\tag{49}$$

Tu je sicer vrednost manjša kot pri uvozu (glej izraz (41)), a še vedno zadovoljivo visoka. Torej tudi tukaj relativno visoka vrednost Likelihood ratio R2L potrди, da se dobljeni logit model dobro prilagodi podatkom. Za Waldovo signifikantnost regresorjev velja (glej (34)):

$$w_{p \text{ krit}}^* = z_{\frac{\alpha}{2}} = \mathbf{1.96}$$

$$\left| \frac{\hat{b}^*(0)}{\sqrt{S_{1,1}^*}} \right| = \left| \frac{\hat{b}_0^*}{\sqrt{S_{1,1}^*}} \right| = |Z_0^*| = |w_0^*| = \mathbf{|-4.8394| > 1.96}$$

$$\left| \frac{\hat{b}^*(1)}{\sqrt{S_{2,2}^*}} \right| = \left| \frac{\hat{b}_1^*}{\sqrt{S_{2,2}^*}} \right| = |Z_1^*| = |w_1^*| = \mathbf{|-2.3871| > 1.96} \quad (50)$$

$$\left| \frac{\hat{b}^*(2)}{\sqrt{S_{3,3}^*}} \right| = \left| \frac{\hat{b}_2^*}{\sqrt{S_{3,3}^*}} \right| = |Z_2^*| = |w_2^*| = \mathbf{|21.648| > 1.96}$$

$$\left| \frac{\hat{b}^*(3)}{\sqrt{S_{4,4}^*}} \right| = \left| \frac{\hat{b}_3^*}{\sqrt{S_{4,4}^*}} \right| = |Z_3^*| = |w_3^*| = \mathbf{|-2.1731| > 1.96}$$

Tudi Waldovo testiranje signifikantnosti regresorjev potrди ustrezno veljavnost dobljenega logit modela.

Nazadnje si pogledjmo še, kako so se izračunali nekateri najbolj tipični residualni testi (glej izraz (35)):

$$e(i, j^*) = \begin{cases} p(i) - \hat{p}(i, j^*) = e_1 \\ \hat{\log it}(i) - \log it(i, j^*) = e_2 \\ e_{pearson}(i, j^*) = e_3 = r_{pi} = \frac{\tilde{p}_i - \hat{p}_i^*}{\sqrt{\hat{p}_i^* \cdot (1 - \hat{p}_i^*)}} = \frac{y_i - \hat{p}_i^*}{m_i \sqrt{\hat{p}_i^* \cdot (1 - \hat{p}_i^*)}}, \quad i = 1, \dots, 237 \end{cases}$$

$$KL_p(j^*) = \sum_{i=1}^{237} \left[p(i) - \hat{p}(i, j^*) \right] \cdot \left[\log[p(i)] - \log[\hat{p}(i, j^*)] \right] = \mathbf{0.0829}$$

$$S_1(j^*) = \sum_{i=1}^{237} \left[|p(i) - \hat{p}(i, j^*)| > 0.1 \right] = \mathbf{0.48945}; \quad S_2(j^*) = \sum_{i=1}^{237} \left[|p(i) - \hat{p}(i, j^*)| < 0.1 \right] = \mathbf{0.51055}$$

$$RAZ_{DOBRI/SLABI}(j^*) = R_L^2(j^*) \cdot N_{SIGN_REGR}(j^*) \cdot n(j^*) \cdot \frac{S_2(j^*)}{S_1(j^*)} = \mathbf{0.775} \times \mathbf{4} \times \mathbf{237} \times \frac{\mathbf{0.51055}}{\mathbf{0.48945}} = \mathbf{766.3}$$

$$CORR(p(i), \hat{p}(i, j^*)) = \frac{COV(p(i), \hat{p}(i, j^*))}{\sqrt{VAR(p(i))} \cdot \sqrt{VAR(\hat{p}(i, j^*))}} = \mathbf{0.77196}$$

$$CORR\left(\hat{\log it}(i), \log it(i, j^*)\right) = \frac{COV\left(\hat{\log it}(i), \log it(i, j^*)\right)}{\sqrt{VAR(\log it(i))} \cdot \sqrt{VAR\left(\hat{\log it}(i, j^*)\right)}} = \mathbf{0.75955} \quad (51)$$

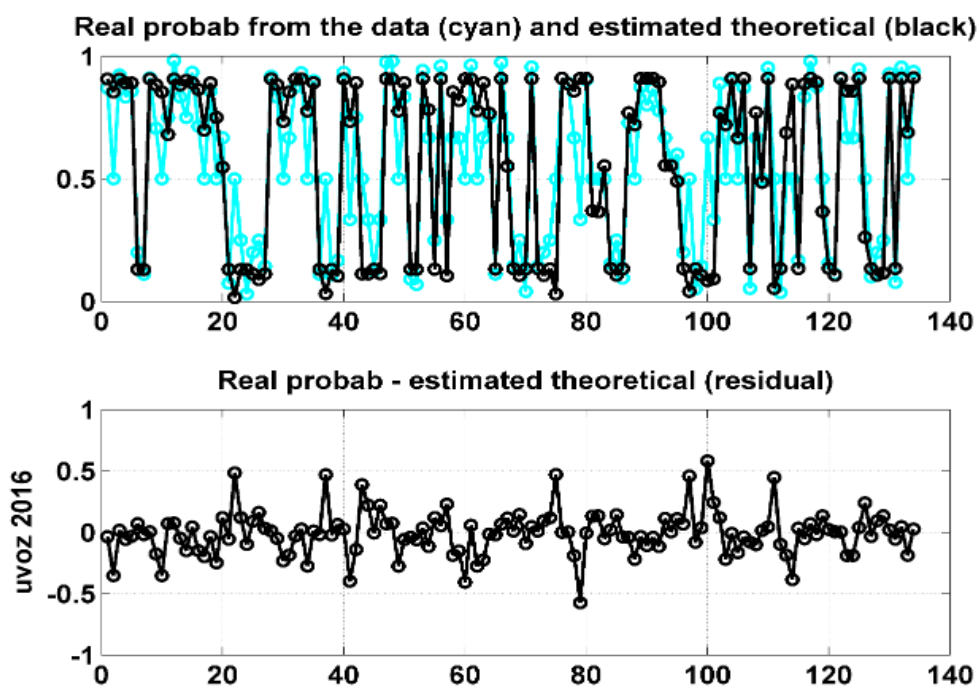
$$C_{true} = \mathbf{153}, \quad C_{false} = \mathbf{16}, \quad C_{uncert} = \mathbf{68}, \quad C_{hits} = \mathbf{0.8954}$$

Tudi prikazani residualni testi potrdijo dobro kvaliteto modela, kar se tiče residualov oz. pogreškov prileganja realnim podatkom. Enako velja za tiste rezultate tesotv, ki niso prikazani. KL metrika je zelo majhna (0.0829), normirano razmerje med dobrimi in slabimi izidi (pravilen ali napačen razred – 0 ali 1) napovedi modela je visoko (766.3), oba korelacijska koeficienta zadovoljivo visoka (0.77196 oz. 0.75955), model pa pravilno zadane napoved kar 153-krat, 16 krat zgreši, 68 je pa nedoločenih situacij (ko so realne verjetnosti točno 0.5). Tako procentualno zagotovi kar 89.54% zadetkov izida razreda (0 ali 1).

6. KVALITETA PRILEGANJA MODELA REALNIM PODATKOM

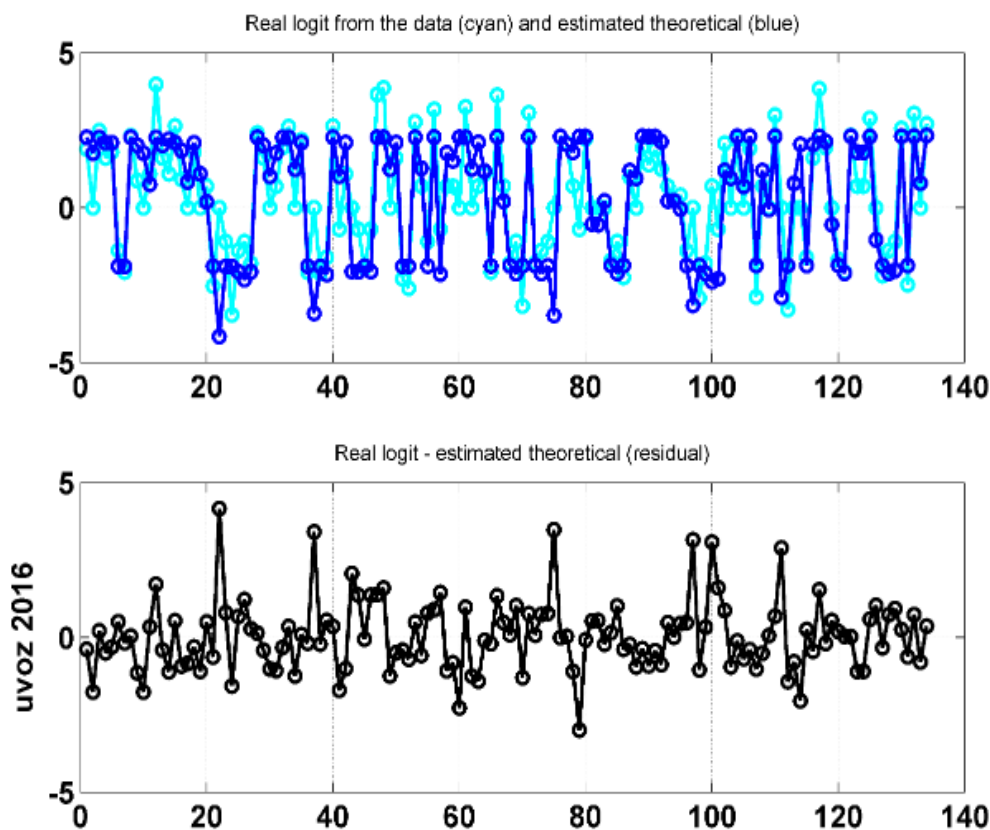
6.1. Kvaliteta Prileganja za Logit model za uvoz 2016

Slika 15 prikazuje dejansko in ocenjeno verjetnost glede na indeks danih 134 podatkov. Kot lahko vidimo, se ocenjene verjetnosti $\hat{p}(i, j^*)$ za status Incoterma v razredu 0 oz. 1, dokaj dobro prilegajo dejanskim verjetnostim $\tilde{p}_i = \frac{y_i}{m_i}$, kjer je y_i število nahajanj (število enk) v razredu 1 pri dani kombinaciji significantnih vhodnih spremenljivk $X(i) = [X_1(i), X_2(i), X_9(i), X_{10}(i)]$, $i = 1, \dots, 134$, in spreminjajočim se kombinacijam vrednosti ostalih spremenljivk $X_{ost}(i) = [X_3(i), X_4(i), \dots, X_8(i), X_{11}(i)]$, $i = 1, \dots, 134$, ki se m_i krat ponovijo. Kot je razvidno na sliki 15, je spodaj prikazan tudi residual, to je pogrešek modela $p(i) - \hat{p}(i, j^*) = e_1$. Slednji je v 103 primerih pod mejo 0.1 in da model pravo oceno razreda (glej $C_{true} = 103$ v izrazu (35)), štirikrat je napačna ocena razreda ($C_{false} = 4$), 27 je pa nedoločenih situacij (ko so realne verjetnosti točno 0.5).



Slika 15: Dejanska ($\tilde{p}_i = \frac{y_i}{m_i}$) in ocenjena verjetnost $\hat{p}(i, j^*)$ glede na indeks danih 134 podatkov, ter residual $p(i) - \hat{p}(i, j^*) = e_1$.

Slika 16 prikazuje dejanske in ocenjene vrednosti logit funkcije glede na indeks danih 134 podatkov. Kot lahko vidimo, se tudi ocenjene vrednosti logit funkcije $\log it(\hat{p}_i^*)$ dokaj dobro prilegajo dejanskim vrednostim logit funkcije $\log it\left(\tilde{p}_i = \frac{y_i}{m_i}\right)$.



Slika 16: Ocenjene vrednosti logit funkcije $\log it(\hat{p}_i^*)$ in dejanske vrednosti logit funkcije

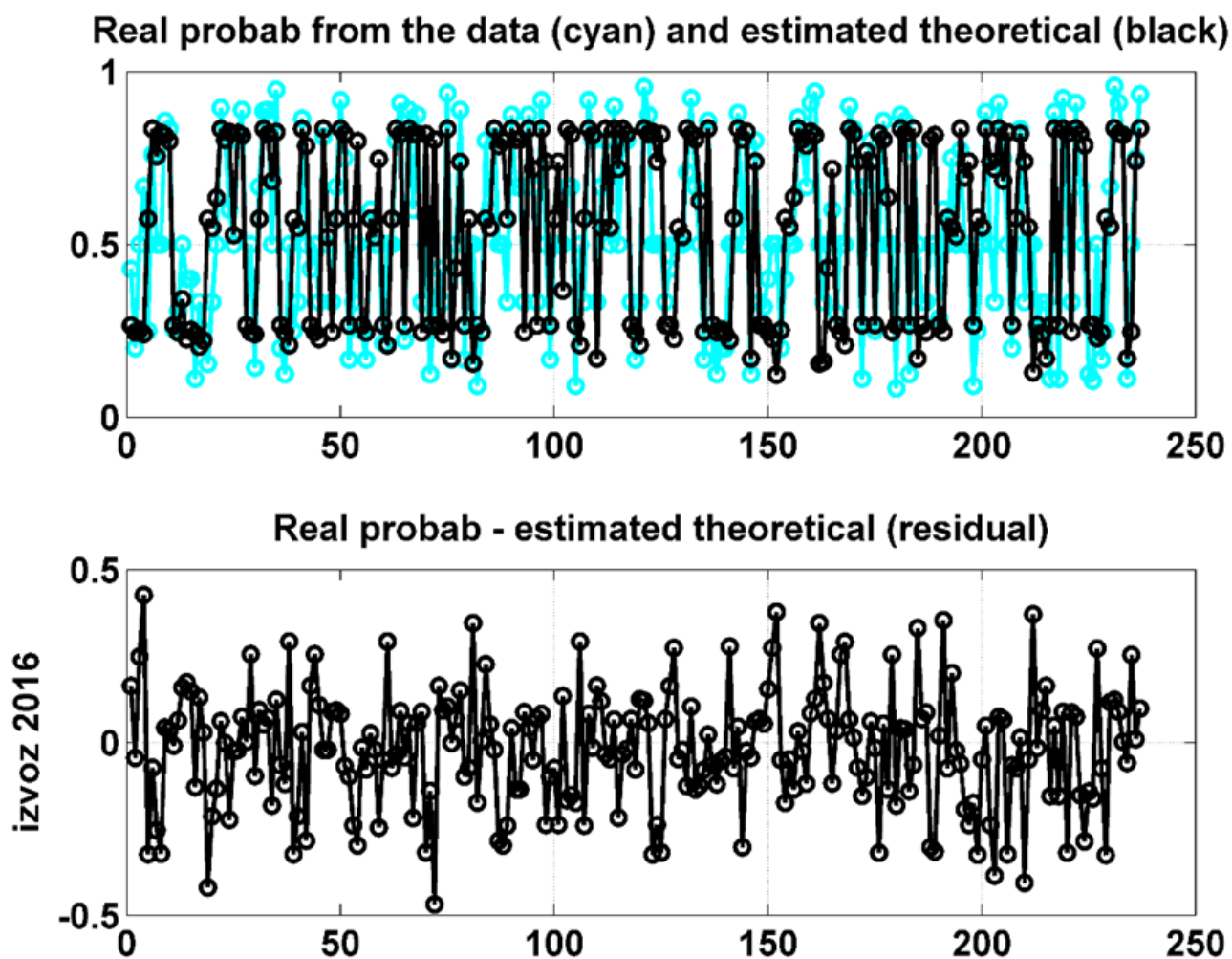
$\log it\left(\tilde{p}_i = \frac{y_i}{m_i}\right)$ glede na indeks danih 134 podatkov, ter residual $\log it(i) - \log it(i, j^*) = e_2$

(pri tem veljata enakosti: $\log it\left(\tilde{p}_i = \frac{y_i}{m_i}\right) = \log it(i)$; $\log it(\hat{p}_i^*) = \log it(i, j^*)$).

6.2. Kvaliteteta Prileganja za Logit model za izvoz 2016

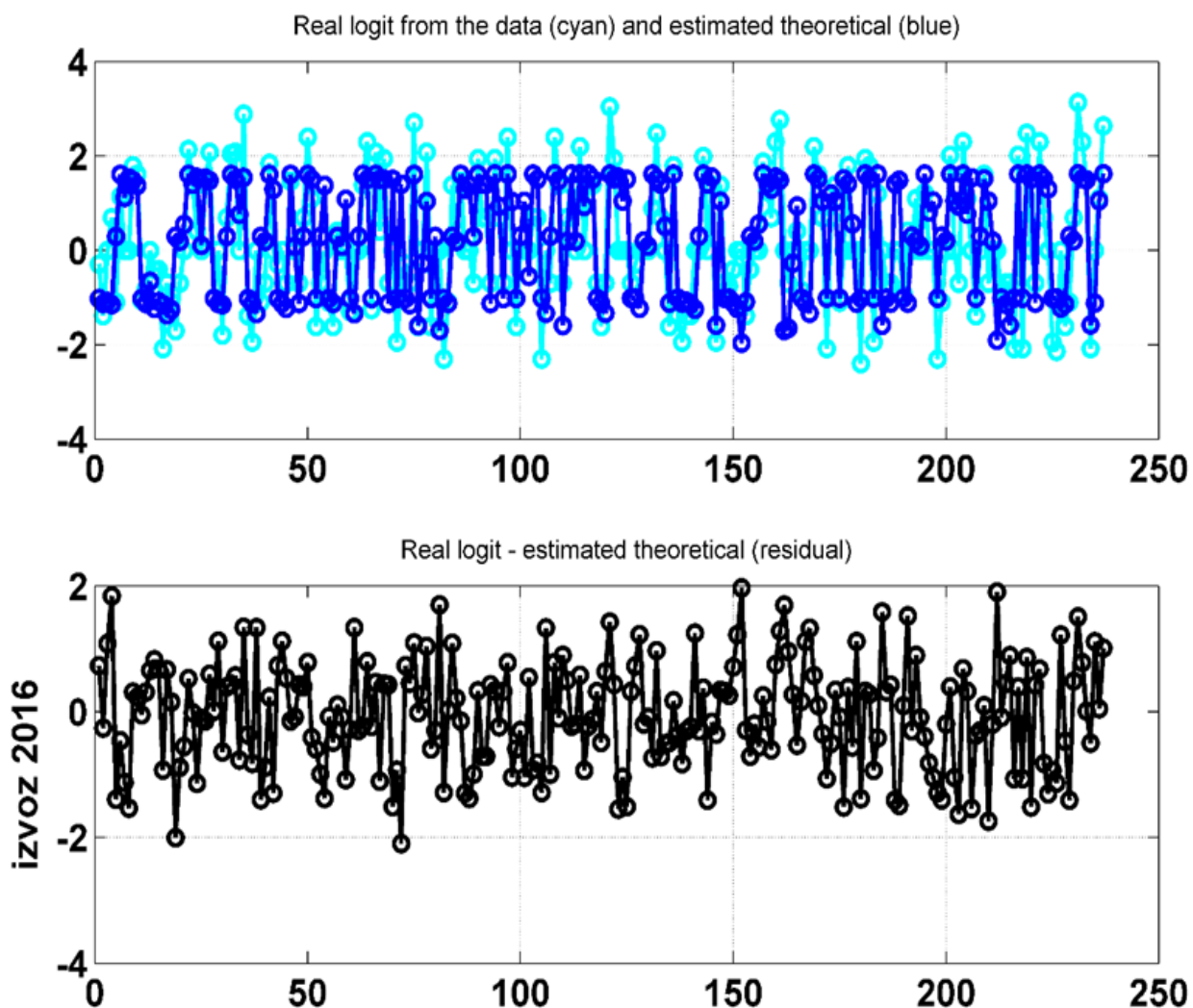
Slika 17 prikazuje dejansko in ocenjeno verjetnost glede na indeks danih 237 podatkov. Kot lahko vidimo, se ocenjene verjetnosti $\hat{p}(i, j^*)$ za status Incoterma v razredu 0 oz. 1, tudi v

primeru izvoza solidno prilegajo dejanskim verjetnostim $\tilde{p}_i = \frac{y_i}{m_i}$, kjer je y_i število nahajanj (število enk) v razredu 1 pri dani kombinaciji signifikantnih vhodnih spremenljivk $X(i) = [X_1(i), X_2(i), X_9(i)]$, $i = 1, \dots, 237$, in spreminjajočim se kombinacijam vrednosti ostalih spremenljivk $X_{ost}(i) = [X_3(i), X_4(i), \dots, X_8(i), X_{10}(i), X_{11}(i)]$, $i = 1, \dots, 237$, ki se m_i krat ponovijo. Kot je razvidno na sliki 17, je spodaj prikazan tudi residual, to je pogrešek modela $p(i) - \hat{p}(i, j^*) = e_1$. Slednji je v 153 primerih pod mejo 0.1 in da model pravo oceno razreda (glej $C_{true} = 153$ v izrazu (51)), 16 krat je napačna ocena razreda ($C_{false} = 16$), 68 je pa nedoločenih situacij (ko so realne verjetnosti točno 0.5).



Slika 17: Dejanska ($\tilde{p}_i = \frac{y_i}{m_i}$) in ocenjena verjetnost $\hat{p}(i, j^*)$ glede na indeks danih 237 podatkov, ter residual $p(i) - \hat{p}(i, j^*) = e_1$.

Slika 18 prikazuje dejanske in ocenjene vrednosti logit funkcije glede na indeks danih 237 podatkov. Kot lahko vidimo, se tudi ocenjene vrednosti logit funkcije $\log it(\hat{p}_i^*)$ dokaj dobro prilegajo dejanskim vrednostim logit funkcije $\log it\left(\tilde{p}_i = \frac{y_i}{m_i}\right)$.



Slika 18: Ocenjene vrednosti logit funkcije $\log it(\hat{p}_i^*)$ in dejanske vrednosti logit funkcije $\log it\left(\tilde{p}_i = \frac{y_i}{m_i}\right)$ glede na indeks danih 237 podatkov, ter residual $\log it(i) - \log it(i, j^*) = e_2$.

7. VIRI IN LITERATURA

Dejan Dragan, 2016: *Logistična regresija s programskim orodjem Matlab*, Fakulteta za logistiko UM.

V Celju, 06.02.2019

Poročilo izdelal: izr. prof. dr. Dejan Dragan