

INTRODUCERE

Societățile de asigurare folosesc metode actuariale pentru a determina aceste variabile utilizându-se în general o modelare de tipul **generalized linear models (GLM)**. Mai jos este o prezentare generală a acestui model actuarial.

MODELUL GLM

Dauna medie și frecvența daunelor pentru categoria autoturismelor deținute de persoane fizice atât pentru riscul de daune materiale cât și pentru cel de vătămări corporale s-au estimat utilizându-se o modelare de tipul Generalized Linear Models.

Modelul GLM reprezintă o extensie a modelului tradițional liniar (care cuprinde regresii simple sau multiple, analize ANOVA) și asociază unei variabile pe care vrem să o previzionăm (numită variabilă „răspuns” sau „variabilă dependentă”) mai multe variabile (numite „factori” sau „predictori” sau „variabile independente”) despre care avem anumite informații.

Scopul modelelor liniare este redarea relației dintre variabila răspuns Y și un număr de variabile predictive, X . Atât modelul clasic linear cât și cel generalizat privesc observațiile Y_i ca fiind realizări ale variabilei aleatoare Y .

Modelele GLM sunt compuse dintr-o gamă largă de modele care includ modelele liniare ca și caz particular. Ipotezele restrictive ale modelului liniar de Normalitate, varianta constantă și aditivitatea efectelor sunt eliminate. În schimb, distribuția variabilei răspuns este membră a familiei exponențiale, iar dispersiei ei este permis să varieze cu media distribuției. Membre ale familiei exponențiale sunt: distribuția Normală, distribuția Poisson, distribuția Binomială, distribuția Gamma, distribuția Invers Gaussiană.

Modelul tipic pentru modelarea numărului de daune sau a frecvenței daunei este Poisson multiplicativ. În afara de faptul că este universal acceptată și cel mai des utilizată în

modelarea numărului daunelor, distribuția Poisson deține și o proprietate particulară care o face și intuitiv adecvată și anume aceea că este invariantă la măsurile de timp.

Similar, un model standard pentru modelarea severității daunelor datorită formei ei generale, este Gamma multiplicativ.

Forma GLM:

$$\bar{y} = X\bar{\beta} + \bar{\varepsilon}$$

\bar{y} - vector $n \times 1$ al valorilor actuale observate, variabile aleatoare independente, normal distribuite.

X - matrice $n \times p$ a variabilelor independente (factori)

$\bar{\beta}$ - vector $p \times 1$ al parametrilor

$\bar{\varepsilon}$ - reprezintă termenul „eroare”, este un vector $n \times 1$ de variabile aleatoare independente, normal distribuite, cu media 0 și dispersia σ^2 .

Relația dintre variabila dependentă și variabilele independente este definită prin $E(y)$.

$$E(\bar{y}) = b'(\theta) = \bar{\mu}$$

O altă componentă a modelului GLM este funcția de legătură („link function”), diferentiabilă, inversabilă, care descrie modul în care media variabilei y_i , μ_i relaționează cu predictorul liniar $\eta_i = x_i' \bar{\beta}$:

$$g(\mu_i) = x_i' \bar{\beta}$$

$$E(\bar{y}) = g^{-1}(\bar{\eta}) = \bar{\mu}$$

Deci, o funcție de medie și nu media, este modelată liniar.

Sintetizând, un model GLM este structurat astfel:

$$\mu_i = E(Y_i) = g^{-1}\left(\sum_j X_{ij}\beta_j + \varepsilon_i\right)$$

$$Var(Y_i) = \frac{\phi V(\mu_i)}{\omega_i}$$

Unde: Y_i - vectorul variabilei răspuns, $g(x)$ - funcția link, inversabilă, care face legătura între răspunsul așteptat și combinația liniară a factorilor observați, X_{ij} - matricea factorilor, β_j - vectorul parametrilor modelului (care trebuie estimați), ε_i - vectorul eroare („offset”), ϕ - parametru, $V(x)$ - funcția variație, ω_i - coeficient care atribuie o credibilitate sau o pondere fiecărei observații.

Tabelul de mai jos sintetizează câteva modele tipice:

\bar{y}	Frecventa	Nr. daune	Dauna medie	Probabilitate (de ex. de reînnoire)
Link $g(x)$	$\ln(x)$	$\ln(x)$	$\ln(x)$	$\ln(x/(1-x))$
Eroare	Poisson	Poisson	Gamma	Binomial
Parametru \emptyset	1	1	Estimat	1
Variația $V(x)$	x	x	x^2	$x(1-x)$
Ponderile ω	Expunere	1	Daune	1
Offset ε	0	$\ln(\text{expunere})$	0	0

Având definit modelul prin X , $g(x)$, $V(x)$, \emptyset , ω , ε , și dat fiind setul de observații y , componentele lui $\bar{\beta}$ se determina prin maximizarea funcției de verosimilitate (sau echivalent, logaritmul funcției de verosimilitate). În esență, aceasta metoda caută să găsească parametrii care, odată aplicați modelului, produc datele observate cu probabilitatea cea mai mare. Verosimilitatea este definită ca fiind produsul probabilităților (sau a funcțiilor densitate de probabilitate în cazul distribuțiilor continue) de observare a fiecărei valori y .

În practică, când există mulți factori cum ar fi vârsta, capacitatea cilindrică, județul și durata poliței, fiecare conținând mai multe niveluri, este mai util să se parametrizeze GLM considerând, pe lângă factorii observați, un termen „intercept” (constanta GLM), care reprezintă un parametru ce se aplica tuturor observațiilor.

O analiză de tip GLM a daunelor necesită un anumit volum de date.

Cu toate că GLM sunt metode multivariate, este utilă și benefică efectuarea unor analize univariate sau bivariate ale datelor înainte de modelarea GLM.

În primul rând, analiza univariată a expunerilor și a daunelor pentru toți factorii indică dacă o variabilă conține suficiente informații pentru a fi inclusă în model (de exemplu, dacă 99% din expunerea unei variabile se află într-un singur nivel, atunci aceasta nu este potrivită pentru modelare).

În al doilea rând, presupunând că există distribuții viabile ale factorilor pe niveluri, atenție suplimentară trebuie acordată nivelurilor individuale care conțin expunere și număr de daune reduse. Dacă aceste niveluri nu sunt în cele din urmă combinate cu alte niveluri, algoritmul de verosimilitate maximă a GLM s-ar putea să nu convergă (dacă un nivel al unui

factor conține zero daune și si se modelează printr-un model multiplicativ, coeficientul teoretic corect pentru acel nivel va fi apropiat de zero, și parametrul estimat corespunzător logaritmului acelui coeficient s-ar putea sa fie atât de mare și negativ încât algoritmul numeric corespunzător verosimilității maxime nu va converge).

Pe lângă investigarea expunerilor și daunelor, analiza unifactoriala a frecvenței și daunei medii furnizează indicații preliminare a efectului produs de fiecare factor în parte.

Deoarece în analiza se utilizează factori categoriali, este necesara o atenție sporita cu privire la modul în care acești factori sunt categorizati astfel încât GLM sa beneficieze de factorii care afectează sistematic experiența și sa se excludă factorii care nu au un efect sistematic. Pentru a distinge între factorii cu efect sistematic și cei cu efect aleator (și deci neprobabil sa se repete în viitor) exista mai multe criterii, printre care: erorile standard ale estimațiilor parametrilor, teste de devianta (teste de tipul III), consistenta în timp și nu în ultimul rând intuiția și simțul realității.

Pe lângă faptul ca furnizează estimațiile de verosimilitate maxima a parametrilor, GLM furnizează și informații suplimentare ce indica certitudinea estimărilor. Un astfel de diagnostic util este reprezentat de erorile standard ale estimațiilor, acestea fiind definite ca fiind rădăcina pătrata a elementului diagonal al lui $-H^{-1}$ unde H (Hessian) este matricea derivatelor de ordin 2 a log verosimilității.

În general, se presupune ca estimațiile parametrilor sunt distribuite asimptotic Normal. În consecință, se poate efectua un simplu test statistic asupra estimațiilor parametrilor pentru a vedea daca efectul fiecărui nivel al factorului este semnificativ diferit fata de nivelul de baza al acelui factor. Pentru aceasta se utilizează de obicei testul χ^2 prin care se compara rădăcina pătrata a estimației parametrului împărțită la dispersia ei cu distribuția χ^2

În practica, deseori interpretarea grafica a estimărilor și a erorilor standard sunt mai utile în testarea adecvării unui model.

Pentru compararea masurilor de devianta a doua modele se utilizează teste de “tipul III” (χ^2 or F-tests) pentru a determina semnificația teoretica a factorilor individuali. Devianta măsoară cu cat diferă valorile previzionate de cele observate.

Pentru validarea modelului, după ce s-a stabilit cat de semnificativi sunt factorii modelați, trebuie sa se investigheze și alte elemente care include: reziduurile (testarea adecvării erorilor), levierul (identifica observațiile care au o influenta excesiva asupra modelului), transformarea Box-Cox (testarea adecvării funcției legătură).

După ce au fost definiți factorii categoriali, este utila o stabilire a corelațiilor dintre factori (de exemplu prin utilizarea statisticii de corelație Cramer's V).

Chiar dacă nu sunt utilizate direct în procesul GLM, o analiză a corelațiilor din portofoliu este utilă în interpretarea rezultatelor produse de către GLM. În particular, acestea pot explica de ce rezultatele multivariate pentru un anumit factor diferă față de rezultatele univariate și pot indica ce factori pot fi afectați de eliminarea sau includerea altor factori în GLM.