

ІНФОРМАЦІЙНІ ТЕХНОЛОГІЇ, СИСТЕМНИЙ АНАЛІЗ ТА КЕРУВАННЯ

УДК 519.246.8

П.І. Бідюк, С.В. Трухан

ПРОГНОЗУВАННЯ АКТУАРНИХ ПРОЦЕСІВ ЗА ДОПОМОГОЮ УЗАГАЛЬНЕНИХ ЛІНІЙНИХ МОДЕЛЕЙ

The method for statistical data analysis in insurance based on application of generalized linear models is studied. These models are extension of linear regression when distribution of random variable can differ from normal however belongs to the class of elliptical distributions. The model constructed can be linear or non-linear (for example, logit or probit). For parameters estimation of the models proposed the generalized least squares (GLS) or the Markov chain Monte Carlo methods are used. The main advantage of GLS is conversion of iterative algorithm which provides maximum likelihood parameter evaluations. The statistical values of losses in auto insurance are used to create the forecasting model for actuarial process selected. The model with Poisson distribution and exponential link function is acceptable for further use because it has minimum value of observational error and reliable value of risk approved by experts. Normal model with identity link function allows to find a result in one iteration with small value of observational error, but it showed "weak" predicted value of losses and impermissible risk assessment.

Keywords: actuarial processes, statistical data, generalized linear models, exponential distributions, losses forecasting.

Вступ

Реформування економічної системи України, її стрімкий перехід до ринкової економіки сформували середовище, в якому виникла нагальна потреба у створенні та використанні сучасної теорії менеджменту ризиків. Насамперед це стосується системи фінансової діяльності підприємств, яка функціонує в умовах впливу множини випадкових збурень, таких як значні коливання курсів валют, затримки платежів між підприємствами, нерациональне використання прибутків, досить часті випадки шахрайства та загальна слабкість економіки перехідного періоду [1].

Ризик – це передумова виникнення страхових відносин, яка визначає межі страхового захисту і є подією з негативними, особливо не-вигідними економічними наслідками, що можуть виникнути у майбутньому в будь-який момент у невідомих масштабах. Власне фактор ризику і необхідність покриття можливої шкоди в результаті його прояву зумовлюють потребу в страхуванні. Завдяки правильному підходу до побудови ефективної системи страхування будь-яка людська діяльність буде захищеною від випадковостей, а така суттєва складова фінансової системи сприятиме стабілізації економіки [2, 3].

На сьогодні теоретичне обґрунтування поняття фінансового ризику, розкриття природи виникнення ризиків у страхуванні, класифікація та системи заходів попередження або мінімізації негативних наслідків не тільки дістали подальший розвиток, а й стали важливим інструментом впливу на розвиток економіки. Ве-

лику роль при цьому відіграють комп'ютерні аналітичні системи, які дають можливість використовувати наявні математичні моделі й створювати нові на основі статистичних даних.

Задача оцінювання фінансових ризиків у сфері страхування є одним із найважливіших етапів фінансового аналізу для підприємств різних форм власності, оскільки для розв'язання задачі менеджменту ризику його потрібно уміти докладно проаналізувати, потім оцінити, використовуючи для цього новітні методики, математичні методи та інформаційні технології. Саме тому проблема моделювання та прогнозування актуарних процесів є надзвичайно актуальною та потребує поглибленого дослідження із використанням сучасних економіко-математичних методів, критеріїв і технологій [3, 4].

Основною метою цієї роботи є дослідження можливості практичного застосування узагальнених лінійних моделей для розв'язання задачі прогнозування величини грошових збитків як оцінки фінансового ризику на основі фактичних статистичних даних страхової компанії.

Постановка задачі

Мета роботи – проаналізувати характеристики узагальнених лінійних моделей з метою їх подальшого використання для прогнозування розвитку актуарних процесів; встановити особливості оцінювання параметрів узагальнених лінійних моделей та вибрати метод оцінювання; побудувати математичну модель для прогнозування збитків у сфері страхування автомобілів на основі фактичних статистичних да-

них; для вибору кращої моделі застосувати множину критеріїв якості.

Узагальнені лінійні моделі

Узагальнена лінійна модель (УЛМ) розглядається як розширення лінійної множинної регресії для випадку однієї залежної змінної та поняття “множинної регресійної моделі”. Головна задача множинної регресії – визначення взаємозв’язку між кількома незалежними змінними (предикторами) та залежною змінною.

Узагальнені лінійні моделі (*GLM – Generalized Linear Models*) – універсальний метод побудови регресійних моделей, що дає змогу враховувати взаємодію між факторами, вид розподілу залежної змінної та припущення про характер розподілу залежної змінної.

УЛМ складається з трьох основних компонент [5]: стохастичної, систематичної, функції зв’язку, та має вигляд

$$\mu_i = E[y_i] = g^{-1}(\sum_j X_{ij} \beta_j + \xi_i), \quad (1)$$

$$\text{Var}[y_i] = \frac{\phi V(\mu_i)}{\omega_i},$$

де y_i – вектор значень залежної змінної; $g(x)$ – функція зв’язку; X_{ij} – матриця, утворена за значеннями факторів; β_j – вектор оцінених факторів моделі; ξ_i – вектор “параметрів зсуву” (або залишки); ϕ – параметр масштабування функції $V(x)$; ω_i – апіорні вагові коефіцієнти ступеня довіри.

Тобто узагальнені лінійні моделі характеризуються такими елементами: законом розподілу залежної змінної Y ; характеристиками та параметрами функції зв’язку $g(\cdot)$; характеристиками лінійного предиктора $\eta = X\beta$.

У процесі побудови УЛМ доречно зробити такі припущення:

– припущення стосовно випадковості: всі компоненти реакції Y – незалежні, їх розподіл належить до сімейства експоненціальних розподілів;

– припущення стосовно систематичності: p предикторів об’єднуються в один “лінійний предиктор” η ;

– наявність функції зв’язку: взаємозалежність між припущенням випадковості та систематичності виражається функцією зв’язку, яка

Таблиця 1. Закони розподілу і їх параметри

Елементи моделі	Нормальний розподіл	Розподіл Пуассона	Біноміальний розподіл	Гамма-розподіл	Обернений гауссів розподіл
Позначення розподілу	$N(\mu, \sigma^2)$	$P(\mu)$	$\frac{B(n, \pi)}{n}$	$G(\mu, \nu)$	$IG(\mu, \sigma^2)$
Дисперсійний параметр ϕ	$\phi = \sigma^2$	1	$\frac{1}{n}$	$\phi = \nu^{-1}$	$\phi = \sigma^2$
Функція кумулянти, $b(\theta)$	$\frac{\theta^2}{2}$	$\exp(\theta)$	$\log(1 + e^\theta)$	$-\log(-\theta)$	$-(-2\theta)^{1/2}$
$c(y, \phi)$	$-\frac{1}{2} \left\{ \frac{y^2}{\sigma^2} + \log(2\pi\sigma^2) \right\}$	$-\log(y!)$	$\log \frac{n}{y}$	$\nu \log(\nu y) - \log(y) - \log(\Gamma(\nu))$	$-\frac{1}{2} \left\{ \log(2\pi\phi y^3) + \frac{1}{\phi y} \right\}$
$\mu(\theta) = E(Y, \theta)$	θ	$\exp(\theta)$	$\frac{e^\theta}{1 + e^\theta}$	$-\frac{1}{\theta}$	$(-2\theta)^{1/2}$
Канонічний зв’язок $\theta(\mu)$	identity	log	logit	reciprocal	$\frac{1}{\mu^2}$
Дисперсійна функція $V(\mu)$	1	μ	$\mu(1 - \mu)$	μ^2	μ^3
$\text{Var}(\mu)$	σ^2	$n\mu(1 - \mu)$	μ	$\frac{\mu^2}{\nu}$	$\frac{\sigma^2}{\mu^3}$

є диференційованою та монотонною, такою, що:

$$E[y] = \mu = g^{-1}(\eta).$$

Множина законів і параметри розподілів залежної змінної наведені в табл. 1.

Функція зв'язку зв'язує лінійний предиктор η із значенням оцінки μ величини Y . У класичній лінійній моделі середнє значення та лінійний предиктор є ідентичними і зв'язок ідентичності правдоподібний до вибору η і μ – довільним чином, але із множини дійсних чисел [5].

Узагальнені лінійні моделі відрізняються від загальної лінійної моделі, окремими випадками якої є множинна регресія, двома моментами:

– розподіл залежної змінної чи змінної реакції може бути негауссовим та не обов'язково неперервним, наприклад біноміальним;

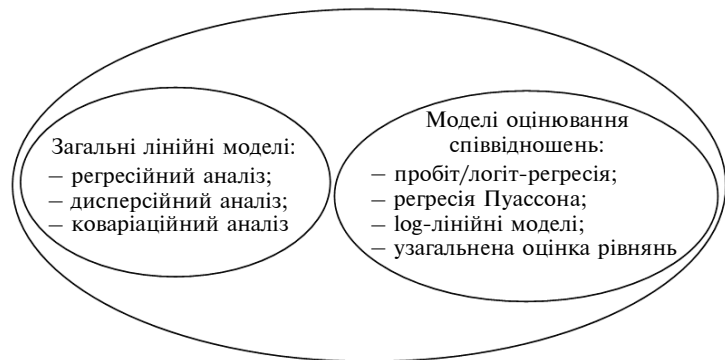
– прогнози значення залежної змінної отримують як лінійну комбінацію предикторів, які “пов'язані” із залежною змінною через функцію зв'язку.

Види УЛМ залежно від закону розподілу залежної змінної та вигляду функції зв'язку наведено в табл. 2.

Таблиця 2. Види УЛМ

Модель	Функція зв'язку	Розподіл залежної змінної
Узагальнена лінійна модель	$g(\mu) = \mu$	Нормальний розподіл
log-лінійна модель	$g(\mu) = \ln(\mu)$	Розподіл Пуассона
Логістична модель	$g(\mu) = \ln\left(\frac{\mu}{1-\mu}\right)$	Біноміальний розподіл
Пробіт-аналіз	$g(\mu) = \Phi^{-1}\mu$	Біноміальний розподіл
Аналіз “виживання”	$g(\mu) = \mu^{-1}$	Гамма-розподіл, експоненціальний розподіл

Таким чином, УЛМ – це узагальнений клас математичних моделей, які включають лінійну регресію, дисперсійний та коваріаційний аналіз, log-лінійні моделі для аналізу випадкових таблиць, пробіт/логіт моделі, регресію Пуассона та деякі інші. На рисунку графічно



Загальна структура УЛМ

зображено структуру узагальнених лінійних моделей.

Кожний розподіл має особливу функцію зв'язку, для якої існує обґрунтована статистика рівності у вимірі лінійного предиктора $\eta = \sum x_j \beta_j$. Цей канонічний зв'язок виникає у випадку, коли $\theta = \eta$, де θ – канонічний параметр, який визначено при введенні функції правдоподібності. Вигляд кожного конкретного розподілу наведено у табл. 2. Узагальнений вигляд розподілу такий:

$$f(y, \theta, \varphi) = \exp\left[\frac{y\theta - b(\theta)}{a(\varphi)} + c(y, \varphi)\right], \quad (3)$$

де a , b , c – функції, які відповідають певному закону розподілу; y – залежна змінна; θ – канонічний параметр або функція деякого параметра певного розподілу; φ – дисперсійний параметр. Функції $b(\cdot)$ надається особливе значення в узагальнених лінійних моделях, оскільки вона описує відношення між середнім значенням та дисперсією у розподілі. Якщо φ відоме, то це експоненціальна модель з канонічним параметром θ . Крім того, експоненціальний розподіл може бути двопараметричним, якщо φ – невідоме. Щільність нормального розподілу можна подати в узагальненому вигляді таким чином:

$$f(y; \mu) = \exp\left[y \frac{\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) - \frac{y^2}{2\sigma^2}\right].$$

Тобто нормальний розподіл належить до сімейства експоненціальних з натуральним параметром $b(\theta) = \mu/\sigma^2$.

Таким чином, канонічний зв'язок для множини експоненціальних розподілів має вигляд, показаний у табл. 3.

Таблиця 3. Канонічні зв'язки розподілів

Розподіл	Канонічний зв'язок
Нормальний розподіл	$\eta = \mu$
Розподіл Пуассона	$\eta = \log(\mu)$
Біноміальний розподіл	$n = \log\left\{\frac{\pi}{1-\pi}\right\}$
Гамма-розподіл	$\eta = \mu^{-1}$
Обернений гауссів розподіл	$\eta = \mu^{-2}$

Для канонічних зв'язків обґрунтованою статистикою є вектор $\mathbf{X}^T \mathbf{y}$:

$$\sum_i x_{ij} y_i, j = 1, \dots, p. \quad (4)$$

Слід зазначити, що канонічний зв'язок приводить до набуття моделлю бажаних властивостей. Так, при використанні вибірок малих розмірів немає ніякої апріорної причини виникнення сукупних систематичних ефектів у моделі без наявності такого зв'язку. В результаті попередньо виконаного аналізу статистичних критеріїв та методів оцінювання якості/адекватності моделі було вибрано функцію правдоподібності та інформаційний критерій Акайке [6].

Особливості оцінювання параметрів УЛМ

Для оцінювання параметрів УЛМ доцільно використовувати узагальнений зважений метод найменших квадратів (УЗМНК) або метод Монте-Карло для марковських ланцюгів (МКМЛ). Застосування УЗМНК дає можливість отримати оцінки параметрів за умови максимізації функції правдоподібності [7], тобто оцінка набуває найбільш вірогідного значення у просторі допустимих значень. МКМЛ характеризується вищим ступенем узагальненості застосування, але обчислювальні витрати на його реалізацію можуть бути значно вищими, ніж для УЗМНК.

У багатьох випадках задача обчислення оцінок параметрів моделей зводиться до пошуку моди. Для оцінювання нелінійних моделей, як правило, застосовують ітераційні оптимізаційні методи. За функціонал якості оцінок можна вибрати, наприклад, такий:

$$J = \frac{\partial}{\partial \theta} \left[\frac{\mathbf{y} \theta - b(\theta)}{a(\psi)} + c(\mathbf{y}, \psi) \right] = \frac{1}{a(\psi)} \left[\mathbf{y} - \frac{\partial}{\partial \theta} b(\theta) \right],$$

де \mathbf{y} – вектор вимірів основної змінної; θ – вектор параметрів; $a(\psi)$ – масштабний множник; $b(\theta)$ – нормуюча константа або “кумулянтна функція”. У такому випадку основне рівняння оцінювання параметрів має вигляд [7]

$$\theta_{k+1} = \theta_k - \frac{\partial}{\partial \theta} l(\theta_k | \mathbf{y}) \left(\frac{\partial^2}{\partial \theta \partial \theta^T} l(\theta_k | \mathbf{y}) \right)^{-1},$$

де θ_k – оцінка параметрів на k -й ітерації алгоритму оцінювання; $l(\cdot)$ – функція правдоподібності.

Ітераційний зважений метод найменших квадратів реалізується за таким алгоритмом:

1) задати початкові значення ваговим коефіцієнтам $v^{(i)}$, наприклад так: $1/v_0^{(i)} = \omega_{ii}$; і побудувати діагональну матрицю $\Omega_0 = \text{diag}[1/v_0^{(1)} \dots 1/v_0^{(p)}]$, де $p = \dim[\theta]$;

2) обчислити оцінку вектора параметрів за виразом для УЗМНК: $\hat{\theta}_k = (\mathbf{X}^T \Omega_k \mathbf{X})^{-1} \mathbf{X}^T \Omega_k \mathbf{y}$;

3) модифікувати значення вагових коефіцієнтів: $1/v_{k+1}^{(i)} = \text{var}(\mu_i)$, де μ_i – i -та компонента функції середнього $\mu = g^{-1}(\mathbf{X}\theta)$;

4) повторювати кроки 2 і 3 до настання збіжності, тобто до виконання умови $(\mathbf{X}\hat{\theta}_k - \mathbf{X}\hat{\theta}_{k+1}) \rightarrow 0$.

При використанні експоненціального сімейства розподілів і виконанні відомих загальних умов до даних наведена процедура забезпечує пошук моди функції правдоподібності, тобто обчислення оцінок параметрів за умови максимізації правдоподібності.

Оцінювання якості узагальнених лінійних моделей

Для аналізу якості моделей і встановлення найкращої моделі для розв'язання певної задачі використовують декілька критеріїв для оцінювання адекватності моделей [6, 7]: загальну точність моделі; помилки I та II-го роду; ROC-криву та індекс GINI.

Загальна точність моделі (CA – common assuagacy) визначається так:

$$CA = \frac{\text{Correct Forecast}}{N}, \quad (5)$$

де Correct Forecast – кількість вірно спрогнозованих випадків, а N – загальна кількість випадків. Загальна точність моделі є дещо суб'єктивною оцінкою, оскільки вона залежить від частки дефолтів у моделі, а також від порога відсікання. Для різних значень порога точність моделі також буде набувати різних значень.

ROC-крива (Receiver Operation Characteristic – робоча характеристика приймача) показує залежність кількості вірно класифікованих позитивних прикладів від кількості невірно класифікованих негативних прикладів. Перші називають істинно позитивними, а інші – негативними множинами. При цьому припускається, що у класифікатора існує певний параметр, варіюючи який, можна отримати певне розбиття на класи. Цей параметр називають порогом або точкою відсікання (cut-off), залежно від його значення будуть отримані різні величини помилок *I-go* та *II-go* роду.

Для аналізу якостей моделі найчастіше використовують такі відносні показники у процентах:

– частка істинно позитивних прикладів (True Positives Rate):

$$TPR = \frac{TP}{TP + FN};$$

– частка хибно позитивних прикладів (False Positives Rate):

$$FPR = \frac{FP}{TN + FP}.$$

Зазвичай для аналізу якості моделей використовують ще дві характеристики: чутливість і специфічність.

Чутливість моделі – це частка істинно позитивних випадків, тобто має місце формула:

$$Se = TPR = \frac{TP}{TP + FN}.$$

Специфічність моделі – це частка істинно негативних випадків, які були вірно класифіковані моделлю:

$$Sp = \frac{TN}{TN + FP}.$$

Зрозуміло, що має місце таке перетворення:

$$Sp = \frac{TN - FP + FP}{TN + FP} = 1 - \frac{FP}{TN + FP} = 1 - FPR.$$

Модель з високою чутливістю надає істинний результат за наявності позитивних випадків (виявляє позитивні приклади).

Навпаки, модель із високою специфічністю найчастіше дає істинний результат за наявності негативних випадків (виявляє негативні приклади). Для побудови графіка ROC-кривої по осі Y відкладаються значення чутливості Se , а по осі X – частка хибно позитивних випадків FPR або $(1 - Sp)$.

Графік ідеального класифікатора ROC-кривої проходить через верхній лівий кут, де частка істинно позитивних випадків становить 1 (ідеальна чутливість), а частка хибно позитивних прикладів дорівнює нулю. Тому чим ближче крива до верхнього лівого кута, тим кращою є здатність моделі передбачувати. Діагональна лінія відповідає класифікатору, який не здатний розпізнати ці два класи.

Оскільки візуальне порівняння ROC-кривих не завжди дає змогу визначити ефективнішу модель, застосовують оцінку площі під кривою AUC (Area Under Curve) обчислюється, наприклад, за методом трапецій:

$$AUC = \int f(x)dx = \sum_i \left[\frac{X_{i+1} + X_i}{2} \right] \cdot (Y_{i+1} - Y_i).$$

Ще одним альтернативним показником оцінювання якості моделі є індекс GINI.

Індекс GINI – це площа області між діагоналлю і кривою Лоренца, поділена на площу усієї області під діагоналлю. Індекс GINI широко використовується для аналізу роздільної здатності системи оцінювання при управлінні кредитними ризиками, тобто оцінки здатності моделі розділяти клієнтів на схильних та нехильних до дефолту. Якщо модель здатна оцінити клієнтів за ймовірністю дефолту, то більшість клієнтів, схильних до дефолту, мають отримати більшу ймовірність дефолту.

Прогнозування актуарного процесу за допомогою УЛМ

За актуарний процес, який вибрано для прикладу застосування УЛМ, взято задачу прогнозування величини грошових збитків у сфері страхування. Експериментальні дані: одна залежна змінна – “Збитки”, тобто розмір випла-

ченої страховки серед автомобілів трьох брендів (ВАЗ, Mitsubishi, Toyota); регіон продажу полісу (Київ, АР Крим, Одеса); рік випуску автомобіля (починаючи з 2006 р.). Загальний розмір вибірки – 9546 значень. Результати побудови УЛМ, отримані на основі припущень щодо законів розподілу залежної змінної та функції зв'язку, подано у табл. 4–6.

Таблиця 4. Порівняльна таблиця результатів

Модель	
Характер розподілу початкової змінної	Функція зв'язку
Гамма	LOG
Нормальний	LOG
Пуассона	LOG
Нормальний	Тотожна

Із табл. 6 видно, що величина ризику для побудованих моделей у середньому коливалась від 40–60 %, що є гранично допустимою величиною та все ж вимагає вжиття додаткових заходів щодо його мінімізації. Ризик втрат характеризується коефіцієнтом варіації, тобто відношенням стандартного відхилення до величини середнього.

При порівнянні моделей з нормальним розподілом та логарифмічною і тотожною функцією зв'язку було виявлено, що інформаційний критерій Акайке набуває приблизно однакового значення 20,66, тому вибирати модель доречно, виходячи із сумарного результату прогнозу величини збитків.

Отже, адекватною, прийнятною для практичного використання є модель із законом роз-

поділу Пуассона та експоненціальною функцією зв'язку через мінімальну величину похибки, показники значущості моделі, максимальне наближення до реальних даних прогнозних значень та достовірну оцінку величини ризику. Нормальна модель з тотожною функцією зв'язку дає змогу отримати результат за одну ітерацію з незначним значенням відносної похибки 1,62 %, але зі “слабкими” прогнозними значеннями збитків та хибною оцінкою ризику.

Висновки

Установлено, що для оцінювання параметрів узагальнених лінійних моделей доцільно застосовувати узагальнений зважений метод найменших квадратів, який у цьому випадку забезпечує отримання незміщених ефективних оцінок. Альтернативою є метод Монте-Карло для марковських ланцюгів.

На основі фактичних даних стосовно виплат за полісами автомобільного страхування побудовано узагальнені лінійні моделі для прогнозування величини можливих втрат. Модель побудовано, виходячи з припущень стосовно відомого закону розподілу. Прийнятною для подальшого використання виявилась модель із законом розподілу Пуассона та експоненціальною функцією зв'язку. Цей факт пояснюється мінімальною величиною похибки, показником адекватності моделі, максимальним наближенням прогнозних значень до реальних даних, а також достовірною оцінкою величини фінансового ризику. Нормальна модель з тотожною функцією зв'язку дає можливість отримати результат за одну ітерацію з відносною похибкою

Таблиця 5. Основні характеристики побудованої моделі

Сумарне значення збитків	Середнє	Стандартне відхилення	Мінімум	Максимум	Стандартна похибка	Дисперсія %
102008320,905	11805,690	15358,118	6273,867	18549,819	0,075	130,091
18111231,380	1897,457	939,910	4010,978	634,054	0,120	49,535
17921032,574	1877,531	1027,567	4234,951	558,354	0,176	54,730
17921032,589	1877,531	999,302	3535,396	118,004	0,188	53,224

Таблиця 6. Результуючі параметри оцінювання моделі

Сумарне значення збитків	Логарифм правдоподібності	Реальні сумарні значення збитків	Відхилення	Ризик втрат
102008320,905	-15742,754	17921032,581	84087288,32	0,984
18111231,380	-98700,167		190198,799	0,495
17921032,574	-42173677,24		0,007	0,547
17921032,589	-98700,167		0,009	0,532

1,62 %, але зі “слабкими” прогнозними значеннями збитків та некоректною оцінкою ризику.

Застосування запропонованої моделі для прогнозування процесів у сфері страхування за допомогою УЛМ гарантує високу точність оцінок прогнозів досліджуваної величини. Використання таких моделей дає можливість підвищити стійкість функціонування страхових компаній завдяки підвищенню якості рішень, які приймаються на основі оцінок прогнозів.

Список літератури

1. *Тэнман Л.Н.* Риски в экономике. – М.: ЮНИТИ, 2002. – 382 с.
2. *Y.Y. Haimes and J.R. Santos*, Risk modeling, assessment and management. New York: John Wiley & Sons, 2009, 447 p.
3. *R. Kaas et al.*, Modern actuarial risk theory. New York: Kluwer Academic Publishers, 2002, 318 p.
4. *J.F. Bouchaud and M. Potter*, Theory of financial risk management. Cambridge: Cambridge University Press, 2001, 218 p.
5. *P. McCullagh and J.A. Nelder*, Generalized Linear Models. New York: Chapman & Hall, 1989, 526 p.
6. *D. Collett*, Modeling Binary Data. New York: Chapman & Hall, 2002, 388 p.
7. *J. Gill*, Generalized linear models – a unified approach. London: Sage Publications, 2001, 101 p.

Рекомендована Радою
Навчально-наукового комплексу
Інститут прикладного системного
аналізу НТУУ “КПІ”

Надійшла до редакції
10 січня 2014 року