

Беспала О.М.

Національний технічний університет України «Київський політехнічний інститут імені Ігоря Сікорського», Київ

ОЦІНКА СТУПЕНЯ ВПЛИВУ МІЖ ПАРАМЕТРАМИ СТРУКТУРНО-ПРИЧИННОЇ МОДЕЛІ

Анотація: Причинно-наслідкові зв'язки лежать в основі сучасної науки та здатності навчати комп'ютери міркуванню подібному до людського. Дослідження причинності у поєднанні з сучасними досягненнями машинного навчання може значно прискорити науковий прогрес. Сучасні наукові досягнення підкреслюють важливість не лише встановлення причинності, а й оцінки ступеня впливу між причинно-наслідковими зв'язками.

В роботі поставлено інтервенційне та контрфактичне запитання: «Як зміняться значення параметрів моделі внаслідок впливу досліджуваного фактора?» та «На який фактор необхідно вплинути для найшвидшої зміни значення досліджуваного параметра моделі?». Відповівши на ці питання, можна сформулювати більш точні гіпотези та прискорити їх перевірку.

У статті запропоновано метод оцінки ступеня впливу між параметрами структурно-причинної моделі, який вирішує проблему виявлення найбільш впливових факторів, а також удосконалює дослідження причинно-наслідкових зв'язків між параметрами, що дозволить покращити прогнозування та управління досліджуваною моделлю. Запропонований підхід представлення даних у структурно-причинній моделі дозволяє робити висновки про те, що є причиною, а що – наслідком і враховувати вплив як прямих, так і непрямих причинно-наслідкових зв'язків. Оцінка ступеня впливу визначає найбільш значущі параметри та найменш впливові фактори. Отже, прогнозування змін значень параметрів моделі стає більш передбачуваним і контрольованим. В процесі оптимізації моделі можна виконати умовне ігнорування найменш впливових параметрів.

Ключові слова: причинно-структурна модель, причинно-наслідкові зв'язки, машинне навчання.

Bespala O.M.

National Technical University of Ukraine «Igor Sikorsky Kyiv Polytechnic Institute», Kyiv

ASSESSING THE DEGREE OF INFLUENCE BETWEEN THE PARAMETERS OF THE STRUCTURAL-CAUSAL MODEL

Abstract: Cause-and-effect relationships underpin modern science and the ability to teach computers human-like reasoning. The study of causality combined with modern advances in machine learning can greatly accelerate scientific progress. Modern scientific advancements emphasize the importance of not only establishing causality, but also assessing the degree of influence between cause-and-effect relationships.

The work poses an interventionist and counterfactual question: «How will the values of the model parameters change due to the influence of the studied factor?» and «Which factor needs to be influenced for the fastest change in the value of the studied model parameter?». Formulating more precise hypotheses and accelerating their testing can be achieved by answering these questions.

This article proposes a method for estimating the degree of influence between the parameters of the structural-causal model, which solves the problem of identifying the most influential factors,

and also improves the study of causal relationships between parameters, which will allow improving the prediction and management of the model under study. The proposed approach of presenting data in a structural-causal model allows you to draw conclusions about what is the cause and what is the effect and take into account the influence of both direct and indirect cause-and-effect relationships. The assessment of the degree of influence determines the most significant parameters and the least influential factors. Consequently, predicting changes in model parameter values becomes more predictable and controllable. Conditional ignoring of the least influential parameters can be used in the model optimization process.

Keywords: *causal-structural model, cause-and-effect relationships, machine learning*

Постановка проблеми

В багатьох галузях нерідко постає задача виявлення причинно-наслідкових зв'язків між параметрами заданої моделі з метою подальшого управління, прогнозування та виконання інших маніпуляцій. Наприклад, мають місце задачі в області медицини, які полягають в дослідженні та виявленні факторів впливу на здоров'я людини [1]. В області екології встановлюються та вивчаються фактори впливу на різні екосистеми, а також їхній вплив та взаємозв'язок з іншими системами, включаючи економічні, політичні та інші галузі. Так зміни глобального клімату вимагають корективів, або і повної реформації в галузі енергетики, хімічної промисловості, тощо. Також тінь таких задач відображається і на бізнес-рішеннях, які мають реагувати на швидку зміну умов.

З розвитком комп'ютерних наук, встановлення причинно-наслідкових зв'язків стає все більш актуальним питанням та має успіхи у вирішенні поставлених деяких задач [2]. Проте виявлення причинного зв'язку може бути недостатньо у випадку, коли необхідно здійснювати об'єктивне управління системою, оскільки потрібні знання про ступінь взаємозв'язку між компонентами та як зміна одного факту вплине на інший і на систему в цілому [3]. Тому постає потреба у вивченні та оцінці взаємозв'язків між параметрами системи та їх ступінь впливу як між собою так і на систему в цілому.

Аналіз останніх досліджень і публікацій

Огляд сучасної наукової літератури свідчить про поєднання машинного навчання з розумінням причинно-наслідкових зв'язків [4]. Використання глибокого навчання для виявлення причинності сприяє науковому прогресу в створенні машин з міркуванням подібним до людського [5, 6]. Також в роботі [7] розглянуто обернену задачу, а саме вивчення причинних факторів, які призводять до певного набору спостережень.

Світова наукова спільнота вже неодноразово стверджувала, що статистичні моделі значно слабші в порівнянні з причинно-наслідковими моделями та й кореляційні тісні зв'язки не завжди вказують на присутність причинно-наслідкового зв'язку. В роботі [8] наведено чітку залежність між народжуваністю людей у Європі та популяцією лелек, що є хибно позитивним результатом, оскільки вплив на популяцію лелек не змінить народжуваності людей.

Моделювання причинно-наслідкових зв'язків за допомогою машинного навчання (нейронних мереж) потребує великого набору високоякісних даних. Втім сучасні методи виявлення причинності можуть виявляти взаємозв'язки і за присутності латентних змінних або неспостережуваних факторів. Але у випадку зміни умов модель потрібно частково або навіть повністю перенавчати. Відповідно, навіть точні причинні моделі не достатні для прийняття рішень чи управління.

Мета і задачі дослідження

Метою дослідження є знаходження оцінки ступеню впливу між причинно-наслідковими зв'язками та виявлення найбільш впливових змінних в системі. Для досягнення цієї мети пропонується вивчення не простих зв'язків між змінними, виявлення тісноти зв'язків та оцінки факторів впливу в причинній моделі.

Результати дослідження

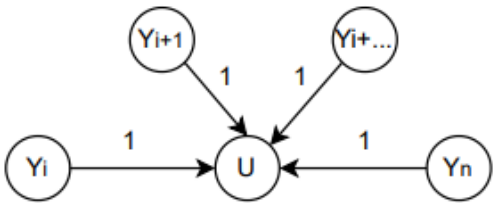
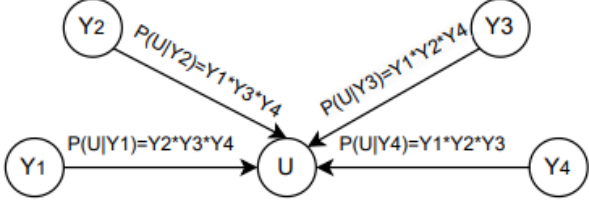

Побудова причинно-наслідкової моделі починається з графічного представлення усіх параметрів, в тому числі прихованих або не вимірних, які є системі та можуть мати вплив. Спрямованій ациклічній графі (DAG) можна віднести до таких графічних моделей, коли вершини графа є параметри моделі, а стрілки між ними вказують на прямий або непрямий вплив між ними. Одним із популярних методів моделювання причинно-наслідкових зв'язків є моделювання структурними рівняннями (SEM). В роботі [9] запропоновано використання структурно-причинної моделі SCM (Structural Causal Model), в якій поєднано SEM та графічні моделі з метою відобразити більш реальну структуру причинності.

Для подальшої роботи будемо вважати, якщо в структурно причинній моделі між двома параметрами існує причинно-наслідковий зв'язок, то між ним встановлюється стрілка і такий зв'язок вказує на прямий вплив. Непрямим впливом, відповідно вважається наявність послідовних зв'язків між досліджуваними параметрами. Якщо на один параметр діє два і більше прямих впливи, то така змінна в моделі є колайдером U. Якщо змінна Y не залежить ні від однієї змінної в моделі, то змінна Y є екзогенною. До екзогенної змінної не входить жодна стрілка. Вхідні змінні моделі є екзогенними змінними. Якщо змінна U залежить принаймні від однієї з інших змінних моделі, то змінна U є ендогенною. До ендогенної змінної входить принаймні одна стрілка. Колайдер є ендогенною змінною. Якщо ендогенна змінна не має жодного впливу, тоді вона є вершиною U_v .

Для визначення максимально впливової змінної P_{max} в системі пропонується виконати декомпозицію причинно-структурної моделі на компоненти $G(U | Y_i)$, де U – колайдер, Y_i – екзогенні змінні. Доречно відмітити, що в загальній структурно-причинній моделі колайдери U та екзогенні змінні Y_i можуть бути і ендогенними змінними. За умовою в колайдерах відбувається процес, який є результатом дії певних операцій над екзогенними змінними. В таблиці 1 наведено формули, за якими визначається вплив кожної екзогенної змінної на колайдери.

Таблиця 1

Вплив екзогенних змінних на колайдер

Операція колайдера	Вплив	Графічна модель
1	2	3
$U = \sum_{i=1}^n Y_i$	$P(U Y_i) = 1$	
$U = \prod_{i=1}^n Y_i$	$P(U Y_i) = \frac{\prod Y_i}{Y_i}$	
$U = \frac{Y_{i=2}}{Y_{i=1}}$	$P(U Y_{i=1}) = \frac{1}{Y_{i=2}}$ $P(U Y_{i=2}) = \frac{Y_{i=1}}{(Y_{i=2})^2}$	

Для знаходження змінної P_{max} , яка має максимальний вплив серед усіх Y_i в $G(U | Y_i)$ використовуємо формулу 1:

$$P_{max} = \max(P(Y_i)). \tag{1}$$

Якщо $P(Y_1) = P(U|Y_1) = \dots = P(U|Y_n)$, тоді всі Y_i , які не є екзогенними змінними в причинно-структурній моделі можна вважати єдиним колайдером U та повторити операцію вже з наступними змінними, які мають прямиий вплив на всі Y_i .

Знаходження змінної P_{max} в загальній причинно-структурній моделі починається від колайдера-вершини U_v до екзогенних змінних. Пошук здійснюється шляхом обходу колайдерів, в яких знайдено P_{max} для кожного $G(U | Y_i)$.

Вилучивши знайдену змінну P_{max} в загальній причинно-структурній моделі можна повторити цикл, таким чином знаходимо послідовну ієрархію змінних, які мають вплив на U_v в спадаючому порядку.

Представлення причинно-структурної моделі можна відобразити у вигляді двовимірного масиву, в якому послідовність проіндексованих змінних наведено в таблиці 2.

Таблиця 2

Представлення масиву даних причинно-структурної моделі

	Екзогенні змінні				Ендогенні змінні				
					Колайдери				Вершина
	X_1	X_2	...	X_i	X_{i+1}	X_{i+2}	...	X_{i+k}	$X_{i=n}$
X_1	[1] [1]	[1] [2]	[1] [...]	[1] [i]	[1] [i+1]	[1] [i+2]	[1] [...]	[1] [i+k]	[1] [n]
X_2	[2] [1]	[2] [2]	[2] [...]	[2] [i]	[2] [i+1]	[2] [i+2]	[2] [...]	[2] [i+k]	[2] [n]
...	[...] [1]	[...] [2]	[...] [...]	[...] [i]	[...] [i+1]	[...] [i+2]	[...] [...]	[...] [i+k]	[...] [n]
X_i	[i] [1]	[i] [2]	[i] [...]	[i] [i]	[i] [i+1]	[i] [i+2]	[i] [...]	[i] [i+k]	[i] [n]
X_{i+1}	[i+1] [1]	[i+1] [2]	[i+1] [...]	[i+1] [i]	[i+1] [i+1]	[i+1] [i+2]	[i+1] [...]	[i+1] [i+k]	[i+1] [n]
X_{i+2}	[i+2] [1]	[i+2] [2]	[i+2] [...]	[i+2] [i]	[i+2] [i+1]	[i+2] [i+2]	[i+2] [...]	[i+2] [i+k]	[i+2] [n]
...	[...] [1]	[...] [2]	[...] [...]	[...] [i]	[...] [i+1]	[...] [i+2]	[...] [...]	[...] [i+k]	[...] [n]
X_{i+k}	[i+k] [1]	[i+k] [2]	[i+k] [...]	[i+k] [i]	[i+k] [i+1]	[i+k] [i+2]	[i+k] [...]	[i+k] [i+k]	[i+k] [n]
$X_{i=n}$	[n] [1]	[n] [2]	[n] [...]	[n] [i]	[n] [i+1]	[n] [i+2]	[n] [...]	[n] [i+k]	[n] [n]

Відповідну структуру масиву можна використати для обчислень впливу між змінними. P_{max} буде знайдено, коли Y_i буде екзогенна змінна для $\max(P(U|Y_i))$, тобто $Y[i][j=m]$, а $Y[i=m][j]$ порожні.

Таблиця 3

Представлення масиву даних впливу між елементами причинно-структурної моделі

	X_{i+1}	X_{i+2}	...	X_{i+k}	$X_{i=n}$
X_1	$P(X_{i+1} X_1)$	$P(X_{i+2} X_1)$...	$P(X_{i+k} X_1)$	$P(X_n X_1)$
X_2	$P(X_{i+1} X_2)$	$P(X_{i+2} X_2)$...	$P(X_{i+k} X_2)$	$P(X_n X_2)$
...	$P(X_{i+1} \dots)$	$P(X_n \dots)$
X_i	$P(X_{i+1} X_i)$	$P(X_{i+2} X_i)$...	$P(X_{i+k} X_i)$	$P(X_n X_i)$
X_{i+1}	$P(X_{i+1} X_{i+1})$	$P(X_{i+2} X_{i+1})$...	$P(X_{i+k} X_{i+1})$	$P(X_n X_{i+1})$
X_{i+2}	$P(X_{i+1} X_{i+2})$	$P(X_{i+2} X_{i+2})$...	$P(X_{i+k} X_{i+2})$	$P(X_n X_{i+2})$
...
X_{i+k}	$P(X_{i+1} X_{i+k})$	$P(X_{i+2} X_{i+k})$...	$P(X_{i+k} X_{i+k})$	$P(X_n X_{i+k})$
$X_{i=n}$	$P(X_{i+1} X_n)=0$	$P(X_{i+2} X_n)=0$...	$P(X_{i+k} X_n)=0$	$P(X_n X_n)=0$

Алгоритм пошуку змінної P_{max} показано на рисунку 1.

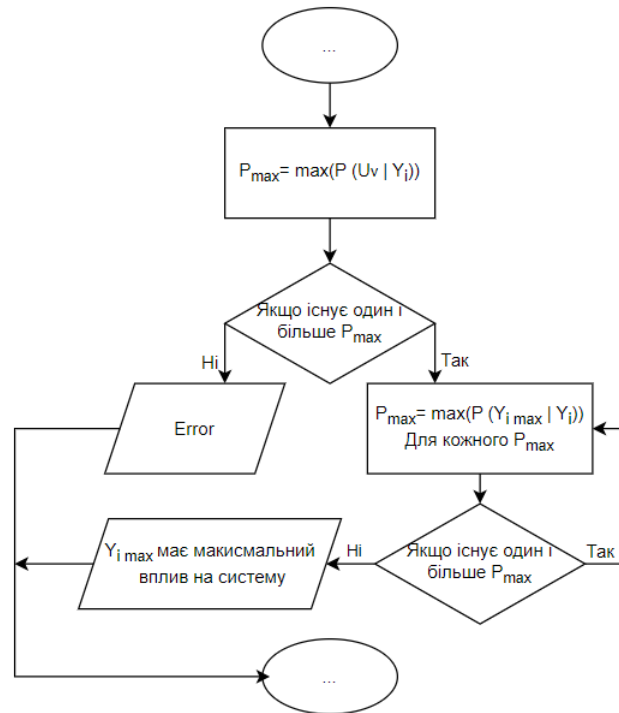


Рис 1. Алгоритм пошуку змінної з максимальним впливом

Запропонований алгоритм було застосовано на практиці в якості визначення найбільш вагомого фактору, який впливає на ризик захворювання людей від забрудненого атмосферного повітря. Для цього було використано розрахунок середньої добової дози впливу бенз(а)пірену на населення міста. Відповідно до цих методичних рекомендацій [10] середня добова доза речовини, *мг/кг-доба*, розраховується за формулою 2:

$$\frac{ADD}{LADD} = \frac{[(Ca * Tout * Vout) + (Ch * Tin * Vin)] * EF * ED}{(BW * AT * D)}. \#(2)$$

Таблиця 4

Вхідні дані. Значення екзогенних змінних

Параметр	Одиниця вимірювання	Характеристика	Стандартне значення
<i>Ca</i>	мг/куб.м	концентрація речовини в атмосферному повітрі	$0,95 \cdot 10^{-6}$
<i>Tout</i>	год/доба	час, що проводиться поза приміщенням	8
<i>Vout</i>	куб.м/год	швидкість дихання поза приміщенням	1,4
<i>Ch</i>	мг/куб.м	концентрація речовини у повітрі приміщення	$0,95 \cdot 10^{-6}$
<i>Tin</i>	год/доба	час, що проводиться у приміщенні	16
<i>Vin</i>	куб.м/год	швидкість дихання у приміщенні	0,63
<i>EF</i>	днів/рік	частота впливу	365
<i>ED</i>	Роки	тривалість впливу	30
<i>BW</i>	Кг	маса тіла	70
<i>AT</i>	Роки	період осереднення експозиції	70
<i>D</i>	Дні	число днів у році	365
<i>SF</i>	мг/(кг*доба)	фактор нахилу	3,1
<i>POP</i>	чол.	чисельність популяції, що підпадає під вплив даного фактора	300000

Індивідуальний канцерогенний та популяційний ризики розраховуються відповідно за формулами 3 та 4:

$$CR = LADD * SF, \#(4)$$

$$PCR = CR * POP. \#(5)$$

Вхідні дані, які використовувались в розрахунках наведені в таблиці 4.

Відповідно до вхідних значень, в таблиці 5 відображено масив з врахуванням послідовностей змінних в причинно-структурній моделі.

Таблиця 5

Представлення масиву даних причинно-структурної моделі

	C a	T o u t	V o u t	C h	T i n	V i n	E F	E D	B W	A T	D	P O P	K 1	K 2	K3	K4	K 5	ADD	CR	PCR		
Ca													0,95									
Tout													8									
Vout													1,4									
Ch													0,95									
Tin													16									
Vin													0,63									
EF																350						
ED																30						
BW																	70					
AT																	70					
D																	365					
POP																					300000	
SF																					3,1	
K1															10,64							
K2															9,576							
K3																20,216						
K4																		212268				
K5																		1788500				
ADD																					0,118685	
CR																						0,367923
PCR																						

Використовуючи формули з колонки 2 таблиці 1 було обчислено значення впливу змінних та за формулою 1 визначено P_{max} в кожному колайдері.

Таблиця 6

Представлення масиву даних впливу між змінними причинно-структурної моделі

	C a	T o u t	V o u t	C h	T i n	V i n	E F	E D	B W	A T	D	P O P	K 1	K2	K 3	K4	K5	ADD	CR	PCR		
Ca													11,2									
Tout													0,133									
Vout													0,76									
Ch													10,08									
Tin													0,599									
Vin													0,152									
EF																606,48						
ED																7075,6						
BW																	25550					
AT																	25550					
D																	4900					
POP																						0,367923

SF																	0,118685	
K1								1										
K2								1										
K3								10500										
K4													0,559127761					
K5													0,06636004					
ADD																	3,1	
CR																		300000
PCR																		

Для візуалізації отриманих результатів було побудовано графічну причинно-структурну модель, на якій відмічено послідовність роботи алгоритму пошуку змінної, яка має найбільший вплив на систему.

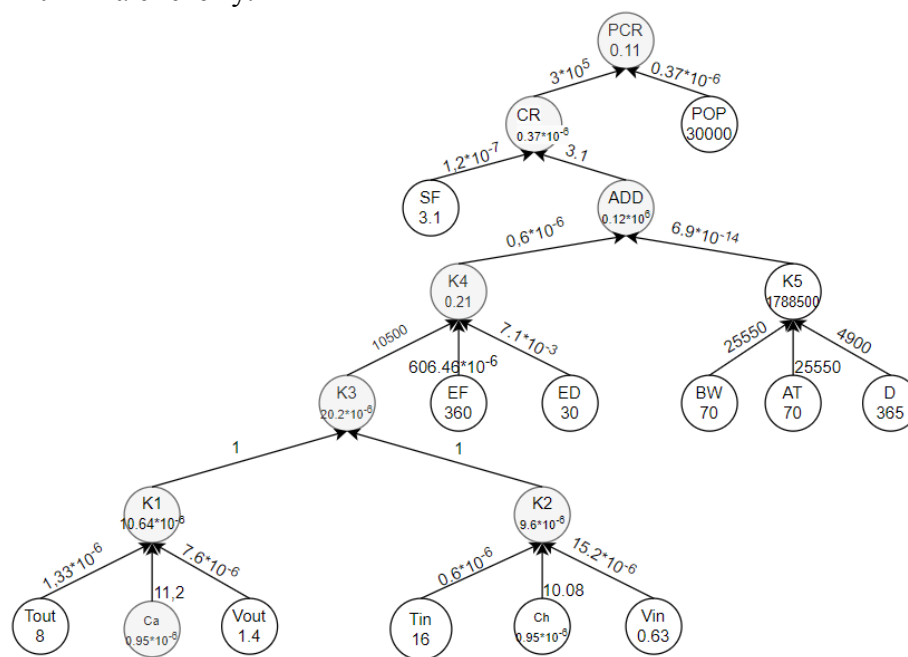


Рис. 2 Графічна модель

Оскільки задача полягала у визначенні фактору, який має найбільший вплив на значення популяційного ризику PCR, тому і початок алгоритму починався саме з цього значення, оскільки алгоритм її визначив як вершину. В результаті, можна стверджувати, що фактор Ca, а саме концентрація речовини в атмосферному повітрі має найбільший вплив на популяційний ризик. Якщо виключити з причинно-структурної моделі змінну Ca та повторивши алгоритм, можна визначити послідовну ієрархію значень, які впливають на змінну PCR. В якості тестування було зроблено однаково рівний вплив $\delta = 0,1$ на кожну з екзогенних змінних, що призвело до максимального збільшення показника PCR= 6115,179 при дії на фактор Ca.

Висновки

Використовуючи парадигми причинного моделювання та сучасні комп’ютерні засоби і методи можна досягти наближення до створення систем з подібним до людського міркуванням. Зокрема, в цій роботі розглянуто метод встановлення впливу зв’язку між змінними в причинно-структурній моделі, що має в перспективі покращити управління системою. Маючи інформацію не тільки про наявність причинно-наслідкового зв’язку між компонентами, а й про їх ступінь впливу (або ж тісноту зв’язку) система може бути більш прогнозованою та гнучкішою в управлінні.

Список літератури:

1. von Kügelgen, Julius, Luigi Gresele, and Bernhard Schölkopf. "Simpson's paradox in Covid-19 case fatality rates: a mediation analysis of age-related causal effects." *IEEE Transactions on Artificial Intelligence*. 2021. Vol. 2, № 1. P. 18-27.
2. Беспала О. М. Інструментарій причинно-наслідкового висновку: огляд та перспективи. *Control systems and computers*. 2020. № 5. С. 52 – 63. ISSN 2706-8145 DOI: 10.15407 / csc.2020.05.052.
3. Беспала О. М. Метод пошуку та оцінки впливу причинно-наслідкових зв'язків в системах прийняття рішень. *Вісник НТУ "ХПІ". Серія: Інформатика та моделювання*. 2020. № 1 (3). С. 59 – 72. ISSN 2079-0031 (Print), ISSN 2411-0558 (Online).
4. Manta, D. C., Hu, E. J., & Bengio, Y. GFlowNets for Causal Discovery: an Overview. In *ICML 2023 Workshop on Structured Probabilistic Inference & Generative Modeling*. 2023, July.
5. Morandé, S., & Tewari, V. Causality in Machine Learning: Innovating Model Generalization through Inference of Causal Relationships from Observational Data. *Qeios*. 2023.
6. Wang, H., Fu, T., Du, Y., Gao, W., Huang, K., Liu, Z., ... & Zitnik, M. Scientific discovery in the age of artificial intelligence. *Nature*. 2023. № 620(7972). P. 47-60.
7. Schölkopf, B., Locatello, F., Bauer, S., Ke, N. R., Kalchbrenner, N., Goyal, A., Bengio, Y. Toward causal representation learning. *Proceedings of the IEEE*. 2021. Vol. 109, № 5. P. 612-634.
8. Matthews, Robert. «Storks deliver babies (p= 0.008).» *Teaching Statistics*. 2000. Vol. 22.2 P. 36-38.
9. Pearl, J. *Causality*. Cambridge university press. 2009. P. 478.
10. Про затвердження методичних рекомендацій «Оцінка ризику для здоров'я населення від забруднення атмосферного повітря»: Наказ Міністерства охорони здоров'я України від 13.04.2007 N 184.

References:

1. von Kügelgen, Julius, Luigi Gresele, and Bernhard Schölkopf. "Simpson's paradox in Covid-19 case fatality rates: a mediation analysis of age-related causal effects." *IEEE Transactions on Artificial Intelligence*. 2021. Vol. 2, No. 1. P. 18-27.
2. Besspala O. M. Toolkit of causal inference: review and perspectives. *Control systems and computers*. 2020. No. 5. P. 52 - 63. ISSN 2706-8145 DOI: 10.15407 / csc.2020.05.052.
3. Besspala O. M. The method of finding and evaluating the influence of causal relationships in decision-making systems. *Bulletin of NTU "Khpi". Series: Informatics and modeling*. 2020. No. 1 (3). P. 59 – 72. ISSN 2079-0031 (Print), ISSN 2411-0558 (Online).
4. Manta, D. C., Hu, E. J., & Bengio, Y. GFlowNets for Causal Discovery: an Overview. In *ICML 2023 Workshop on Structured Probabilistic Inference & Generative Modeling*. 2023, July.
5. Morandé, S., & Tewari, V. Causality in Machine Learning: Innovating Model Generalization through Inference of Causal Relationships from Observational Data. *Qeios*. 2023.
6. Wang, H., Fu, T., Du, Y., Gao, W., Huang, K., Liu, Z., ... & Zitnik, M. Scientific discovery in the age of artificial intelligence. *Nature*. 2023. No. 620(7972). P. 47-60.
7. Schölkopf, B., Locatello, F., Bauer, S., Ke, N.R., Kalchbrenner, N., Goyal, A., Bengio, Y. Toward causal representation learning. *Proceedings of the IEEE*. 2021. Vol. 109, No. 5. P. 612-634.
8. Matthews, Robert. "Storks deliver babies (p= 0.008)." *Teaching Statistics*. 2000. Vol. 22.2 P. 36-38.
9. Pearl, J. *Causality*. Cambridge university press. 2009. P. 478.
10. On the approval of methodological recommendations "Risk assessment for public health from atmospheric air pollution": Order of the Ministry of Health of Ukraine dated April 13, 2007 N 184.