

# РОЗПІЗНАВАННЯ ТА КОРЕКЦІЯ ЗОБРАЖЕНЬ ДЕСЯТКОВОГО СКЛАДАННЯ У СТОВПЧИК ЗА ДОПОМОГОЮ CNN ТА АЛГОРИТМУ КОКА-ЯНГЕРА-КАСАМІ

П. П. Пеліх<sup>1, а</sup>, В. М. Кригін<sup>2, б</sup>

<sup>1</sup>Національний технічний університет України  
«Київський політехнічний інститут імені Ігоря Сікорського»,  
НН Фізико-технічний інститут

<sup>2</sup>Міжнародний науково-навчальний центр інформаційних технологій та систем  
НАН України та МОН України

## Анотація

У даній роботі розглянуто метод розпізнавання та корекції зображень складання десяткових виразів у стовпчик. Метод використовує згорткову нейронну мережу та двовимірне узагальнення алгоритму Кока-Янгера-Касамі для зважених двовимірних контекстно-вільних грамастик. Описано процес побудови мережі та її навчання, правила граматики для розпізнавання складання десяткових чисел у стовпчик, а також опис комп'ютерної програми, де реалізовано запропонований метод.

*Ключові слова:* згорткові нейронні мережі, алгоритм Кока-Янгера-Касамі, аналіз зображень математичних виразів

## Вступ

Структурне розпізнавання візуальних образів — це сукупність методів, які аналізують зв'язки між об'єктами на зображенні для їх розпізнавання. Воно знайшло своє застосування в обробці документів, а також різноманітних рукописних виразів, написання яких чітко визначено деякими правилами, до яких відносяться й математичні вирази, що є об'єктом нашого дослідження.

У цій роботі продемонстровано метод, що використовує класифікатор рукописних цифр та узагальнення алгоритму Кока-Янгера-Касамі. Основна відмінність від роботи [1] полягає у використанні згорткової нейронної мережі, що є більш сучасним підходом для класифікації об'єктів на зображенні.

Структура роботи наступна: у першому розділі міститься високорівневий опис методу, що містить три основні етапи, наступні два розділи містять більш детальний опис цих етапів, а в останньому розділі наведено результати експериментів на модельних та реальних даних.

## 1. Опис програми

### 1.1. Вхідні дані

На вхід програмі подається одноканальне 8-бітне зображення  $\Pi$  розміру  $H \times W$ , на якому міститься складання десяткових чисел у стовпчик без знаку «+» та без риски під доданками. Приклад вхідного зображення наведений на рис. 1.

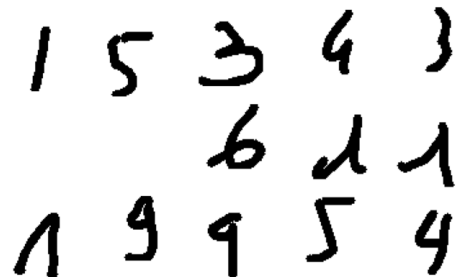


Рис. 1. Приклад вхідного зображення

### 1.2. Кроки програми

Щоб визначити, чи є коректним вираз на вхідному малюнку, зображення проходить через етапи, що описано у даному підрозділі.

**Пошук та розпізнавання цифр.** На першому кроці виконується пошук цифр. Зображенням проходить ковзне вікно, вміст якого подається на вхід класифікатора. Результат класифікації ставиться у відповідність поточному ковзному вікну.

**Розпізнавання дробового виразу.** Наступним етапом є структурний аналіз зображення, який проводиться узагальненим алгоритмом Кока-Янгера-Касамі для роботи зі зваженими двовимірними контекстно-вільними грамастиками, результатом якого стане чисельний показник некоректності виразу на зображенні.

**Корекція вхідного зображення.** Корекцією ми називаємо таку мінімально можливу заміну цифр на малюнку, щоб вираз став коректним. Зображення цифр, на які будуть замінятися вхідні, беруться з деякого набору еталонів.

<sup>а</sup>pavpel-ipt22@iit.kpi.ua

<sup>б</sup>valeriy.krygin@gmail.com

Якщо на зображенні розпізнано коректний вираз, воно не модифікується. Якщо ж помічена помилка, алгоритм корегує зображення та віддає його користувачу програми.

## 2. Пошук та розпізнавання цифр

### 2.1. Прохід по зображенню ковзним вікном

Як уже було зазначено, пошук цифр на зображенні відбувається за допомогою ковзного вікна. Вміст кожної позиції вікна проходить через класифікатор, який для кожної цифри вказує ступінь упевненості у тому, що в поточній позиції знаходиться певна цифра. На цьому етапі потрібно обрати

- розмір ковзного вікна,
- розмір горизонтального кроку,
- розмір вертикального кроку.

Зазначимо, що прохід може бути декілька і, можливо, з різними параметрами, бо цифри на зображенні можуть мати різні розміри, що потребує різних розмірів ковзного вікна.

### 2.2. Класифікація вмісту ковзного вікна

Кожному ковзному вікну для кожного можливого класу треба надати чисельну характеристику схожості вмісту на відповідну цифру або щось інше. Для цього було використано класифікатор, який повертає числа з проміжку  $[0, 1]$  для кожної цифри від 0 до 9 і для відсутності цифри. У якості останнього класу виступають підзображення, які містять або фон, або тільки частину цифри, що не може бути віднесено до класів цифр. Останній клас будемо позначати символом  $\epsilon$ . Сума характеристик схожості для усіх цифр та відсутності цифри у кожному віконці дорівнює одиниці. Надалі цю характеристику будемо скорочено називати ймовірністю. У якості класифікатора було вирішено використовувати згорткові нейронні мережі (англ. Convolutional Neural Networks, CNN).

### 2.3. Підготування даних

Для тренування нейронної мережі були вибрані зображення цифр з набору математичних виразів CROHME [2]. Для доповнення (англ. augmentation) набору були взяті наступні послідовно застосовані перетворення до зображень з оригінального набору:

**Масштабування.** Розмір вхідного символу випадковим чином масштабується з коефіцієнтом масштабування у діапазоні  $[0.66, 1]$ .

**Зсув.** Відмасштабований вхідний символ зсувається випадковим чином по вертикалі та горизонталі необмеженим чином, тобто можливі ситуації, коли частина символу або навіть він увесь виходить за межі зображення.

**Зашумлення.** До зображення випадковим чином додається на вибір шум Гаусса або шум «солі та перцю» (англ. salt and pepper noise) різних інтенсивностей.

Перше перетворення обумовлено тим, що на реальному зображенні можуть зустрічатися цифри різних

розмірів. Друге перетворення створює не тільки такі навчальні приклади, у яких цифри відцентровані на зображенні, але й такі, на яких присутні тільки частини від цифр, які у даній задачі мають клас  $\epsilon$ , тобто не класифіковані як цифри. Третє перетворення має створити наближення до реальних умов, коли на зображенні присутній шум різного походження.

Отримані на попередньому етапі зображення нормуються за допомогою перетворень, що описано нижче.

**Нормалізація.** Значення у кожному пікселі нормалізуються таким чином, щоб усі вони лежали у діапазоні  $[0, 1]$ .

**Інверсія.** Значення у кожному пікселі змінюється відповідно до формули

$$p_{new} = 1 - p_{old},$$

де  $p_{new}$  — нове значення у пікселі,  $p_{old}$  — старе значення.

### 2.4. Архітектура згорткової нейромережі

Штучна нейромережа побудована з декількох складних блоків, кожний з яких складається з

- 1) шару пакетної нормалізації (англ. Batch Normalization); у [3] було показано, що додавання цього шару суттєво прискорює навчання моделі та підвищує її точність;
- 2) шару згортки;
- 3) активація ReLU.

Два таких складних блоки ідуть підряд, маючи різні параметри для згорткового шару, а після кожної пари іде шар Max Pooling з розміром вікна  $2 \times 2$  та кроком 2. Після декількох таких послідовностей іде повнозв'язний шар з функцією *softmax* в кінці. Архітектура CNN наведена на рис. 2.

На вхід мережа приймає одноканальні зображення розміру  $50 \times 50$ , причому значення у кожному пікселі лежать у діапазоні  $[0, 1]$ .

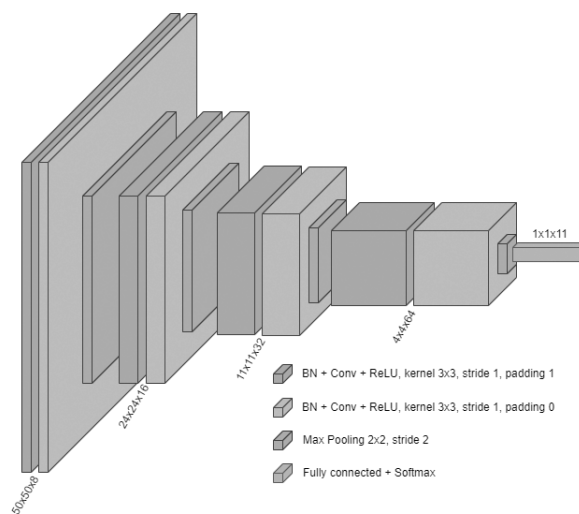


Рис. 2. Архітектура обраної згорткової нейромережі

## 2.5. Обробка отриманих результатів

Результати роботи класифікатора для кожного ковзного вікна зберігаються у матриці  $P(x_1, y_1, x_2, y_2, \ell)$ , де

- $(x_1, y_1)$  — координати верхнього лівого кута ковзного вікна;
- $(x_2, y_2)$  — координати нижнього правого кута ковзного вікна;
- $\ell$  — клас, для якого зберігається значення у матриці,  $\ell \in \{0, \dots, 9, \varepsilon\}$ .

Надалі буде використовуватись скорочений запис  $\mathbf{p}_i = (x_i, y_i)$ , де  $i$  залежатиме від контексту. Для кожного ковзного вікна  $w(\mathbf{p}_1, \mathbf{p}_2)$  у клітинці  $P(\mathbf{p}_1, \mathbf{p}_2, \ell)$  зберігається  $1 - P_\ell(w)$ , де  $P_\ell(w)$  — ймовірність того, що вікно  $w$  містить об'єкт класу  $\ell$ . Інші значення матриці заповнюються значеннями  $+\infty$ .

## 3. Розпізнавання виразу

Подальший етап проводить структурний аналіз зображення та надає певну оцінку коректності зображення.

### 3.1. Узагальнений алгоритм Кока-Янгера-Касамі

Алгоритм Кока-Янгера-Касамі (англ. Cocker-Younger-Kasami algorithm, CYK algorithm) — це алгоритм, який дозволяє встановити, чи можна у мові, породженою даною контекстно-вільною граматикою, вивести задане слово. Якщо так, можна також вивести спосіб отримання цього слова з набору правил граматики. Цей алгоритм було відкрито незалежно різними вченими у 60-х роках минулого століття, за що він і отримав таку назву.

Наведений для одновимірних контекстно-вільних грамастик алгоритм природнім чином узагальнюється і для роботи з двовимірними, а також до більших розмірностей. Ідея узагальнення одновимірного алгоритму на двовимірний випадок була наведена у [4], де також наведена його складність —  $O(m^2 n^2 (m+n))$ . Окрім узагальнення на більші розмірності грамастик він також може бути модифікований для роботи зі зваженими грамастиками, де правила не бінарні, а вказують ступінь можливості конкатенації символів цими правилами для отримання нетермінальних символів.

Нехай  $X = \{1, \dots, W\} \times \{1, \dots, H\} \times \{1, \dots, W\} \times \{1, \dots, H\}$ , а  $G = \langle T, N, g^V, g^H, g, P \rangle$  — граMATика, де  $T$  — множина термінальних символів,  $N$  — множина нетермінальних символів,  $g^V : N^3 \rightarrow \{0, 1\}$  — правила вертикальної конкатенації,  $g^H : N^3 \rightarrow \{0, 1\}$  — правила горизонтальної конкатенації,  $g : N \times T \rightarrow \{0, 1\}$  — правила перейменування нетермінальних символів на термінальні,  $P : X \times T \rightarrow \mathbb{R}$  — матриця ваг термінальних символів. Модифікація алгоритму для зважених двовимірних контекстно-вільних грамастик наведена у алг. 1.

### Algorithm 1 Узагальнений алгоритм Кока-Янгера-Касамі

Введемо допоміжні функції. Функція розрахунку штрафу за горизонтальне склеювання двох частин зображення

$$S(\mathbf{p}_L, \mathbf{p}_R, n_U, n_D, y) = \mathbf{1}(g^H(n_U, n_D, n) = 0) \cdot \infty + P(\mathbf{p}_L, x_R, y, n_U) + P(x_L, y + 1, \mathbf{p}_R, n_D). \quad (1)$$

Функція розрахунку штрафу найкращого правила горизонтального склеювання двох частин зображення

$$H(\mathbf{p}_L, \mathbf{p}_R) = \min_{\substack{n_U \in N \\ n_D \in N \\ y_L \leq y < y_R}} S(\mathbf{p}_L, \mathbf{p}_R, n_U, n_D, y).$$

Функція розрахунку штрафу за вертикальне склеювання двох частин зображення

$$S^V(\mathbf{p}_L, \mathbf{p}_R, n_L, n_R, y) = \mathbf{1}(g^V(n_L, n_R, n) = 0) \cdot \infty + P(\mathbf{p}_L, x_R, y, n_L) + P(x_L, y + 1, \mathbf{p}_R, n_R). \quad (2)$$

Функція розрахунку штрафу найкращого правила вертикального склеювання двох частин зображення

$$V(\mathbf{p}_L, \mathbf{p}_R) = \min_{\substack{n_L \in N \\ n_R \in N \\ x_L \leq x < x_R}} S^V(\mathbf{p}_L, \mathbf{p}_R, n_L, n_R, y).$$

Перейдемо до описання самого алгоритму.

**Вхід:** ГраMATика  $G = \langle T, N, g^V, g^H, g, P \rangle$  з обчисленими вагами для термінальних символів.

**Результат:** Матриця  $P'$  — обчислені ваги для кожного підзображення і зображення в цілому.

```

for  $t \in T$  do
  for  $x_L = \overline{1}, \overline{W-1}, y_L = \overline{1}, \overline{H-1}$  do
    for  $x_R = \overline{x_L+1}, \overline{W}, y_R = \overline{y_L+1}, \overline{H}$  do
       $P'(\mathbf{p}_L, \mathbf{p}_R, t) = P(\mathbf{p}_L, \mathbf{p}_R, t)$ 
    end for
  end for
end for
for  $s = \overline{1}, \overline{H \cdot W}$  do
  for  $x_L = \overline{1}, \overline{W}, y_L = \overline{1}, \overline{H}$  do
    for  $\mathbf{p}_R : (y_R + 1 - y_L) \cdot (x_R + 1 - x_L) = s$  do
      for  $n \in T \cup N$  do
         $h \leftarrow H(\mathbf{p}_L, \mathbf{p}_R, n)$ 
         $v \leftarrow V(\mathbf{p}_L, \mathbf{p}_R, n)$ 
         $r \leftarrow \min_{t \in T} (P'(\mathbf{p}_L, \mathbf{p}_R, t) + \mathbf{1}(g(t, n) = 1) \cdot \infty)$ 
         $P'(\mathbf{p}_L, \mathbf{p}_R, n) \leftarrow \min(P'(\mathbf{p}_L, \mathbf{p}_R, n), h, v, r)$ 
      end for
    end for
  end for
end for
return  $P'$ 

```

### 3.2. Етап корекції

Після того, як отримали матрицю  $P'$ , яка містить вагу кожного символа, а також для кожного підзображення та зображення в цілому, запускаємо алг. 2.

**Algorithm 2** Алгоритм корекції зображення

Нехай  $I \in N$  — нетермінальний символ, який від-  
повідає коректному зображенню,  $\mathbb{E}$  — множина ета-  
лонних зображень. Даний алгоритм використовує  
функції (1) та (2). Введемо функцію пошуку накра-  
пцю горизонтального розбиття підзображення, на  
якому міститься символ  $n$ ,

$$\hat{H}(\mathbf{p}_L, \mathbf{p}_R, n) = \underset{\substack{n_U \in N \\ n_D \in N \\ y_L \leq y < y_R}}{\operatorname{argmin}} S(\mathbf{p}_L, \mathbf{p}_R, n_U, n_D, y),$$

та функцію пошуку найкращого вертикального роз-  
биття підзображення, на якому міститься символ  $n$ ,

$$\hat{V}(\mathbf{p}_L, \mathbf{p}_R, n) = \underset{\substack{n_L \in N \\ n_R \in N \\ x_L \leq x < x_R}}{\operatorname{argmin}} S^V(\mathbf{p}_L, \mathbf{p}_R, n_L, n_R, y).$$

Перейдемо до описання самого алгоритму.

**Вхід:** Матриця  $P'$ , зображення  $\mathbb{I}$ .

**Результат:** Зображення з коректним виразом, яке  
якомога більше схоже на  $\mathbb{I}$ .

$\mathbf{p}_L \leftarrow (1, 1)$

$\mathbf{p}_R \leftarrow (W, H)$

$s \leftarrow I$

$w \leftarrow P'(\mathbf{p}_L, \mathbf{p}_R, s)$

**for**  $s \in N$  **do**

$n_L, n_R, y \leftarrow \hat{H}(\mathbf{p}_L, \mathbf{p}_R, n)$

$n_U, n_D, x \leftarrow \hat{V}(\mathbf{p}_L, \mathbf{p}_R, n)$

$n_R \leftarrow \underset{t \in T}{\operatorname{argmin}} (P'(\mathbf{p}_L, \mathbf{p}_R, t) + \mathbb{1}(g(t, n) = 1) \cdot \infty)$

$D_0 \leftarrow \{(\mathbf{p}_L, \mathbf{p}_R, n_R)\}$

$D_x \leftarrow \{(\mathbf{p}_L, x, y_R, n_L), (\mathbf{p}_L, x + 1, y_R, n_R)\}$

$D_y \leftarrow \{(x_L, y, \mathbf{p}_R, n_U), (x_L, y, \mathbf{p}_R, n_D)\}$

$D \leftarrow D_0 \cup D_x \cup D_y$

$\mathbf{p}_L, \mathbf{p}_R, s \leftarrow \underset{(\mathbf{p}'_L, \mathbf{p}'_R, s') \in D}{\operatorname{argmax}} P'(\mathbf{p}'_L, \mathbf{p}'_R, s')$

**end for**

**if**  $\min_{t \in T} P'(\mathbf{p}_L, \mathbf{p}_R, t) \neq P'(\mathbf{p}_L, \mathbf{p}_R, s)$  **then**

Підзображення в  $(\mathbf{p}_L, \mathbf{p}_R)$  замінити на еталонне  
з  $\mathbb{E}$ , відповідне до  $s$ , на зображенні  $\mathbb{I}$ .

**end if**

**return**  $\mathbb{I}$

**3.3. Набір правил складання у стовпчик**

Для того, щоб можна було застосувати алгоритм  
Кока-Янгера-Касамі для аналізу зображення, треба  
побудувати граматику, з якою він буде працювати.  
Умовно правила можна розділити на три частини:  
розпізнавання цифр, розпізнавання стовпчиків та їх  
конкатенацію.

Знайти розташування цифр на зображенні недоста-  
точно. Через те, що вони «стрибають» на зображенні,  
їх треба конкатенувати з частинами фону

$$\begin{aligned} 1 &\rightarrow \varepsilon | 1 & \dots & 9 &\rightarrow \varepsilon | 9, \\ 1 &\rightarrow 1 | \varepsilon & \dots & 9 &\rightarrow 9 | \varepsilon, \\ 1 &\rightarrow \frac{1}{\varepsilon} & \dots & 9 &\rightarrow \frac{9}{\varepsilon}, \\ 1 &\rightarrow \frac{\varepsilon}{1} & \dots & 9 &\rightarrow \frac{\varepsilon}{9}. \end{aligned}$$

Опишемо правила для двох цифр, які стоять в  
одному стовпчику. Нам потрібно знати суму цих цифр

за модулем 10, а також знати, чи ця сума є більшою  
за 9

$$\begin{aligned} A_{(m+n) \bmod 10, (m+n) \operatorname{div} 10} &\rightarrow \frac{m}{n}, \quad m, n = \overline{0, 9}, \\ A_{(m+n) \bmod 10, (m+n) \operatorname{div} 10} &\rightarrow \frac{n}{m}, \quad m, n = \overline{0, 9}. \end{aligned}$$

Важливо нагадати, що зазвичай у складанні в  
стовпчик приймають участь два числа необов'язково  
однакової довжини. У таких випадках перед коро-  
тшим числом можна поставити стільки нулів, скіль-  
ки потрібно для компенсації його «недостатньої» дов-  
жини, але замість них на зображенні маємо фон

$$\begin{aligned} A_{m,0} &\rightarrow \frac{m}{\varepsilon}, \quad m = \overline{0, 9}, \\ A_{m,0} &\rightarrow \frac{\varepsilon}{m}, \quad m = \overline{0, 9}. \end{aligned}$$

Стовпчики можна поділяти за комбінацією двох  
ознак: породження зайвого десятку та поглинання  
зайвого десятку. Ці типи будемо позначати символом  
 $V_{i,j}$ ,  $i, j \in \{0, 1\}$ , де перший індекс дорівнює одиниці,  
якщо стовпчик поглинає десяток, та нулю у іншому  
випадку. Аналогічно, одиниця або нуль для другого  
індексу, якщо стовпчик генерує або не генерує деся-  
ток, відповідно. Опишемо ці правила формулами

$$\begin{aligned} V_{0,\ell} &\rightarrow \frac{A_{\ell,k}}{k}, \quad k = \overline{0, 9}, \\ V_{1,\ell} &\rightarrow \frac{A_{\ell,k}}{k+1}, \quad k = \overline{1, 9}, \\ V_{1,1} &\rightarrow \frac{A_{0,9}}{0}. \end{aligned}$$

Правила конкатенації самих стовпчиків вигляда-  
ють наступним чином:

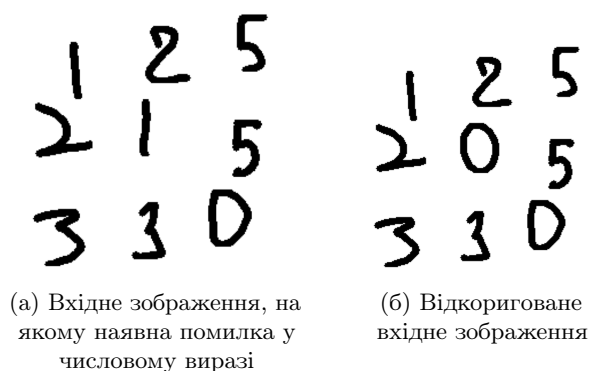
$$\begin{aligned} I &\rightarrow V_{0,0}, \\ I &\rightarrow V_{0,0} | V_{0,0}, \\ I &\rightarrow V_{1,0} | V_{0,1}, \\ I &\rightarrow I | V_{0,0}, \\ I &\rightarrow I_0 | V_{0,1}, \\ I_0 &\rightarrow V_{1,0}, \\ I_0 &\rightarrow I | V_{1,0}, \\ I_0 &\rightarrow I_0 | V_{1,1}. \end{aligned}$$

Символ  $I$  відповідає коректному виразу,  $I_0$  — конка-  
тенація стовпчиків, за якої поглинається десяток.

**4. Опис результатів**

У цьому підрозділі представлено результати експе-  
риментів, що було проведено за допомогою програми,  
яка реалізує описаний метод.

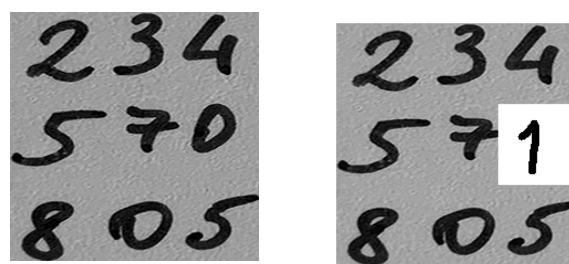
Спершу розглянемо тривіальний випадок, коли на  
вхідному зображенні представлено коректний вираз  
(рис. 1). Для символу  $I$  (коректний вираз) програма  
надає ступінь несхожості приблизно 0.05, що, вра-  
ховуючи похибку у розпізнаванні цифр від класи-  
фікатора, інтерпретується як коректне зображення.  
Після етапу корекції було отримано те ж саме зобра-  
ження, що підтвержує те, що на зображенні дійсно  
коректний вираз.



(а) Вхідне зображення, на якому наявна помилка у числовому виразі

(б) Відкориговане вхідне зображення

Рис. 3. Синтетичний приклад некоректного виразу



(а) Вхідне зображення

(б) Відкориговане зображення

Рис. 4. Реальне зображення некоректного виразу

Розглянемо випадок, коли на зображенні є одна помилка (рис. 3(а)). Цього разу символ  $I$  має ступінь несхожості біля 1.01, що вже вказує на наявність помилки. Після етапу корекції на зображенні була замінена одиниця, що знаходиться у центрі (рис. 3(б)), що є однією з можливих корекцій зображення.

На рис. 4(а) представлено реальне зображення з некоректним виразом, яке було також подано на вхід програмі. На виході було отримано відкориговане зображення (рис. 4(б)).

## Висновки

У даній роботі розглянуто сумісне використання згорткової нейромережі та методу структурного аналізу зображень. На прикладі зображень складання у стовпчик продемонстрована ефективність сукупності цих методів. Слід зазначити, що цей підхід не обмежується розпізнаванням таких зображень, і його можна використовувати для аналізу інших структурованих зображень, якщо відповідним чином замінити класифікатор і граматику. Сучасні обчислювальні можливості відкривають нові перспективи використання подібних методів.

## Перелік використаних джерел

1. Průša D., Hlavac V. 2D context-free grammars: Mathematical formulae recognition // Proceedings of the Prague Stringology Conference. — 2006. — 01. — P. 77–89.
2. ICDAR 2019 CROHME + TFD: Competition on Recognition of Handwritten Mathematical Expressions and Typeset Formula Detection / Mahdavi M., Zanibbi R., Mouchere H., Viard-Gaudin C., and Garain U. // 2019 International Conference on Document Analysis and Recognition (ICDAR). — 2019. — sep. — P. 1533–1538.
3. Ioffe S., Szegedy C. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. — 2015. — Vol. 37. — P. 448–456. — Access mode: <https://proceedings.mlr.press/v37/ioffe15.html>.
4. Шлезингер М. И., Главач В. Десять лекций по статистическому и структурному распознаванию образов. — Наукова думка, 2004. — 545 с.