

18. Parascript. Medical Imaging. AccuDetect. Available at: <http://www.parascript.com/medical-imaging/> (accessed 25 November 2013).

19. Sadykov S. S., Bulanova Y. A. Algorithm of localization of breast cancer in the background of mastopathy. *11th International Conference on Pattern Recognition and Image Analysis: New Information Technologies (PRLA-11-2013)*, 2013, vol. 2, pp. 717–721.

УДК 004.93'12

МЕТОД КОРРЕКЦИИ ОШИБОК КЛАССИФИКАЦИИ РАСПОЗНАННЫХ СИМВОЛОВ

Статья поступила в редакцию 30.01.2014, в окончательном варианте 16.02.2014.

Брейман Александр Давидович, кандидат технических наук, доцент, Национальный исследовательский университет «Высшая школа экономики», 101000, Российская Федерация, г. Москва, ул. Мясницкая, 20, e-mail: abreyman@hse.ru

Яковлев Илья Александрович, аспирант, Московский государственный университет приборостроения и информатики, 107996, Российская Федерация, г. Москва, ул. Стромынка, 20, e-mail: krofes@gmail.com

Процесс распознавания текстовых документов неизбежно связан с возникновением ошибок распознавания, для выявления и исправления которых используют методы пост-обработки, как правило, опирающиеся на словарный поиск. Использование словарей позволяет достичь приемлемого качества распознавания для латиницы, кириллицы и других фонетических алфавитов, однако мало-пригодно для языков, в которых выделение отдельных слов в письме нехарактерно или обязательно (китайский, японский, корейский, вьетнамский и прочие языки). В статье рассмотрены существующие методы, направленные на решение данной проблемы, а также описан новый подход к исправлению некоторых видов ошибок, основанный на применении ансамблей нейронных сетей (по нейронной сети на каждый возможный символ), позволяющий сократить количество ошибок в результате распознавания иероглифического письма, а для фонетических алфавитов – снизить зависимость от качества словарей.

Ключевые слова: оптическое распознавание символов, ошибки распознавания символов, пост-обработка ошибок распознавания, система верификации распознавания, система коррекции ошибок распознавания без словаря, распознавание иероглифов, нейронные сети, нейросетевые ансамбли

OPTICAL CHARACTER RECOGNITION ERRORS CORRECTION METHOD

Breyman Aleksandr D., Ph.D. (Engineering), Associate Professor, National Research University “Higher School of Economics”, 20 Myasnitskaya St., Moscow 101000, Russian Federation, e-mail: abreyman@hse.ru

Yakovlev Ilya A., post-graduate student, Moscow State University of Instrument Engineering and Computer Science, 20 Stromynka St., Moscow, 107996, Russian Federation, e-mail: krofes@gmail.com

Optical recognition of text documents is inevitably error-prone process. To identify and correct that errors systems use post-processing techniques that are usually based on dictionary search. Using dictionaries can bring an acceptable quality of recognition for Latin, Cyrillic and other phonetic alphabets, but of little use for the languages in which the selection of individual words is untypical or optional (Chinese, Japanese, Korean, Vietnamese and other languages). This paper discusses known methods to address this problem, and proposes a new approach to correcting certain types of errors, based on the application of neural networks ensembles (containing distinct neural network for each possible character), which allows to reduce the num-

ber of hieroglyphic recognition errors and to reduce dependence on the quality of dictionaries while recognizing texts in phonetic alphabets.

Keywords: optical character recognition, character recognition errors, post-processing of recognition errors, the verification system for recognition results, dictionary-less recognition error correction system, hieroglyph recognition, neural networks, neural network ensembles

Введение. Оптическое распознавание символов (ОРС) представляет собой процесс преобразования растровых изображений машинно-печатного или печатного текста в цифровой формат, доступный для редактирования на ЭВМ, содержащий распознанный текст и, возможно, его разметку.

В ходе распознавания могут происходить ошибки, приводящие к искажению его результатов. Большинство используемых методов пост-обработки опираются на словарный поиск, достигая приемлемых результатов для фонетических алфавитов, но не для иероглифического идеографического либо морфемного письма (например, для современного китайского, корейского, японского языков). В то же время важность международного информационного обмена со странами Восточной Азии постоянно возрастает, и повышение качества распознавания символов иероглифического письма может способствовать его интенсификации. Поэтому *целью данной статьи* является аналитический обзор ряда существующих методов для исправления ошибок распознавания символов, а также представление нового метода коррекции некоторых ошибок, ориентированного на иероглифическое письмо.

Методы повышения качества распознавания. Идея повышения качества распознавания текста с помощью различных манипуляций с данными возникла достаточно давно. Условно все методы этой группы можно разделить на два типа: методы, определяющие свойства исходного изображения с текстом и выполняющие его коррекцию для повышения точности распознавания; методы, работающие с полученным результатом распознавания, т.е. выполняющие пост-обработку полученного текста с исправлением ошибок, уже допущенных системой распознавания.

К первой группе можно отнести метод, описанный в работе [2]. Ее авторы для повышения результатов качества распознавания оцифрованных документов машинно-печатного текста проводят анализ исходного документа (изображения), оценивая характеристики изображения по пяти метрикам, нормированным в [0; 1].

Small Speckle Factor – мера, оценивающая количество черных пятен на документе.

White Speckle Factor – мера, оценивающая количество белых «частиц» с замкнутым черным контуром на изображении.

Touching Character Factor – мера, оценивающая близость расположения символов на изображении относительно друг друга.

Broken Character Factor – мера, оценивающая степень того, как часто в тексте встречаются символы с разорванным контуром.

Font Size Factor – мера, оценивающая размер шрифта на изображении.

После проведения анализа документа, в зависимости от значений полученных параметров, выбирается и выполняется один из 14 алгоритмов преобразования изображения. Результат работы выбранного алгоритма передается на вход системы распознавания.

В [2] использовано два способа оценки качества распознавания – посимвольный и пословный. Для тестирования использовался корпус (тестовый материал) из 41 документа Министерства энергетики США, отсканированных с разрешением 300 dpi с бумажных оригиналов, полученных с телетайпа, а также напечатанных на печатной машинке, в том числе под копиру, и распознанных с использованием Caere OmniPage Pro v8.0 со средней точностью распознавания отдельных символов 65,72 % (для каждого документа точность распо-

знавания символов, т.е. доля верно распознанных символов в документе, находилась в интервале от 50 до 80 %) и средней точностью распознавания слов 49,01 %. В результате предобработки изображений было достигнуто повышение точности распознавания символов на 24 %, а слов – на 30 %.

Метод, описанный в работе [9], также относится к первой группе. Он основан на алгоритме восстановления слабых сигналов FastICA, до этого не применявшемся к анализу и восстановлению изображений.

Ко второй группе (применение пост-обработки результатов распознавания) можно отнести методы, развиваемые в работах [4, 7, 8], основанные на лингвистическом анализе результирующего текста. В случае отсутствия распознанного слова в словаре языка выполняется его коррекция. Схематично работа таких методов представлена на рис. 1.

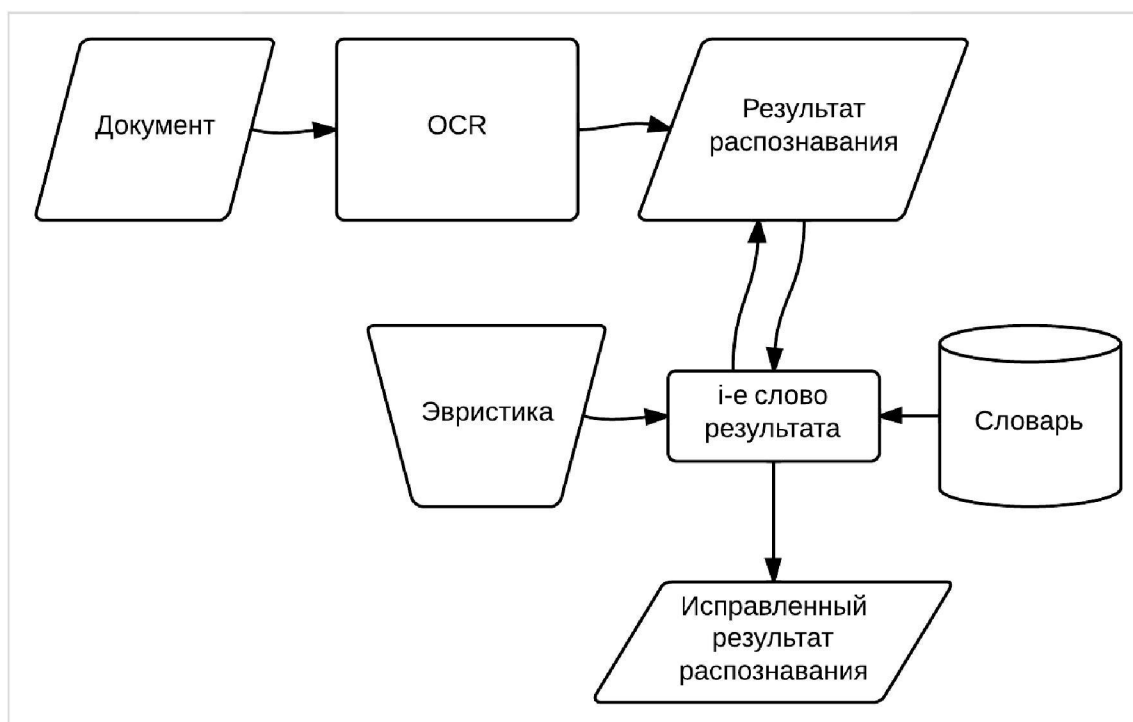


Рис. 1. Общий принцип работы систем коррекции ошибок результатов распознавания

Самой известной из систем, реализующих методы второй группы, является система OCRSpell [8], содержащая точные и эвристические алгоритмы, которые выполняют коррекцию распознанного символа, используя многокритериальную оценку, с дополнительной возможностью исправления текста пользователем. Если слово не найдено в словаре языка, для него создается два списка кандидатов на замену.

Первый список создается на основе статического анализа языка текста и генерации на основе такого анализа возможных вариантов исходного слова. Для каждого сгенерированного слова рассчитывается байесовская функция ранжирования, определяющая вероятность соответствия распознанного слова варианту исходного.

Второй список кандидатов создается на основе команд графического интерфейса пользователя, позволяющего указать ошибки слияния и разбиения символов в распознанном тексте. Указанные пользователем ошибки также используются для генерации возможных исходных слов.

В случаях, когда статическая и динамическая генерация слов не дают нужных результатов, для распознанного слова применяется эвристический алгоритм, пытающийся найти наиболее похожее в словаре слово – с попыткой привести распознанное слово в нормальную форму (например, форму единственного числа, именительного падежа для существительных).

К полученному после применения алгоритмов слову «возвращается» предполагаемое число и форма распознанного слова, и результат передается оператору.

Хотя подходы, относящиеся к первой группе, универсальны и применимы для любых документов, они демонстрируют лучшие результаты лишь при коррекции искажений определенного типа (например, клякс, или «просветов» символов с оборотной части листа).

Недостатком подходов второй группы является необходимость выделения отдельных слов результата распознавания (токенов) и применения словарей для проверки корректности результатов. Выделение токенов в тексте, к примеру, невозможно в ряде иероглифических языков, таких как японский, китайский, корейский и др., где использование пробелов для отделения слов либо не обязательно, либо вовсе не используется.

Метод коррекции ошибок с применением нейронных сетей. Распространенный подход к решению задач, связанных с распознаванием символов при помощи нейронных сетей [3], – построение одной нейронной сети (НС), вектор входных данных которой представляет преобразованное в нужный формат изображение символа, а нейроны выходного слоя служат для отнесения исходного изображения к определенному классу.

Решение задачи при таком подходе сводится к синтезу и обучению НС, выполняющей задачу классификации лучше, чем это делает система распознавания, ошибки которой необходимо корректировать. При этом работа модуля распознавания соответствует поиску областей с изображениями символов на документе, а задача классификации символов осуществляется НС.

Очевидным недостатком такого подхода является неизбежный рост сложности решаемой задачи при увеличении количества классов, для различения которых используется НС, а также необходимость в полном переобучении всей сети при добавлении новых обучающих примеров или при расширении количества символов, которые необходимо распознавать с помощью сети.

Процедура синтеза НС зависит от размера и типа обучающей выборки. Топология сети подбирается в ходе обучения на основе вычисления ошибки обучения сети на множестве обучающих примеров. Для останова процесса обучения используется ограничение, накладываемое на количество эпох обучения сети. Для выбора слоя, на котором происходит приращение количества нейронов, используется мера количества связей, возникающих при добавлении очередного нейрона. Для обучения НС в зависимости от размеров обучающей выборки используются алгоритмы обратного распространения ошибки [6] и эластичного распространения ошибки [5].

Вместо синтеза НС, которая занималась бы верификацией и коррекцией любого предъявляемого символа, авторами настоящей статьи предлагается формировать нейросетевой ансамбль – множество НС, каждая из которых выполняет верификацию и коррекцию одного конкретного символа некоторого алфавита, а первоначальную гипотезу о принадлежности символа определенному классу выдвигает система распознавания.

Такой подход упрощает процедуру обучения, так как каждая НС, входящая в ансамбль, должна различать всего два класса изображений: первый класс соответствует символу алфавита, за который отвечает конкретная сеть, второй класс не соответствует этому символу (на вход сети подан образ, содержащий символ, не соответствующий символу, за который «отвечает» сеть). Такая архитектура позволяет легко добавлять новые классы, которые требуется различать (каждому новому классу соответствует новая сеть). Кроме того,

при изменении обучающей выборки не требуется перечислять все сети, а достаточно перечислить только ту сеть, в обучающей выборке для которой произошли изменения.

Для коррекции результата распознавания из ансамбля выбирается НС с максимальным выходным значением (если НС, продемонстрировавших одинаковое максимальное значение, несколько, выбор одной из них осуществляется случайным образом). Тот символ алфавита, за распознавание которого «отвечает» выбранная НС, принимается за результат коррекции.

Помимо нейросетевого ансамбля, система коррекции ошибок содержит модуль формирования данных для работы сети, модуль обучения НС и модуль, управляющий процессом выполнения коррекции ошибки классификации.

Как видно из графического представления метода (рис. 2), нейросетевой ансамбль представляет собой дополнительный модуль, преобразующий результаты распознавания системы ОРС в итоговые результаты.

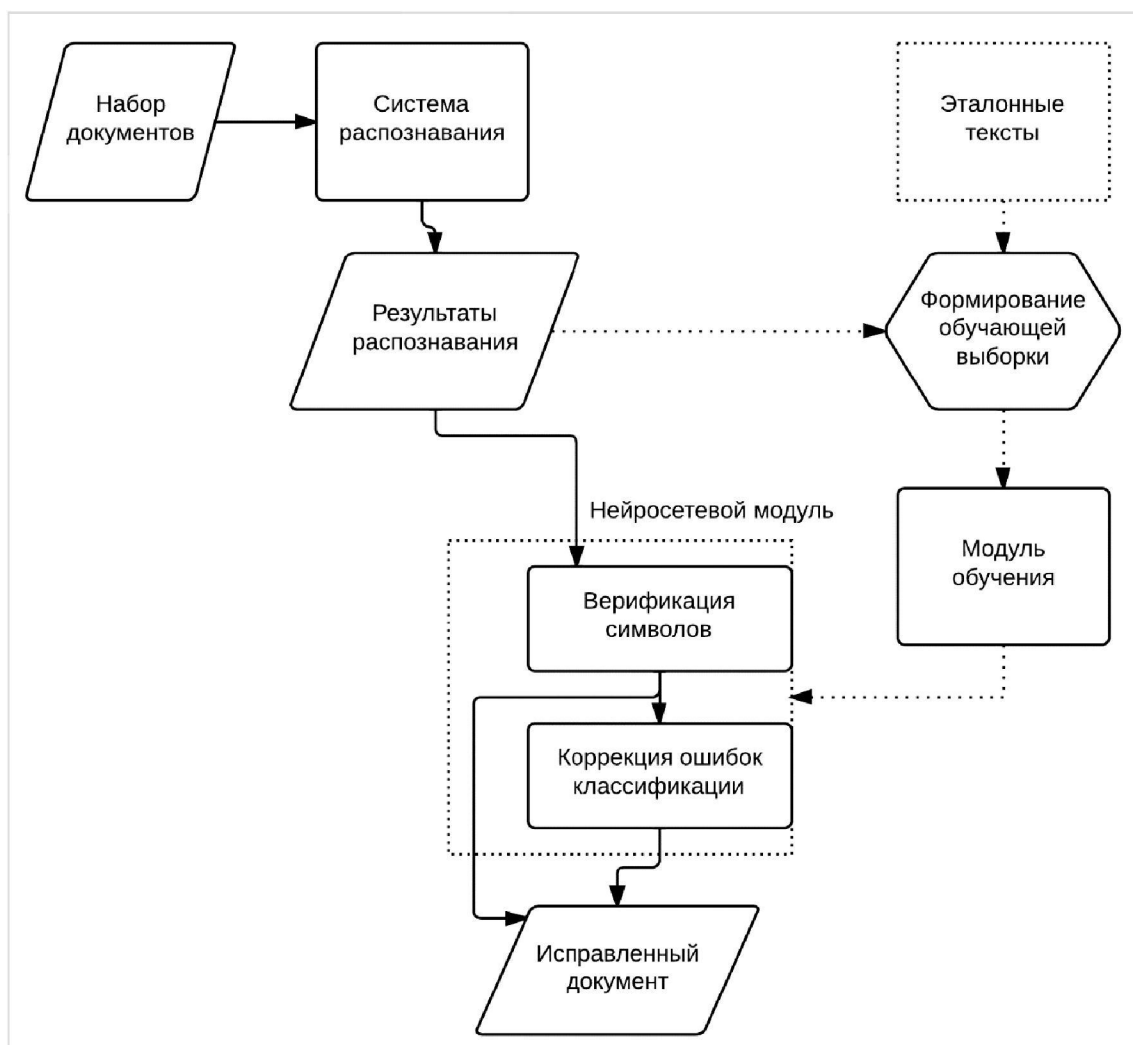


Рис. 2. Метод поиска и коррекции ошибок классификации

К достоинствам предложенного метода можно отнести следующее:

- независимость от словарей естественного языка (метод применим для коррекции результатов распознавания любого естественного языка);

- сокращение времени обучения и времени работы в режиме поиска и коррекции ошибок за счет использования ансамбля НС вместо одной большой НС;
- простота в дообучении нейросетевого ансамбля при расширении алфавита языка.

Тем не менее, необходимость в начальном обучении нейросетевого ансамбля требует предъявления заранее подготовленного и выверенного набора пар (Изображение; Эталонный текст), причем качество результата зависит от полноты предъявленного обучающего набора данных. Кроме того, вычислительная сложность метода прямо пропорциональна количеству различаемых классов (символов). Эта проблема может быть частично решена использованием эвристического метода снижения количества проверок символа на этапе коррекции путем учета статистики ошибок классификации символов.

Результаты экспериментов. Система коррекции ошибок распознавания (СКОР) была реализована авторами на языке C#3.0 в виде программной библиотеки. Выполненная разработка содержит модуль верификации и коррекции ошибок классификации, обучающий модуль и модуль получения данных. В качестве внешней системы распознавания в экспериментах использовалась ABBYY Mobile OCR Engine [1].

Оценка качества работы СКОР состоит из двух частей. В первой части оценивается качество верификации распознанных символов VQ. Значение VQ лежит в интервале [0; 1], где «0» означает всегда неверную верификацию, а «1» – всегда верную. Для её оценки выделяется четыре состояния, в которые может попасть символ после процедуры верификации (табл. 1). Верификация называется положительной, если она оценивает результат распознавания как верный, и негативной – в противном случае.

Таблица 1

Возможные состояния символа после распознавания и верификации

| | Верное распознавание | Неверное распознавание |
|---------------------------|----------------------|------------------------|
| Положительная верификация | VSC | VSU |
| Негативная верификация | VFU | VFC |

Тогда

$$VQ = 1 - \sum_{i=0}^n (|VSU_i| + |VFU_i|) / \sum_{i=0}^n |\{Ch_i | Ch_i \in Alphabet\}|, \quad (1)$$

где Alphabet – множество символов распознаваемого языка; Ch_i – множество символов, полученных в ходе распознавания i-го изображения; VSU_i – множество ошибочно положительно верифицированных неверно распознанных символов i-го изображения (ошибки первого рода); VFU_i – множество ошибочно негативно верифицированных верно распознанных символов i-го изображения (ошибки второго рода).

Следующим этапом является вычисление оценки качества работы модуля коррекции CQ. Поскольку коррекция выполняется только после негативной верификации, то состояния VSC и VSU не влияют на CQ. Для его оценки также используется четыре метрики (табл. 2).

Таблица 2

Возможные состояния символа после негативной верификации

| | Верное распознавание | Неверное распознавание |
|----------------------|----------------------|------------------------|
| Успешная коррекция | CSVF | CSVS |
| Неуспешная коррекция | CUVF | CUVS |

В этом случае:

$$CQ = 1 - \sum_{i=0}^n (|CUVS_i| + |CUVF_i|) / \sum_{i=0}^n (|VFC_i| + |VFU_i|), \quad (2)$$

где $CUVS_i$ – множество символов, верно верифицированных как ошибочно распознанные, но которые не удалось исправить на верный символ; $CUVF_i$ – множество символов i -го изображения, ошибочно верифицированных как неправильно распознанные, и которые были исправлены на другой, заведомо неправильный символ; VFC_i – символы i -го изображения, верно помеченные модулем верификации как ошибочно распознанные.

Параметры набора изображений, использованного для обучения нейросетевого ансамбля, приведены в табл. 3.

Таблица 3

Параметры обучающего набора изображений

| Описание параметра | Параметр |
|------------------------------------|---|
| Количество изображений | 10 |
| Язык текста на изображениях | Корейский (хангыль, ханчча) |
| Общая длина эталонного текста | 8063 символа |
| Размер алфавита | 726 символов |
| Источник изображений | Сканирующее устройство |
| Свойства представления текста | Простой документ |
| Свойства шрифта текста | Стандартное начертание Начертание курсивом (8 страница) Шрифт MyeongJo-Medium |
| Цветовые свойства исходного текста | Простой |
| Наличие искажений на изображении | Нет |

Оценка качества работы модулей верификации и коррекции выполнялась на тестовых данных, описанных в табл. 4.

После обучения нейросетевого ансамбля была проведена серия экспериментов на тестовых данных. Результаты экспериментов представлены в табл. 5, а результаты коррекции ошибок для набора тестовых данных представлены в табл. 6.

Таблица 4

Описание набора тестовых данных

| Описание параметра | Параметр |
|------------------------------------|---|
| Количество изображений | 30 |
| Язык текста на изображениях | Корейский |
| Общая длина эталонного текста | 18 761 символ, входящий в алфавит; 873 символа отфильтрованы, так как не входили в алфавит |
| Источник изображений | Сканирующее устройство (1–18 изображение) Фотоаппарат Canon 650D (19–22 изображение) Камера смартфона Samsung Galaxy Note (23–30 изображение) |
| Свойства представления текста | Простой документ (1–22) Документ со сложной версткой (22–30) |
| Свойства шрифта текста | Стандартное начертание (1–25, 27–30) Начертание курсивом (26) Шрифт MyeongJo-Medium (1–22) Шрифт GoThic-Medium (23–25, 27–30) Шрифт Microsoft Batang (26) |
| Цветовые свойства исходного текста | Простой (1–22, 26–30) Сложный (23–25) |
| Наличие искажений на изображении | Нет (1–22, 30) Блики от вспышки (25) Размытие от движения (22–24, 26–29) |

Таблица 5

Основные показатели работы системы поиска и коррекции ошибок

| Пакет | Положительно верифицировано | | Исправлено | | Не исправлено | |
|---------|-----------------------------|----------|------------|----------|---------------|----------|
| | VCS _{средн.} | % от Екл | CSVS | % от Екл | CUVS | % от Екл |
| Скан. | 37,39 | 67,17 % | 30,17 | 54,19 % | 7,22 | 12,97 % |
| Фото | 67,25 | 77,75 % | 51,75 | 59,83 % | 15,50 | 17,92 % |
| Телефон | 85,75 | 71,61 % | 51,88 | 43,32 % | 33,88 | 28,29 % |
| Общий | 54,27 | 70,6 % | 38,83 | 50,52 % | 15,43 | 20,08 % |
| | Негативно верифицировано | | Исправлено | | Не исправлено | |
| | VFU _{средн.} | % от VCS | CUVF | % от VFU | CSVF | % от VFU |
| Скан. | 18,11 | 48,44 % | 6,06 | 33,4 % | 12,06 | 66,56 % |
| Фото | 18,25 | 27,14 % | 5,25 | 28,8 % | 13,00 | 71,23 % |
| Телефон | 79,25 | 92,42 % | 68,50 | 86,4 % | 10,75 | 13,56 % |
| Общий | 34,43 | 63,45 % | 22,60 | 65,6 % | 11,83 | 34,37 % |

Таблица 6

Результаты работы обученной системы

| Пакет | Кол-во изображений | Кол-во символов | Верно классифицировано | Ошибок классификации | Ошибок после коррекции | Δ , % | VQ | CQ |
|---------|--------------------|-----------------|------------------------|----------------------|------------------------|--------------|------|------|
| Скан. | 18 | 427,94 | 372,28 | 55,67 | 31,56 | -43,31 % | 0,91 | 0,76 |
| Фото | 4 | 423,75 | 337,25 | 86,50 | 40,00 | -53,76 % | 0,91 | 0,76 |
| Телефон | 8 | 325,38 | 205,63 | 119,75 | 136,38 | 13,88 % | 0,65 | 0,38 |
| Все | 30 | 400,03 | 323,17 | 76,87 | 60,63 | -21,12 % | 0,86 | 0,57 |

Как видно из табл. 6, в среднем для всего набора тестовых изображений произошло снижение количества ошибок классификации на 21 %. Однако, если обратить внимание на увеличение количества ошибок классификации для изображений, снятых на камеру мобильного устройства, и в то же время близкое к 50 % снижение количество ошибок классификации для остальных изображений, то можно сделать предположение о том, что влияние на качество коррекции оказывают исходные свойства изображений. На рис. 3 для сравнения представлены три изображения текста, полученные с помощью сканирующего устройства (верхние иероглифы), фотокамеры Canon 650D (средние) и фотокамеры смартфона Samsung Galaxy Note (нижние иероглифы).



Рис. 3. Сравнение строк текста, полученных с помощью сканирующего устройства, фотокамеры и смартфона

Как видно из рис. 3, изображения, полученные с камеры смартфона, ввиду низкого качества приобретают «рваную» структуру, что сказывается на результатах как распознавания, так и коррекции ошибок.

На рис. 4 представлены основные показатели работы системы на исходном пакете тестовых данных, охарактеризованных в табл. 4.

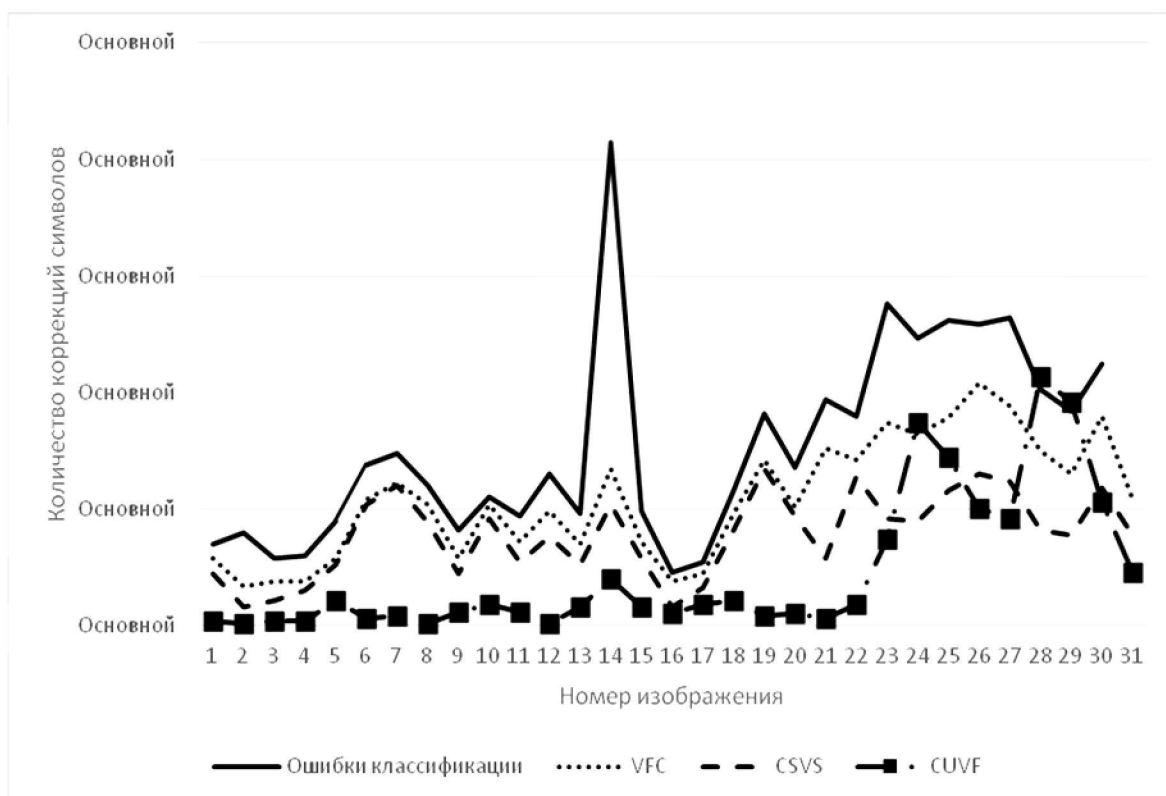


Рис. 4. Основные показатели работы системы коррекции ошибок на тестовом наборе данных

Как видно из рис. 4, для изображений, полученных со сканирующего устройства (1–18) и с фотокамеры (19–22), характерны значения CUVF, близкие к нулю – это отражает высокое качество распознавания символов. В то же время для изображений, полученных с камеры смартфона (23–30), характерны значения CUVF, составляющие десятки единиц, что связано с низким качеством изображений.

Поэтому был проведен эксперимент, в котором в обучающую выборку были добавлены изображения, полученные с помощью камеры смартфона. В табл. 7 и на рис. 5 представлены результаты работы системы, полученные при добавлении к обучающей выборке нейросетевого ансамбля шести изображений, сделанных с помощью камеры смартфона Samsung Galaxy Note. Для получения изображений были взяты те же исходные текстовые документы, которые были использованы (отсканированы) для составления первоначальной обучающей выборки. Нейросетевой ансамбль при этом работал в режиме дообучения, т.е. использовались уже созданные НС, к которым применялась процедура обучения на дополненных данных.

Таблица 7

Результаты работы системы после дообучения нейросетевого ансамбля

| Пакет | Кол-во изображений | Кол-во символов | Верно классифицировано | Ошибок классификации | Ошибок после коррекции | Δ , % | VQ | CQ |
|---------|--------------------|-----------------|------------------------|----------------------|------------------------|--------------|------|------|
| Скан. | 18 | 427,94 | 372,28 | 55,67 | 38,94 | -30,04 % | 0,90 | 0,66 |
| Фото | 4 | 423,75 | 337,25 | 86,50 | 51,00 | -41,04 % | 0,88 | 0,68 |
| Телефон | 8 | 325,38 | 205,63 | 119,75 | 52,38 | -56,26 % | 0,86 | 0,72 |
| Все | 30 | 400,03 | 323,17 | 76,87 | 44,13 | -42,58 % | 0,89 | 0,69 |

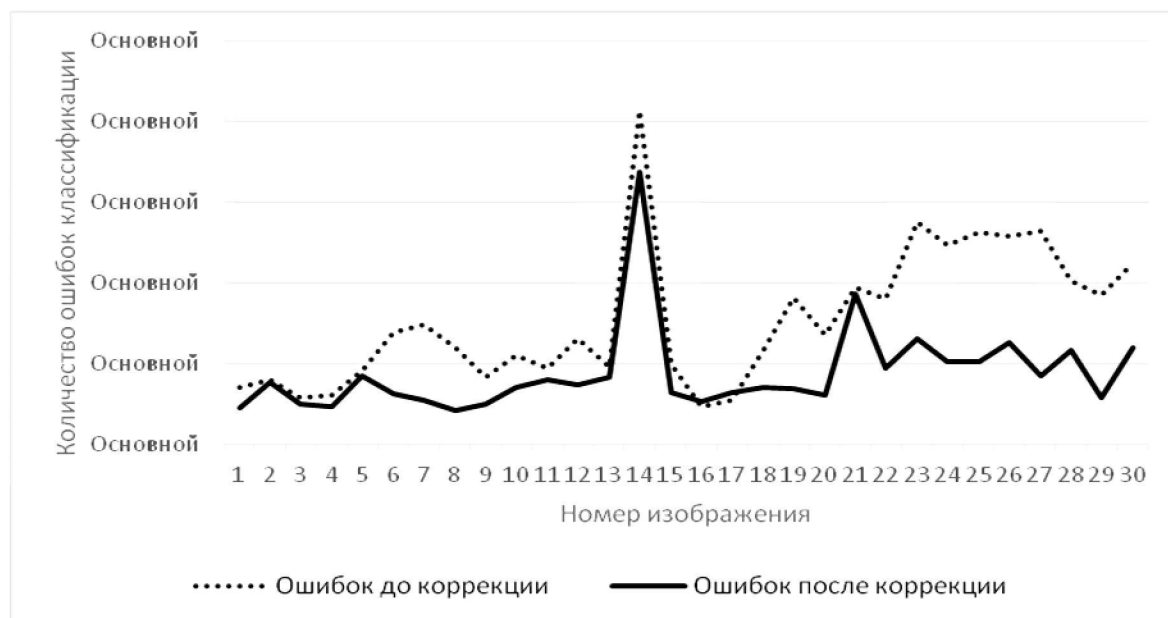


Рис. 5. Количество ошибок до и после коррекции нейросетевого ансамбля

Как видно из табл. 7 и рис. 5, при добавлении в обучающую выборку изображений, полученных с помощью камеры смартфона, на части тестовой выборки, составленной из фотографий со смартфона, произошло уменьшение итогового количества ошибок классификации в среднем на 20 %. В то же время качество работы системы на отсканированных и сфотографированных документах упало на 13,3 % и 12,7 % соответственно. Возможным путем решения этой проблемы может быть наращивание нейросетевого ансамбля с выделением отдельных сетей для распознавания каждого символа алфавита, полученного из источников изображений с разными характеристиками.

Заключение. Полученные в работе результаты позволяют сделать вывод об эффективности предложенного метода коррекции ошибок результата распознавания ОРС при помощи нейросетевого аппарата. Достоинством такого подхода является работа с отдельными распознанными символами, а не с отдельными словами результата распознавания. Это позволяет использовать данный метод при коррекции результатов распознавания языков, в которых токенизация текста отсутствует или необязательна.

Очевидные ограничения такого подхода по сравнению с подходами, не использующими НС, связаны с необходимостью обучения нейросетевого ансамбля.

Список литературы

1. ABBYY Mobile OCR Engine. – Available at: <http://www.abby.ru/mobile-ocr-engine> (accessed 06.11.2013).
2. Cannon M. Quality Assessment and Restoration of Typewritten Document Images / M. Cannon, J. Hochberg, P. Kelly // *International Journal on Document Analysis and Recognition*. – 1999. – Vol. 2. – P. 80–89.
3. Ni D. X. Application of Neural Networks to Character Recognition / D. X. Ni // *Proceedings of Students/Faculty Research Day, CSIS, Pace University*. – 2007.
4. Reynaert M. Non-Interactive OCR Post-Correction for Giga-Scale Digitization Projects / M. Reynaert // *Proceedings of the 9th International Conference on Computational Linguistics and Intelligent Text Processing*. – 2008. – P. 617–630.
5. Riedmiller M. A direct adaptive method for faster backpropagation learning: the RPROP algorithm / M. Riedmiller, H. Braun // *Proceedings of the International Conference on Neural Networks*. – IEEE Press, 1993. – P. 586–591.
6. Rumelhart D. E. Learning representations by back-propagating errors / D. E. Rumelhart, G. E. Hinton, R. J. Williams // *Nature*. – 1986. – № 323 (6088), 8 October. – P. 533–536.
7. Strohmaier C. A visual and interactive tool for optimizing lexical postcorrection of OCR results / C. Strohmaier, C. Ringlstetter, K. U. Schulz, S. Mihov // *Proceedings of the IEEE Workshop of Document Image Analysis and Recognition, DIAR'03*. – 2003.
8. Taghva K. OCRSpell: An Interactive Spelling Correction System for OCR Errors in Text / K. Taghva, E. Stofsky // *International Journal on Document Analysis and Recognition*. – 2001. – Vol. 3. – P. 125–137.
9. Tonazzini A., Bedini L., Salerno E. Independent component analysis for document restoration / A. Tonazzini, L. Bedini, E. Salerno // *International Journal on Document Analysis and Recognition*. – 2004. – Vol. 7, issue 1. – P. 17–27.

References

1. ABBYY Mobile OCR Engine. Available at: <http://www.abby.ru/mobile-ocr-engine> (accessed 6 November 2013).
2. Cannon M., Hochberg J., Kelly P. Quality Assessment and Restoration of Typewritten Document Images. *International Journal on Document Analysis and Recognition*, 1999, vol. 2, pp. 80–89.
3. Ni D. X. Application of Neural Networks to Character Recognition. *Proceedings of Students/Faculty Research Day, CSIS, Pace University*, 2007.
4. Reynaert M. Non-Interactive OCR Post-Correction for Giga-Scale Digitization Projects. *Proceedings of the 9th International Conference on Computational Linguistics and Intelligent Text Processing*, 2008, pp. 617–630.
5. Riedmiller M., Braun H. A direct adaptive method for faster backpropagation learning: the RPROP algorithm. *Proceedings of the International Conference on Neural Networks*. IEEE Press, 1993, pp. 586–591.
6. Rumelhart D. E., Hinton G. E., Williams R. J. Learning representations by back-propagating errors. *Nature*, 1986, no. 323 (6088), 8 October, pp. 533–536.
7. Strohmaier C., Ringlstetter C., Schulz K. U., Mihov S. A visual and interactive tool for optimizing lexical postcorrection of OCR results. *Proceedings of the IEEE Workshop of Document Image Analysis and Recognition, DIAR'03*, 2003.
8. Taghva K., Stofsky E. OCRSpell: An Interactive Spelling Correction System for OCR Errors in Text. *International Journal on Document Analysis and Recognition*, 2001, vol. 3, pp. 125–137.
9. Tonazzini A., Bedini L., Salerno E. Independent component analysis for document restoration. *International Journal on Document Analysis and Recognition*, 2004, vol. 7, issue 1, pp. 17–27.