

Jak stroić OCR w programie eSZOK

i uzyskać wyższą skuteczność rozpoznawania dokumentów przez program ?



Dokument jest przeznaczony wyłącznie dla osób, które mają przynajmniej podstawową wiedzę z zakresu obsługi oprogramowania eSZOK, którego producentem jest Centrum Technologii Informatycznej.

Dokument przygotował : mgr inż. Zygmunt Wilder – specjalista z zakresu wykorzystania sztucznej inteligencji w procesach zachodzących w księgowości zarówno biur rachunkowych jak i w firmach.

Data utworzenia : 2023.02.10

Wprowadzenie.

Oprogramowanie eSZOK jest zbiorem kilku programów, które wspierają polskie biura rachunkowe we wprowadzaniu danych księgowych do systemu księgowego ERP OPTIMA firmy Comarch SA.

Całe oprogramowanie zostało wymyślone w 2017 roku przez Zygmunta Wilder i napisane przez specjalistów z Centrum Technologii Informatycznej. Jest ono ciągle rozwijane pod potrzeby klientów i zmiany przepisów podatkowych.

Jednym z ważnych elementów oprogramowania jest funkcjonalność pozwalająca na automatyzację wprowadzania dokumentów kosztowych bezpośrednio ze skanów faktur. Polega to na tym, że dostarczając do programu skan np.: faktury VAT za paliwo system samodzielnie odczytuje dane z takiego skanu i wprowadza je do programu księgowego. Nie ma tu potrzeby ręcznego przepisywania danych z faktury do księgowości.

Oczywiście nie wyeliminujemy tu pracy księgowych. Taki dokument trzeba sprawdzić, czy został on prawidłowo rozpoznany przez oprogramowanie OCR. OCR to skrót od Optical Character Recognition, czyli Technologii Rozpoznawania Znaków. Jest to proces, w którym specjalne oprogramowanie analizuje obraz lub skan dokumentu i przekształca znaki w nim zawarte na tekst, który może być edytowany i przeszukiwany. OCR jest często używany do digitalizacji dokumentów papierowych, takich jak : książki, artykuły, faktury i inne, aby uczynić je bardziej dostępnymi i łatwiejszymi do zarządzania. Technologia ta jest szczególnie przydatna w przypadku dokumentów o dużej objętości lub archiwizowania zasobów cyfrowych.

Po rozpoznaniu obrazu kolejnym etapem jest dopasowanie informacji do tego co oczekuje księgowy. W przypadku oprogramowania eSZOK jest to odpowiednie przypisanie np.: numeru faktur, daty wystawienia faktury, kwot z faktury itp. To kolejny etap uzupełniania/dopasowywania informacji do tego, co oczekuje użytkownik programu. Za to odpowiedzialny jest system sztucznej inteligencji.

System sztucznej inteligencji AI (artificial intelligence) to komputerowy system, który jest zaprojektowany do wykonywania zadań takich jak : rozpoznawanie mowy, rozumienie języka naturalnego, uczenie się, rozwiązywanie problemów i wiele innych. AI opiera się na algorytmach i modelach matematycznych, które pozwalają mu na uczenie się i doskonalenie swoich działań na podstawie analizy danych. W zależności od rodzaju i złożoności systemu AI, jego zastosowanie może być bardzo szerokie i obejmować wiele różnych dziedzin, takich jak medycyna, finanse, transport, marketing i wiele innych.

W ostatnich latach sztuczna inteligencja stała się jednym z najszybciej rozwijających się obszarów technologii i ma znaczący wpływ na wiele aspektów naszego życia. Niemniej jednak, istnieją też poważne obawy dotyczące jej potencjalnego wpływu na rynek pracy, bezpieczeństwo i prywatność, a także etyczne i filozoficzne aspekty jej rozwoju i zastosowania.

Każdy system OCR dla księgowości wymaga jednak strojenia (dopasowania do potrzeb użytkownika). Wiele firm, które tworzą tego typu rozwiązania wykonuje to za pośrednictwem zespołów wielu osób, które nad tym pracują, jednocześnie oczekując za to odpowiedniej opłaty. System eSZOK daje możliwość samodzielnego strojenia, dopasowanego do potrzeb konkretnego biura rachunkowego lub firmy.

System eSZOK uwzględnia rozwiązania AI do wspierania pracy działów księgowości i wykorzystuje aż trzy poziomy sztucznej inteligencji :

- 1.0. Kodowanie maszynowe, które polega na tym że programista lub użytkownik określił algorytmy (instrukcje) podpowiadające co program ma rozpoznawać i jak to interpretować.
- 2.0. Uczenie maszynowe (uczenie z nauczycielem), gdzie AI na podstawie wskazanych przez użytkownika informacji zapamiętuje je i w przyszłości podpowiada odpowiednią informację.
- 3.0. Uczenie ze wzmocnieniem, które polega na tym, że system AI dostaje informacje o nieprawidłowości i samodzielnie stosuje inne rozwiązanie, prowadząc do prawidłowego wyniku poprzez zmianę podejmowania decyzji.
- 4.0. Jest to obecnie najwyższy poziom AI i tutaj systemy sztucznej inteligencji współpracują z innymi systemami AI. To rozwiązanie nie jest stosowane w oprogramowaniu eSZOK ze względu na ograniczenia obowiązującego obecnie w Polsce prawa np.: RODO.

Rozdział 0.

Jak zacząć ?

Jeśli chcesz rozpocząć strojenie oprogramowania eSZOK musisz zacząć od zrozumienia kiedy to strojenie ma sens i jak to zweryfikować.

Tak jak wspomniałem na początku potoczne określenie OCR w oprogramowaniu eSZOK składa się z dwóch etapów :

1. Rozpoznanie obrazu - jest to faktyczny proces OCR.
2. Dopasowanie treści do oczekiwań działu księgowego - jest to zastosowanie AI.

Jeśli zeskanowany dokument nie jest czytelny ze względu na to, że jest mały kontrast tekstu, jest pokreślony, pomięty, napisany odręcznie, do góry nogami itp. to uczenie AI nie ma sensu, lub jest niemożliwe. Proponuję zapoznać się z innym moim dokumentem, jakim jest : „Poradnik korzystania z aplikacji OCR by CTI”.


Żeby to sprawdzić należy po wejściu do programu Kancelaria wyświetlić szczegóły interesującego nas dokumentu naciskając na liście dokumentów przycisk VAT, lub ED i przejść następnie na zakładkę „Treść OCR”.

Nauka ta polega na zastosowaniu odpowiedniego REGEXU. Wyrażenia regularne (regular expressions, regex) są używane przez wiele programów wyszukujących określone ciągi znaków w tekstach. Technicznie wyrażenie regularne jest zbiorem zasad, które ciąg znaków powinien spełnić. Jeśli ciąg znaków spełnia podane zasady, mówimy, że pasuje do wyrażenia regularnego (ang. to match).

Funkcjonalność strojenia pozwala na zapisanie ręczne wzoru REGEX i jeśli wzór ten występuje na dokumencie to wprowadzenie odpowiedniej informacji do formularza z danymi rozpoznanymi przez OCR.

Przykładów zastosowania takiego rozwiązania jest wiele. Natomiast potrzeba prawidłowego rozpoznawania numeru faktur pokazała jak ułatwia to pracę użytkownika programu. To właśnie skłoniło nas do realizacji tego pomysłu i napisania takiego rozwiązania. Często zdarzało się że oprogramowanie rozpoznające obraz zamiast podawać w numerze faktury znaki „/” podpowiadało coś innego np.: „1”, „l”, „L”, „\” itp.

Przykładowo numer „FK/165/12/2022/PO” został przeczytany jako „FK1165/12120221PO”. Jak widać na tym przykładzie znak „/” program błędnie zinterpretował jako „1”. W tym konkretnym przypadku wszelkie inne zaawansowane metody uczenia nie są skuteczne i należy zastosować ręczne dostrojenie systemu.

SCHIESSL POLSKA Sp. z o.o. ul. Karczkowska 46 02-871 Warszawa tel. + 48 22 750 42 94, fax + 48 22 750 42 96 e-mail: schiessl@schiessl.pl www: www.schiessl.pl NIP: 951-16-35-342 Nr rejestrowy BDO: 000000027	
Data sprzedaży: 2022-12-20 Data wystawienia: 2022-12-20	
Faktura VAT FK/165/12/2022/PO Oryginał Na podstawie zamówienia: ZO 178434 (Of: 0)	

Takie dostrojenie wykonuje się poprzez dodanie pozycji szablonu REGEX w sekcji zamienniki używając przycisku „dodaj wiersz”.

Zamienniki						
NIP	Regex	Priorytet	Regex do podmiany	Regex tekst podmiany	Typ	
Dodaj wiersz		Zapisz		Usuń zaznaczony		Usuń wszystkie

Program w pozycji NIP wpisze samoczynnie numer NIP z dokumentu na podstawie którego chcemy stroić. W polu Regex musimy wprowadzić odpowiedni ciąg, który chcemy wyszukać na dokumencie. Treść wzorcowa, do której się odnosimy jest widoczna na zakładce Treść OCR. W polu priorytet możemy wprowadzić informację, jaki algorytm program ma stosować jako nadrzędny (jeśli jest ich kilka równoważnych). Kolejność jest stosowana od najmniejszego do największego. Następnie

wprowadzamy jaki znak chcemy zamienić na jaki i wybieramy typ danych, które chcemy wprowadzać do formatki z danymi rozpoznany przez OCR.

Możemy tego typu metodę uczenia stosować do nauki : numeru NIP, numeru dokumentu, wartość brutto, numer rachunku bankowego, data wystawienia dokumentu, data sprzedaży, forma zapłaty, termin płatności, numer dokumentu korygowanego, waluta, kurs waluty, data kursu waluty.

Do tworzenia zapytań REGEX można znaleźć w internecie wiele przydatnego oprogramowania wspierającego w konstruowaniu odpowiedniego wzoru. Przydatne tu będą poniższe linki :

<https://regexr.com>

<https://chat.openai.com>

<https://regex101.com>

Przykład 1

Jeśli chcielibyśmy dostroić OCR tak by dokument płatny kartą płatniczą był traktowany jako płatny gotówką to możemy to zrobić stosując odpowiednio REGEX. Jeśli np.: na naszym dokumencie jest słowo „Visa”, lub „Karta” to możemy je zastąpić słowem „gotówka”. Tworzymy wówczas dwie pozycje zamienników. W polu „Regex” wpisujemy „Visa”. W polu „do podmiany” ciąg znaków, które chcemy zamienić w naszym wypadku „Visa”, a w polu „Regex tekst podmiany” słowo na które chcemy zamienić nasze szukane słowo na „gotówka”. W ostatniej kolumnie musimy wybrać typ danych „forma płatności”.

W drugim przypadku, gdy na dokumencie znajduje się informacja o płatności jako napis „Karta” to odpowiedni wpisujemy „Karta”, „Karta”, „gotówka” i wybieramy typ danych „forma płatności”. Istotna jest tu wielkość liter.

Zamienniki						
	NIP	Regex	Priorytet	Regex do podmiany	Regex tekst podmiany	Typ
→	9720865431	Visa	0	Visa	gotówka	Forma Płatności
	9720865431	Karta	0	Karta	gotówka	Numer Dokumentu

Dodaj wiersz Zapisz Usuń zaznaczony Usuń wszystkie

Przykład 2

Innym przykładem użycia REGEX jest sytuacja, gdy wiemy, że przed informującą nas informacją jest konkretny napis. W naszym przykładzie nauczymy program rozpoznawania daty wystawienia dokumentu.

Jest to przypadek faktury obcojęzycznej i przed datą wystawienia jest słowo „crédito:”, a następnie konkretna data w formacie 31/01/2023. W tym przypadku użyjemy REGEX, który będzie składał się z dwóch warunków :

Pierwszy to informacja, że przed datą musi być konkretne słowo (?<=crédito:). Nawiasy są tu istotne.

Drugi warunek mówi o formacie daty $d\{2}\backslashd\{2}\backslash20\d\{2}$. Ten drugi warunek jest nam potrzebny, ponieważ słowo „crédito:” występuje we wzorcu tego dokumentu jeszcze w innych miejscach. Kompletny REGEX wygląda tak: $(?<=crédito:)\d\{2}\backslashd\{2}\backslash20\d\{2}$.

Zamienniki						
	NIP	Regex	Priorytet	Regex do podmiany	Regex tekst podmiany	Typ
→	19647148	$(?<=crédito:)\d\{2}\backslashd\{2}\backslash20\d\{2}$	0			Data Wystawienia

Dodaj wiersz Zapisz Usuń zaznaczony Usuń wszystkie

Przykład 3

Najczęściej występującym problemem jak już wspominałem jest błędne czytanie przez system rozpoznawania obrazu ukośników „/” w numerze dokumentu. Zazwyczaj ukośnik jest błędnie interpretowany jako „1”.

Np.: numer faktury VAT „0351/21/FVF” jest podpowiadany jako „03511211FVF”. Oczywiście możemy sobie z tym poradzić stosując opcję zamienników. Pierwszym etapem jest odnalezienie w tekście ciągu który odszukał system rozpoznawania obrazu. Możemy w naszym przypadku zastosować regex:

$(?<=Nr)\d\{4}\backslashd\{2}\backslashFVF$

Nasz regex wyszukuje na początku tekst „Nr „ (ten ciąg jest pomijany w wyniku). Następnie szuka czterech cyfr. Następnie szuka dowolnego znaku. Kolejnym wynikiem jest ciąg dwóch cyfr za którym jest dowolny znak i napis FVF.

W polu „Regex do podmiany” wpisujemy informację, która mówi o tym by program zamienił znak, który jest przed napisem FVF: $.(?=FVF)$. Kropka przed nawiasem jest istotna i nawias również.

W polu „Regex tekst podmiany” wstawiamy „/” (ukośnik).

W ten sposób otrzymamy wynik: „0351121/FVF”. Uwaga to nie jest gotowe rozwiązanie. Brakuje jeszcze zamiany piątej z kolei jedynki na ukośnik.

Zamienniki						
	NIP	Regex	Priorytet	Regex do podmiany	Regex tekst podmiany	Typ
→	6312667237	$(?<=Nr)\d\{4}\backslashd\{2}\backslashFVF$	1	$.(?=FVF)$	/	Numer Dokumentu

Dodaj wiersz Zapisz Usuń zaznaczony Usuń wszystkie

Rozdział 2.

Strojenia AI 2.0.

Inną z metod strojenia OCR jest nauka z nauczycielem, gdzie to operator wskazuje miejsce na skanie z którego system ma pobierać konkretną informację. Będąc na danych szczegółowych dokumentu możemy rozpocząć uczenie metodą AI 2.0. Naciskamy na klawiaturze przycisk Shift i jednocześnie przy wciśniętym lewym przycisku myszki zaznaczamy ramką treść, którą chcemy by system nam pobrał do strojenia danych OCR.

W ten sposób możemy nauczyć program lokalizacji na dokumencie, gdzie znajduje się treść, którą chcemy rozpoznawać. W ten sposób wskazujemy w jakim miejscu na dokumencie program ma szukać konkretnej informacji (współrzędne x, y) w przyszłości. Oczywiście oprócz samej lokalizacji uwzględniany jest też prawidłowy format zaznaczonego fragmentu tekstu (np. DD.MM.RRRR). W ten sposób możemy pomóc odnajdywać: numer dokumentu, daty wystawienia, daty sprzedaży, termin płatności.

The screenshot shows a software application for processing VAT documents. The interface includes several panels:

- Dane poddawane dokumentu:** Fields for document ID (631221431), date of issue (13-02-2023), date of sale (31-07-2021), and other document details.
- Dane szczegółowe kontrahenta:** Fields for contractor name (REMONDIS GŁIWICE SPÓŁKA Z OGRANICZONĄ ODPOWIEDZIALNOŚCIĄ), address (Kaszubka 2, Gliwice), and other identification numbers.
- Podgląd dokumentu:** A preview of the scanned document. A red box highlights a field containing '31.07.2021'. A context menu is open over this field, showing options like 'NIP', 'Numer Dokumentu', 'Wartość Brutto', 'Data Wystawienia', 'Data Sprzedaży', 'Forma Płatności', 'Termin Płatności', 'Numer Dok. Korygowanego', 'Waluta', 'Kurs Waluty', and 'Data Kursu'.
- Table:** A table with columns: Stawka, Netto, VAT, Brutto, Kategoria 1 opis, Odlicze., and Podstaw. The table contains several rows of data, including VAT rates and net amounts.

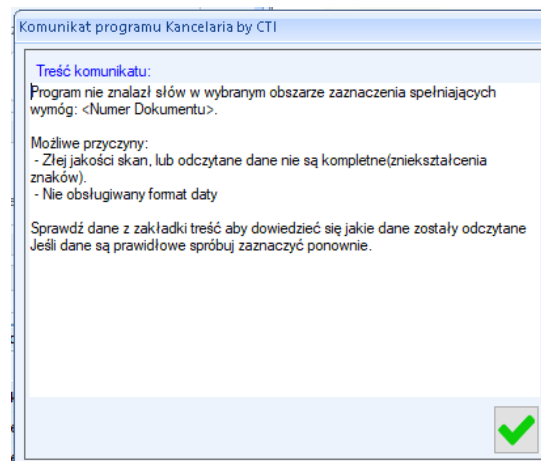
Po zwolnieniu przycisku myszki pojawi się okienko, które może wskazać nam kilka wyników (ciągów znaków), które pobrał z zaznaczonego fragmentu dokumentu. Im bardziej dokładnie zaznaczymy interesujący fragment, tym dokładniej zostanie wskazany koordynat x, y na dokumencie. Ma to duże znaczenie w jakości rozpoznawania kolejnych dokumentów od tego dostawcy. Jeśli w okienku pojawią się różne wyniki alternatywne to powinniśmy wybrać ten prawidłowy. Jest to związane z tym, że stosujemy różne programy rozpoznające obraz i w zależności od rodzaju wydruku każdy z nich sprawdza się lepiej w różnych przypadkach. Po naciśnięciu przycisku „zatwierdź” informacja zostanie przeniesiona do formularza z danymi faktury i jednocześnie nasz wybór zostanie zapisany jako wynik strojenia OCR. Jeśli żaden z wyników nam nie odpowiada to oczywiście możemy wycofać się z uczenia wybierając opcję „anuluj”.

631-22-11-431
: PL 6312211431

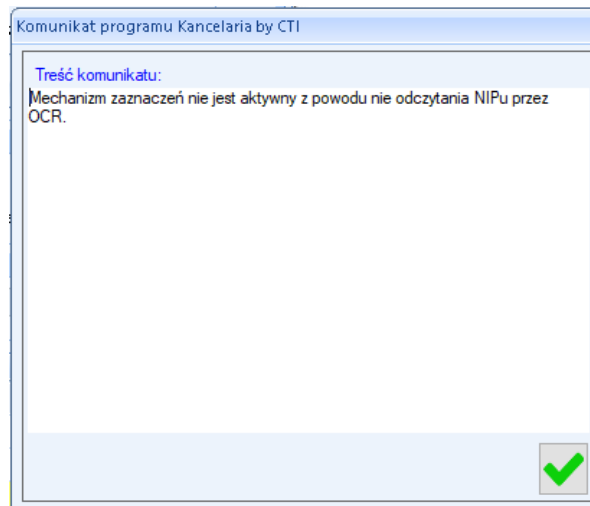
NIP: 6311907277

Mr. Faktury		Mr. Klienta	Data wyst. faktury		
		1000203	31.07.2021		
A VAT przedaży /towaru		Nazwa Usługi	a Netto	Wartość Netto	Stawka VAT
ŚWIĘTOJAŃSKA 35 44100 GLIWICE					
21	Wywóz odpadów komuna Pojemnik - 2401		56,41	56,41	8%
21	Wywóz odpadów komunalnych Pojemnik - 2401	1,000 szt	56,41	56,41	8%
21	Wywóz odpadów komunalnych	1,000 szt	56,41	56,41	8%

Jeśli jednak po zaznaczeniu wybranego fragmentu dokumentu otrzymamy komunikat o treści „Program nie znalazł słów w wybranym obszarze zaznaczenia” to oznacza, że jakość skanu dokumentu nie jest odpowiednia do tego by rozpoznać dane. Wówczas strojenie OCR nie jest możliwe.



Strojenie nie jest też możliwe, gdy nie udało się rozpoznać numeru NIP sprzedawcy na dokumencie. Wówczas otrzymamy również odpowiedni komunikat.



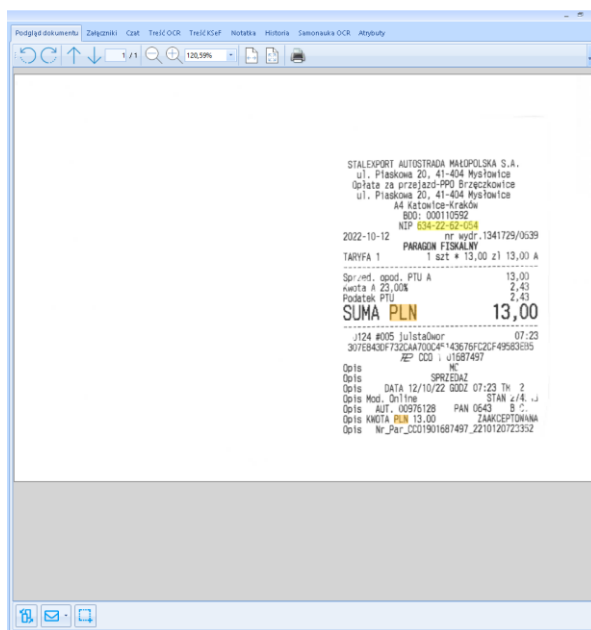
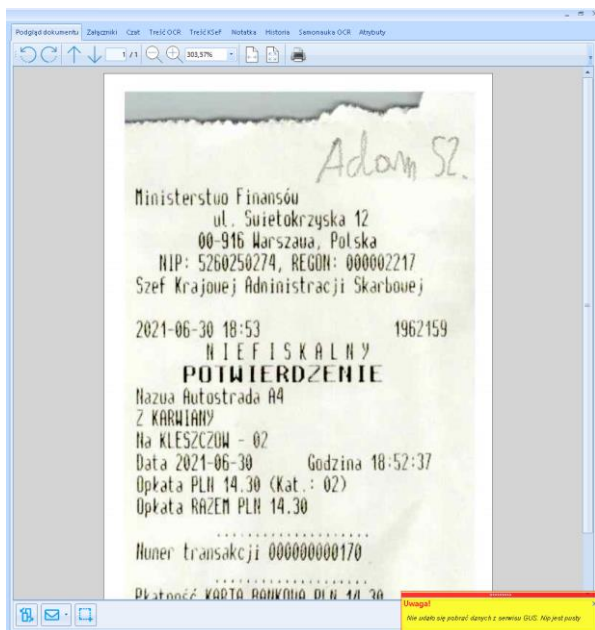
Stosując ten sam skrót klawiszy Shift i lewy przycisk myszy na wyświetlonej liście mamy jeszcze możliwość wyboru numeru NIP. Opcja ta jednak nie służy do uczenia rozpoznawania numeru NIP, a jedynie do tego by za jej pomocą skopiować ze skanu tę informację i przenieść ją do formatki z danymi. Uczenie, gdzie w dokumencie znajduje się numer NIP nie ma sensu, gdyż wszystkie metody uczenia są ograniczone do szablonów dokumentów od dostawców z konkretnym numerem NIP i jeśli program nie rozpoznał wcześniej samodzielnie NIP to nie ma informacji do jakiego wzoru się odnieść. Sama funkcjonalność kopiowania NIP za pomocą tej kombinacji klawiszy powstała w celu przenoszenie do formatki z danymi numery NIP np.: z poza naszego kraju i UE, których nie rozpoznajemy na dzień dzisiejszy.

Uczenie nie funkcjonuje dla : wartości brutto, formy płatności, numeru dokumentu korygowanego, waluty, kursu waluty, daty kursu waluty. Jest tak dlatego, że informacje te zazwyczaj na dokumencie (nawet przy tym samym szablonie) znajdują się w różnych pozycjach x, y i uczenie pozycji w której ma być wyszukiwana informacja doprowadziłoby do błędnego strojenia OCR.

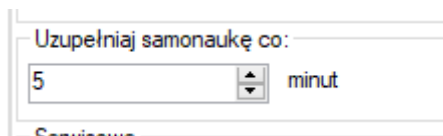
UWAGA !

Strojenie tego typu ma sens jedynie przy dokumentach, dla których jesteśmy pewni, że pozycja przez nas zaznaczona nie zmienia swojego położenia x, y. Najczęściej stosujemy ją dla dokumentów o wymiarze A4. Przy dokumentach o innych wymiarach np.: kwitek za autostradę może się okazać, że zeskanowany dokument za każdym razem ma interesujące nas informacje w różnych miejscach x, y na skanie. Można tu oczywiście te dokumenty ujednoczyć podczas samego wykonywania skanu. Jest to funkcja, która w większości współczesnych skanerów jest dostępna i należy w tym przypadku zapoznać się z instrukcją obsługi konkretnego urządzenia. Funkcjonalność ta może być również możliwa do ustawienia na wielu współczesnych smartfonach za pomocą których klienci wykonują zdjęcia faktur przesyłanych do biura rachunkowego.

Na poniższych przykładach pokazane jest jak różne mogą być to skany od tego samego dostawcy (wystawcy dokumentu).



Odświeżanie wiedzy systemu OCR używającego koordynatów odbywa się w zależności od ustawionego czasu w konfiguracji samego programu OCR. Parametr ten znajduje się na zakładce „Konfiguracja połączeń” w kontrolce „Uzupełnij samonaukę co :”



Informacja o tym czy dokument od tego dostawcy jest objęty strojeniem pokazuje się w zakładce „Samonauka OCR” w sekcji „Wzorce zaznaczenia”. W sekcji tej możemy usunąć koordynaty (zaznaczenia pozycji x, y) używając przycisku „usuń ...”. Czasami takie usunięcie jest potrzebne, gdy użytkownik błędnie nauczył program pozycji x, y.

Wzorce zaznaczenia		
NIP	Typ	Ile razy użyt.
9720865431	Data wystawienia	0
9720865431	Numer dokumentu	0

Jeśli chcemy sprawdzić jak użytkownik oznaczył pozycje x, y danych na skanie to możemy to obejrzeć na zakładce „podgląd dokumentu” naciskając przycisk pod skanem „pokaż/ukryj koordynatory samonauki”. Jest to trzecia ikonka od lewej. Jeśli ikonka ta nie jest widoczna to trzeba zmienić kontrolkę podglądu naciskając przycisk pierwszy od lewej.



Wzorce				
	NIP	Typ	Wzorzec	Ile razy użyt.
→	7790001083	Numer dokumentu	[0-9]	85

Usuń zaznaczony
Usuń wszystkie

Taki wzorzec można też usunąć, jeśli po jakimś czasie użytkownika okaże się, że program pomylił się przy strojeniu, lub firma ta zmieniała sposób numeracji i program podpowiada numer dokumentu zgodny z wzorcem, który kiedyś był prawidłowy, a obecnie jest błędny. Usunięcie takiego wzorca może mieć też znaczenie, gdy użytkownik programu kancelaria nauczył program nieprawidłowego zachowania i np.: wpisywał przez jakiś czas numer faktury jako ciąg znaków „NR : 00907321070356100”, a obecnie chciał by była to jedynie liczba „00907321070356100”.

Automatyczne strojenie dotyczy jedynie numeru dokumentu.

Program uczy się wzorców dla wszystkich baz danych (firm biura rachunkowego). Jeśli strojenie jeden z użytkowników wykona na jednym z klientów biura rachunkowego to na wszystkich innych firmach będzie widoczny efekt.

Podsumowanie.

Wszystkie czynności związane ze strojeniem programu OCR prowadzą do znacznego zwiększenia skuteczności rozpoznawania danych.

Kolejność strojenia powinna być wykonywana począwszy od samonauki systemu (Wzorce AI 3.0), uczenie poprzez zaznaczanie treści na dokumencie (Wzorce zaznaczeń AI 2.0). W przypadku gdy poprzednie metody uczenia zaawansowanego nie przyniosły wystarczającego efektu to można bazę dostroić za pomocą REGEX (Zamienniki AI 1.0).

mgr inż. Zygmunt Wilder

Centrum Technologii Informatycznej ZYGMUNT WILDER

ul. Świętojańska 35

44-100 Gliwice

NIP : 6311907277