

# Algorytmy i struktury danych z językiem Python

## Algorytm Aho-Corasick

Angela Czubak

## 1 Teoria

### 1.1 Wprowadzenie

Algorytm Aho-Corasick został opracowany przez Alfreda V. Aho oraz Margaret J. Corasick. Jego celem jest znalezienie wzorców  $\mathcal{P} = \{P_0, \dots, P_k\}$  pochodzących z pewnego słownika w tekście.

Cechą charakterystyczną tego algorytmu jest to, że szukanie wystąpień zadanych słów następuje "na raz", dzięki czemu złożoność obliczeniowa tego algorytmu wynosi  $O(m + z + n)$ , gdzie  $m$  - długość tekstu, w którym wyszukujemy,  $z$  - liczba wystąpień wzorców w zadanym tekście oraz  $n = \sum_{i=0}^k |P_i|$  - sumy długości tychże.

Algorytm ten jest stosowany np. w komendzie UNIX-a - **fgrep**.

Idea algorytmu opiera się na drzewach trie i automatach, których tworzenie omówię dalej.

### 1.2 Drzewo trie

**Definicja 1.** *Drzewem trie dla zbioru wzorców  $\mathcal{P}$  nazywamy takie ukorzenione drzewo  $\mathcal{K}$ , że:*

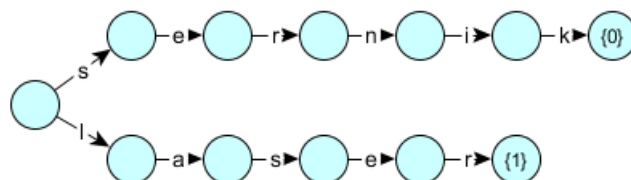
1. *Każda krawędź drzewa  $\mathcal{K}$  jest etykietowana jakimś znakiem*
2. *Każde dwie krawędzie wychodzące z jednego wierzchołka mają różne etykiety*

**Definicja 2.** *Etykietą węzła  $v$  nazywamy konkatencję etykiet krawędzi znajdujących się na ścieżce z korzenia do  $v$ . Oznaczamy ją jako  $\mathcal{L}(v)$ .*

3. *Dla każdego  $P \in \mathcal{P}$  istnieje wierzchołek  $v$  taki, że  $\mathcal{L}(v) = P$ , oraz*
4. *Etykieta  $\mathcal{L}(v)$  jakiegokolwiek liścia  $v$  jest równa jakiemuś  $P \in \mathcal{P}$*

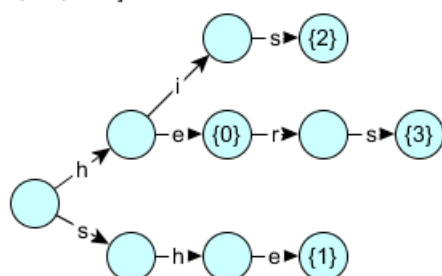
Przykładowe rysunki drzew trie znajdują się na następnej stronie. Numer w wierzchołku oznacza indeks słowa należącego do słownika, które jest etykietą tego wierzchołka.

$P=[\text{semik, laser}]$



Rysunek 1: Drzewo trie dla słów  $\mathcal{P} = \{\text{semik, laser}\}$

$P=[\text{he, she, his, hers}]$



Rysunek 2: Drzewo trie dla słów  $\mathcal{P} = \{\text{he, she, his, hers}\}$

### 1.3 Konstrukcja drzewa trie

Jak budować drzewo trie dla  $\mathcal{P} = \{P_0, \dots, P_k\}$ ? Procedura jest następująca:

1. Rozpocznij od stworzenia korzenia
2. Umieszczaj kolejne wzorce jeden po drugim według poniższych kroków:
  - (a) Poczynając od korzenia, podążaj ścieżką etykietowaną kolejnymi znakami wzorca  $P_i$
  - (b) Jeśli ścieżka kończy się przed  $P_i$ , to dodawaj nowe krawędzie i węzły dla pozostałych znaków  $P_i$
  - (c) Umieść identyfikator  $i$  wzorca  $P_i$  w ostatnim wierzchołku ścieżki

Jak łatwo zauważyć, konstrukcja drzewa zajmuje  $O(|P_0| + \dots + |P_k|) = O(n)$ .

### 1.4 Wyszukiwanie wzorca w drzewie

Wyszukiwanie wzorca  $P$  odbywa się następująco:

Tak długo, jak to możliwe, podążaj ścieżką etykietowaną kolejnymi znakami  $P$

1. Jeśli ścieżka prowadzi to wierzchołka z pewnym identyfikatorem,  $P$  jest słowem w naszym słowniku  $\mathcal{P}$
2. Jeśli ścieżka kończy się przed  $P$ , to słowa nie ma w słowniku
3. Jeśli ścieżka kończy się w wierzchołku bez identyfikatora, to słowa nie ma w słowniku

Wyszukiwanie zajmuje więc  $O(|P|)$ .

Naiwnie postępując, moglibyśmy chcieć wyszukiwać wzorce w tekście tak, by dla każdego znaku tekstu próbować iść wzdłuż krawędzi odpowiadającym kolejnym znakom - jeśli po drodze przejdziemy przez wierzchołki z identyfikatorami, to znaleźliśmy słowa im odpowiadające. Gdy już nie ma krawędzi, którą moglibyśmy przejść, zaczynamy wyszukiwanie dla kolejnego znaku tekstu. Jednak takie wyszukiwanie zajęłoby  $O(nm)$  czasu, gdzie  $m$  - długość tekstu,  $n$  - suma długości wzorców.

By przyspieszyć wyszukiwanie wzorców, rozszerzamy drzewo trie do **automatu**.

## 1.5 Automat

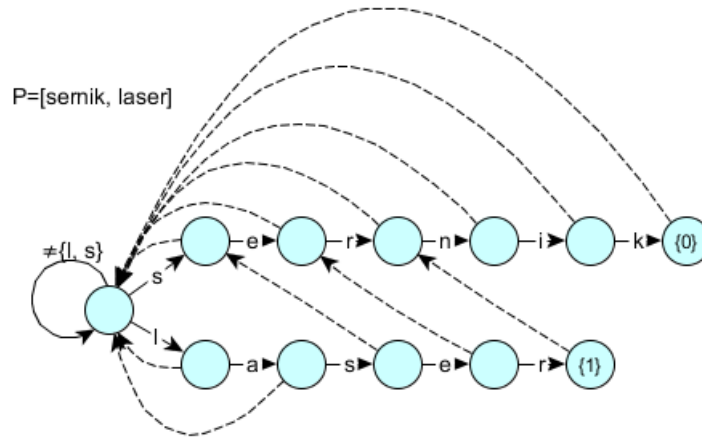
**Definicja 3.** *Automatem deterministycznym nazywamy piątkę uporządkowaną  $(\Sigma, Q, q, \delta, F)$ , gdzie*

1.  $\Sigma$  - skończony alfabet
2.  $Q$  - skończony zbiór stanów
3.  $q \in Q$  - stan początkowy
4.  $\delta : Q \times \Sigma \rightarrow Q$  - funkcja przejścia, przypisującą parze  $(q, a)$  nowy stan  $p$ , w którym znajdzie się automat po przeczytaniu symbolu  $a$  w stanie  $q$
5.  $F \subset Q$  - zbiór stanów końcowych

W naszym przypadku automat, w związku z wprowadzeniem dodatkowych pojęć, nie będzie ściśle deterministyczny. Automat będziemy budować na podstawie drzewa trie, zatem uściślam: zbiór stanów będą stanowiący węzły drzewa, a stanem początkowym będzie korzeń, do którego będę się czasami odnosiła jako 0. Wprowadzamy następujące funkcje:

1. Funkcja **goto**, oznaczana jako  $g(q, a)$  - jest to odpowiednik funkcji przejścia  $\delta$ , odpowiada ona krawędziom w drzewie, ponadto zachodzą jeszcze pewne własności; w skrócie:
  - jeśli krawędź  $(q, v)$  jest etykietowana przez  $a$ , to  $g(q, a) = v$
  - $g(0, a) = 0$  dla każdego  $a$  nie będącego etykietą krawędzi wychodzącej z korzenia - automat ma pozostać w stanie początkowym, jeśli nie można znaleźć dopasowania
  - w przeciwnym przypadku  $g(q, a) = \emptyset$  - brak przejścia w automacie
2. Funkcja **failure**, oznaczana jako  $f(q)$ , dla każdego stanu różnego od początkowego ( $q \neq 0$ ) zwraca stan, do którego powinniśmy się udać w przypadku niemożności zastosowania funkcji  $g(q, a)$  - nie istnieje krawędź wychodząca z  $q$ , etykietowana przez  $a$ . Stanem tym jest węzeł odpowiadający najdłuższemu właściwemu sufiksowi  $L(q)$  (najdłuższemu podsłowu

$y$ , krótszemu niż samo słowo  $s$ , takiemu, że istnieje słowo  $t$  o niezerowej długości, że  $s = ty$ ). Chodzi o to, by nie przegapić żadnego potencjalnego dopasowania wzorca - np. biorąc słowa *laser*, *sernik* i szukając w tekście *lasernik*, zaczniemy od dopasowania do słowa *laser*, a powinniśmy mieć jeszcze możliwość przejścia do gałęzi odpowiadającej słowu *sernik*, by także je odnaleźć. Funkcje przejścia w tym przypadku przedstawiono na rysunku (3). Funkcja  $f(q)$  jest zawsze dobrze zdefiniowana, gdyż  $\mathcal{L}(0) = \epsilon$  jest sufiksem każdego słowa.



Rysunek 3: Automat dla słów  $\mathcal{P} = \{\textit{sernik}, \textit{laser}\}$ . Przerywaną linią oznaczono krawędzie odpowiadające funkcji  $f(q)$

3. Funkcja **wyjścia**, oznaczana jako  $out(v)$ , zwraca indeksy wzorców, do których znajdujemy dopasowanie w stanie  $q$ .

## 1.6 Przeszukiwanie tekstu

Załóżmy, że mamy do dyspozycji gotowy automat oraz tekst  $T[1 \dots]$ , w którym szukamy wzorców. Procedura ta prezentuje się w następujący sposób:

**Algorithm 1.1:** SEARCH( $T[0 \dots m - 1]$ )

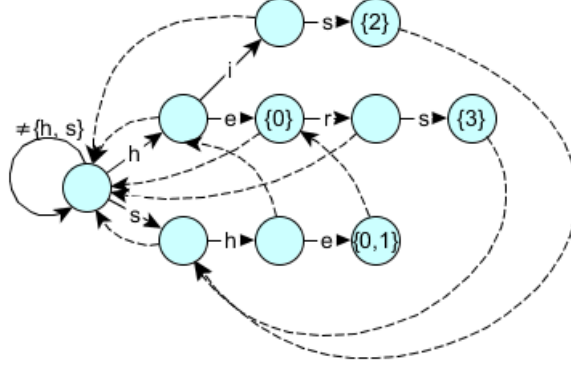
```

 $q \leftarrow 0$ 
for  $i \leftarrow 0$  to  $m - 1$ 
do
  while  $g(q, T[i]) = \emptyset$ 
  do
     $q \leftarrow f(q)$ 
    comment: podążaj za funkcją failure, aż znajdziesz dopasowanie
   $q \leftarrow g(q, T[i])$ 
  comment: przejdź do dopasowanego stanu
  if  $out(q) \neq \emptyset$ 
  then output  $(i), out(q)$ 

```

Weźmy automat jak na rysunku (4):

$P=[he, she, his, hers]$



Rysunek 4: Automat dla słów  $\mathcal{P} = \{he, she, his, hers\}$ . Przerywaną linią oznaczono krawędzie odpowiadające funkcji  $f(q)$

Przeszukamy przy jego pomocy tekst *ushers*:

1. Czytamy znak  $u$  - zostajemy w korzeniu
2. Czytamy znak  $s$  - przechodzimy po odpowiedniej krawędzi, czyli idziemy do węzła etykietowanego  $s$
3. Czytamy znak  $h$  - przechodzimy po odpowiedniej krawędzi, czyli idziemy do węzła etykietowanego  $sh$
4. Czytamy znak  $e$  - przechodzimy po odpowiedniej krawędzi, czyli idziemy do węzła etykietowanego  $she$ ; wypisujemy, że znaleźliśmy słowa o indeksach 0 i 1 na pozycji 3
5. Czytamy znak  $r$  - korzystamy z krawędzi **failure**, następnie przechodzimy po odpowiedniej krawędzi, czyli idziemy do węzła etykietowanego  $her$
6. Czytamy znak  $s$  - przechodzimy po odpowiedniej krawędzi, czyli idziemy do węzła etykietowanego  $hers$ ; wypisujemy, że znaleźliśmy słowo o indeksie 3 na pozycji 5

Złożoność czasowa takiej procedury jest  $O(m + z)$ , gdzie  $m$  - długość tekstu, w którym wyszukujemy, a  $z$  - liczba wystąpień wzorca w tekście.

Wynika to z faktu, że liczba wywołań funkcji  $f(q)$  jest ograniczona z góry przez  $m$  - w danym momencie możemy wywołać funkcję  $f(q)$  co najwyżej tyle razy, ile znaków zdążyliśmy przeczytać z  $T$ , a w sumie możemy ich przeczytać  $m$  (po wywołaniu tej funkcji "przesuwamy" początek sufiksu na pewną literę z ciągu  $T$  - przesunąć ten początek możemy maksymalnie  $m$  razy).

Podobnie, funkcja  $g(q, a)$  jest wywoływana dokładnie raz dla każdego  $a \in T$  - zostanie ona wywołana  $m$  razy. Wystąpienie wzorca możemy zgłaszać w czasie stałym, stąd zgłoszenie wszystkich zajmie  $O(z)$ .

## 1.7 Budowa automatu

W konstrukcji automatu możemy wyróżnić dwie fazy.

### 1.7.1 Faza I

1. Tworzymy drzewo trie dla słownika  $\mathcal{P}$   
 item Dla każdego  $P_i \in \mathcal{P}$  ustawiamy  $out(v) = i$  dla wierzchołka  $v$  etykietowanego przez  $P_i$
2. Uzupełnij funkcje przejść dla korzenia

$$g(0, a) = 0$$

dla każdego znaku  $a$  (należącego do alfabetu  $\Sigma$ ) takiego, że nie etykietuje on żadnej krawędzi wychodzącej z korzenia

Złożoność czasowa takiej procedury, dla pewnego niezmiennego alfabetu, wynosi  $O(n)$ , gdzie  $n = \sum_{i=0}^k |P_i|$ .

### 1.7.2 Faza II

Przedstawię ją w postaci pseudokodu:

**Algorithm 1.2:** PHASE2(void)

```

Q ← QUEUE()
for a ∈ Σ
  do { if q ← g(0, a) ≠ 0
      then { f(q) ← 0
           Q.ENQUEUE(q)
    while !Q.ISEMPTY()
      do { r ← Q.DEQUEUE()
          for a ∈ Σ
            do if u ← g(r, a) ≠ ∅
                then { Q.ENQUEUE(u)
                     v ← f(r)
                     while g(v, a) = ∅
                       do v ← f(v) // (*)
                     f(u) ← g(v, a)
                     out(u) ← out(u) ∪ out(f(u)) // (**)

```

Jak widać funkcje  $f$  i  $out$  są wyliczane dla wierzchołków w kolejności BFS. Dzięki temu wierzchołki znajdujące się bliżej korzenia zostały już obsłużone, gdy zachodzi potrzeba skorzystania z odpowiednich funkcji na nich wykonywanych.

Rozważmy wierzchołki  $r$  i  $u = g(r, a)$ , w takim przypadku  $r$  jest rodzicem  $u$ . Co więcej,  $\mathcal{L}(u) = \mathcal{L}(r)a$ . Jakie więc powinno być  $f(u)$ ? Przypomnijmy, że  $f(u)$  powinno wskazywać na najdłuższy właściwy sufix  $\mathcal{L}(u)$ . Z tego wynika, że powinniśmy spróbować dopasować  $f(u) = g(f(r), a)$ , bo  $\mathcal{L}(f(u))$  może być sufiksem  $\mathcal{L}(g(f(r), a))$ , o ile taka krawędź istnieje. Jeśli jej nie ma, to próbujemy

$f(u) = g(f(f(r)), a)$ , itd., aż znajdziemy odpowiedni wierzchołek (pesymistycznie może to być korzeń). W linii oznaczonej (\*) wykonujemy właśnie te czynności.

Czynności oznaczone (\*\*) wykonujemy, gdyż wzorce rozpoznawane w stanie  $f(u)$  (i jedynie te) są właściwymi sufiksami  $\mathcal{L}(u)$  i dlatego powinny być rozpoznawane także w stanie  $u$ .

Jaka jest złożoność powyższej procedury? Zauważmy, że jest ona podobna do BFS-a - stąd przechodzenie po drzewie, pomijając linię oznaczoną (\*), zajmie czas proporcjonalny do rozmiaru drzewa - tj.  $O(n)$ . A jak określić, ile razy wykonamy linię (\*)?

Rozważmy ścieżkę złożoną z wierzchołków  $u_1, \dots, u_l$ , która jest tworzona podczas dodawania wzorca  $a_1 \dots a_l$ . Oznaczmy dodatkowo  $df(u)$  jako głębokość w drzewie wierzchołka  $f(u)$ , zatem  $df(u_1), \dots, df(u_l)$  to ciąg głębokości dla wierzchołków z rozważanej ścieżki, wszystkie są  $\geq 0$ .

Zauważmy ponadto, że głębokość kolejnego wierzchołka może wzrosnąć co najwyżej o 1, czyli  $df(u_{i+1}) \leq df(u_i) + 1$ , zatem wartości  $df$  wzrastają sumarycznie co najwyżej o  $l$  podczas przechodzenia tej ścieżki.

Kiedy wyliczamy położenie  $f(u_{i+1})$ , każde wywołanie linii (\*) przybliży  $v$  do korzenia, a stąd wartość  $df(u_{i+1})$  będzie mniejsza od  $df(u_i)$  co najmniej o jeden. W związku z ograniczeniem od dołu, możemy zmniejszać kolejne wartości  $df(u_i)$  co najwyżej tyle, ile razy zostały one zwiększone, czyli linia (\*) zostanie wykonana  $\leq l$  razy dla pewnego wzorca o długości  $l$ .

Sumaryczna długość wzorców wynosi  $n$ , dlatego podczas budowy automatu linia (\*) zostanie wykonana co najwyżej  $n$  razy.

Zastanówmy się jeszcze, ile czasu zajmuje wykonanie linii (\*\*). Zauważmy, że przed wykonaniem przypisania,  $out(u) = \emptyset$  albo  $out(u)$  jest równy indeksowi  $\mathcal{L}(u)$ . Jakielwiek wzorce znajdują się w  $out(f(u))$ , są one na pewno krótsze niż  $\mathcal{L}(u)$ , bo  $f(u)$  jest bliżej korzenia - zatem zbiory te są rozłączne. Możemy więc przyjąć, że reprezentujemy je przez listy, do których da się dołączać drugą w stałym czasie.

Zatem złożoność czasowa drugiej fazy wynosi  $O(n)$ .

## 1.8 Determinizacja automatu

Ze względu na występowanie w automacie funkcji  $f(u)$  jest on nie deterministyczny - nie znamy przejść ze wszystkich stanów dla wszystkich możliwych znaków (np.  $a$ ), będziemy musieli więc czasem wykonać wiele przejść z użyciem funkcji  $f(u)$ , aż dojdziemy do stanu  $v$ , w którym istnieje dobrze określone przejście  $g(v, a)$ . Można jednak podejść inaczej o budowy automatu - mianowicie wprowadzić funkcję  $next(u, a)$ , którą wyliczamy w następujący sposób:

**Algorithm 1.3:** PHASE3(*void*)

```

Q ← QUEUE()
for a ∈ Σ
  do { if q ← g(0, a) ≠ 0
        then Q.ENQUEUE(q)
        next(0, a) ← g(0, a)
  while !Q.ISEMPTY()
    do { r ← Q.DEQUEUE()
          for a ∈ Σ
            do if u ← g(r, a) = ∅
                  then { Q.ENQUEUE(u)
                          v ← f(r)
                          while g(v, a) = ∅
                            do v ← f(v)
                          next(r, a) ← g(v, a)
                        else { next(r, a) ← g(r, a)
                              Q.ENQUEUE(g(r, a))

```

Procedurę tę wykonujemy po wyliczeniu funkcji  $f(u)$ . Wtedy wyszukiwanie upraszcza się do:

**Algorithm 1.4:** SEARCHDETERMINISTIC( $T[0 \dots m-1]$ )

```

q ← 0
for i ← 0 to m-1
  do { q ← next(q, T[i])
        comment: przejście jest deterministyczne
        if out(q) ≠ ∅
          then output (i), out(q)

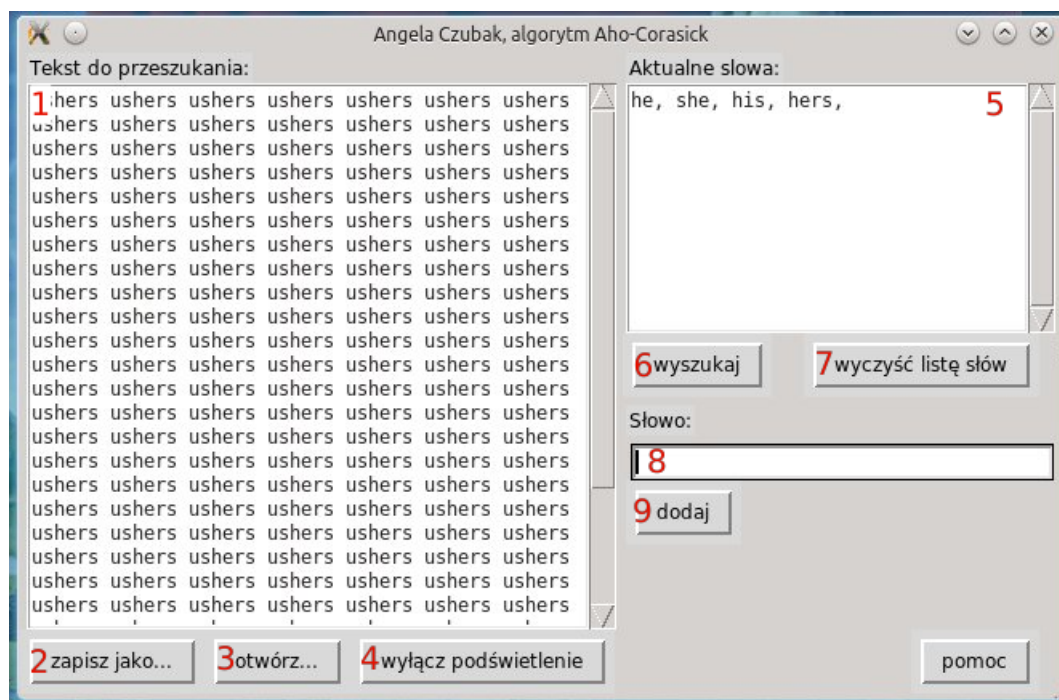
```

Jednak wprowadzenie takich przejść wiąże się ze znacznym obciążeniem pamięciowym, dlatego ja w swojej implementacji pominę te kroki.



## 2 Opis interfejsu

Zaimplementowano interfejs graficzny ułatwiający korzystanie z napisanego kodu. Został on przedstawiony na rysunku (5).



Rysunek 5: Interfejs graficzny programu zaliczeniowego

Poniżej znajduje się opis poszczególnych elementów interfejsu:

### 1. TEKST DO PRZESZUKANIA

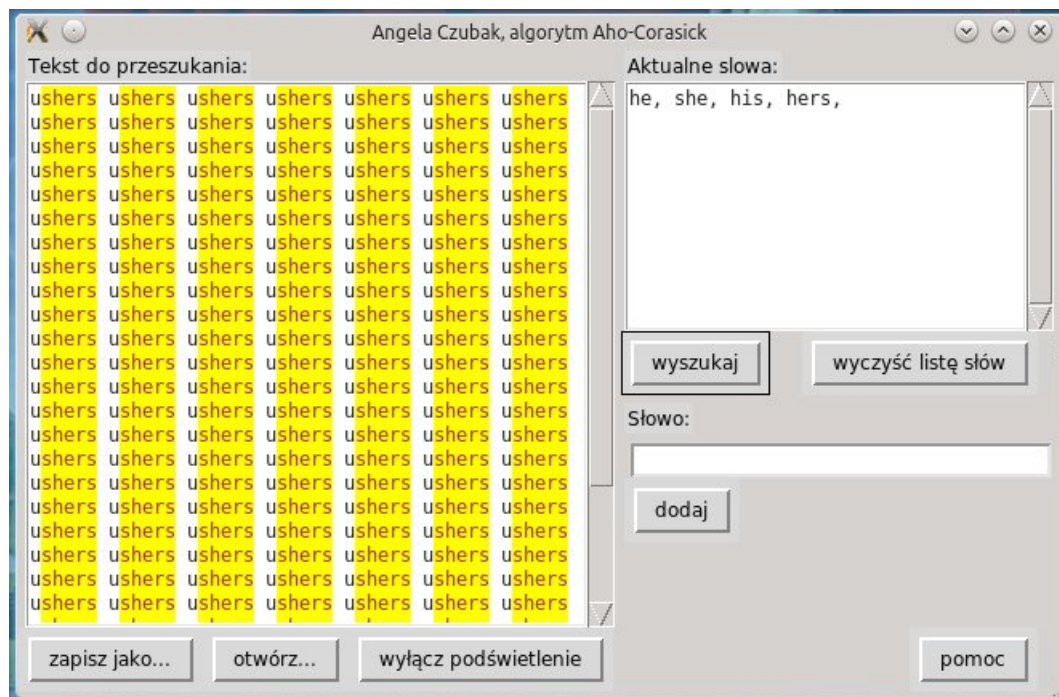
Jest to pole, w którym umieszczamy tekst, w którym będziemy wyszukiwać wzorce. Możemy tam wprowadzać tekst wprost z klawiatury lub wczytać tekst z pliku. W tym drugim przypadku należy kliknąć przycisk **otwórz...**, a następnie wybrać plik z użyciem okna dialogowego. Tekst, który wprowadzimy do tego pola możemy zapisać. By to zrobić, należy kliknąć przycisk **zapisz jako...**, a następnie wybrać odpowiednią nazwę pliku. Po wykonaniu wyszukania znalezione wzorce zostaną podświetlone na żółto (patrz rysunek (6)). By wyłączyć to podświetlenie, należy kliknąć na przycisk **wyłącz podświetlenie**.

### 2. ZAPISZ JAKO...

Tekst, który wprowadzimy do tego pola możemy zapisać. By to zrobić, należy kliknąć przycisk **zapisz jako...**, a następnie wybrać odpowiednią nazwę pliku.

### 3. OTWÓRZ...

Zamiast wpisywać tekst, możemy otworzyć gotowy plik tekstowy. W tym celu należy kliknąć przycisk **otwórz...**, a następnie wybrać plik z użyciem okna dialogowego.



Rysunek 6: Interfejs graficzny programu zaliczeniowego, działanie podświetlenia

### 4. WYŁĄCZ PODŚWIETLENIE

Po wykonaniu wyszukiwania znalezione wzorce zostaną podświetlone na żółto (patrz rysunek (6)). By wyłączyć to podświetlenie, należy kliknąć na przycisk **wyłącz podświetlenie**.

### 5. AKTUALNE SŁOWA

W tym polu znajdują się słowa (wzorce), które będą wyszukiwane w tekście przy użyciu algorytmu Aho-Corasick. By wyszukać wzorce, należy kliknąć przycisk **wyszukaj**, wtedy znalezione wystąpienia z pola **Tekst do przeszukania** zostaną podświetlone na żółto (patrz rysunek (6)). By wyczyścić listę słów, należy kliknąć przycisk **wyczyść listę słów**. By dodać słowo należy umieścić wymyślony przez nas wzorzec w polu **Słowo**, a następnie wcisnąć ENTER na klawiaturze lub przycisk **dodaj**.

### 6. WYSZUKAJ

By wyszukać wzorce, należy kliknąć przycisk **wyszukaj**, wtedy znalezione wystąpienia z pola **Tekst do przeszukania** zostaną podświetlone na żółto (patrz rysunek (6)).

## 7. WYCZYŚĆ LISTĘ SŁÓW

By wyczyścić listę słów, należy kliknąć przycisk **wyczyść listę słów**.

## 8. SŁOWO

W tym polu wpisujemy wzorec, który chcemy wyszukiwać w tekście znajdującym się w polu **Tekst do wyszukania**. Następnie należy nacisnąć ENTER na klawiaturze lub przycisk **dodaj**.

## 9. DODAJ

Po wpisaniu słowa w polu **Słowo**, które chcemy wyszukać w tekście znajdującym się w polu **Tekst do wyszukania**, można nacisnąć przycisk **dodaj**, by dodać słowo do listy wzorców.

# 3 Kod źródłowy

Cały kod oraz opis zmian można podejrzeć na <https://github.com/cosmia/pythonProject>.

## 3.1 Klasa MyList

Ponieważ złożoność dołączania jednej listy do drugiej w *Pythonie* jest zależna od długości tej drugiej, postanowiono napisać własną implementację listy tak, aby łącznie następowało w czasie stałym.

### 3.1.1 MyListException

Wyjątek związany z operacjami na obiektach klasy *MyList*

### 3.1.2 Element

Klasa opisująca element listy.

Posiada ona następujące **poła**:

1. *arg* - zawartość tego elementu listy
2. *follow* - następny element na liście

**Metody:**

1. konstruktor `__init__(self, arg=None, follow=None)`- *arg* - element znajdujący się w liście, *follow* - następny element na liście
2. `setData(self, arg)` - ustawienie zawartości elementu listy na *arg*
3. `setNext(self, follow)` -ustawienie następnego elementu na liście na *follow*, rzuca wyjątkiem *MyListException*, jeśli *follow* nie jest klasy *Element*
4. `getData(self)` - zwraca zawartość elementu listy
5. `getNext(self)` - zwraca następny element na liście

### 3.1.3 MyList

Klasa opisująca moją wersję listy.

Posiada ona cztery **poła**:

1. *first* - pierwszy element listy, jeśli lista jest pusta, jest to *None*
2. *last* - ostatni element listy, jeśli jest to lista pusta, jest to *None*
3. *current* - jest to zmienna pomocnicza, używana przy iterowaniu listy, początkowo ustawiona na *None*
4. *length* - długość listy

**Metody:**

1. konstruktor bezargumentowy
2. *add(self, argument)* - dodaje element argument do listy na ostatniej pozycji
3. *\_\_add\_\_(self, other)\_\_* - dodaje do siebie dwa obiekty *MyList*, zmienia pierwszy obiekt, zwraca wskaźnik na pierwszy obiekt; rzuca *MyListException*, jeśli drugi argument nie jest obiektem *MyList*
4. *\_\_iter\_\_(self)\_\_* - zwraca iterator dla listy *MyList*
5. *next(self)* - zwraca następny element na liście, metoda dla iteratora
6. *\_\_len\_\_(self)\_\_* - zwraca długość listy
7. *\_\_eq\_\_(self, other)\_\_* - metoda porównująca listy, zwraca *True*, jeśli listy równe, *False* wpp
8. *\_\_ne\_\_(self, other)\_\_* - metoda sprawdzająca, czy listy są różne, zwraca *True*, jeśli tak; *False* wpp
9. *\_\_contains\_\_(self, other)\_\_* - metoda sprawdzająca, czy lista zawiera *other*, zwraca *True*, jeśli tak; *False* wpp

### 3.1.4 Kod

```
1  #!/usr/bin/python
2  # -*- coding: utf-8 -*-
3
4  class MyListException(Exception):
5      '''wyjatek dla klasy MyList'''
6      def __init__(self, mes):
7          '''konstruktor, argumentem tresc przy rzucaniu wyjatku'''
8          self.value = mes
9      def __str__(self):
10         '''podaje tresc wyjatku'''
11         return self.value
12
13 class Element:
14     '''klasa opisujaca element MyList'''
15     def __init__(self, arg=None, follow=None):
16         '''konstruktor; arg - element znajdujacy sie w liscie ,
```

```

17         follow — następny element na liście'''
18     if follow is not None and not isinstance(follow, Element):
19         raise MyListException("argument is not an Element")
20     self.arg = arg
21     self.follow = follow
22     def setData(self, arg):
23         '''ustawienie zawartosci elementu listy na arg'''
24     self.arg = arg
25     def setNext(self, follow):
26         '''ustawienie nastepnego elementu na liście na follow'''
27     if not isinstance(follow, Element):
28         raise MyListException("argument is not an Element")
29     self.follow = follow
30     def getData(self):
31         '''zwraca zawartosc elementu listy'''
32     return self.arg
33     def getNext(self):
34         '''zwraca następny element na liście'''
35     return self.follow
36
37 class MyList:
38     '''lista, która będzie można łączyć z drugą w czasie stałym
39     jest to uproszczona lista, nie zawiera np. usuwania elementów,
40     gdyż nie wydaje się to potrzebne'''
41     def __init__(self):
42         '''konstruktor, tworzy pustą listę'''
43     self.first = None
44     self.last = None
45     self.current = None
46     self.length = 0
47     def add(self, argument):
48         '''dodaje argument do listy na ostatniej pozycji'''
49     if self.first is None:
50         self.first = Element(argument, None)
51         self.last = self.first
52     else:
53         tmp = Element(argument, None)
54         self.last.setNext(tmp)
55         self.last = tmp
56     self.length += 1
57     def __add__(self, other):
58         '''dodaje do siebie dwa obiekty MyList
59         zmienia pierwszy obiekt, zwraca wskaźnik na pierwszy obiekt'''
60     if other is None or other.first is None:
61         return self
62     if not isinstance(other, MyList):
63         raise MyListException("the other argument is not a MyList")
64     if self.first is None:
65         self.first = other.first
66         self.length = other.length
67         self.last = other.last
68         return self
69     #print other.first
70     self.last.setNext(other.first)
71     self.length += other.length
72     return self
73     def __iter__(self):
74         '''metoda zwracająca iterator'''
75     self.current = self.first
76     return self
77     def next(self):
78         '''zwraca następny element na liście'''

```

```

79     if self.current is None:
80         raise StopIteration
81     else:
82         tmp = self.current.getData()
83         self.current = self.current.getNext()
84         return tmp
85     def __len__(self):
86         '''metoda zwracajaca dlugosc listy'''
87     return self.length
88     def __eq__(self, other):
89         '''metoda porownujaca listy
90         zwraca True, jesli listy rowne, False wpp'''
91     if not isinstance(other, MyList):
92         return False
93     dl = len(other)
94     if dl != len(self):
95         return False
96     iter1 = iter(self)
97     iter2 = iter(other)
98     for i in range(dl):
99         e1 = iter1.next()
100        e2 = iter2.next()
101        if e2 != e1:
102            return False
103    return True
104    def __ne__(self, other):
105        '''metoda sprawdzajaca, czy listy sa rozne
106        zwraca True, jesli tak; False wpp'''
107    return not self == other
108    def __contains__(self, other):
109        '''metoda sprawdzajaca, czy lista zawiera other
110        zwraca True, jesli tak; False wpp'''
111    for i in self:
112        if other == i:
113            return True
114    return False

```

myList.py

## 3.2 Klasa Node

### 3.2.1 Kod

```

#!/usr/bin/python
2 # -*- coding: utf-8 -*-

4 from myList import *

6 class NodeException(Exception):
7     '''wyjatek dla klasy Node'''
8     def __init__(self, value):
9         '''konstruktor, argumentem tresc przy rzucaniu wyjatku'''
10    self.napis = value
11    def __str__(self):
12        '''podaje powod wyjatku'''
13    return repr(self.napis)
14    __repr__ = __str__

16 class Node:
17     '''klasa opisujaca wezel/stan w automacie/drzewie Trie'''
18    def __init__(self):

```

```

20     '''konstruktor bezargumentowy
        accept = MyList() - pusta lista ,
        fail = None
22     edges = {}'''
    self.accept = MyList()
24     self.edges = {}
    self.fail = None
26     def labelCorrect(self, label):
        '''sprawdza, czy label jest poprawna etykieta krawedzi
28         jesli nie, rzuca NodeException'''
        if not isinstance(label, (str, unicode)):
30             raise NodeException("label must be a character")
        if len(label) != 1:
32             raise NodeException("label must be exactly one character long
                ")
        def nodeCorrect(self, node, n=1):
34            '''sprawdza, czy node jest obiektem klasy Node
                jesli nie, to rzuca NodeException
36            n to numer argumentu, którym jest node w jakiej funkcji
                sluzy uszczegolowieniu, ktory argument jest bledny'''
            if not isinstance(node, Node):
38                lan = "the "
                if n != 1: lan += "second "
                lan += "argument must be a node"
40                raise NodeException(lan)
            def getAccept(self):
42                '''zwraca liste indeksow slow, ktore akceptuje ten wezel,
                    lista jest pusta, jesli ten wezel niczego nie akceptuje'''
44                return self.accept
            def getLabels(self):
46                '''zwraca liste etykiet dla krawedzi wychodzacych z tego wezla'''
48                return self.edges.keys()
            def getAim(self, label):
50                '''zwraca wezel, do ktorego prowadzi krawedz z etykieta label
                    jesli brak takiej krawedzi, zwraca None'''
52                self.labelCorrect(label)
            if label not in self.edges:
54                if self.fail is None:
                    return self
                    else:
56                    return None
                else:
58                    return self.edges[label]
            def getFail(self):
60                '''zwraca wezel odpowiadajacy najdluzszemu wlasciwemu sufiksowi
                    slowa, ktore do ktorego probowalismy znalezc dopasowanie'''
62                return self.fail
            def setAccept(self, number):
64                '''ustalamy, ze ten wezel akceptuje slowo o indeksie number lub
                    slowa o
                    indeksach z MyList number
66                rzuca NodeException, jesli number nie jest calkowita liczba
                    nieujemna
                    albo obiektem MyList calkowitych liczb nieujemnych'''
68                if isinstance(number, MyList):
                    for i in number:
70                        if not isinstance(i, (long, int)) or i < 0:
72                            mes = "argument should be a non-negative integer or long or
                                a set of those"
                                raise NodeException(mes)
                                #print number
74                            self.accept + number
76

```

```

78         return
if not isinstance(number, (long, int)):
    raise NodeException("argument should be an integer or long or
    a set of those")
80 if number < 0:
    raise NodeException("argument must be non-negative")
82 self.accept.add(number)
    def setAim(self, label, node):
84         '''ustalamy, ze z tego wezla bedzie wychodzic krawedz
            etykietowana label i bedzie ona prowadzic do node
86         rzuca NodeException, jesli label niepoprawna lub node nie jest
            wezlem'''
        self.labelCorrect(label)
88        self.nodeCorrect(node, 2)
        self.edges[label] = node
90        def setFail(self, node):
            '''ustalamy, ze najdluzszy sufixs slowa, do ktorego probowalismy
92            dopisowac w tym wezle odpowiada wezlowi node
            rzuca wyjatkiem, jesli node nie jest wezlem'''
94        self.nodeCorrect(node)
        self.fail = node

```

node.py

### 3.3 Klasa AhoCorasick

#### 3.3.1 Kod

```

1  #!/usr/bin/python
   # -*- coding: utf-8 -*-
3
   from node import *
5   from Queue import *
7   class AhoCorasickException(Exception):
        '''wyjatek dla klasy Node'''
9       def __init__(self, value):
        '''konstruktor, argumentem tresc przy rzucaniu wyjatku'''
11      self.napis = value
        def __str__(self):
13          '''podaj powod wyjatku'''
          return repr(self.napis)
15
   class AhoCorasick:
17       '''klasa opisujaca drzewo Trie / automat, sluzacy wyszukiwaniu
            wzorcow'''
        def __init__(self):
19            '''konstruktor bezargumentowy
            n - korzen drzewa, pusty
21            words - pusta liczba slow zakodowanych w drzewie'''
            self.n = Node()
23            self.words = []
            self.built = False
25            def addWord(self, word):
                '''dodaje slowo word do automatu/drzewa
27                rzuca AhoCorasickException, jesli word nie jest stringiem
                    lub zbudowano juz automat'''
29            if self.built:
                raise AhoCorasickException("automaton has been built already")
31            if not isinstance(word, (str, unicode)):

```



```

33         raise AhoCorasickException("argument is not a string")
34     dl = len(word)
35     if dl == 0: return #nie dodajemy pustego slowa
36     wezel = self.n
37     i = 0
38     #idziemy dopoki mozemy po istniejacych wezlach
39     while i < dl:
40         litera = word[i]
41         labels = wezel.getLabels()
42         if litera in labels:
43             wezel = wezel.getAim(litera)
44         else:
45             break
46         i += 1
47     #a teraz tworzymy nowe, jesli taka potrzeba
48     while i < dl:
49         litera = word[i]
50         wezel.setAim(litera, Node())
51         wezel = wezel.getAim(litera)
52         #wezel.setFail(self.n) #na poczatku najdluzszy wlasciwy
53         #sufiks to slowo puste
54         #mozna to w sumie robic przy budowaniu automatu...
55         i += 1
56     #jesli jeszcze nie dodalismy tego slowa
57     if wezel.getAccept() == MyList():
58         ktore = len(self.words)
59         wezel.setAccept(ktore)
60         self.words.append(word)
61     def lookUp(self, word):
62         '''sprawdza, czy dane slowo wystepuje w drzewie
63         zwraca True, jesli tak; False wpp
64         rzuca AhoCorasickException, jesli word nie jest strigiem'''
65     if not isinstance(word, str):
66         raise AhoCorasickException("argument is not a string")
67     if word == "": return False
68     i = 0
69     dl = len(word)
70     wezel = self.n
71     while i < dl:
72         litera = word[i]
73         labels = wezel.getLabels()
74         if litera not in labels:
75             return False
76         wezel = wezel.getAim(litera)
77         i += 1
78     if wezel.getAccept() != MyList():
79         return True
80     return False
81     def build(self):
82         '''konstruuje automat skonczony na podstawie drzewa, ktore
83         powstaje podczas dodawania slow metoda addWord'''
84     q = Queue()
85     for i in self.n.getLabels():
86         s = self.n.getAim(i)
87         s.setFail(self.n)
88         q.put(s)
89     while not q.empty():
90         r = q.get()
91         for a in r.getLabels():
92             u = r.getAim(a)
93             q.put(u)
94             v = r.getFail()

```

```

93     while v.getAim(a) is None: #jesli stad nie ma tego przejścia
94         v = v.getFail() #to szukaj krotszego dopasowania
95     u.setFail(v.getAim(a))
96     u.setAccept(u.getFail().getAccept()) #dodaj nowe akceptowane
97     slowa
98 self.built = True
99     def makeTree(self, wordList):
100     '''konstruuje drzewo i automat na podstawie listy slow wordList
101     takze dodaje do istniejacego automatu wzorce z wordList
102     rzuca AhoCorasickException, jesli wordList nie jest lista
103     stringow
104     lub zbudowano juz automat'''
105 if self.built:
106     raise AhoCorasickException("automaton has been built already")
107 )
108 if not isinstance(wordList, list):
109     raise AhoCorasickException("argument is not a list")
110 for i in wordList:
111     if not isinstance(i, (str, unicode)):
112         raise AhoCorasickException("element of the list is not a string")
113 for i in wordList:
114     self.addWord(i)
115 self.build()
116     def clear(self):
117     '''czysci automat i drzewo'''
118 self.words = []
119 self.n = Node()
120 self.built = False
121     def search(self, tekst, returnSet=False):
122     '''wyszukuje wzorce w zmiennej tekst, zwraca string z wiadomoscia
123     o wynikach
124     domyslne argument returnSet mowi o formacie zwracanej wartosci
125     jesli returnSet jest False, to zwracamy string z informacjami
126     jesli returnSet jest True, to zwracamy zbior krotek o dlugosci
127     dwa, krotka
128     zawiera pozycje, na ktorej znalazla slowo, oraz indeks
129     slowa
130     jesli tekst nie jest zmienna string, to rzuca
131     AhoCorasickException'''
132 if not isinstance(tekst, (str, unicode)):
133     raise AhoCorasickException("argument is not a string")
134 node = self.n
135 if not returnSet: message = ""
136 else: message = set()
137 dl = len(tekst)
138 for i in range(dl):
139     while not node.getAim(tekst[i]):
140     node = node.getFail()
141     node = node.getAim(tekst[i])
142     if node.getAccept() != set():
143     zbior = node.getAccept()
144     for j in zbior:
145         if returnSet:
146             message.add((i, j))
147         else:
148             message += "Found \" + self.words[j] + \" in position \" + str(i) +
149             "\n"
150 if not returnSet and message == "":
151     message = "Nothing found\n"
152 if not returnSet: message = message[:len(message)-1]
153 return message

```

## 4 Testy

## 5 Podsumowanie

## 6 Bibliografia

- [http://pl.wikipedia.org/wiki/Algorytm\\_Aho-Corasick](http://pl.wikipedia.org/wiki/Algorytm_Aho-Corasick)
- <http://www.cs.uku.fi/~kilpelai/BSA05/lectures/slides04.pdf>
- [http://pl.wikipedia.org/wiki/Deterministyczny\\_automat\\_skończony](http://pl.wikipedia.org/wiki/Deterministyczny_automat_skończony)
- <http://pl.wikipedia.org/wiki/Pods\OT4\lowo>
- <https://wiki.python.org/moin/TimeComplexity>