Automat Aho-Corasik — omówienie —

Piotr Mikołajczyk

1 Wprowadzenie

Niech $\mathcal{W} = \{w_1, ..., w_k\}$ będzie zbiorem wzorców (niepustych słów nad alfabetem Σ). Otrzymawszy słowo $T \in \Sigma^m$ chcemy wyznaczyć wszystkie indeksy z [m], na których w T zaczyna się pewien wzorzec z \mathcal{W} . Naturalnym rozwiązaniem jest k-krotne użycie jednego z wydajnych (liniowych) algorytmów do szukania pojedynczego wzorca. Podejście takie oczywiście jest poprawne, niestety skutkuje złożonością czasową $\Theta(km)$. Okazuje się, że na podstawie \mathcal{W} jesteśmy w stanie skonstruować w czasie $\Theta(\sum_{i \in [k]} |w_i|)$ odpowiednią strukturę, która będzie potrafiła znaleźć naraz wszystkie wystąpienia wzorców oglądając T wyłącznie raz. Przed nami: automat Aho-Corasick.

2 Idea

Pomysł polega na zbudowaniu odpowiedniego deterministycznego automatu skończonego \mathcal{A} . Czytając T będziemy podróżować zgodnie z krawędziami \mathcal{A} . W każdym stanie s będziemy mieć zapisaną listę tych wzorców z \mathcal{W} , które kończą się w s – tzn. aby znaleźć się w stanie s musieliśmy przejść skądś krawędziami etykietowanymi kolejnymi literami pewnego wzorca. Spacerowanie po \mathcal{A} wzdłuż T możemy interpretować jako poruszanie się jednocześnie po wszystkich słowach z \mathcal{W} .

Zaczniemy konstrukcję od utworzenia drzewa trie na podstawie \mathcal{W} . Następnie dodamy do każdego stanu informacje o wzorcach które się w nim kończą i podniesiemy funkcję przejścia z częściowej do totalnej. Na koniec uprościmy automat. Po tych zabiegach otrzymamy strukturę, dzięki której wyszukiwanie wielu wzorców okaże się bardzo proste:

```
def find_occurrences(self, text, n):
    state = self._root
    for i in range(1, n + 1):
        state = state.next(text[i])
        for keyword, keyword_len in state.output():
        yield i - keyword_len + 1, keyword
```

Warto zwrócić uwagę na fakt, że jeśli dany zbiór wzorców chcemy wyszukać w wielu różnych tekstach $T_1, ..., T_t$, to wystarczy tylko raz skonstruować automat Aho-Corasick i użyć go niezależnie na każdym z T_j .

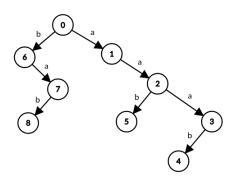
3 Konstrukcja

Uwaga: sednem naszych rozważań jest opis sposobu, w jaki można przechodzić pomiędzy stanami automatu. Wygodnie nam będzie utożsamić ze sobą różne podejścia realizacji – równoważne dla nas będą pojęcia zbioru etykietowanych krawędzi (skierowanych) pomiędzy dwoma stanami, funkcji mapujących pary (stan, symbol) w stan oraz lokalne w stanach tablice routingu.

3.1 Drzewo trie

Mając zbiór słów W nad alfabetem Σ możemy w prosty sposób zbudować tzw. drzewo trie, ozn. \mathcal{T}_{W} . Jest to skierowane, ukorzenione drzewo, którego krawędzie etykietowane są symbolami z Σ . Każde słowo $w \in W$

ma odpowiadającą mu ścieżkę z korzenia, etykietowaną kolejnymi literami w. Jeśli $w \in \mathcal{W}$ nie jest prefiksem żadnego innego słowa z \mathcal{W} , to jego ścieżka kończy się w liściu. Inaczej ujmując, istnieje bijekcja pomiędzy ścieżkami w $\mathcal{T}_{\mathcal{W}}$ a wszystkimi (różnymi) prefiksami słów z \mathcal{W} . Krawędź etykietowaną symbolem a pomiędzy s_1 a s_2 możemy kodować jako trójkę (s_1, a, s_2) .



Rysunek 1: $\mathcal{T}_{\{aaab,aab,bab,ba\}}$

3.2 Plan

Zacznijmy temat konstrukcji od zdefiniowania funkcji, które będą nam potrzebne. Niech W oznacza zbiór wzorców, \mathcal{T} oznacza drzewo trie dla W (skrót dla \mathcal{T}_{W}), \mathfrak{r} oznacza korzeń \mathcal{T} , \mathcal{S} oznacza zbiór wierzchołków \mathcal{T} (czvli docelowo zbiór stanów \mathcal{A}).

• goto :: $S \times \Sigma \mapsto S$ – ta funkcja w zasadzie reprezentuje zbiór krawędzi \mathcal{T} ; dodatkowo jednak, dla każdego symbolu w Σ , który nie etykietuje żadnej krawędzi wychodzącej z \mathfrak{r} , dodajemy pętlę w \mathfrak{r} ; czyli:

$$\mathtt{goto}(s_1,a) = \begin{cases} s_2 & \text{jeśli } (s_1,a,s_2) \in E[\mathcal{T}] \\ \mathfrak{r} & \text{jeśli } s_1 = \mathfrak{r} \ \land \ \forall_s \ (\mathfrak{r},a,s) \notin E[\mathcal{T}] \end{cases}$$
 None

- output :: $S \mapsto \wp(W)$ jak zostało już wspomniane, chcemy w każdym stanie s przechowywać listę wzorców, których wystąpienie w T będzie poświadczone odwiedzeniem stanu s; w sporym zakresie możemy to obliczyć już przy konstrukcji T (np. wracając do rysunku 1, output(1) = \varnothing , output(5) = $\{aab\}$, output(7) = $\{ba\}$); będziemy musieli jeszcze jednak obsłużyć przypadki, gdy jeden wzorzec będzie sufiksem drugiego (np. output(4) = $\{aaab, aab\}$)
- fail :: S → S ostatecznie A będzie skonstruowane na wierzchołkach T więc już teraz możemy spróbować zinterpretować obecność w stanie s po przeczytaniu j-tego symbolu T otóż etykiety na ścieżce pomiędzy r a s tworzą najdłuższy możliwy prefiks któregoś ze wzorców z W, który da się dopasować w T kończąc na pozycji j; funkcja goto jest częściowa, a więc dla niektórych stanów i niektórych symboli nie mamy zdefiniowanego przejścia (na razie otrzymamy None); aby nie utknąć i zachować poprawność interpretacji, to fail(s) musi wskazywać na taki stan s', żeby słowo powstałe ze ścieżki r → s' było najdłuższym możliwym sufiksem słowa odpowiadającemu r → s; w dalszej części utożsamiać będziemy ścieżki ze słowami;

Chodzenie krawędziami fail opiera się na identycznym pomyśle co rekurencyjne używanie tablic Border w algorytmie Morrisa-Pratta. Warto również zauważyć, że fail nie ma sensu definiować dla \mathfrak{r} .

next :: S×Σ → S - trzy zdefiniowane funkcje powyżej zupełnie wystarczają, aby nasz docelowy automat robił to, czego pragniemy; jak jednak łatwo zauważyć, próba przeczytania jednego znaku z T może wymagać więcej niż jednego przejścia krawędzią korzystając z goto oraz fail - możemy musieć schodzić wzdłuż fail wielokrotnie; aby uniknąć tego, pod sam koniec utworzymy funkcję next, która skompresuje nam tego typu ścieżki i zastąpi goto oraz fail

3.3 Obliczanie goto, output, fail

Jak już wiemy co chcemy policzyć, to możemy przejść do tego jak to zrobić. Funkcję goto oraz część output, jak już zauważyliśmy, możemy obliczyć konstruując $\mathcal{T}_{\mathcal{W}}$:

```
def _construct_goto(self, keywords, alphabet):
     for k, k_len in keywords:
       self._enter(k, k_len)
3
     for a in alphabet:
5
       if self._root.goto(a) is None:
          self._root.update_goto(a, self._root)
   def _enter(self, keyword, keyword_len):
     current_state = self._root
10
11
12
     while j < keyword_len and current_state.goto(keyword[j]) is not None:
13
       current_state = current_state.goto(keyword[j])
14
       j += 1
15
16
     for a in keyword[j:keyword_len + 1]:
17
       next_state = AhoCorasickAutomaton.Node()
18
       current_state.update_goto(a, next_state)
       current_state = next_state
20
21
     current_state.append_outputs([(keyword, keyword_len)])
22
```

Funkcja def _enter(...) uzupełnia częściowo skonstruowane drzewo trie o ścieżkę dla podanego wzorca, w razie potrzeby tworząc nowe węzły. W pętli w wierszu 5. uzupełniamy definicję goto dla korzenia.

Funkcję fail powinniśmy obliczać dla coraz bardziej odległych od r węzłów – podobnie jak w algorytmie Morrisa-Pratta długości najdłuższych borderów obliczaliśmy na podstawie wartości dla krótszych słów:

```
def _construct_fail(self, alphabet):
     q = Queue()
     for s in (self._root.goto(a) for a in alphabet):
3
       if s != self._root:
          q.put(s)
5
          s.update_fail(self._root)
     while not q.empty():
       current = q.get()
       for a, child in ((a, current.goto(a)) for a in alphabet):
10
          if child is not None:
11
            q.put(child)
12
13
            fallback = current.fail()
14
            while fallback.goto(a) is None:
              fallback = fallback.fail()
16
            child_fallback = fallback.goto(a)
18
            child.update_fail(child_fallback)
19
            child.append_outputs(child_fallback.output())
20
```

Przechodzimy zatem po \mathcal{T} za pomocą algorytmu BFS. Będąc w wierzchołku current, widząc krawędź etykietowaną przez a do wierzchołka child, chcemy obliczyć wartość fail(child). Szukamy więc najdłuż-

szego sufiksu słowa generowanego przez $\mathfrak{r} \leadsto \mathtt{child}$ obecnego w \mathcal{T} . Będzie to najdłuższy sufiks ścieżki $\mathfrak{r} \leadsto \mathtt{current}$, który można przedłużyć o symbol a.

Znalazłszy odpowiedni wierzchołek – child_fallback – powiększamy zbiór output(child) o wartości z output(child_fallback) – z oczywistej indukcji wynika, że output(child_fallback) zawiera już wszystkie wzorce, które są sufiksami $\mathfrak{r} \leadsto \text{child}_f$ allback.

3.4 Obliczanie next

Jeśli goto była określona dla $s \in \mathcal{S}$ oraz $a \in \Sigma$, to $\mathtt{next}(s, a) = \mathtt{goto}(s, a)$. W przeciwnym wypadku będąc w s i chcąc przejść symbolem a, powinniśmy schodzić krawędziami fail tak długo, aż znajdziemy się w s', dla którego $\mathtt{goto}(s', a) \neq \mathtt{None}$. Jest to podobny zabieg co poprawka Knutha do algorytmu Morrisa-Pratta.

```
def _construct_next(self, alphabet):
     q = Queue()
2
     for a in alphabet:
3
       a_child = self._root.goto(a)
       self._root.update_next(a, a_child)
       if a_child != self._root:
          q.put(a_child)
     self._root.use_only_next()
     while not q.empty():
10
       current = q.get()
11
       for a, child in ((a, current.goto(a)) for a in alphabet):
12
          if child is not None:
13
            q.put(child)
14
            current.update_next(a, child)
15
          else:
16
            fallback = current.fail()
            current.update_next(a, fallback.next(a))
18
       current.use_only_next()
```

Instrukcje w wierszach 8. oraz 19. pozwalają zwolnić pamięć zapominając o wartościach funkcji goto i fail. Formalnie dopiero teraz, struktura $\mathcal A$ ze zbiorem stanów $\mathcal S$ oraz funkcją przejścia next jest deterministycznym automatem skończonym.

4 Analiza

4.1 Analiza poprawności

Z tego, co już zostało przedstawione, powinna wynikać poprawność zastosowanej metody. Aby jednak formalnie jej dowieść, wystarczy postawić trzy proste lematy:

Lemat. Niech $s, t \in \mathcal{S}$ i niech u, v będą słowami generowanymi przez ścieżki $\mathfrak{r} \leadsto s$ i $\mathfrak{r} \leadsto t$. Wówczas $fail(s) = t \iff v$ jest najdłuższym właściwym sufiksem u, który jest prefiksem pewnego $w_i \in \mathcal{W}$.

Lemat. Stowo k należy do output $(s) \iff w \in W$ oraz k jest sufiksem stowa generowanego przez $\mathfrak{r} \leadsto s$.

Lemat. Po przeczytaniu j znaków z T znajdujemy się w stanie $s \iff slowo$ generowane przez $\mathfrak{r} \rightsquigarrow s$ jest najdłuższym sufiksem T[1..j], który jest prefiksem pewnego $w \in \mathcal{W}$.

Ich uzasadnienia są prostymi dowodami rekurencyjnymi zapisanymi we fragmentach kodu powyżej.

4.2 Analiza złożoności

Konstrukcja całego automatu polega na sekwencyjnym wywołaniu procedur obliczających funkcje przejścia:

```
class AhoCorasickAutomaton:
def __init__(self, keywords, alphabet):
self._root = AhoCorasickAutomaton.Node()
self._construct_goto(keywords, alphabet)
self._construct_fail(alphabet)
self._construct_next(alphabet)
```

Bardzo łatwo zauważyć, że każdą z nich wykonać można w czasie $\Theta(\sum_{i \in [k]} |w_i|)$ – za każdym razem przechodzimy po całym drzewie \mathcal{T} odwiedzając każdy wierzchołek dokładnie raz oraz wykonując $\Theta(1)$ operacji w każdym z nich. Zakładamy również, że łączenie zbiorów output wykonuje się w czasie stałym (jak np. listy wiązane). Warto dodać, że moc Σ wlicza się w zapotrzebowanie (zarówno pamięciowe jak i czasowe) liniowo – rozmiar grafu jest z góry ograniczony przez $|\Sigma| \cdot \sum_{i \in [k]} |w_i|$ (z każdego wierzchołka może wychodzi $|\Sigma|$ krawędzi). Konstrukcja $\mathcal A$ oraz wyszukiwanie w T odbywają się w czasie liniowo proporcjonalnym do wielkości $\mathcal T$.

Pozostaje zastanowić się, ile czasu zajmie przeczytanie tekstu T (|T|=m) i wypisanie wszystkich wystąpień wzorców. Moc Σ traktujemy jak stałą.

Najpierw rozważmy sytuację, w której nie zgłaszamy żadnego wystąpienia wzorca (możemy np. tylko zaznaczyć pozycje T o których wiemy, że na nich coś się kończy):

- używając funkcji goto oraz fail wykonamy co najwyżej 2m przejść w A po przeczytaniu dowolnego symbolu przejdziemy dokładnie raz według funkcji goto i być może wielokrotnie krawędziami fail; łatwo jednak zaobserwować, że każde przejście fail skraca odległość z obecnego wierzchołka do r, a z r zawsze można wyjść używając goto; zatem nie możemy przejść krawędziami fail więcej razy niż przeszliśmy goto;
 - można również spojrzeć na to w ten sposób, że przechodząc automatem \mathcal{A} wzdłuż T, pamiętamy dwa wskaźniki w T słowo pomiędzy nimi to właśnie to, które jest generowane przez ścieżkę z korzenia do obecnego wierzchołka; przejście goto przesuwa 'prawy' wskaźnik o 1 w prawo, przejście fail przesuwa 'lewy' wskaźnik o dodatnią liczbę pozycji w prawą, ale na tyle małą by nie przeskoczyć prawego;
- \bullet używając funkcji next, przy każdym przeczytanym symbolu zTwykonujemy zawsze jedno przejście mruchów

Problem może się pojawić, gdy będziemy chcieli wypisywać wszystkie wystąpienia wzorców. Pesymistycznym przykładem jest $\mathcal{W} = \{a, a^2, a^3, ..., a^k\}$ oraz $T = a^m$. Liczba wystąpień będzie $\Theta(km)$. Jednakże, nie jest to wada zastosowanego algorytmu – każdy inny algorytm zgłaszający wszystkie wystąpienia wzorców musiałby wykonać co najmniej tyle samo pracy.

Zatem możemy określić złożoność czasową wyszukiwania wielu wzorców w tekście za pomocą automatu Aho-Corasick jako liniową w stosunku do sumy długości wszystkich wzorców oraz ilości ich wystąpień.

5 Zastosowanie do wzorców 2D

Okazuje się, że bardzo prosto można użyć automatu Aho-Corasick do wyszukiwania wzorców dwuwymiarowych. W skrócie:

- 1. niech Tbędzie kwadratową tablicą rozmiaru $m^2,\ P$ kwadratową tablicą rozmiaru n^2 o różnych kolumnach, n < m
- 2. niech $\mathcal W$ będzie zbiorem kolumn P
- 3. automatem Aho-Corasick szukamy W w każdej z kolumn T; jeśli fragment j-tej kolumny od indeksu i pokrywa się z k-tą kolumną P, to niech $T_P[i][j] = k$
- 4. algorytmem KMP szukamy w każdym wierszu T_P ciagu (1, 2, ..., n)

Całość działa w czasie liniowym od wielkości wejścia i rozmiaru zbioru z którego pochodzą elementy w T oraz P.