

Programming Technology: Xử lý chuỗi với phép tách từ

Tran Viet Khoa
tvkhoa.husc@gmail.com

Hue University— Ngày 30 tháng 3 năm 2020

Mục lục

1	Giới thiệu	2
2	Cơ bản về lớp string	2
2.1	Khai báo	2
2.2	Nhập dữ liệu từ <code>stdin</code> , <code>stream</code>	2
2.3	Độ dài của chuỗi	3
2.4	Phép ghép chuỗi	3
2.5	Phép truy cập và duyệt	3
3	Bài toán tách từ	4
3.1	Xử lý chuỗi bằng các hàm <code>algorithm</code>	4
3.2	Xử lý chuỗi bằng các hàm <code>regex</code>	6

1 Giới thiệu

Bài toán tách từ (Tokenizing a string) là một bài toán cơ bản và được ứng dụng nhiều trong các bài toán thực tế. Ví dụ như:

1. Xử lý ngôn ngữ tự nhiên.
2. Xử lý văn bản, tìm kiếm và nhận dạng từ.
3. xử lý cú pháp lệnh của trình biên dịch.

Những ngôn ngữ tiên tiến như Python, Java cung cấp các thư viện rất mạnh mẽ cho xử lý trên. Đối với ngôn ngữ C các hàm `strtok()` và `strtok_r()` của thư viện `<string.h>` xử lý cho công việc này. Riêng đối với lớp chứa `string` của STL C++ thì ta cũng có thể sử dụng các hàm trên bằng cách chuyển qua chuỗi của ngôn ngữ C và gọi hàm trên.

Bài hướng dẫn này trình bày cách xây dựng hàm xử lý bài tách từ của xâu thuộc lớp chứa `string` bằng cách dùng các hàm của thư viện `<algorithm>` và thư viện `<regex>` để cài đặt. Các cài đặt này rất hiệu quả và có tính công nghệ, tức là giúp bạn có thể tiếp cận hàm tương tự như vậy trên các ngôn ngữ lập trình khác.

Nội dung bao gồm:

- Cơ bản về lớp `string`.
- Cài đặt các bài toán tách từ bằng thư viện `algorithm`.
- Cài đặt các bài toán tách từ bằng thư viện `regex`.



Thông tin

+ Một số bài toán thực hành có thể dễ dàng cài đặt với mảng ký tự và với thuật toán đơn giản, tuy nhiên tính tổng quát là không cao.

2 Cơ bản về lớp `string`

2.1 Khai báo

- Khai báo thư viện:

```
1 #include <string>
2 using namespace std;
```

Hoặc thay thế bằng `#include <bits/stdc++.h>`

- Khai báo đối tượng.

```
1 // Khai báo đối tượng str.
2 string str;
3 // Khai báo đối tượng str bằng hàm dựng
4 string str("abc");
```

2.2 Nhập dữ liệu từ `stdin`, `stream`

- Sử dụng toán tử tải bội « với đối tượng `cin`, Ví dụ:

```
1 string firstName;
2 cin >> firstName; // lấy dữ liệu từ bộ đệm bàn phím stdin.
3 cout << firstName;
```

Nhược điểm của toán tử này nó coi ký tự khoảng như (tab `\t`, space `' '`) là ký tự kết thúc cho nên ta không quét được hết dữ liệu.

- Sử dụng hàm `getline()`, ví dụ:

```

1 string fullName;
2 getline (cin, fullName);
3 cout << fullName;

```

2.3 Độ dài của chuỗi

- Sử dụng hàm `length()`, Ví dụ:

```

1 string str = "ABCDEFGHIJKLMNOPQRSTUVWXYZ";
2 cout << str.length();

```



Thông tin

+ Vì lớp `string` là vector các ký tự cho nên hàm `size()` được xem như là bí danh (alias) của `length()`. Bạn có thể dùng hai hàm trên.

2.4 Phép ghép xâu

- Sử dụng toán tử tải bội `+`, ví dụ:

```

1 string firstName = "Khoa";
2 string lastName = "Tran Viet";
3 string fullName = firstName + " " + lastName;
4 cout << fullName;

```

- Sử dụng hàm `append()`, ví dụ:

```

1 string firstName = "Khoa";
2 string lastName = "Tran Viet";
3 string fullName = firstName.append(lastName);
4 cout << fullName;

```

2.5 Phép truy cập và duyệt

- Sử dụng toán tử tải bội `[]` hoặc hàm `at()`, ví dụ:

```

1 // Khai báo và dựng (hàm dựng) đối tượng.
2 string str ("Programming is Fun");
3 // Duyệt in bằng chỉ mục với toán tử [].
4 // tương tự mảng chỉ số i=0..str.length()-1.
5 for (unsigned i=0; i<str.length(); ++i)
6     cout << str[i];
7 // Duyệt in bằng chỉ mục với hàm at().
8 for (unsigned i=0; i<str.length(); ++i)
9     cout << str.at(i);

```

- Sử dụng con trỏ `iterator` để duyệt xâu với các hàm `begin()`, `end()`, `rbegin()`, `rend()`, ví dụ:

```

1 // Khai báo và khởi tạo chuỗi.
2 string str = "Dogma I am God";
3 // Khai báo biến lặp.
4 string::iterator it;
5 // Khai báo biến lặp ngược.
6 string::reverse_iterator it1;
7 // Duyệt in các ký tự của chuỗi nhờ biến lặp.
8 cout << "The string using forward iterators is : ";
9 for (it=str.begin(); it!=str.end(); it++)
10     cout << *it;
11 cout << endl;
12 // Duyệt chiều ngược lại bằng con trỏ ngược.

```

```

13 cout << "The reverse string using reverse iterators is : ";
14 for (it1=str.rbegin(); it1!=str.rend(); it1++)
15     cout << *it1;
16 cout << endl;

```

3 Bài toán tách từ

Tách từ từ chuỗi (string tokenization) là phép tách một chuỗi thành các phần gọi là token, ví dụ str="Programming is fun" sẽ được tách thành "Programming", "is" và "fun". Phép toán này được dùng rất nhiều trong chương trình dịch nhằm phân tích một câu cú pháp lệnh để biết được đâu là biến, hàm, hằng số, phép toán.

Nếu một xâu được tiền xử lý và đẹp như xâu ví dụ trên thì bài toán này rất đơn giản. Trong thực tế, ta gặp nhiều bài toán xử lý chia xâu thành các token phức tạp hơn, ví dụ xét xâu: str=" *Good friends, good books, and a sleepy conscience: this is the ideal life. — Mark Twain*". Nhận thấy xâu này có các ký tự dấu phẩy, chấm, hai chấm, gạch ngang để phân tách từ ngoài ký tự trắng thông dụng. Nếu câu viết không đúng cú pháp thì việc tách càng khó, ví dụ: str=" hello again..... sea you late...! ".

Thông thường đối với những xâu như trên, người ta sẽ mất một công sức cho công việc gọi là tiền xử lý làm đẹp xâu, ví dụ xóa ký tự trắng thừa, xóa ký tự sai cú pháp. Và như vậy bài toán trở nên rất phức tạp.

Bài toán đặt ra cho các bạn là: Cho một xâu *str*, việc phân cách từ trong xâu bằng cách sử dụng các ký tự đặc biệt (người ta gọi là delimiter) có thể được liệt kê trước như dấu chấm, phẩy, hai chấm, chấm phẩy, than, hoa thị,.. và cả dấu cách (space). Hãy tách các từ trên.

3.1 Xử lý xâu bằng các hàm algorithm

Bài tập 1: Tìm từ dài nhất

Cho một xâu ký tự gồm các ký tự Alphabet ở dạng chữ hoa và chữ thường, ngoài ra còn có các ký tự đặc biệt sau: apostrophe,full stop,comma,semicolon,colon, ký tự space và newline. Trong các ký tự trên chỉ có ký tự apostrophe là dùng để nối từ, các ký tự còn lại nếu đúng cú pháp là để phân biệt từ. Hãy lập trình tìm từ có chiều dài lớn nhất trong xâu cho trước.

Input

+ Gồm nhiều dòng, mỗi dòng là một xâu ký tự có tối đa 80 ký tự.

Output

+ Ứng với mỗi dòng in chiều dài lớn nhất của các từ tương ứng trong các xâu.

Samples

input	output
If you're smart you will solve this right.	6

Tham khảo:

Bài tập PT074, <http://oj.hueuni.edu.vn/practice/problem/174/details>

Bài toán này được giải bởi hai bước chính sau:

1. Đọc xâu *str* và tách ra từng từ dựa vào việc nhận diện ký tự phân tách. Kết quả lưu vào vector.
2. Quá đơn giản cho việc tìm xâu có độ dài lớn nhất trong một vector chứa chúng.

Bài toán cần các hàm sau để giải.

1. Hàm `find_first_not_of()` - Hàm này dùng tìm kiếm sự vắng mặt (absence) tức là không có ký tự ở trong chuỗi. Tìm kiếm chuỗi cho ký tự đầu tiên không khớp với bất kỳ ký tự nào được chỉ định trong các đối số của nó. Khi pos được chỉ định, tìm kiếm chỉ bao gồm các ký tự tại hoặc sau vị trí pos, bỏ qua các ký tự trước vị trí pos đó.

Kết quả trả về là vị trí đầu tiên không phù hợp, nếu không tìm thấy trả về hằng `string::npos`. Kiểu của trị trả về là `size_t` (số nguyên không âm).

Xét ví dụ mã nguồn sau:

```
1 string str ("look for non-alphabetic characters...");
2 size_t found = str.find_first_not_of("abcdefghijklmnopqrstuvwxyz ");
3 if (found!=string::npos){
4     cout << "The first non-alphabetic character is " << str[found];
5     cout << " at position " << found << '\n';
6 }
```

7

8 Output:

9 The first non-alphabetic character is - at position 12.

Như vậy với hàm trên ta lặp đi lặp lại việc tìm các ký tự phân cách, xác định được vị trí của chúng để sao chép các chuỗi con lưu vào vector nhưng cần thêm một địa chỉ nữa và được biết đến với hàm thứ hai.

2. Hàm thứ hai là `find_first_of()` - Hàm tìm ký tự trong chuỗi. Tìm kiếm chuỗi cho ký tự đầu tiên khớp với bất kỳ ký tự nào được chỉ định trong các đối số của nó.

Kết quả trả về là vị trí đầu tiên tìm thấy, nếu không tìm thấy trả về hằng `string::npos`. Kiểu của trị trả về là `size_t` (số nguyên không âm).

Xét ví dụ mã nguồn sau:

```
1 string str ("Please, replace the vowels in this sentence by asterisks.");
2 // Tìm các nguyên âm.
3 size_t found = str.find_first_of("aeiou");
4 while (found!=string::npos){
5     // Tìm thấy và thay thế bằng dấu *.
6     str[found]='*';
7     //Tìm tiếp cho vị trí tiếp theo. lặp cho đến khi gặp ::npos.
8     found=str.find_first_of("aeiou",found+1);
9 }
```

10 cout << str << endl;

11

12 Output:

13 Pl**s*, r*pl*c* th* v*w*ls *n th*s s*nt*nc* by *st*r*sks.

3. Cuối cùng là hàm copy xâu con `string substr (size_t pos = 0, size_t len = npos) const;` - Gọi là hàm tạo xâu con. Trả về một đối tượng chuỗi được xây dựng mới với giá trị của nó được khởi tạo từ một bản sao (copy) của một chuỗi con của đối tượng chuỗi gọi hàm.

Chuỗi con là một phần của đối tượng bắt đầu tại vị trí ký tự `pos` và kéo dài các ký tự với chiều dài `len` (hoặc cho đến khi kết thúc chuỗi, tùy theo điều kiện nào đến trước).

Xét ví dụ mã nguồn sau:

```
1 string str="We think in generalities, but we live in details.";
2 string str2 = str.substr (3,5); // Cắt được xâu con "think"
3 std::size_t pos = str.find("live"); // Trả về địa chỉ của xâu "live" trong str.
4 string str3 = str.substr (pos); // cắt được xâu con "live" cho đến cuối.
5 cout << str2 << ' ' << str3 << '\n';
```

6 Output:

7 think live in details.

Với sự kết hợp của ba hàm trên, thuật toán cho bài toán trình bày (accept) với mã nguồn:

a.cpp

```
1 #include <bits/stdc++.h>
2 using namespace std;
3 using VS = vector<string>;
4
5 VS split(const string& text, const string& delim){
6     VS tokens;
7     size_t start = text.find_first_not_of(delim), end = 0;
8     while((end = text.find_first_of(delim, start)) != string::npos){
9         tokens.push_back(text.substr(start, end - start));
10        start = text.find_first_not_of(delim, end);
11    }
12    if(start != string::npos) tokens.push_back(text.substr(start));
13    return tokens;
14 }
15 int main()
16 {
17     // Khai báo các ký tự phân tách.
18     string delims="; :!, ";
19     string s;
20     //Vòng lặp đọc lần lượt các xâu.
21     while (getline(cin, s)){
22         // Phân tách và lưu vào vector.
23         VS tokens = split(s, delims);
24         //Tìm max.
25         unsigned int max=0;
26         for (auto &s: tokens) {
27             if (s.size() > max) max =s.size();
28         }
29         cout<< max <<endl;
30         //Xóa vector dọn chỗ cho lần lặp tiếp.
31         tokens.clear();
32     }
33     return 0;
34 }
```



Nhận xét

+ Với bài toán này, nhiều bạn nghĩ đơn giản hơn nhiều và đã accept được, ví dụ duyệt qua các ký tự xâu và kiểm tra xem có phải ký tự alphabet (isalpha) hoặc ký tự mã #39, nếu thỏa thì tăng biến đếm. so sánh với max để in kết quả. Nhưng nếu xâu không phải alphabet thì bài toán không giải quyết được.

+ Ý thứ hai quan trọng là bài toán lưu được các từ để dùng xử lý sau này.

3.2 Xử lý xâu bằng các hàm regex

Bài toán thứ hai hoàn toàn tương tự, nội dung như sau:

Bài tập 2: Số từ và từ dài nhất

Một chuỗi ký tự có thể bao gồm các ký tự Alphabet, ký tự trắng, dấu phẩy (,), dấu chấm (.), dấu hai chấm(:). Một từ là một chuỗi con chỉ chứa các ký tự Alphabet, ngoài ra hai ký tự liền trước và liền sau của chuỗi con đó (nếu có) sẽ không phải là ký tự Alphabet. Hãy lập trình tìm từ có chiều dài lớn nhất trong chuỗi và số từ của chuỗi đó.

Input

+ Gồm một dòng duy nhất là một chuỗi ký tự có tối đa 200 ký tự.

Output

+ Dòng thứ nhất in chiều dài lớn nhất của các từ trong chuỗi.
+ Dòng thứ hai in số từ của chuỗi.

Samples

input	output
Programming Java in oj	11 6

Tham khảo:

Bài tập PT075, <http://oj.hueuni.edu.vn/practice/problem/175/details>



Nhận xét

- + Quá dễ cho bài toán này khi nội dung hoàn toàn tương tự, thậm chí các ký tự phân cách lại ít hơn.
- + Accept thôi, chỉnh lại tí nộp.

Mã nguồn:

```
bNotAc.cpp
1 #include <bits/stdc++.h>
2 using namespace std;
3 using VS = vector<string>;
4 //Sử dụng hàm đã ac ở câu a.
5 VS split(const string& text, const string& delim){...
6 }
7 int main()
8 {
9     string delims = " ,.:";
10    string s;
11    getline(cin, s);
12    tokens = split(s, delims);
13    unsigned int max=0;
14    for (auto &s: tokens) {
15        if (s.size() > max) max =s.size();
16    }
17    cout<< max <<endl;
18    cout<<tokens.size();
19    return 0;
20
21 }
```



Thách thức

+ Chiều nộp hơn 10 lần vẫn WA, mới phát hiện ra testcase bài 174 yếu hơn, thiếu những trường hợp là một loạt ký tự phân cách dính liền nhau và debug mãi không được. Thầy nhờ mấy em debug tiếp, ví dụ sau là sai: "Hoc, hoc nua, hoc mai... hehe", kết quả ra: 4, 6 mà thuật mình ra 6, 6. Sửa một hồi ra 5, 6 gần đúng, toát mồ hôi. Ở đây cũng gặp trường hợp chạy máy mình báo đúng, không dùng DevC mà dùng command Line, bực quá chạy trên IDEone lúc đúng, lúc sai nhưng nộp là sai.

+ Cuối cùng nghĩ phương án khác, hồi trước có làm về xử lý chuỗi nhớ có sử dụng cái là regex kiểu như compiler trên shell linux. Thế là nhờ Giáo sư Goggle, tìm và thử cuối cùng cũng accept được câu b. Có điều giải thích nó cũng hơi mệt. mà hôm nay ngồi nguyên một ngày. Nhờ các bạn đọc hộ và phân biện.

Mã nguồn:

```
bAccept.cpp
1 #include <bits/stdc++.h>
2 using namespace std;
3
4 int main()
5 {
6     string s;
7     getline(cin,s);
8     regex words_regex("[^\\s.,;!?]+");
9     auto words_begin = sregex_iterator(s.begin(), s.end(), words_regex);
10    auto words_end = sregex_iterator();
11
12    unsigned int max=0,c=0;
13    for (sregex_iterator i = words_begin; i != words_end; ++i){
14        // Nhắc quá cũng không lưu vector, chú (*i).str() chính là các từ được
    tách ra.
15        if ((*i).str().length() >max)
16            max = (*i).str().length();
17        c++;
18    }
19    cout<<max<<endl;
20    cout<<c<<endl;
21    return 0;
22 }
```