

< 研究ノート >

正規表現について

その 4 ツールエンジンの効率化

坂本 義行・江戸 浩幸

Tool Efficiency on Regular Expression

Yoshiyuki SAKAMOTO and Hiroyuki EDO

概要

前三編において、Perl, sed, grep, Tcl, Python, Expect, Emacs など、よく使われるツールに実装されている正規表現エンジンの長所を探ると同時に、その短所を避けるための工夫について検討をしてきた。

今回はその第 4 段階として、さらに「詳説正規表現」を読み進め、ツールエンジンが従来型 NFA か POSIX NFA を調べるテスト、「ループ展開」の技術について、どのような正規表現を構築すれば、検索において効率化が図られるのか。C 言語で用いられるコメント文中で最高効率の検索を行うための正規表現について、実験を含めて詳細な検討が加えられた技術についての知見を報告する。

キーワード：正規表現、ツールエンジン、ループ展開、C コメント、特殊要素

1 エンジンの種類のテスト

ツールの正規表現エンジンをテストする作業には、いくつかの手順がある。

1.1 NFA と DFA の基本テスト

理論的には、エンジンの基本種類を調べるためには、永久に終わらないマッチについて調べればよい。

awx についてのテスト

```
echo =XX==...== | awx "/X(.+)*X/
{ print }
```

もし =XX==... の行が即座に出力され

ば、その awx は (大半のものと同じく) DFA を使っていることになる。もし手持ちの awx が mawx や Mortice Kern Systems 社の awx であれば、NFA が使われている。

1.2 従来型 NFA か POSIX NFA かを調べるテスト

従来型 NFA はマッチを見つけた瞬間に止まるため、永遠に続くような正規表現も、マッチがあれば通常すぐに終了する。

従来型か POSIX かのテストは、これまで述べて最適化を考えると難しくなる場合が多い。POSIX で即座に成功マッチを返すの

は、これが必要性認識の最適化を採用しているからである。つまり正規表現を「マッチするかしないか」のように二者択一的に使うと、マッチの範囲は重要ではなくなり、わざわざ最長マッチ（この場合は特定マッチ）を見つける必要もなくなる。

2 ループ展開

「ループ展開」と呼ぶ技術がある。これはある一般的な表現の速度を上げるのに力を発揮する。対象となるループは「(this|that|...)*」のようなパターンの表現で使われるスターによるものである。前出の永久マッチ「(\\.[^\\]+)*」がこのパターンに当てはまる。マッチ不成立を報告するのに半永久的な時間がかかることを考えれば、これは速度向上を試すのに適した例といえる。

2つの異なる方法で、この技術に到達することができる。

1. さまざまなマッチテスト例の中で、実際に「(\\.[^\\]+)*」のどの部分が成功したかを検討し、使われた部分パターンの流れを記録しておく。そしてそこから浮かび上がったパターンに基づいて効率的な表現を再構築する。
2. もう一つの方法は、マッチしたい構文を、さらに高い次元から見る方法である。

2.1 方法1：これまでの経験に基づいて正規表現を構築する

「(\\.[^\\]+)*」を分析する上で、いくつかのマッチ対象文字列を調べて、全体マッチの中で厳密にどの部分パターンが使われているかを知る。

次の場合は、「[^\\]+ \\.[^\\]+ \\.[^\\]+」となる。

"he said \\hi there" and left"

この例と表1では、パターンが分かりやすいように表現をハイライトしてある。

最初の4つの例を考えてみよう。クォート間の表現が「[^\\]+」で始まり、次にいくつかの「\\.[^\\]+」という並びが続く。これ全体を正規表現として表現し直すと、「[^\\]+(\\.[^\\]+)*」になる。

汎用「ループ」展開パターンの構築

ダブルクォート文字列をマッチする場合、「通常要素+(特殊要素 通常要素+)*」という汎用パターンに当てはまる。これに開きクォートと閉じクォートを加えると、

「[^\\]+(\\.[^\\]+)*」ができる。この表現の中にある2つの「[^\\]+」が文字列の先頭、および特殊な文字の後ろで通常の文字を要求している点にある。

では「[^\\]+*(\\.[^\\]+)*」をよく見てみよう。必要条件である開きクォートとその後続くかもしれない通常文字をマッチする。

表1 ループ展開表現例

対象文字列	実際に使われる表現
"hi there"	" [^" \\]+"
"just one \" here"	" [^" \\]+ \\.[^" \\]+"
"some \" quoted \" things"	" [^" \\]+ \\.[^" \\]+ \\.[^" \\]+"
"with \" a\" and \" b\"."	" [^" \\]+ \\.[^" \\]+ \\.[^" \\]+ \\.[^" \\]+"
" \" ok \" \"n"	" \\.[^" \\]+ \\.[^" \\]+"
"empty \" \" quote"	" [^" \\]+ \\.[^" \\]+"

2.2 本格的な「ループ展開」パターン

この種類の表現に用いる汎用パターン：
「開始 通常要素*(特殊要素 通常要素*)
* 終了」

永久マッチの回避

「`[^\\]*(\\[^\n]*)*`」が永久マッチにならないのは、非常に重要な以下の3つの点
が守られているからである。

特殊要素と通常要素の開始位置を決して交差させない

特殊要素と通常要素が絶対に同じ位置で
マッチしないようにすれば、対象文字列に
マッチできる特殊要素の「並び」が必ず一つ
だけ存在するはずである。この一つの並びを
テストする方が、何億もの並びテストするよ
りはるかに速く、永久マッチを回避すること
ができる。通常要素を「`[^\\]`」とし、特殊要素
「`\\`」とするここでの例は、この条件を満
たしている。これらの要素は決して同じ文字
でマッチを開始することはない。つまり、後
者は先頭にバックスラッシュを要求するの
に対し、前者は明示的にバックスラッシュを許
していない。

特殊要素は無とマッチしてはいけない

特殊要素はマッチするときには必ず最低1
個の文字とマッチしなければならない。
一つの特特殊要素でマッチできるテストは、複
数の特殊要素でマッチできてはならない。

Pascalの`{...}`型コメントとスペースからな
る文字列をマッチする場合、コメント部分と
マッチする正規表現は「`{[^\n]*}`」であり、
全体表現は「`([^\n]*)*`」となる。

特殊要素	通常要素
「 <code>\\</code> 」	「 <code>[^\n]</code> 」

これを今まで作っておいた「通常要素*
(特殊要素 通常要素*)*」というパターン
に組み込むと、「`([^\n]*)*(\\[^\n]*)*`」
になる。

注意すべき一般事項

「`(...)*`」のような繰り返し制御文字が複
数レベルで存在しているのは、重要な危険信
号であるが、こうした表現であってままた
く問題ないものも多い。例えば次のような
ものである。

- ・「`(Re:*)*`」 任意数の'Re:'をもつシーケ
ンスとマッチする。
- ・「`(*)\{0-9\}+`」 (スペースによって区
切られた)ドル金額とマッチする。
- ・「`(.*\n)+`」 一つ以上の行とマッチする。

2.3 方法2：高次元の視点で見する方法

「`\\`」は、時々現れるエスケープされた要
素の処理に用いられるだけである。いずれか
を許す選択を用いれば便利な正規表現ができ
るが、エスケープされた数少ない(一つもな
い場合の方が多い)要素のために、マッチ全
体の効率性を考えると無駄である。

「`^[^\n]+`」が終わる度に、閉じクォートが
別のエスケープのマッチかという前と同じ状
況がやってくるのである。

最終的に、「`^[^\n]+(\\[^\n]*)*`」とな
る。

2.4 方法3：クォート付きのインター ネットホスト名

ループ展開の技術を実現する方法を2つ述
べたが、3つ目ともいえる似た方法を紹介す
る。

これは、`prez.whitehouse.gov` や `www.yahoo.com`
といったドメイン名をマッチする正規表
現を構築するときである。これを「`[a-z]+(\\[a-z]+)*`」
と書くと、文字列が通常の「`^[^\n]`」
の並びで、これが「`\\`」で区切られ、全体が
「...」で囲まれていると考えれば、ループ展
開パターン「`^[^\n]+(\\[^\n]*)*`」を作る
ことができる。

2.5 考察

このダブルクォートの文字列の例についてまとめると、「`[^\\]*(\\.[^\\]*)*`」の表現には利点が多く、しかもほとんど欠点がないことがわかる。

欠点

- ・効率のために読みやすさを多少犠牲にしている（「`^(^\\|\\.)*`」）と比較

- ・効率のために維持性を多少犠牲にしている

利点

- ・速度の点で、新しい正規表現は、マッチが不可能な場合でも破綻することはない

3 C コメントの展開

C 言語では、コメントは `/*` で始まり、`*/` で終了し、複数行にまたがることもある。

3.1 複雑な正規表現

単に「`/*[^\n]**/`」とするだけではうまくいかない。これでは内部に有効な `*` を含む `/* * コメント部分 * */` などをマッチできない。対象となるコメントは `/*...*/` ではなく、`/x...x/` として考える。「`/*[^\n]**/`」は「`/x[^\n]*x/`」のように、ずっと読みやすくなる。

3.2 厳密でない考え方

1. 開きデリミタをマッチする
2. テキスト本体をマッチする：実際には「終了デリミタでない任意のもの」をマッチする
3. 閉じデリミタをマッチする

この場合、終了デリミタが単独の文字でない場合、成功しない。これを解決する一つの方法は、`x` を終了デリミタの始まりとして解釈する方法である。「終了デリミタでない任意のもの」は、

- ・`x` ではない任意のもの、すなわち「`[^\n]`」

- ・`x` で、その後ろにスラッシュが続かないもの、すなわち「`^\n/`」

のいずれかになる。すなわち、

「`/x[^\n]^\n/*x/`」ができる。

もう一つの方法は、`x` が先行する場合に限り、スラッシュを終了デリミタとみなす方法である。「終了デリミタでない任意のもの」は、

- ・スラッシュでない任意のもの、すなわち「`^\n/`」

- ・前に `x` を伴わないスラッシュ 1 個、すなわち「`^\n/`」

これより、「`/x[^\n][^\n]/*x/`」ができる。

しかし、いずれの方法もうまくいかない。前の方法は、`/xx_foo_xx/` の例で、後の例は、`/x/_foo_/x/` の例でうまくいかないことが分かる。

正しい表現に修正する

解決策は、「本当に必要なものを書く」という視点に立ち返ることである。

すなわち、「`/x[^\n]x+[^\n]/*x+/`」が得られる。

実際のコメント（`x` の代わりに `*` が使われる）では次のようにしなければならない。

「`/*([^\n]*\n+)[^\n]*/`」

となる。正規表現を読むことは簡単なことではない。肝心なのは、十分に気を配りながら、複雑な表現を頭の中で慎重に区切って考えることだ。

3.3 C コメント表現のループ展開

効率の点から、正規表現の展開について考えてみよう。表 2 は、ループ展開のパターンの中に挿入できる表現を示したものである。通常要素の並びが、常に終了デリミタの最初の文字で終わり、次に続く文字が末尾を構成しない場合に限り、特殊要素として取り扱いたい。

これらの要素を汎用展開パターンに挿入すると次のようになる。

表2 C言語コメント用のループ展開要素

「開き 通常*(特殊 通常*)* 閉じ」		
要素	対象	正規表現
開き	コメント冒頭	/x
通常*	コメントテキストから1個以上の'x'を含むところまで	[^x]*x+
特殊	終了スラッシュ以外 (かつ'x'ではない) のもの	[^/x]
閉じ	後続のスラッシュ	/

```
「/\x[\^x]*x+([\^x][\^x]*x+)**/」
```

現実を考える

「/\x[\^x]*x+([\^x][\^x]*x+)**/」は、すぐにそのまま使えるというわけではない。当然ながら、コメントは/x...x/ではなく/*...*/である。これは各xを\<*と置き換えれば簡単に直る(文字クラスであれば各xは*に置き換わる)。

```
「/\*[\^*]*\*+([\^*][\^*]*\*+)**/」
```

これを用いる際に関係してくる問題は、コメントがしばしば複数行にまたがることである。

4 スムーズに流れる正規表現

文字列定数内のコメントに似た要素がマッチされないようにするにはどうしたらよいかという問題が残されている。

4.1 マッチするように導くには

Tclの例:

```
set COMMENT {/\*[\^*]*\*+([\^*][\^*]*\*+)**/}
# コメントをマッチする正規表現
set DOUBLE {"\.[^\\" ]*" }
# ダブルクォート文字列をマッチする正規表現
regsub -all "$DOUBLE|$COMMENT" $text
{ } text
```

これでは、不十分で、次のように変更すべきである。

```
set COMMENT {/\*[\^*]*\*+([\^*][\^*]*\*+)**/}
# コメントをマッチする正規表現
```

```
set DOUBLE {"\.[^\\" ]*" }
# ダブルクォート文字列をマッチする正規表現
```

```
regsub -all {"$DOUBLE"}|$COMMENT" $text
{\1 } text
```

この変更点は

- ・マッチが文字列選択肢であった場合、\1 (Perlの\$1に対応するTclの置換文字列)を埋める丸括弧を加えた点。マッチがコメント選択肢の場合には\1は空のままになる。

- ・置換文字列を上に出てきた\1にした点。これにより、ダブルクォート文字列がマッチされると、置換は同一のダブルクォート文字列になる。

最後に、\tのようにシングルクォートで囲まれたCの文字定数を処理する必要がある。これには、単に丸括弧の中にもう一つの選択肢を加えれば済む。

```
set COMMENT
{/\*[\^*]*\*+([\^*][\^*]*\*+)**/}
# コメントをマッチする正規表現
set COMMENT2 {"//[\^n]*"}
# C++の//コメントをマッチする正規表現
set DOUBLE {"\.[^\\" ]*" }
# ダブルクォート文字列をマッチする正規表現
set SINGLE {"\.[^\\" ]*" }
# シングルクォート文字列をマッチする正規表現
```

表現

```
regsub -all
"($DOUBLE)|($SINGLE)|$COMMENT|
$COMMENT2" $text {\1} text
```

この基本的な構造はかなり単純である。これにより効率化が図られた。

4.2 よく制御された正規表現は速い

正規表現エンジンが進む流れを制御して、もっと速くマッチさせることができる。

例えば、選択肢がマッチするためには、先頭の文字はスラッシュ、シングルクォートまたはダブルクォートでなければならない。時間と労力のかかる方法でエンジンに模索させるのではなく、「[^'"]」を選択肢として加えて、エンジンに直接指示を出す。実際には、これらが連続している場合は一気に捕まえられるので、代わりに「[^'"]+」使うことにする。実験ではこれにより、1桁も短縮できた。

もっとすばやくマッチさせる方法がある。

- ・一般的な選択肢は「\$OTHER+」なので、これを丸括弧の中で最初に持ってくる
- ・人間がより上位の視点に立って局所的な最適化を行うことで、正規表現エンジンをマッ

チに導いているからである。

4.3 まとめ

2つの文字列部分パターンを次のように置き換える。

```
set DOUBLE { "[^"]*(\\ "[^"])*" }
set SINGLE { "[^']**(\\ '[^']*)*" }
```

さらに、25%短縮できた。図1はこれまでの短縮時間の成果をすべて示したものである。

5 エンジンの最高効率化

汎用の最適化よりも重要なこととして、NFA エンジンから最高効率を引き出すには、自分の目標を達成するために必要なものと不要なものを見極めると同時に、それを実現するために、NFA エンジンがどう動けばよいのかを把握することが鍵となる。

5.1 最適化の紆余曲折

```
"\bchar\b|\bconst\b|\bdouble\b...
\bsigned\b|\bunsigned\b|\bwhile\b"
```

上の文字列を使ってC言語の予約語を含む行を見つける場合を考える。

具体的な例として、Perlのソースプログラムからこれらの単語を含む行を数える短い

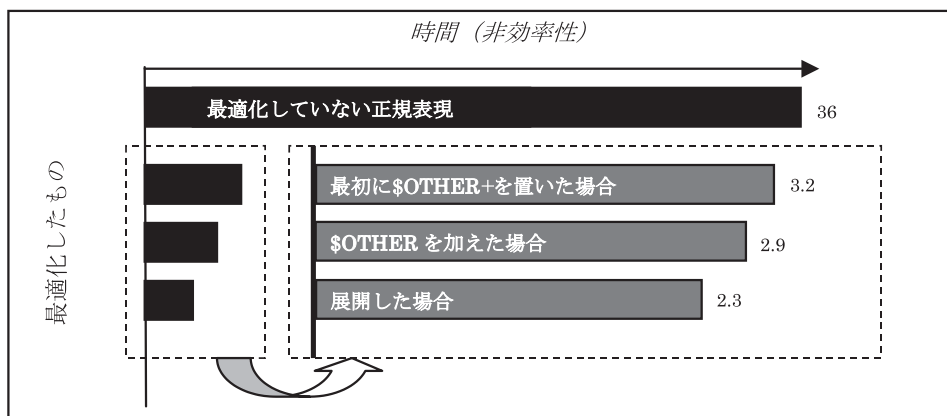


図1 十分に制御された最適化による短縮時間効果

Perl スクリプトを作った。その一つを次に示す。

```
$count = 0;
while ( < > )
{
    $count++ if m < \bchar\b
        | \bconst\b
        | \bdouble\b
        .
        .
        .
        | \bsigned\b
        | \bunsigned\b
        | \bwhile\b > x;
}
```

ここでは Perl の /x 機能が使われている。これは正規表現におけるクラスの部分で自由形式のスペースを許す。長く複数行にわたる「m <...> x」の内部では、空白文字を除くすべてが（この例では多くの選択肢を含む）一つの正規表現となる。次のような別のスクリプトも試してみた。

```
$count = 0;
while ( < > ) {
    $count++ next if m/\bchar\b/;
    $count++ next if m\bconst\b/;
    $count++ next if m\bdouble\b/;
    .
    .
    .
    $count++ next if m\bsigned\b/;
    $count++ next if m\bunsigned\b/;
    $count++ next if m\bwhile\b/;
}
```

どちらからも同じ結果が得られたが、2番目の方が6倍以上も速い。それはバックトラックにかかる処理負担の多くが取り除かれたからである。正規表現が「. *」など、無差別の文字で始まるマッチを許してしまう構造で始まる場合は、最適化が無効になるが、この例

では適用可能である。

他の最適化を考える

これをさらに推し進め、手作業で各選択肢に共通した先頭要素を選択の冒頭に持つてくることができる。この作業は、

「this|that|the_other」を

「th(is|at|e_other)」に変更するのと似ている。こうした変更を行うことで、比較的速度の遅かった選択を、頭文字識別を使わなくても「th」がマッチした後でのみ実行することができる。

総じて、プログラムの効率化は、アイデアとロジックを用いることで完全ではないにせよ大方実現することができる。筆者の責任はこうした高速化方式の中で最速の方法を教えることである。Perlをはじめ、筆者がこれらのテストに用いたツールにしても、次のバージョンでは内部最適化速度の順位が変わるかもしれない。

あとがき

今回も前回に引き続き、各エンジンによって、実際の検索効率が大きく異なることに着目して、文字列により、エンジンの種類でどのような検索が行われるかについて詳細な検討を加えた。とくに、「ループ展開」のパターンで効率化を図るにはどのような正規表現が有効であるかを検討した。さらにエンジンの種類によってCコメント表現のループ展開において、スムーズに流れる正規表現の記述の仕方、実際のエンジンでの表現の違いによる速度差を実験により確かめられている。次回はPerlについて、その機能の実態を詳細に検討し、実際の応用プログラムに適用した場合の効率的な正規表現について検証を加えていく予定である。

参考文献

- 1) Jeffrey E.F. Friedl 著、歌代和正監訳、1999、オ

ライリー・ジャパン発行

- 2) 江戸浩幸・坂本義行、「電子メール自動集計システム - I」、東京家政学院筑波女子大学紀要、第4集、2000.
- 3) 江戸浩幸・坂本義行、「電子メール自動集計システム II」、東京家政学院筑波女子大学紀要、第5集、2001.
- 4) 坂本義行・江戸浩幸、「正規表現について - そ

の1 どう読むか - 」、東京家政学院筑波女子大学紀要、第5集、2001.

- 5) 坂本義行・江戸浩幸、「正規表現について - その2 処理メカニズム - 」、東京家政学院筑波女子大学紀要、第6集、2002.
- 6) 坂本義行・江戸浩幸、「正規表現について - その3 正規表現を工夫する - 」、東京家政学院筑波女子大学紀要、第7集、2003.