

CS 61A/CS 98-52

Mehrdad Niknami

University of California, Berkeley

Motivation

How would you find a substring inside a string?

Something like this? (Is this good?)

```
def find(string, pattern):
    n = len(string)
    m = len(pattern)
    for i in range(n - m + 1):
        is_match = True
        for j in range(m):
            if pattern[j] != string[i + j]:
                is_match = False
                break
        if is_match:
            return i
```

What if you were looking for a *pattern*? Like an email address?

Background

Text processing has been at the heart of computer science since the **1950s**

- Regular languages: 1950s (Kleene)
- Context-free languages (CFLs): 1950s (Chomsky)
- Regular expressions (regexes) & automata: 1960s (Thompson)
- LR parsing (left-to-right, rightmost-derivation): 1960s (Knuth)
- Context-free parsers: 1960s (Earley)
- String searching (Knuth-Morris-Pratt, Boyer-Moore, etc.): 1970s
- Periods & critical factorizations: 1970s (Cesari-Vincent)
- [...] Critical factorizations in linear complexity: **2016** (Kosolobov)

Research is still ongoing ...apparently more in Europe?

Background

Most of you will probably graduate without learning string processing. Instead, you'll learn how to process images and Big Data™.

Which makes me sad. :(You should know how to solve solved problems!

Learn & use 100%-accurate algorithms before 85%-accurate ones!

$O(mn)$ -time `str.find(substring)` is **bad!** You can do *much* better:

- Good algorithms finish in $O(m + n)$ time & space (e.g. [Z algorithm](#))
- The best/coolest finish in $O(m + n)$ time but $O(1)$ **space!!!**

So, today, I'll teach a bit about string processing. :)

You can learn more in CS 164, CS 176, etc. (Have fun!)

In formal language theory:

- **Alphabet:** any set (usually a *character set*, like English or ASCII)
→ Often denoted by Σ
- **Letter:** an element in the given *alphabet*, e.g. “x”
- **String** (or **word**): finite sequence of *letters*, e.g. “hi”
- **Language:** a set of *strings*, e.g. {“a”, “aa”, “aaa”, ...}
→ Often denoted by L

We might omit the quotes/braces, so we'll use the following denotations:

- ε : empty string (i.e., “”)
- \emptyset : empty language (i.e., empty set $\{\}$)

Formal Grammars

Languages can be infinite, so we can't always list all the strings in them. We therefore use **grammars** to describe languages.

For example, this grammar describes $L = \{ "", "hi", "hihi", \dots \}$:

$$S \rightarrow T$$

$$T \rightarrow \varepsilon$$

$$T \rightarrow T \text{ "h" "i"}$$

We call S a **nonterminal** symbol and $"h"$ a **terminal** symbol (i.e., letter). Each line is a **production rule**, producing a **sentential form** on the right.

To make life easier, we'll denote these by uppercase and lowercase respectively, omitting quotes and spaces when convenient.

We then merge and simplify rules via the pipe (OR) symbol:

$$S \rightarrow S \text{ hi} \mid \varepsilon$$

Regular Languages

The following are **regular languages** over the alphabet Σ :

- \emptyset
- $\{\varepsilon\}$
- $\{\sigma\} \quad \forall \sigma \in \Sigma$
- The union $A \cup B$ of any regular languages A and B over Σ
- The concatenation AB of any regular languages A and B over Σ
- The repetition (Kleene star) A^* of any regular language A over Σ
 $A^* = \{\varepsilon\} \cup A \cup AA \cup AAA \cup \dots$

Notice that all finite languages are regular, but not all infinite languages.

Regular languages do not allow arbitrary “nesting” (e.g. parens).

Regular Grammars

A **regular grammar** is a grammar in which all productions have at most one nonterminal symbol, all of which appear on either the left or the right.

In other words, this is a regular grammar:

$$\begin{aligned} S &\rightarrow A b c \\ A &\rightarrow S a \mid \varepsilon \end{aligned}$$

This is not a regular grammar (but it is *linear* and *context-free*):

$$\begin{aligned} S &\rightarrow A b c \\ A &\rightarrow a S \mid \varepsilon \end{aligned}$$

and neither is this (it is *context-sensitive*):

$$\begin{aligned} S &\rightarrow S s \mid \varepsilon \\ S s &\rightarrow S t \end{aligned}$$

A language is regular iff it can be described by a regular grammar.

Regular Expressions

A **regular expression** is an easier way to describe a regular language. It's essentially a pattern for describing a regular language.

For example, in $[abcw-z]^*(1+2|3)?4\?$, we have:

- $[abcw-z]$ (a *character set*) means “either a, b, c, w, x, y, or z”.
- Asterisk (a.k.a. “Kleene star”, a *quantifier*) means “zero or more”
- Plus (another quantifier) means “one or more”
- Question mark (another quantifier) means “at most one”
- Backslash (“escape”) before a special character means *that character*
- Pipe (the OR symbol $|$) means “either”, and parentheses group

So this matches zero or more of a, b, c, w, x, y, z, followed by either nothing or by 3 or by 1's followed by 2, followed by 4 and a question mark.

Regular Expressions

Regular expressions (**regexes**) are equivalent to regular grammars¹, e.g.

$$\underbrace{\underbrace{[abcw-z]^*}_{X} (1+2|3)? 4 \setminus ?}_{Z}$$

is equivalent to

$$S \rightarrow Z 4 ?$$

$$Z \rightarrow Y 2 \mid X 3 \mid \epsilon$$

$$Y \rightarrow Y 1 \mid X 1$$

$$X \rightarrow X a \mid X b \mid X c \mid X w \mid X x \mid X y \mid X z \mid \epsilon$$

Here, the regex is more compact. Sometimes, the grammar is smaller.

¹If you've seen backreferences: those are **not** technically valid in regexes. 

Regular Expressions

Python has a regex engine to find text matching a regex:

```
>>> import re
>>> m = re.match('.* ([a-z0-9._-]+)@([a-z0-9._-]+)',
                 'hello cs61a@berkeley.edu cs98-52')
>>> m
<re.Match object; span=(0, 24),
                 match='hello cs61a@berkeley.edu'>
>>> m.groups()
('cs61a', 'berkeley.edu')
```

Notice that these could all be handled by `re.match`:

- Substring search (`str.find`)
- Subsequence search (`re.match(".*b.*b", "abbc")`)

The `grep` tool (from `ed`'s `g/re/p = global/regex/print`) does this for files.

Million-dollar question:

How do you find text matching a regex?

Two steps:

- 1 Parse the regex (pattern) to “understand” its structure
- 2 Use the regex to parse the actual text (corpus)

It turns out that:

- 1 Step 1 is theoretically harder, but practically easier.
(This can be done similarly to how you parsed Scheme.)
- 2 Step 2 is theoretically easier, but practically harder.

This is because we need parsing the *corpus* to be **fast**.

Regular Expressions

How do you solve each step?

Both steps are often done using “recursive-descent”—similarly to how your Scheme parser parsed its input.

Basically: try every possibility recursively.
“Backtrack” on failure to try something else.

Problem: Recursive-descent can take **exponential time!**

Example (where “a{3}” is shorthand for “aaa”):

```
>>> re.match("(a?){25}a{25}", "a" * 25)
```

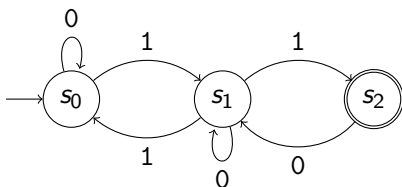
Can we hope to parse corpora in time **linear** to their lengths?

Yes, using finite automata.

Finite Automata

A finite automaton (FA) consists of the following (example below)²:

- An *input alphabet* Σ ($\{0, 1\}$ here)
- A finite set of *states* S ($\{s_0, s_1, s_2\}$ here)
- An *initial state* $s_0 \in S$ (s_0 here)
- A set of *accepting* (or *final*) *states* $F \subset S$ ($\{s_2\}$ here)
- A *transition function* $\delta : S \times \Sigma \rightarrow 2^S$ (the arrows here)



²Note that an FA is *not* quite the same thing as a finite-state machine (FSM).

Notice the transition function δ outputs a *subset* of states.

In a *deterministic* finite automaton (DFA), the transition function always outputs a set with **exactly** one state (a *singleton*).

i.e., in a DFA, the next state is **determined** by the input & current state. (i.e., every state has exactly 1 arrow leaving it for each possible input.)

In a *nondeterministic* finite automaton (NFA), the above is not true.

Finite Automata

Finite automata are **language recognizers**: you feed a string as an input, and if it accepts the input string, the string is in its language.³


In particular:

⇒ Finite automata **recognize *regular* languages**, and *nothing else*!

Therefore, we can:

- 1 Convert regex pattern to FA
- 2 Feed corpus to FA in **linear time**!
- 3 ...
- 4 Profit!

But how can we do this?

³*Pumping lemma*: A long-enough input must contain a repeatable substring. (Why?) 

Finite Automata from Regular Expressions

Consider: $(a|b)^*(1+2|3)$. Ask: **Where in the pattern can we be?**

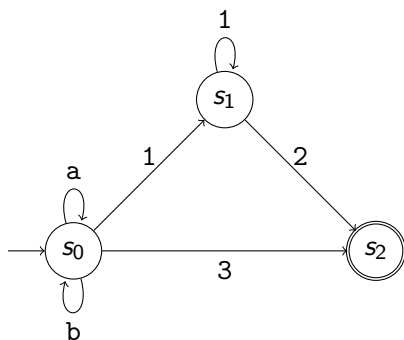
$$\begin{aligned} \textcircled{1} \quad s_0 &= \bullet(a|b)^*(1+2|3) \\ &= \bullet(\bullet a | \bullet b)^*(1+2|3) \\ &= \bullet(\bullet a | \bullet b)^* \bullet(1+2|3) \\ &= \bullet(\bullet a | \bullet b)^* \bullet(\bullet 1+2 | \bullet 3) \end{aligned}$$

$$\textcircled{2} \quad s_1 = (a|b)^*(\bullet 1+\bullet 2|3)$$

$$\begin{aligned} \textcircled{3} \quad s_2 &= (a|b)^*(1+2\bullet|3) \\ &= (a|b)^*(1+2\bullet|3)\bullet \end{aligned}$$

$$\begin{aligned} s_2 &= (a|b)^*(1+2|3\bullet) \\ &= (a|b)^*(1+2|3\bullet)\bullet \end{aligned}$$

$$s_2 = (a|b)^*(1+2\bullet|3\bullet)\bullet\bullet$$



(Expanding a state to its equivalents is a mathematical *closure* operation.)

Finite Automata from Regular Expressions

We created a *deterministic finite automaton* (DFA) from a regex!

It can find regular patterns (substrings, subsequences, etc.) in **linear** time.

However: *there is no such thing as a free lunch.*

What is the caveat?

Caveat:

The number of states $|S|$ can be **exponential** in the size of the pattern m . (This is sometimes referred to as the size of the *state space*.)

Why? Because we compute **subsets** of locations in the pattern, and we could encounter around 2^m subsets for a pattern of length m .

Solution?

DFA minimization: Always merge states that behave identically.

→ We already did this. It often works well.

Boring: Use an NFA instead of a DFA. (Track state items separately.)

→ Guarantees $O(m)$ memory usage, but running time is $O(mn)$.

Clever: Build the DFA lazily as needed by input.

→ Memory usage becomes $O(m + n)$, but running time approaches $O(n)$.

Conclusion

Automata are extremely powerful! They can do many other cool things:

- *Levenshtein automata* can recognize corpora that are k “edit distances” (insertions, deletions, or mutations) away from a pattern.
- When given a stack, LR automata can parse *context-free languages* (like many programming languages) in linear time (CS 164).
- Büchi automata, which allow *infinite-length* input strings, are used for formal verification of computer programs.

Finite-state machines (very similar) are widely used in digital design:

- Used in engineering to prove digital systems work as intended
- Used to optimize power consumption, logic circuitry, etc.

Conclusion

This is just the tip of the iceberg for string algorithms (and automata). Languages, grammars, and automata are also used in computational linguistics, computational biology/genomics (DNA alignment/matching)...

It is extremely easy to graduate and avoid languages & automata. But they provide the keys for solving many otherwise difficult problems. You can see more in EE/CS 144, 149, 151, 164, 172, 176, 219C, 291E...

~ Related Words of Wisdom ~

Minimizing the number of states in your design (e.g. factoring out duplicate data) helps keep designs clean & bug-free.

One reason: single source of truth. If it *can't* be wrong, it won't be.

Kleeneliness is next to Gödeliness.

Bonus: What language does S describe?

$$S \rightarrow S a \mid \varepsilon$$

Hmm, union and concatenation sure look like addition & multiplication...

$$S = Sa + \varepsilon$$

$$S\varepsilon = Sa + \varepsilon$$

$$S(\varepsilon - a) = \varepsilon$$

$$S = \frac{\varepsilon}{\varepsilon - a}$$

...wait, what?

Oh, right—Taylor series...

$$S = \varepsilon + a + a^2 + \dots$$

:-) Please don't ask...

Thank you!