Formal Languages

We'll use the *English* language as a running example.

Definitions.

Examples.

- A *string* is a *finite set* of *symbols*, where each *symbol* belongs to an *alphabet* denoted by Σ .
- The set of all strings that can be constructed from an alphabet Σ is Σ*.
- If *x*, *y* are two strings of *lengths* |*x*| and |*y*|, then:
 - xy or $x \circ y$ is the *concatenation* of x and y, so the length, |xy| = |x| + |y|
 - $(x)^R$ is the *reversal* of x
 - the k^{th} -power of x is

$$x^{k} = \begin{cases} \epsilon & \text{if } k = 0\\ x^{k-1} \circ x, & \text{if } k > 0 \end{cases}$$

- equal, substring, prefix, suffix are defined in the expected ways.
- Note that the language \emptyset is *not* the same language as ϵ .

Operations on Languages

Suppose that L_E is the *English* language and that L_F is the *French* language over an alphabet Σ .

- Complementation: $\overline{L} = \Sigma^* L$ \overline{L}_E is the set of all words
- Union: $L_1 \cup L_2 = \{x : x \in L_1 \text{ or } x \in L_2\}$

 $L_E \cup L_F$ is the set

• Intersection: $L_1 \cap L_2 = \{x : x \in L_1 \text{ and } x \in L_2 \}$

 $L_E \cap L_F$ is the set

Concatenation: L₁ ◦ L₂ is the set of all strings xy such that x ∈ L₁ and y ∈ L₂

Q: What is an example of a string in $L_E \circ L_F$?

Q: What if L_E or L_F is \emptyset ? What is $L_E \circ L_F$?

• Kleene star: L*. Also called the Kleene Closure of L and is the concatenation of zero or more strings in L.

Recursive Definition

- Base Case: $\epsilon \in L$
- Induction Step: If $x \in L^*$ and $y \in L$ then $xy \in L^*$
- Language Exponentiation Repeated *concatenation* of a language *L*.

$$L^{k} = \begin{cases} \{\epsilon\} & \text{if } k = 0\\ L^{k-1} \circ L, & \text{if } k > 0 \end{cases}$$

• **Reversal** The language *Rev(L)* is the language that results from *reversing all strings in L*.

Q: How do we *define* the strings that belong to a *language* such as *English, French, Java, arithmetic,* etc.

Example: For the *language of arithmetic*, \mathcal{LA} :

Define $\Sigma = \{\mathbb{N}\} \cup \{+, -, =, (,)\}$ then

")((2(+4(= "
$$\in \Sigma^*$$

but

")((2(+4(="
$$\notin \mathcal{LA}$$
.

1

Regular Expressions

A regular expression over an alphabet Σ consists of

- 1. Symbols in the alphabet
- 2. The symbols {+, (,),* } where + means OR and * means *zero* or *more times*.

Recursive Definition.

Let the set \mathcal{RE} of ALL regular expressions, be the smallest set such that:

- **Basis**: $\emptyset, \epsilon, a \in RE, \forall a \in \Sigma$
- Inductive Step: if R and S are regular expressions $\in \mathcal{RE}$, then so are: $(R + S), (RS), R^*$

Examples: Let $\Sigma = \{0, 1\}$:

Regular Expression	Corresponding Language
$(0+1)^*$	all binary strings
$((0+1)(0+1)^*)$	all non-empty binary str.
$((0+1)(0+1))^*$	all even lengthed binary
$\epsilon + 0 + 0(0 + 1)^*0$	au strings not starting or
$11(0+11)^*$	all strings with I's in
-	pairs and starting with

Relating Regular Expressions to Languages

Let $\mathcal{L}(\mathcal{R})$ represent the *language* constructed by the *regular expression* R.

We define $\mathcal{L}(\mathcal{R})$ *inductively* as follows:

Base Case:

- $\mathcal{L}(\emptyset) = \emptyset$
- $\mathcal{L}(\epsilon) = \{\epsilon\}$
- For any $a \in \Sigma$, $\mathcal{L}(a) = \{a\}$

Induction Step: If R is a *regular expression*, then by definition of R,

- R = ST, or
- R = S + T, or
- $R = S^*$

where *S* and *T* are *regular expressions* and by *induction*, $\mathcal{L}(S)$ and $\mathcal{L}(T)$ have been defined.

We can define the language denoted by R, ie., $\mathcal{L}(\mathcal{R})$ as follows:

- $\mathcal{L}((S+T)) = \mathcal{L}(S) Union \mathcal{L}(T)$
- $\mathcal{L}((ST)) = \angle (S) \circ \angle (T)$
- $\mathcal{L}(\mathcal{S}^*) = \left(\mathcal{L}(\mathcal{S}) \right)^*$

Q: Why is this definition important? Use these definitions when proving languages are equivalent to Example each other.

Q: What is a *regular expression* R_A to denote the language of strings consisting of only an *even number* of *a*'s?

e.g., aa, aaaa, aaaaaaaa etc.

(ha) to

Q: What is a regular expression, \mathcal{R}_{AB} , for the language of strings consisting of an *even number* of *a*'s *sandwiched* between 1 or more *triples* of *b*?

eg., bbbaabbb, or bbbaaaaaabbb

RR RARD

Equivalence. We say that two regular expressions R and S are *equivalent* if they *describe* the *same language*.

In other words, if $\mathcal{L}(\mathcal{R}) = \mathcal{L}(\mathcal{S})$ for two regular expressions R and S then R = S.

Examples.

• Are *R* and *S* equivalent?

Regular Expression Equivalences

There exist *equivalence axioms* for *regular expressions* that are very similar to those for *predicate/propositional logic*.

Equivalences for Regular Expressions $R_1S_7 T a re R_2E$

- Commutativity of union: $\mathcal{R} + S \equiv S + \mathcal{R}$
- Associativity of union: (R+5)+T = R+(S+T)
- Associativity of concatenation: R(ST) = (RS) T
- Left distributivity: R(S+T) = RS + RT
- Right distributivity: $(S+T)R \equiv SR + TR$
- Identity of Union: $R + \phi \equiv R$
- Identity of Concatenation: $R \mathcal{E} \subseteq R$
- Annihilator for concatenation: $\mathcal{R} \not \!\!\! D \equiv \not \!\! Q \equiv \not \!\!\! Q \mathcal{R}$
- Idempotence of Kleene star: $\mathcal{R}^{*} = \mathcal{R}^{*}$

Theorem (Substitution) If two substrings R and R' are equivalent then if R is a substring of S then replacing R by R' constructs a new regular expression equivalent to S.

Equivalent Regular Expressions

Q: How can we determine whether two *regular expressions* denote the *same language*?

Examples.

Prove that

1

$$(0110 + 01)(10)^* \equiv 01(10)^*$$

Proof.

$$(0110+01)(10)^* \equiv (0110+01\varepsilon)(10)^* \text{ by id. of contrat.} \\ \equiv (01(10+\varepsilon))(10)^* \text{ left. dist.} \\ \equiv (01)((10+\varepsilon)(10)^*) \text{ assoc. of contrat.} \\ \equiv 01(10(10)^*+\varepsilon(10)^*) \text{ right dist.} \\ \equiv 01(10(10)^*+(10)^*) \text{ id of contrat.} \\ \equiv 01(10)^* \text{ by substitution of equiv.} \\ \equiv R.E. \\ Notice \text{ that } L(10(10)^*) \leq L(10)^*) \end{aligned}$$

Another Example.

Prove that R denotes the *language* L of all strings that contain an even number of 0s.

 $R = 1^* (01^* 01^*)^*$

Equivalently,

$$x \in L \Leftrightarrow x \in \mathcal{L}(R)$$

Proof.



- Let $x \in \mathcal{L}(R)$.
- Then $x \in \mathcal{L}\left(\int^{\mathfrak{p}} \left(0 \int^{\mathfrak{q}} \left(0 \int^{\mathfrak{q}} \right)^{\mathfrak{q}} \right) \stackrel{\mathfrak{q}}{=} \mathcal{L}\left(\int^{\mathfrak{q}} \left(0 \int^{\mathfrak{q}} \right) \mathcal{L}\left(0 \int^{\mathfrak{q}} \right) \right)$
- Let $x = y(zw)^*$ then $y \in \mathcal{L}(\mathbb{I}^*)$, $\mathcal{E} \in \mathcal{L}(\mathbb{O}\mathbb{I}^*)$, $w \in \mathcal{L}(\mathbb{O}\mathbb{I}^*)$
- Therefore, y has $On h \int f$.
- Therefore, w has exactly one 0.
- Therefore, z has -exactly one 0.
- So, $x = y(zw)^*$ has a multiple of two 0s'. X E] 81



- Suppose that *x* is an *arbitrary* string in *L*.
- $\Rightarrow x$ has an *even* number of 0s. Denote by 2k for some $k \in \mathbb{N}$.
- How can we rewrite x consisting of θ s and 1s?

• Let
$$x = y_0, y_1, y_2, \dots, y_k$$
, so $y_0 = | \dots | \mathcal{E}(\mathcal{A})$
each $y_i := (0 | \dots |)(0 | \dots |) = 0 | 0 | \mathcal{E}(0 | \mathcal{A}) / (0 | \mathcal{A})$
 $| \leq i \leq k$
• So $x = y_0 y_1 \dots y_k \in \mathcal{L}(1^*)(\mathcal{L}(01^*01^*))^* = \mathcal{L}(1(01^*01^*)^*).$
Q: Can every possible type of string be represented by a regular
expression?

To answer this, we turn to *Finite State Machines*.

String Matching and Finite State Machines

- Given *source code* (say in Java)
- Find the comments may need to remove comments for software transformations

```
public class QuickSort {
   private static long comparisons = 0;
   private static long exchanges
                             = 0;
  * Quicksort code from Sedgewick 7.1, 7.2.
   public static void quicksort(double[] a) {
       shuffle(a);
                                     // to guard against worst-case
       quicksort(a, 0, a.length - 1);
   }
   public static void quicksort(double[] a, int left, int right) {
      if (right <= left) return;</pre>
       int i = partition(a, left, right);
       quicksort(a, left, i-1);
       quicksort(a, i+1, right);
   }
   private static int partition(double[] a, int left, int right) {
      int i = left - 1;
      int j = right;
      while (true) {
          while (less(a[++i], a[right])) // find item on left to swap
                                         // a[right] acts as sentinel
          while (less(a[right], a[--j]))
                                        // find item on right to swap
             if (j == left) break;
                                         // don't go out-of-bounds
          if (i >= j) break;
                                         // check if pointers cross
          exch(a, i, j);
                                         // swap two elements into place
      }
                                         // swap with partition element
      exch(a, i, right);
      return i;
   }
```

Q. What patterns are we looking for? //text In / text Q. What do we know if we see a / followed by a * In a comment / in a comment text not in a comment, **Q.** What do we know if we see /* followed by a * might be at the end of the comment if next char is a I not at end text not at end of homment. fext, /4 Let's represent these ideas with a *diagram*. /* start) test in Ł tep, ¥ 12 accept 84 ¥

Deterministic Finite State Automata (DFSA or DFA)

A DFA consists of:

- · Q. the set of statis
- ∑. analphabet from with strings com from.
 s∈Q. start state
- F⊆Q final or accepting states.
- S. QXZ -> Q transition function which takes a state and impost and

Comment Example. outputs a state.

- Q={start, /, 11, /*, *, accept}.
- E = {text, \n, ', *]
- s = 5 + a(+
- $F = \{ G \in \mathcal{C} \in \mathbb{P}^+ \}$
- δ

Example cont...



Q: What if we want to know which state the input "**//" ends at if we *begin* at *start*?

Two Options. $A \neq //$ 1. Compute: $S\left(\delta\left(\delta\left(\delta(s + x, t, *), \#), /\right), /\right) = 1/$ 2. Define δ^* . = $\delta^*\left(s + x, t, *, \# / 1, *\right) = 1/$

Formal definition of $\delta^*(q, x)$ (reading left to right):

$$\delta^*(q,x) = \begin{cases} q & \text{if } x = \epsilon \\ \delta(\delta^*(q,z),a) & \text{if } x = za, a \in \Sigma, z \in \Sigma^* \end{cases}$$

Exercise: Rewrite δ^* as for reading a string from right to left.
86

Regular Expressions and DFA

- The set of strings *accepted* by an *automaton* defines a *lan-gauge*.
- For automaton M the language M accepts is $\mathcal{L}(M)$.
- Given regular expression R, find M such that

$$\mathcal{L}(R) = \mathcal{L}(M).$$

Examples.

Let regular expression $R_1 = (1 + 00)^*$.

Q. Which strings belong to $\mathcal{L}(R_1)$? $\mathcal{L}(R_1) = \{ s \in \{0, j\}^{\prime} | a | 0's are in pairs \}$

Q: What is a *DFA* M_1 such that $\mathcal{L}(M_1) = \mathcal{L}(R_1)$?



DFSA Conventions

- Strings ending at a *final state* are *accepted* (if we want to accept/reject).
- Drop *dead* states.
- Group *elements* that go from and to the *same states*.

Examples cont.

Let regular expression $R_2 = 1(1 + (01))^*$.

Q. Which strings belong to $\mathcal{L}(R_2)$? over y d is sand whet is.

Q: What is a *DFA* M_2 such that $\mathcal{L}(M_2) = \mathcal{L}(R_2)$?



$$\delta: \quad \delta(q_0, 0) = \partial \quad \delta(q_0, 1) = \frac{2}{7} \delta(q_1, 0) = \frac{2}{7} \circ$$
$$\delta(q_1, 1) = \frac{2}{7} \frac{\delta(q_2, 0)}{\delta(q_2, 0)} = \frac{\delta(q_2, 1)}{\delta(q_2, 1)} = \frac{2}{7} \delta(q_1, 0 \text{ or } 1) = \frac{2}{7}$$

Q: How do we know that our *machine M* is *correct*?

We can show this by proving that $\delta^*(q_0, x)$ only accepts those strings in $\mathcal{L}(R_2)$.

Q: What might be a *good* way to do this?

Proving a DFA is Correct

