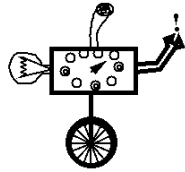


C S 6 0



Inductive Definitions and Grammars

October 31, 2001

Inductive Definitions

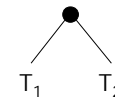
- Inductive definitions are the main “constructive” way to define infinite sets.
- We will need infinite sets in much of what follows.

Inductive Definitions

- Elements of an inductive definition of a set S :
 - **Basis**: Specifies that certain items are in S .
 - **Induction rule(s)**: Introduce new items in S based on existence of other, usually simpler, items.
 - **Extremal clause**: Says that the only items in S are those derivable by the previous two elements, applied any finite number of times.
 - Often implicit

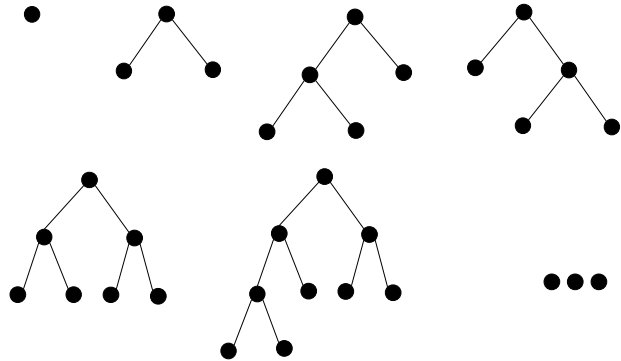
Example: Binary Trees

- \bullet is a binary tree.
- If T_1 and T_2 are binary trees, then so is:



- Extremal clause: The only binary trees are those constructible by a finite number of applications of the above rules.

Examples of Binary Trees



Example: Natural Numbers ω

- **Basis:** 0 is in ω .
- **Induction:** If n is in ω , so is the successor of n (variously denoted n' , $S(n)$, or $n+1$).
- **Extremal:** The only elements in ω are those derivable by applications of the above rules.
- Examples: 0, $0'$, $0''$, $0'''$, $0''''$, ... are all elements of ω .

Notes

- ω is an infinite set.
- ω does *not* contain infinity (∞) as an element.
 - Why?
- ω 's members are all finite.

Decimal Numerals

- We can agree by convention that
 - 1 stands for $0'$,
 - 2 stands for $0''$,
 - ...
 - 9 stands for $0''''''''$.
 - Beyond that, give an algorithm for generating additional numerals:
10, 11, 12, 13, ...

1-adic Numerals

- The only digit is 1.
- The empty string (denoted λ so it is readable) stands for 0.
- 1X (1 followed by X) stands for X' .
- The numerals are:
 $\lambda, 1, 11, 111, 1111, 11111, \dots$
- Could also use lists: $[], [1], [1, 1], [1, 1, 1], \dots$

2-adic Numerals

- The digits are 1 and 2.
- The empty string (denoted λ so it is visible) stands for zero.
- The numerals are:
 $\lambda, 1, 2, 11, 12, 21, 22, 111, 112, \dots$
- Unlike binary numerals, there is no redundancy (1, 01, 001, 0001, ... all mean the same thing in binary).

Roman Numerals

- The digits are I, V, X, L, C, D, M.
- There is no string for 0.
- You know the rest.

Numerals vs. Numbers

- Numbers are abstract.
- Numerals are a concrete representation.
 - There's a reason why we don't call them "Roman numbers"

Strings over an alphabet

- An *alphabet* is a set of arbitrary symbols.
- The set of *all* finite strings over an alphabet Σ is denoted Σ^* .
- Example:
 $\{a, b\}^* =$
 $\{\lambda, a, b, aa, ab, ba, bb, aaa,$
 $aab, aba, \dots\}$

Formal definition:

- Basis:
 λ is in Σ^* .
- Inductive rule:
If $x \in \Sigma^*$ and $\sigma \in \Sigma$, then
 $x\sigma$ (x followed by σ) is in Σ^* .
- Extremal clause.

Languages

- A *language* over Σ^* is any subset of Σ^* .
– A very general definition!
- Examples, where $\Sigma = \{a, b\}$
 - $\{a, b\}^*$ itself
 - $\{ \}$ the empty language
 - $\{ba, baba\}$ maybe your first language
 - $\{\lambda, aa, aaaa, aaaaa, aaaaaa, \dots\}$ the language of an even number of a's.

More Languages

- Still with $\Sigma = \{a, b\}$:
 - $\{\lambda, ab, ba, aabb, abab, baab, bbaa, aaabbb, aababb, \dots\}$ the language in which the number of a's equals the number of b's.
 - $\{a, b, aa, bb, aab, aba, baa, abb, bab, bba, \dots\}$ the language in which the number of a's is *not* equal to the number of b's.
 - $\{\lambda, ab, abab, aabb, aababb, abaabb, ababab, \dots\}$ a slightly less obvious language.

Languages

- There are lots of languages, some very weird.
- To be of computational interest, a language needs to be defined **inductively**.
- Given a language, need a way of telling whether a given string is in the language or not (called *parsing* the string).

Non-Trivial Language Defined Inductively

- $L = \{\lambda, ab, abab, aabb, aababb, \dots\}$
- Basis: λ is in L .
- Inductive rules:
 - If X is in L , so is aXb .
 - If X_1 and X_2 are in L , so is X_1X_2 .

Grammars: A Shorthand

- Spelling everything out with these inductive definitions is laborious.
- We need a shorthand, especially for more complex languages.

Grammatical Definition

- Pick a symbol s *not* in the alphabet of interest.
- \rightarrow is a symbol meaning "can be rewritten as".
- Specify a collection of *grammar rules*:
 - $s \rightarrow ab$
 - $s \rightarrow aSb$
 - $s \rightarrow SS$
- Starting with s , what strings can we generate by replacing left-sides with right-sides of rules?
 - A sequence of such replacements is called a *derivation*.
 - Defines a language: the strings we can generate which *don't* contain s .

Using the Grammar Rules

- Given the grammar:

$S \rightarrow ab$
 $S \rightarrow aSb$
 $S \rightarrow SS$

- Example derivations of strings in the language:

$S \Rightarrow ab$
 $S \Rightarrow aSb \Rightarrow aabb$
 $S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aaabbb$
 $S \Rightarrow SS \Rightarrow abS \Rightarrow abab$
 $S \Rightarrow SS \Rightarrow SSS \Rightarrow ababab$
 $S \Rightarrow SS \Rightarrow aSbS \Rightarrow aabbS \Rightarrow aabbaSb \Rightarrow aabbaabb$

More General Grammars

- Instead of just S , allow multiple symbols, (called *nonterminals*) none of which are in the alphabet of the language.
 - The symbols in the alphabet of the language are called **terminals**.
 - Pick one to be the “**start symbol**”.

Additive Arithmetic Expressions

- The start symbol is A .
- The terminals are $\{a, b, c, +\}$.
- The productions are:

$A \rightarrow V$
 $A \rightarrow V + A$
 $V \rightarrow a$
 $V \rightarrow b$
 $V \rightarrow c$

Phase-Structure Grammars

- In a *phase-structure grammar*, rules can have arbitrary left-hand and right-hand sides involving terminals and nonterminals:

$S \rightarrow abc$
 $S \rightarrow aSQ$
 $bQc \rightarrow bbcc$
 $cQ \rightarrow Qc$

- Arbitrary grammars are too hard to work with
- In CS you're likely to see only *context-free grammars*
 - Left-hand side is always a nonterminal

