

Stochastic grammars and RNA

In a nutshell

- Primary sequence structure is described by hidden Markov models
- Hidden Markov models are stochastic regular grammars
- Context free grammars generalize regular grammars
- Secondary sequence structure (of RNA) is described by stochastic context free grammars

Plan

- Grammars, derivation, parsing, and the Chomsky hierarchy
- Regular grammars and finite state automata
- Context free grammars and reconstruction of derivation trees
- Stochastic regular grammars
- Stochastic context free grammars: the CYK algorithm

KVL

Seminar on computational biology 1999-12-20

Page 1

A grammar $G = (N, T, R, S)$

$$\begin{aligned} S &\rightarrow aS \\ S &\rightarrow bS \\ S &\rightarrow \epsilon \end{aligned}$$

Nonterminal symbols $N = \{S\}$, terminal symbols $T = \{a, b\}$, a start symbol S , and rules R (above).

In biosequence analysis, terminal symbols will be nucleotides or proteins, e.g. $T = \{A, C, G, T\}$.

Rules are sometimes called *productions*.

Rules may be written compactly: $S \rightarrow aS \mid bS \mid \epsilon$ or in Backus-Naur form: $\langle S \rangle ::= a \langle S \rangle \mid b \langle S \rangle \mid \langle empty \rangle$

Derivation, and the language generated by a grammar

Derivation generates a sequence of terminal symbols. Example: $S \Rightarrow aS \Rightarrow abS \Rightarrow abbaS \Rightarrow abba$

We may write $S \Rightarrow^* abba$ for derivation in zero or more steps.

The language generated by a grammar G is $L(G) = \{\alpha \in T^* \mid S \Rightarrow^* \alpha\}$ where S is the start symbol of G .

Example: $L(G) = \{\epsilon, a, b, aa, ab, ba, bb, aaaa, aabb, \dots\}$

In biosequence analysis, a 'language' will be a set of DNA or protein sequences.

KVL

Seminar on computational biology 1999-12-20

Page 2

Grammar classes: the Chomsky hierarchy (ca. 1956)

By restricting the possible form of grammar rules, we get a hierarchy of increasingly powerful grammar classes:

Grammar class	Admissible rules
Regular grammars	$W \rightarrow aW$ and $W \rightarrow a$
Context free grammars	$W \rightarrow \beta$
Context sensitive grammars	$\alpha_1 W \alpha_2 \rightarrow \alpha_1 \beta \alpha_2$, non-contracting
Unrestricted grammars	$\alpha_1 W \alpha_2 \rightarrow \gamma$

Above, α, β etc. stand for arbitrary sequences of terminals and nonterminals.

Notes

- Regular grammars and context free grammars are widely used in computer science
- Programming languages are described by the so-called LR or LL subclasses of context free grammars
- Only part of programming language syntax is described by the LR -grammar (e.g. type rules are not).
- Unrestricted grammars are also called term rewrite systems
- Chomsky's original goal was to describe natural languages. He (and everybody else) failed.

KVL

Seminar on computational biology 1999-12-20

Page 3

Regular grammars and finite state automata

The FMR-1 triplet region contains repeats of the triplets CCG and occasionally AGG, e.g.

GGC CTG, GCG CCG CTG, GCG CCG CCG CTG, GCG CCG CCG CTG, ...

Such regions can be described by the grammar

$$\begin{aligned} S &\rightarrow gW_1 \\ W_1 &\rightarrow cW_2 \\ W_2 &\rightarrow gW_3 \\ W_3 &\rightarrow cW_4 \\ W_4 &\rightarrow gW_5 \\ W_5 &\rightarrow gW_6 \\ W_6 &\rightarrow aW_4 \mid cW_7 \\ W_7 &\rightarrow tW_8 \\ W_8 &\rightarrow g \end{aligned}$$

The language generated by a regular grammar may be recognized by a finite state automaton (and vice versa).

KVL

Seminar on computational biology 1999-12-20

Page 4

Regular grammars and regular expressions

It's customary to use *regular expressions* as a shorthand for regular grammars.

Every regular expression corresponds to a regular grammar and vice versa.

The FMR-1 grammar can be written as the regular expression:

$$gcgcg((a+c)gg)^*ctg$$

In UNIX grep or emacs notation, that is:

$$gcgcg([ac]gg)^*ctg$$

More regular expressions: PROSITE patterns

Regular expressions are used to describe 'signature' conserved protein sequences and their variants.

Brackets [RK] indicate choice, braces {EDRKHPCG} choice from the complement, x matches anything.

(Figure 9.3)

A context free grammar for palindromes over $\{a, b\}$

$$S \rightarrow aSa \mid bSb \mid a \mid b$$

Derivation of the palindrome *abaababaa*:

$$S \Rightarrow aSa \Rightarrow aaSaa \Rightarrow aabSbaa \Rightarrow abaababaa$$

Context free grammars and RNA secondary structure

RNA sequences can form 'stem loops' when one part of the sequence matches another part ($a \leftrightarrow u, c \leftrightarrow g$).

Three-base RNA stem loops with 'head' *gaaa* or *gcaa* can be described by:

$$\begin{array}{l} S \quad \rightarrow \quad aW_1a \mid cW_1g \mid gW_1c \mid uW_1a \\ W_1 \quad \rightarrow \quad aW_2u \mid cW_2g \mid gW_2c \mid uW_2a \\ W_2 \quad \rightarrow \quad aW_3u \mid cW_3g \mid gW_3c \mid uW_3a \\ W_3 \quad \rightarrow \quad gaaa \mid gcaa \end{array}$$

Expressive power in the Chomsky hierarchy

Let $\alpha, \beta \in (N \cup T)^*$ denote arbitrary sequences of nonterminals and terminals.

- Context sensitive grammars can describe copy languages, such as:

$$\{\alpha\alpha \mid \alpha \in (N \cup T)^*\}$$

Context free grammars cannot.

- Context free grammars can describe palindrome languages (where α^{-1} is α reversed), such as:

$$\{\alpha\alpha^{-1} \mid \alpha \in (N \cup T)^*\}$$

(Context free grammars can describe proper parenthetical nesting, as used in programming languages).

Regular grammars cannot.

- Regular grammars can describe languages with uncorrelated repeats, such as:

$$\{a^n \alpha a^m \mid \alpha \in (N \cup T)^*\}$$

But they cannot encode the additional requirement $n = m$.

Derivation trees and parsing

The derivation from a context free grammar may be shown as a tree.

(The derivation from a regular grammar is a degenerate – linear – tree, not very interesting).

Parsing: find a derivation tree for a given sequence, if any

Given a grammar G and a sequence α ,

- is sequence α derivable from G ?
- if so, what derivation trees would produce sequence α ?

Grammar classes and 'parsers' (recognizers)

For each grammar class there is a characteristic 'machine type' that solves the parsing problem for that class. More expressive grammar classes require more powerful recognizers (and more time and space).

Let n be the length of the sequence being parsed.

Grammar class	Recognizer	Time complexity	Space complexity
Regular grammars	Finite state automata	Linear	Constant
Context free grammars	Pushdown automata	$O(n^3)$	$O(n^2)$
Context sensitive grammars	Linear bounded automata	NP-complete	PSPACE-complete
Unrestricted grammars	Turing machines	Undecidable	Unbounded

KVL

Stochastic context free grammars

For each non-terminal symbol W , assign a distribution to the set of rules $W \rightarrow \alpha$.

	Simple (deterministic)	Stochastic
Regular	PROSITE patterns	Primary structure, probabilistically (HMM)
Context free	RNA secondary structure	RNA secondary structure, probabilistically

Chomsky normal form

A grammar is on Chomsky normal form if all rules have form $W \rightarrow W_1 W_2$ or $W \rightarrow a$.

Every context free grammar can be transformed to Chomsky normal form.

KVL

Stochastic regular grammars

For each non-terminal symbol W , assign probabilities to rules of form $W \rightarrow a W_1$ or $W \rightarrow a$.

Think of W as a state and a as an emitted symbol.

The stochastic grammar emits symbols on state transitions.

A hidden Markov model (HMM) emits symbols without state transitions, and emits nothing on state transitions.

The two kinds of machines are interconvertible (by introducing extra states and transitions).

Thus alignment of a biosequence α using an HMM is parsing of sequence α using a stochastic regular grammar.

KVL

Finding the most probable parse tree: the CYK algorithm

Dynamic programming (tabulation). Analogous to the Viterbi algorithm, which finds the most probable alignment.

A probabilistic version of the Cocke-Younger-Kasarni (CYK) algorithm for context free grammar parsing (ca. 1968).

Input: A stochastic context free grammar G on Chomsky normal form, and a sequence $x = x_1..L$.

Let the grammar G have terminal symbols $T = \{V_1, \dots, V_M\}$ and start symbol $S = V_1$.

Let $e_v(a)$ be p if $V_a \rightarrow a$ with probability p , and 0 otherwise.

Let $t_v(y, z)$ be p if $V_v \rightarrow V_y V_z$ with probability p , and 0 otherwise.

Output: A table γ such that $\gamma(i, j, v)$ is the probability of the most probable parse tree that derives $x_i..j$ from V_v .

Algorithm: Put $\gamma(i, i, v) = e_v(x_i)$ for $i = 1..L$ and $v = 1..M$.

For $i = 1..(L - 1)$ and $j = (i + 1)..L$ and $v = 1..M$, put

$$\gamma(i, j, v) = \max_{y+z=j} \max_{k=i+1}^{j-1} (\gamma(i, k, y) \cdot \gamma(k+1, j, z) \cdot t_v(y, z))$$

The probability of the most probable derivation of x from G is $\gamma(1, L, 1)$.

Also, a traceback $\tau(i, j, v)$ is built as in the Viterbi algorithm.

KVL