

5 Formal Models of Language (pjb48)

A linguist produces the grammar $G = (\mathcal{N}, \Sigma, S, \mathcal{P})$ where:

$$\begin{aligned} \mathcal{N} &= \{S, X, Y, V, C\} \\ \Sigma &= \{a, \textit{contagious}, \textit{highly}, \textit{virus}\} \\ S &= S \\ \mathcal{P} &= \{S \rightarrow a X, X \rightarrow Y \textit{virus} \mid \textit{virus}, Y \rightarrow V C \mid C, \\ &\quad V \rightarrow \textit{highly} V \mid \textit{highly}, C \rightarrow \textit{contagious} C \mid \textit{contagious}\} \end{aligned}$$

- (a) Draw all the trees with 4 leaves that can be derived from this grammar. [2 marks]
- (b) Based on corpus data the linguist assigns probabilities to each rule in his grammar. Describe how the probability of a string is calculated from the rule probabilities. [2 marks]

A mathematician prefers to generate the strings of a language inductively. She defines a homomorphism: $\{(a, a), (c, \textit{contagious}), (h, \textit{highly}), (v, \textit{virus})\}$. She defines $L \subset \Sigma^*$ where $\Sigma = \{a, c, h, v\}$ using the following axioms and rules:

$$\begin{aligned} &\overline{av} \text{ (a1)} \\ &\frac{u_1 v}{u_1 c v} \text{ (r1) where } u_1 \in \Sigma^* \\ &\frac{a c u_1}{a h c u_1} \text{ (r2) where } u_1 \in \Sigma^* \\ &\frac{u_1 h u_2}{u_1 h h u_2} \text{ (r3) where } u_1, u_2 \in \Sigma^* \end{aligned}$$

- (c) Let $L_i = \{u \in L \mid \textit{length}(u) \leq i\}$. Find all members of L_4 [3 marks]
- (d) Describe L as a regular expression and specify a Deterministic Finite Automaton, M , such that $L(M) = L$. [4 marks]
- (e) Provide an expression for the conditional entropy of X for L_5 , where X is a random variable over Σ . A numerical value is not required. [5 marks]
- (f) Suggest some hypotheses about human language processing that we could test based on the models mentioned in this question. Provide reasons for your hypotheses. [4 marks]