

速習 ⚡ オートマトンと形式言語①

有限オートマトンと正規言語



Mitsuha Asada

Engineering toC Product Unit2

概要

- オートマトンとは何か？ / 性質は？ / 何に役立つの？
- 情報科学系の学部で基礎的な範囲として学ぶ内容を中心に触れていく
 - ▶ 扱わないこと
 - チューリング機械
 - 判定可能性
 - 帰着可能性
 - 計算可能性
 - …
- 参考文献
 - ▶ 計算理論の基礎 原著第 3 版 1. オートマトンと言語
 - Sipser による名著。情報系の学生は全員読んでいる [要出典]

目次

1. 前提知識	5
1.1. 前提知識について	6
1.2. 集合	7
1.3. 列・組	8
1.4. 無向グラフ	9
1.5. 有向グラフ	10
1.6. 文字列	11
2. 正規言語	12
2.1. 有限オートマトン	13
2.2. 言語	15
2.3. 計算	16
2.4. 非決定性	17

2.5. 非決定性有限オートマトン	18
2.6. DFA と NFA の等価性	19
2.7. 正規演算	20
2.8. 非正規言語	21
2.9. ポンピング補題	22

1. 前提知識

前提知識について

- ・ いくつかの前提知識がありますが、基本的な理解があれば十分です
- ・ 本章ではそれぞれの前提知識について簡単に説明します
 - ▶ 時間の都合上、それぞれの厳密な定義や証明は省略します
 - ▶ 勉強会終了後に質問していただくことは大歓迎です

集合

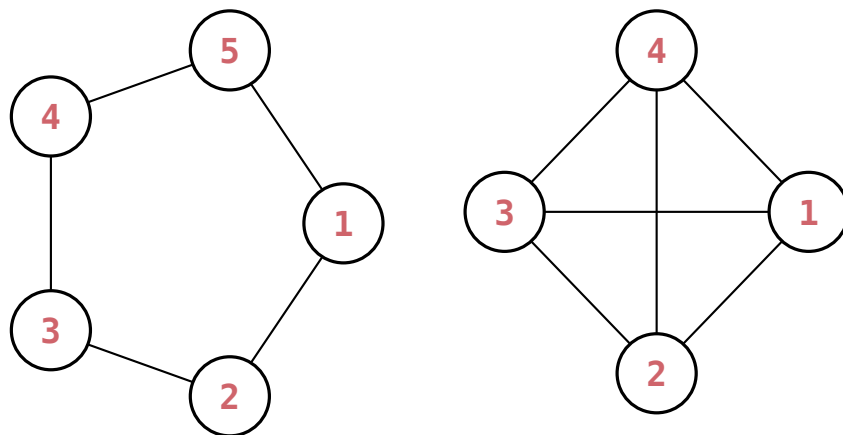
- 集合とは、元/要素の集まりのこと
 - ▶ $S_1 = \{1, 2, 3\}$
 - ▶ $S_2 = \{a, b, c\}$
 - ▶ $S_3 = \{1, a, \{1, 2\}\}$
- 部分集合: A のすべての要素が S に属する $A \subseteq S$
- 真部分集合: A が S の部分集合かつ A と S が異なる $A \subsetneq S$
- 空集合: 要素を 1 つも含まない集合 \emptyset
- 演算
 - ▶ 和集合: A または B の要素をすべて含む集合 $A \cup B$
 - ▶ 積集合: A かつ B の要素をすべて含む集合 $A \cap B$
 - ▶ 補集合: A に含まれていない要素の集合 \overline{A}

列・組

- 列: 順序付けられた要素の集まり $s = (s_1, s_2, \dots, s_n)$
- 組: 有限な列 $t = (t_1, t_2, \dots, t_m)$
 - ▶ k 個の要素からなる列を k 個組と呼ぶ
 - ▶ 2個組を特に順序対と呼ぶ
- 集合 A, B について、 A の要素を1番目の要素、 B の要素を2番目の要素とした順序対の全体集合を 直積 (または Cartesian 積) と呼ぶ $A \times B$
 - ▶ $A = \{1, 2\}, B = \{a, b\}$ のとき、 $A \times B = \{(1, a), (1, b), (2, a), (2, b)\}$

無向グラフ

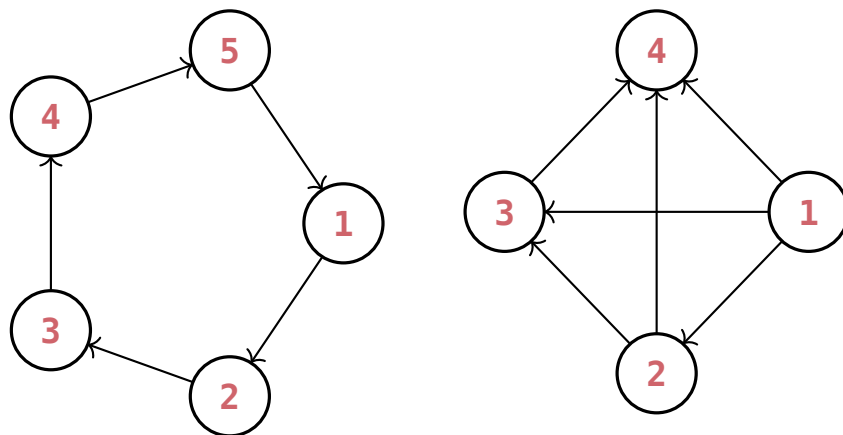
- 無向グラフ (グラフ) : 頂点と辺からなる構造



- 次数: ある頂点に繋がっている辺の個数
- グラフ G の頂点がグラフ H の頂点の部分集合で、かつ G の辺が H の辺の部分集合であるとき、 G は H の部分グラフであるという

有向グラフ

- 有向グラフ: 頂点と有向辺からなる構造



- 次数: ある頂点に繋がっている辺の個数
 - 入次数: 頂点に入ってくる辺の個数
 - 出次数: 頂点から出ていく辺の個数

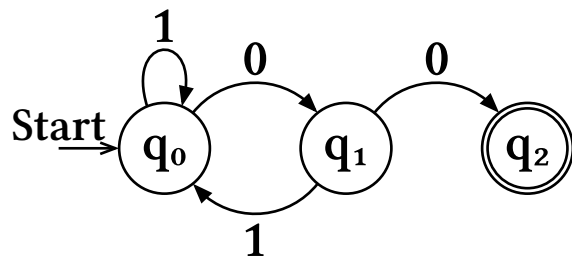
文字列

- 任意の空でない有限集合 Σ をアルファベットと呼ぶ
- アルファベット Σ の元を文字と呼ぶ
- そのアルファベットの文字からなる有限列をアルファベット上の文字列と呼ぶ
 - ▶ $\Sigma_1 = \{a, b\}$ ならば、 $a, b, ab, ba, aa, bb, aba, \dots$ はすべて Σ_1 上の文字列
- w を Σ 上の文字列とするとき、 w に含まれる文字の数を w の長さと呼び、 $|w|$ で表す
 - ▶ 例えば、 $\Sigma_1 = \{a, b\}$ ならば、 $|a| = 1, |b| = 1, |ab| = 2, |aba| = 3$

2. 正規言語

有限オートマトン

- 有限の状態を持ち、入力された文字列を読み取って状態を遷移させるモデルを有限オートマトンと呼ぶ



- アルファベット $\Sigma = \{0, 1\}$ について、00で終わる文字列を受理する
 - 開始状態 (q_0)、受理状態 (q_2) が定められている
 - 状態から状態への移動を遷移と呼ぶ
 - 文字列を読み終えたとき、受理状態にあれば出力は受理、そうでなければ拒否

有限オートマトン

- 有限オートマトンの正式な定義を示す
- 有限オートマトンは、5つ組 $(Q, \Sigma, \delta, q_0, F)$ で定義される
 - ▶ Q は状態と呼ばれる有限集合
 - ▶ Σ はアルファベットと呼ばれる有限集合
 - ▶ $\delta : Q \times \Sigma \rightarrow Q$ は遷移関数
 - ▶ $q_0 \in Q$ は開始状態
 - ▶ $F \subseteq Q$ は受理状態の集合
- Q. $F = \emptyset$ のとき、これは有限オートマトンか？

言語

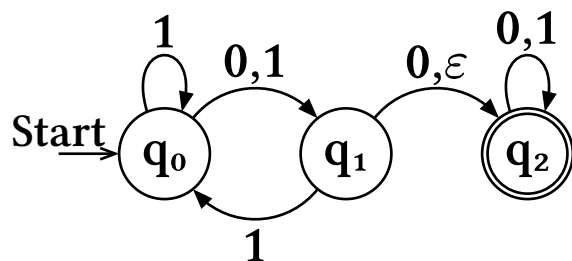
- 機械 M が受理するすべての文字列の集合を A とするとき、 A は機械 M の言語であるという $L(M) = A$
 - ▶ また、 M は A を認識する、という
- 機械は複数の文字列を受理できるが、常に唯一の言語を認識する
- Q. 機械が文字列を 1 つも受理しない場合、その機械の言語は何か？

計算

- $M = (Q, \Sigma, \delta, q_0, F)$ を有限オートマトンとし、 $w = w_1 w_2 \dots w_n$ をアルファベット Σ 上の文字列とする
- 以下の 3 つの条件を満たす状態の列 $r_0, r_1, \dots, r_n \in Q$ が存在するなら、 M は w を受理するという
 - ▶ $r_0 = q_0$
 - ▶ $i = 0, \dots, n - 1$ のとき、 $\delta(r_i, w_{i+1}) = r_{i+1}$
 - ▶ $r_n \in F$
- なお、ある有限オートマトンで認識される言語を 正規言語 と呼ぶ

非決定性

- ・ 機械の状態と次に読み出される文字によって、次の状態が一意に定まるとき、その機械を決定性機械と呼ぶ
- ・ そうでない機械を非決定性機械と呼ぶ
- ・ 非決定性は決定性の一般化なので、任意の決定性機械は非決定性機械
- ・ 以下に非決定性有限オートマトン (NFA) の例を示す



- ・ 取り得る可能性をすべて考慮して、文字列を受理するか拒否するかを決定する

非決定性有限オートマトン

- 以下に、非決定性有限オートマトンの正式な定義を与える
- 非決定性有限オートマトンは、5つ組 $(Q, \Sigma, \delta, q_0, F)$ で定義される
 - ▶ Q は状態と呼ばれる有限集合
 - ▶ Σ はアルファベットと呼ばれる有限集合
 - ▶ $\delta : Q \times (\Sigma_\epsilon) \rightarrow \mathcal{P}(Q)$ は遷移関数
 - $\Sigma_\epsilon = \Sigma \cup \{\epsilon\}$
 - 任意の集合 Q に対して、 $\mathcal{P}(Q)$ を Q の部分集合全体の集合 = 冪集合と呼ぶ
 - ▶ $q_0 \in Q$ は開始状態
 - ▶ $F \subseteq Q$ は受理状態の集合

DFA と NFA の等価性

- DFA と NFA は同じ言語のクラスを認識する (!)
 - ▶ NFA の方が表現力が高いように見えるのに
 - ▶ NFA の記述は DFA の記述よりも簡単である場合が多いので、同じ表現力であるという事実は重要
- 等価とは、2つの機械が同じ言語を認識すること

正規演算

- 2つの言語 A, B について、以下の演算を定義する
 - ▶ 和集合演算: A または B の文字列をすべて含む言語 $A \cup B = \{x \mid x \in A \vee x \in B\}$
 - ▶ 連結演算: A の文字列の後に B の文字列が続くような文字列をすべて含む言語 $A \circ B = \{xy \mid x \in A \wedge y \in B\}$
 - ▶ スター演算: A の文字列が0回以上連結されたような文字列をすべて含む言語 $A^* = \{x_1x_2\dots x_n \mid n \geq 0 \wedge x_i \in A \text{ for } i = 1, \dots, n\}$
- 正規言語は正規演算に関して閉じている
 - ▶ オートマトンを書けば自明に分かる (証明略)

非正規言語

- 正規言語でない言語も存在する
 - ▶ $\Sigma = \{0, 1\}$ とするとき、 $A = \{0^n 1^n \mid n \geq 0\}$ は正規言語ではない
- 直感的には「判定のために状態が無限に必要な言語」は正規言語ではないように思える
 - ▶ しかし、 $B = \{w \mid w \text{ は部分文字列として } 01 \text{ と } 10 \text{ を同じ個数含む}\}$ は正規言語
 - ▶ ホワイトボードにオートマトンを示す
 - これは性質を考えると最初と最後の文字が同じであるかどうかを判定するオートマトンと同じ
- 文字列の数を数えるという構造自体は同じなのに、 A は正規言語でないのに B は正規言語

ポンピング補題

- A が正規言語であるならば、ある数 p が存在して、 s が少なくとも長さ p である A の任意の文字列であるとき、 s は以下の3条件を満たす3つの部分文字列 x, y, z に分割できる
 - ▶ 任意の $i \geq 0$ について、 $xy^iz \in A$
 - ▶ $|xy| \leq p$
 - ▶ $|y| > 0$
- Q. 言語 A が正規言語でないことを示すにはどうすればよいか？