

漢文訓読をあらわす三つのモデルとチョムスキー階層

島野達雄

1. 返り点の働きとリング表示

1-1 百聞不_レ如_二一見_一(百聞は一見に如かず)は、(1)百聞と一見は漢字 2 字を続けて読み、(2)如は一見をさきに読み、不は如_二一見_一をさきに読む、という二つのルールにもとづいて読み下す。(1)のルールは漢字どうしを「順読」し、(2)は漢字 1 字をあとから「返読」する。

ここで漢字 1 字を丸カッコの対(ついで)として()であらわし、順読する漢字 2 字は()()とならべ、返読する場合は、さきに読む漢字列を()で包むことにすると、百聞不_レ如_二一見_一は、()()((()()))とあらわせる。このように返り点つき漢文を丸カッコであらわすことを、丸カッコの対の上部と下部をつなぐとリング(輪)になるので、漢文訓読のリング表示とよぶ(右図)。

【百聞不_レ如_二一見_一のリング表示】

()()((()()))

○○((○○))

○○ (○○)

1-2 漢文訓読は歴史的な産物であり、レ点・一

二点・上下点などの返り点は、自然言語である漢文(文言文)の語順を、自然言語である日本語の語順に置き換える役割をしている。返り点はいわば**自然メタ言語**である。

そのため有朋自遠方来には、有_レ朋自_二遠方_一来(朋有り遠方より来たる)、有_フ朋自_二遠方_一来_上(朋の遠方より来たる有り)の二通りの読み方がある。本稿では、有_レ朋自_二遠方_一来は()()()(), 有_フ朋自_二遠方_一来_上は()()()()と別々のリング表示であらわす。

また、使子路問津焉には、現行の使_二子路_一問_レ津焉、および歴史的な使役形特有の下点と二点を複合させた使_フ子路_二問_上津焉の二通りの返り点のつけ方があるが、どちらも「子路をして津(しん)を問わしむ」という日本語になる。本稿では、このような**異点同順**の例は、どちらも()()((()))()と同じリング表示であらわす。なお、黙字の焉は、日本語にはあらわれないが、リング表示には存在していることを注意しておく。

1-3 レ点は、原則としてレ点をはさむ前後 2 漢字の語順を入れ替える働きがある。以_レ友輔_レ仁(友を以て仁を輔く)のリング表示は()()(), レ点が連続する不_レ踰_レ矩(矩を踰えず)は((())), 匹夫不_レ可_レ奪_レ志也(匹夫も志を奪うべからず)は()()(((())))()になる。

原則として、と述べた理由は、たとえば百聞不_レ如_二一見_一の不のように、如_二一見_一という「漢字と返り点が混じった列」をさきに読むケースがあるからである。

1-4 一二点(一二三…), 上下点(上(中)下), 甲乙点(甲乙丙…), 天地点(天地(人))などは**大返り**とよばれ、一・上・甲・天点がついた漢字を読んだあと、二・下・乙・地点がついた漢字 1 字を読む、という原則がある。

博学_二於文_一-(博く文を学ぶ)は $()(())$, 遠望_三田間有_二一馬_一-(遠く田間に一馬有るを望む)は $()(())(())$ になる.

上下点の中に一二点がある不_二以_一千里_一称_上也(千里を以て称せざる也)は $((())(())())$, 甲乙点のなかに上下点, 一二点がある欲_二得_一備_二学徳_一者_上友_{甲レ}之(学徳を備えし者を得て之を友とせんと欲す)は $((())(())(())(()))$ になる.

文末の「友_{甲レ}之」には甲点とレ点の複合返り点がついており, 友と之の2漢字の語順を入れ替えたのち, 乙点のついた欲に返っている.

2. 返り点だけを用いるモデル

2-1 前章で述べたレ点および一二点, 上下点などの返り点だけを用いる場合の, すべてのリング表示の集合 $\{(), (()), (()), (())(), (())(), \dots\}$ を作り出す仕組みを, 形式言語学の手法を用いて述べてみよう. (注: 漢字数ごとのリング表示の数はカタラン数になる.)

その仕組みは3つの規則からできている.

規則1 $S \rightarrow ab$

規則2 $S \rightarrow aSb$

規則3 $S \rightarrow SS$

矢印 \rightarrow は「書き換えることができる」と読む.

規則1は「 S は ab に書き換えることができる」, 規則2は「 S は aSb に書き換えることができる」, 規則3は「 S は SS に書き換えることができる」と読む.

S は開始記号とよび, この仕組みの出発点を意味するとともに, 変数の役割をしている.

たとえば, 規則1を使うと, ab が導かれる. 規則2と規則1を使うと $aabb$ が導かれる. 規則3と規則1を使うと $abab$ が. 規則1, 2, 3を組み合わせてを使うと, $aabbab$ が導かれる. 最後の例は, $S \Rightarrow SS \Rightarrow aSbab \Rightarrow aabbab$ と導出される.

このようにして生まれる文字列(語(word)とよぶ)の a を左カッコ(に置き換え, b を右カッコ)に置き換えるとリング表示になる.

規則1, 2, 3を繰り返し何度も適用すれば, すべてのリング表示を導くことができる.

2-2 形式言語学では, 集合 $\{a, b\}$ をアルファベット(終端記号ともいう), うえの3つのような書き換え規則をまとめて文法とよび, 文法によって生成されるすべての語(word)の集合を言語とよぶ.

ここでは, 規則1, 2, 3を文法 G_0 とし, G_0 によって生成される言語を $L(G_0)$ と書く.

$L(G_0)$ は同じアルファベット $\{a, b\}$ から次の文法 G_0' によっても生成される.

規則1' $S \rightarrow \varepsilon$

規則2 $S \rightarrow aSb$

規則3 $S \rightarrow SS$

規則1'の ε (イプシロン)は空語(empty word). 「何も並べない」ことを意味している.

規則1',2,3からも $L(G_0)$ が生成される。つまり、一つの言語を生成する文法は、幾つも考えられる。同じ言語を生成する2つの文法を**等価**とよぶ。

3. 返り点・再読をみとめるモデル

3-1 再読とは、未(いま)だ~ず, 当(まさ)に~べし, のように, 漢字を2度読むことをいう。この場合, リング表示を()の対ではなく [] のような別の対にすればよい。

返り点だけを使い, 同時に再読をおこなうことを認める場合は,

規則 1 $S \rightarrow ab$

規則 2 $S \rightarrow aSb$

規則 3 $S \rightarrow SS$

規則 4 $S \rightarrow pSq$

であらわせる。pとqを [と] に置き換えれば, リング表示となる。

規則 1, 2, 3, 4 を文法 P_0 とし, 言語を $L(P_0)$ と書く。アルファベットは $\{a, b, p, q\}$ の4種類。

なお, 再読が入れ子になる例は現実にはほとんどないが, 規則 4 は猶_レ未_レ成(なお未だ成らざるがごとし), つまりリング表示 [[()]] が導出されることを認めている。

4. 返り点・再読・豎点をみとめるモデル

4-1 これまで見たように, 返り点だけを用いるときは漢字 1 字だけを返読する。軽蔑臣下を「臣下を軽蔑す」と訓読するには, 軽_二-蔑_一臣_二下_一のように軽蔑の2漢字の中央に**豎点**(たててん, ハイフン, 漢字連結記号, 連続符号とも)を置く必要がある。

豎点を使うばあいのリング表示は, 次のようにハイフンを使うのが妥当であろう。

豎点 1 個で 2 漢字をつなぐ, 卑_二-下_一之_二- (之を卑下す)は(-(())-), 患_二所_一-以立_二- (立つ所以(ゆえん)を患(うれ)う)は((-(())-)), 此非_二吾所_一-以居_二- 処子_二-也(此れ吾れ子を居処せしむる所以に非ざるなり)は(())(-(())-)-(())。この例は豎点 1 個で 2 漢字をつないだ列が入れ子になって 2 度あらわれている。

豎点 2 個で 3 漢字をつなぐ奴_二-僕_一-視之_二- (之を奴僕視す)は, (-(-(())-)-)。

豎点 3 個で 4 漢字をつなぐ馴_二-致_一-服_二-習_一天下之心_二- (天下の心を馴致服習す)は, (-(-(-(())(()))-)-)-)。

4-2 返り点のほかに豎点 1 個の使用を認める場合は, 文法 G_0 の規則 1, 2, 3 に規則 5 を加えた文法 G_1 で生成される。アルファベットは $\{a, b, t\}$ の3種類。言うまでもなくtをハイフンに置き換えれば, リング表示になる。

規則 1 $S \rightarrow ab$

規則 2 $S \rightarrow aSb$

規則 3 $S \rightarrow SS$

規則 5 $S \rightarrow ataSbtb$

規則 5 は、「 S を(- (と)-)で囲んでもよい」ことを示している。

2 個以上の豎点で 3 漢字以上をつなげるときは、「豎点が 1 つあれば、もう 1 つ追加してもよい」ことを規定する規則 6 を追加する。

規則 6 $taSbt \rightarrow tataSbtbt$

規則 6 は、「 ta と bt のあいだにある S は、 $taSbt$ に書き換えてもよい」ことを示しており、たとえば、漢字 1 字を豎点 2 個でつながれた漢字 3 字の熟語で返読する奴-僕-視之-は、規則 5, 6 をつかって、 $S \Rightarrow ataSbtb \Rightarrow atataSbtbtb \Rightarrow atataabbtbtb$ と導出できる。

豎点 3 個以上のときも同様に規則 5, 6 で導出できる。

規則 1, 2, 3, 5, 6 の、豎点の無制限の使用を認める文法を G_K とし、言語を $L(G_K)$ とする。

なお、規則 6 の代わりに

規則 v $S \rightarrow taSbt$

とすることも考えられるが、これでは $S \Rightarrow taSbt \Rightarrow taabbt$ つまり-(())-というハイフンに始まりハイフンに終わる語 (word) が導出されてしまう。また、変数 X, Y を使って、

規則 w $S \rightarrow aXSYb$

規則 x $X \rightarrow XX$

規則 x' $X \rightarrow ta$

規則 y $Y \rightarrow YY$

規則 y' $Y \rightarrow bt$

を追加しても、 X, Y の繰り返しの回数が同じとは決まっていないので、 $S \Rightarrow aXSYb \Rightarrow aXXSYb \Rightarrow atataSbtb$ という a と b の個数が異なる語 (word) が導出されてしまう。

4-3 さて、規則 6 のかわりに別途、豎点 2 個の使用を認める、

規則 12 $S \rightarrow atataSbtbtb$

を設けることもできる。規則 1, 2, 3, 5, 12 を文法 G_2 とよぶ。

豎点 3 個の使用を認める場合は、さらに次の規則 13 を付け加え、文法 G_3 とする。

規則 13 $S \rightarrow atatataSbtbtbtb$

このようにして文法の列 G_0, G_1, G_2, \dots 、および言語の列 $L(G_0), L(G_1), L(G_2), \dots$ が定義できる。

ここで、豎点の使用を m 個まで認める文法 G_m が生成する言語 $L(G_m)$ は、 m をどんどん大きくすると、「豎点の使用に制限のない、すべてのリング表示」にどんどん近づく。つまり $L(G_m)$ はどんどん $L(G_K)$ に近づき、文法 G_∞ と文法 G_K は等価と言える。

4-4 むろん、返り点・豎点に加えて再読を認めるときは、返り点の規則 1, 2, 3 に再読の規則 4 を加え、さらに豎点の使用を認める規則 5, 6 を加える。

規則 1, 2, 3, 4, 5, 6 を文法 P_K とし、言語を $L(P_K)$ とする。

また、規則 12, 13...を追加してゆき、 m 個までの豎点と再読を認める文法を P_m 、言語を $L(P_m)$ と書く。 m をどんどん大きくすると、 $L(P_m)$ はどんどん $L(P_K)$ に近づき、文法 P_∞ と文法 P_K は等価と言える。

5. チョムスキー階層とオートマトン

5-1 ノーム・チョムスキーは、1956年に Three Models for the Description of Language (言語をあらわす3つのモデル)、1959年に On Certain Formal Properties of Grammars (文法に関するある種の形式的性質について)と題した、形式文法を幾つかの型に分類する論文を発表した。現在ではこの分類を **チョムスキー階層**とよんでいる。

形式言語学における形式文法は、**非終端記号**(開始記号 S をふくむ、変数の役割をする記号)、**終端記号**(アルファベット、空語の ϵ をふくむ)の二種の記号と、これらの記号列の有限個の**書き換え規則**(生成規則)、の三つにより構成される。

チョムスキー以来の伝統にしたがい、非終端記号を A, B, \dots の大文字、終端記号を a, b, \dots の小文字であらわし、非終端記号と終端記号が混在する(片方だけでも良い。また同じ記号が重複しても良い)記号列をギリシア文字の α, β, \dots であらわす(ここでは混在列とよぶ)。

チョムスキー階層は次の四つに分かれている。

- (1) 書き換え規則の矢印の左辺が常に非終端記号ひとつであり、右辺が a 、および aB または Ba 、という形の、最も厳しい条件の書き換え規則だけを認める文法を**正規文法**、生成される言語を**正規言語**とよぶ。正則文法、正則言語とよぶこともある。
- (2) 左辺が非終端記号ひとつで、右辺は混在列、という条件だけの文法を**文脈自由文法**、言語を**文脈自由言語**とよぶ。返り点だけの文法 G_0 、返り点と再読の文法 P_0 、および返り点と堅点 m 個までを認める文法 G_m 、返り点と再読と堅点 m 個までを認める文法 P_m は、この文脈自由文法に属す。これらの規則の左辺は常に開始記号 S だけになっている。
- (3) $\alpha A \beta \rightarrow \alpha \gamma \beta$ という形の規則をもつ文法は、書き換えが混在列 α と混在列 β という前後の文脈に依存するという意味から、**文脈依存文法**とよび、言語を**文脈依存言語**とよぶ。返り点と無制限の堅点を認める文法 G_K 、返り点と再読と無制限の堅点を認める文法 P_K は、規則 $6 \quad taSbt \rightarrow tataSbtbt$ が $\alpha A \beta \rightarrow \alpha \gamma \beta$ の形をしており、この文脈依存文法に属する。
- (4) 最後に、規則にどんな制約もおかない、もっとも緩(ゆる)やかな文法を**句構造文法**、言語を**句構造言語**とよぶ。無制限文法、帰納的可算言語とよぶこともある。

5-2 これらの形式言語は、それぞれ上位の言語の真部分集合であることが知られている。すなわち、

正規言語 \subset 文脈自由言語 \subset 文脈依存言語 \subset 句構造言語

となっており、正規言語は文脈自由言語であり文脈依存言語であり句構造言語でもある。

また、文脈自由言語でない文脈依存言語が存在するが、本稿では、 G_K, P_K が文脈自由文法でなく、 $L(G_K), L(P_K)$ が文脈自由言語でない証明はおこなっていない。

5-3 正規言語は**有限オートマトン**に、言語 $L(G_0), L(P_0), L(G_m), L(P_m)$ が属する文脈自由言語は**プッシュダウンオートマトン**に、言語 $L(G_K), L(P_K)$ が属する文脈依存言語は**線形拘束オートマトン**に、句構造言語は**チューリングマシン**に受理されることが知られている。

ここでのオートマトンとは、読み取った語 (word) を内部で処理し、受理するか拒否するか、判定できる「仮想機械」をいう。たとえば、 $aabb, abaaabbb$ など $L(G_0)$ に属する「正しくカッコの付けられた語 (word)」を受理し、 $aab, baab$ など $L(G_0)$ に属さない語を拒否するオートマトンは、「 $L(G_0)$ を受理する」という。

5-4 最後に堅点でつながれた漢字列をまとめて漢字 1 字とみなし、卑₋下₋之₋(之を卑下す)を $(-())$ つまり $ataabb$ とあらわすモデル(偏リング表示とよぶ)の文法を考えてみよう。変数 T を導入し、

規則 7 $S \rightarrow aTSb$

規則 8 $T \rightarrow ta$

規則 9 $T \rightarrow TT$

としたとき、規則 1, 2, 3, 7, 8, 9 の文法 G_T は、 G_K と同様に無制限に堅点の使用を認めているが、矢印の左辺がすべて非終端記号ひとつだから、 G_K と異なり文脈自由文法になる。

再読をふくむ規則 1, 2, 3, 4, 7, 8, 9 の文法 P_T も同様で、言語 $L(G_T)$ および言語 $L(P_T)$ は文脈自由言語になり、プッシュダウンオートマトンによって受理される。

6.まとめ

6-1 本稿では、返り点・再読・堅点をふくむ漢文が形式文法であらわせ、生成される形式言語がチョムスキー階層に対応し、オートマトンに認識されることを駆け足で説明した。

今後の漢文教育を情報科学に結びつける手段として、この漢文訓読の形式文法、形式言語理論が活用されることを期待する。

5-4 で示したように、漢文訓読を対象とする数理漢文学には、情報科学の一分野(具体例)として発展する余地が大いにある。

参考文献

- (1) 『これならわかる返り点—入門から応用まで—』古田島洋介, 新典社
- (2) 『漢文訓読入門』古田島洋介・湯城吉信, 明治書院, 2011
- (3) 『計算論とオートマトン理論』A. サローマ, 野崎昭弘・町田元・山崎秀記・横森貴共訳, サイエンス社, 1988
- (4) 『オートマトン・言語理論入門』大川知・広瀬貞樹・山本博章, 共立出版, 2012