

ビデオゲーム環境における自然な発話の 教師無し二重分節と強化学習による意味付け

Unsupervised Double Articulation and Semantic Acquisition of Natural Speech in a Video Game Environment

山口 皓太郎^{1*} 岡 夏樹¹ 谷口 忠大² 尾崎 僚²
Kotaro Yamaguchi¹ Natsuki Oka¹ Tadahiro Taniguchi² Ryo Ozaki²

¹ 京都工芸繊維大学

¹ Kyoto Institute of Technology

² 立命館大学

² Ritsumeikan University

Abstract: 乳幼児は母語の音声発話を聞いて教師無しで二重分節（音韻への分節化と単語への分節化）を行い、単語の意味を学習する能力を持っていると考えられる。これまでに、明瞭に発話された小規模な語彙からなる音声発話の二重分節が工学的に実現できることが示されている（Taniguchi, Nagasaka, & Nakashima, 2016）が、本研究では迷路ゲーム環境における比較的自然的な発話の教師無し二重分節を試みた。合わせて、遅れのない報酬だけを扱う単純化した Q 学習を用いて、分節化された単語の意味学習の可能性を試した。その結果、迷路ゲーム環境における方向指示の発話に限定した場合、上下左右の意味を概ね正しく学習することができた。

1 はじめに

人間の言語獲得過程にはいくつかの発達段階がみられる。幼児は音同士の統計的な共起性を用いて連続音声信号を分節化し、単語を抽出できることが示されている [1]。幼児は時系列音声データからの単語分割問題を解いていることになる。つまり、幼児は言語を学習する過程で、言語に関する事前知識を持たずに時系列音声信号を認識し、単語列に分節化した上で、語彙を獲得する必要がある。その獲得には以下の 3 つの過程が必要であると考えられている [2]。

- 音声入力を単語に分節化するため、母語の音声的な特徴を分析する
- 音声入力を単語に分節化する
- 分節化した単語に意味を付与する

谷口らは事前に単語の知識を持たず、日本語の音節のみ認識可能な自律移動ロボットに、その場その場で発話

文により教示を行うことで、場所に関する語彙を獲得させることを検証した [3]。場所概念獲得モデルは、状態をパーティクルで表現する自己位置推定の手法である Monte-Carlo Localization [4] に場所概念を導入した確率的生成モデルを用いた。また音声認識器には、大語彙連続音声認識システム Julius^{*1} を使用し、形態素解析器には、教師なし形態素解析手法 [5] が実装された latticelm^{*2} を使用した。これらを用いた場所概念の学習の実験により、実環境においても多くの場合で目的の場所付近に場所概念がそれぞれ形成されることを確認した。

田口らは、発話からの単語の切り出しと単語の音素系列の学習に焦点を当て、連続音声と指示対象の直示による多様な言い回しでの教示から、単語の正しい分節化とその音素系列、および、単語と指示対象の直接的な対応関係を学習する手法を提案した [6]。与えられた発話と指示対象の関係を MDL (Minimum Description Length) 原理を用いて単語リストを構成し、単語の意味や文法を学習し、学習結果を用いて MDL に基づいた不

* 連絡先：京都工芸繊維大学工学部設計工学域情報工学課程
インタラクティブ知能研究室
〒606-8585 京都市左京区松ヶ崎橋上町 1
E-mail: yamaguchi@ii.is.kit.ac.jp

^{*1} 使用バージョン: dic tation-kit-v4.3.1-linux GMM 版,
<http://julius.sourceforge.jp/>

^{*2} latticelm 0.4, <http://www.phontron.com/latticelm/>

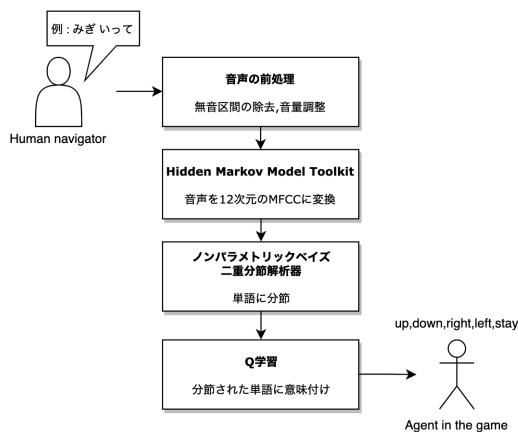


図1 提案システムの構成図

要単語の削除と単語 bigram による単語連結を行い、単語リストを再構築した。単語リストの再構築と、意味と文法の学習を繰り返すことで、正しい音素系列を得るシステムを構築した。提案手法の検証結果から、教示の言い回しや単語の知識を事前に与えることなく、音素系列を平均 83.6 % の音素正解精度で獲得できることを示した。

このように近年、言語獲得の研究が盛んに行われている。本研究の目的は、人の比較的自然的な音声発話を教師無しで二重分節し（音韻への分節化と単語への分節化を行い）、分節化された単語の意味を学習するシステムを構築することである。

2 提案システム

本研究の提案システムの構成図を図1に示す。

2D の仮想世界（迷路ゲーム環境）を設計し、その環境内にエージェントを1体用意する。用意したエージェントがユーザの発話に従い迷路を進んでいくゲームを作成した。作成したゲーム画面を図2に示す。

2.1 分節化

入力音声を単語列に分節化するための前処理として、本研究では HTK (Hidden Markov Model Toolkit)^{*3} を用いて MFCC (Mel-Frequency Cepstrum Coefficients) に変換する。MFCC に変換した後、2.2 章に示すノンパラメトリックベイズ二重分節解析器 [7] を用いて単語を分節化する。

^{*3} 使用バージョン：3.4.1,
<http://htk.eng.cam.ac.uk/>

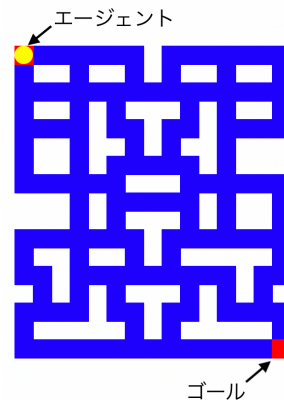


図2 本実験で使用した迷路ゲーム画面

2.2 ノンパラメトリックベイズ二重分節解析器

谷口らは人間の音声言語に含まれる二重分節構造に基づいた教師なし語彙獲得手法である NPB-DAA (Non-parametric Bayesian Double Articulation Analyzer) を提案した [7]。二重分節構造とは、音声学における最小単位である音素を下位層、音素で構成される単語を上位層とする二層の構造のことである。谷口らは、言語モデルと音響モデルを統合し、一つの生成モデルとして表現した階層ディリクレ過程隠れ言語モデルを提案し、これに対してブロック化ギブスサンプリングを行うことで言語モデルと音響モデルを同時に学習することを可能とした。階層ディリクレ過程隠れ言語モデルは Hierarchical Dirichlet Process Hidden Semi-Markov Model [8] を拡張して得られる、潜在的に二重分節構造をもつ時系列データに対する生成モデルである。言語モデルと音響モデルの同時推定により音素認識誤りに対応し、母音のみで構成され、明瞭に発話された人工的な音声データに対して高精度な単語分割ができることを示した。しかし、より自然的な発話音声での単語分割は検証されていない。そこで本研究では人の比較的自由的な発話音声を入力とし、教師無し二重分節を試みる。

2.3 ラベル付け

分節結果に意味付けをするために、各単語列にラベル付けを行う。例えば、入力音声「みぎいって」、「ひだりいって」が「みぎ/いっ/て」、「ひだり/いっ/て」と分節化されたとするとそれぞれ「0 1 2」、「3 1 2」とラベル付けする。ラベルは NPB-DAA の処理ごとに結果が変わるため、「みぎ」という音声が必ずしもラベル 0 にラ

ベル付けされるとは限らない。

2.4 学習

ラベル付けされた単語への意味付けについては、最初ではできるだけ単純な方法を試す方針をとり、遅れのない報酬（即時報酬）だけを扱う単純化した Q 学習を採用する。Q 値は迷路ゲーム開始時に全て 0 で初期化する。学習における状態、エージェントが選択する行動、得られる報酬、Q 値の更新について以下にそれぞれ示す。

2.4.1 状態

エージェントは、入力音声を NPB-DAA で分節化した単語列にラベル付けを行った結果を状態として持つ。例えば、ラベル s1, ラベル s2, ラベル s3 が入力された時、これを状態として持つ。

2.4.2 行動

エージェントの行動は上, 下, 左, 右の 4 通りである。ただしエージェントが選択した行動の先に壁が存在する場合はエージェントはその場に待機する。エージェントの行動選択には、ソフトマックス手法を用いる。ソフトマックス手法では状態 s のときに行動 a を選択する確率を $\pi(s, a)$ として、

$$\pi(s, a) = \frac{\exp\left(\frac{Q(s, a)}{\tau}\right)}{\sum_{a_i \in A} \exp\left(\frac{Q(s, a_i)}{\tau}\right)} \quad (1)$$

と計算される確率で行動を選択する。ただし、 τ は温度パラメータで、温度が大きいほど全ての行動が同じ確率で選択されやすくなり、小さければ、Q 値の大きい行動をとりやすくなる。温度は $\tau = 0.5$ とした。1 発話には複数のラベルが対応しているため、各ラベルから選択された行動のうち、最も多い行動をエージェントは実行する。例えば、ラベル s1, ラベル s2, ラベル s3 を状態として持つ時、各ラベル (s1, s2, s3) の Q 値に従って行動が上, 上, 右と選択されたとすると、エージェントは最も選択された数が多い上に行動する。一方で、行動が上, 下, 右と選択された場合は、Q(s1, 上), Q(s2, 下), Q(s3, 右) の中で値が最も大きい行動を実行する。ただし、Q 値が等しい場合は、その中からランダムに行動を実行する。

2.4.3 報酬

評価は「良い」または「悪い」の 2 種類で評価「良い」には +1, 評価「悪い」には -1 の報酬をそれぞれ与える。

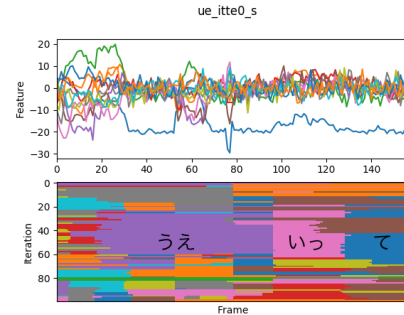


図 3 分節結果 (MFCC12 次元)

2.4.4 Q 値の更新

実際に実行した行動の Q 値を更新する。例えば、ラベル s1, ラベル s2, ラベル s3 を状態として持った時、行動上, 上, 右が選択され、上に行動した場合、Q(s1, 上) と Q(s2, 上) と Q(s3, 上) を更新する。Q 値の更新式は以下ようになる。

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r - Q(s, a)) \quad (2)$$

パラメータは学習率 $\alpha = 0.1$, 割引率 $\gamma = 0$ とする。割引率 $\gamma = 0$ としたのは即時報酬のみを扱うためである。

3 分節化の予備実験

迷路誘導時に発話しそうな音声を第一著者が発声したものを録音し、単語に分節化することが可能かを調査するために予備実験を行った。音声録音時に発話前後の無音区間が入らないようにするために、一定の閾値を設定し、その閾値以上の音声 flowed 時のみ録音するシステムを作成した。この録音システムを用いて（無響室ではない）通常の居室で発話者以外人がいない状態で「うえいって」、「したいって」、「みぎいって」、「ひだりいって」という音声を各 3 個ずつ、合計 12 個録音した。録音した音声を HTK を用いて MFC C12 次元に変換し、NPB-DAA を用いて分節化を行った。その結果例を図 3 に示す。図の上部に示す波形は入力音声の時系列データである。図の下部に示すグラフは音声から推定した単語のまとまりを表す。グラフの横軸は時間 [ms], 縦軸はギブスサンプリングのイテレーション回数を示す。今回の予備実験においてイテレーション回数は 100 回とした。結果から、音声の後半部分にノイズが入っているが単語列に分割できているように考えられる。

次に、ラベル数を 0 ~ 5 の 6 個として、予備実験で用いたデータをラベル付けした結果を表 1 に示す。表 1 か

表 1 分節化の予備実験結果

入力音声	ラベル
うえいって	4 1 0
	1 0
	4 1 0
したいって	5 1 0
	5 0
	5 0
ひだりいって	3 1 0
	3 4 1 0
	4 3 4 1 0
みぎいって	2 3 0
	2 1 0
	2 1 0

ら「うえ」がラベル 4 に、「した」がラベル 5 に、「ひだり」がラベル 3 に、「みぎ」がラベル 2 に概ね対応していることがわかる。さらに、「いって」という音声は多くの場合ラベル 1, ラベル 0 が付与されていると考えられる。以上の結果から、概ね想定通り分節結果にラベル付けができていていることがわかる。

4 実験

本実験の目的は実験参加者の比較的自由的な発話内容を録音し、単語に分節化することに加えて、分節化された単語に遅れない報酬だけを扱う単純化した Q 学習で試験的に意味付けを行うことである。実験では JBL 社の ENDURANCE RUN のマイクを用いて録音をする。手順 1 で実験参加者の自然的な発話の教示データのため込んだ後、手順 2 で 2.4 に示す学習をオフラインで行う。

4.1 手順 1

手順 1 では実験参加者の比較的自由的な発話を録音することを目的とする。実験参加者に図 2 の迷路ゲーム画面を提示し、迷路内にいるエージェントに対して自由的な発話をしてもらう。エージェントは実験参加者の発話にかかわらず 1.5 秒毎に上, 下, 左, 右のいずれかに移動する。実験前に実験参加者に対して以下の指示をする。

- (1) 実験参加者の発話によって迷路画面内のエージェントをゴールに誘導する。
- (2) 幼児に話しかけるように自由に発話をする。
- (3) エージェントは 1.5 秒毎に上, 下, 左, 右のいずれかに行動する。
- (4) 発話を 40 回録音した時点で実験終了とする。

実験終了後に発話データをもとに NPB-DAA を用い

て分節化を行う。

4.2 手順 2

手順 2 では手順 1 で収集した発話データの分節結果に、強化学習で試験的に意味付けを行うことを目的とする。手順 2 の実験手順は次の通りである。

- (1) 図 3 の迷路ゲーム画面内にいるエージェントに対して手順 1 で収集した発話の分節結果（ラベル付け結果）を選択する。例えばエージェントに右に行くように指示するときは、「みぎいって」、「みぎみぎ」の分節結果を選択する。上, 下, 左に行くように指示するときも同様に選択する。各方向の行動指示の選択回数に差が出ないように均等に選択する。今回、同一方向の行動指示の発話には（例えば、「みぎいって」、「みぎみぎ」のように）2 種類あるが、この選択頻度も均等にする。また、収集した発話の頻度に比例した確率で、方向指示以外の発話も選択する。
- (2) 選択された発話の分節結果と Q 値をもとに 2.4.2 の方法に従い、エージェントが上, 下, 左, 右のいずれかに移動する。
- (3) エージェントの行動に対して +1（「良い」）、または -1（「悪い」）のどちらかの評価をする。発話が行動指示であった場合は、それとエージェントが選択した行動が合っていたかどうかにより決めるが、行動指示でない発話に対する報酬は、それがたまたまゴールに向かう行動であれば、+1 を、そうでなければ -1 を与える。
- (4) (1) ~ (3) の流れを 100 回繰り返す。

以上の処理を行うシミュレータを Python で作成した。

5 結果

今回は理工系の大学生 1 名を実験参加者として、実施してみた実験に関して参考までに報告する。

5.1 手順 1 の結果

1 回目、実験参加者に 4.1 に示す指示をしてから実験を行ったところ、録音することができた発話内容は「うえ」、「した」、「ひだり」、「みぎ」などの 1 単語で構成される発話が多かった。結果として、比較的自由的な発話データを収集することができなかつたため、2 回目、次の指示をして再度実験を行った。

- (1) 実験参加者の発話によって迷路画面内のエージェントをゴールに誘導する。
- (2) 発話は「うえいって」、「したいって」、「ひだりいって」、「みぎいって」、「うえうえ」、「したした」、「ひだりひだり」、「みぎみぎ」、「そうそう」、「そうだよ」、「ちがうちがう」、「ちがうよ」の12種類が発話できることにする。ただし、少し言い換えて自由に発話しても構わない。
- (3) 実験中に実験参加者が発話できる内容を確認できるように、上記12種類の発話内容を記載したメモ用紙を置いておく。
- (4) エージェントは1.5秒おきに上下左右のいずれかに行動する。
- (5) 発話を40回録音した時点で実験終了とする。

2回目に録音した発話内容と、ラベル数を8個として分節化を行った結果を表2に示す。

表2の「(不明瞭な発話)」とは、「えー」や「んー」といった発話音声小さい言い淀みや感嘆詞であった。表2より、「いって」はラベル0、ラベル1に対応しており、この部分に関しては多くの発話において正しく分節化できていることがわかる。また「うえ」にラベル7、「みぎ」にラベル5、「ひだり」にラベル2が多く含まれていることがわかる。一方で「した」はラベル1に対応しているように考えられるが、「いって」のラベル0、ラベル1と重複しているため、正しくラベルを割り当てることができなかつたと考えられる。「ひだり」と同様に、上記以外の他の発話もラベルを正確に割り当てることができなかつた。

次に、「(不明瞭な発話)」とデータ数の少なかった「そうそう」、「そうだよ」、「ちがうちがう」、「ちがうよ」を除いた26個のデータを用いてラベル数を7個として分節化を行った。その結果を表3に示す。表3より「いって」はラベル「5 4」、「うえ」はラベル0、「した」はラベル3、「ひだり」はラベル6、「みぎ」はラベル1が割り当てられていることが多いことがわかる。表2と表3から、不明瞭な発話やデータ数が少ない発話を除くことで比較的自然的な発話を単語に分節化できたということを示す結果となった。

5.2 手順2の結果

表2の分節結果を用いて2.4の学習方法で意味付けを試みた。1試行の学習回数を100回として10試行実験を行い、各Q値の平均を算出した。結果を表4に示す。また、Q値(Q(7, 上), Q(1, 下), Q(2, 左), Q(5, 右))の学習経過の平均値(平均値の変化)をグラフにしたも

表2 発話内容と分節結果

発話内容	発話回数	ラベル
うえいって	3	7 0 1
		1 5 0 1
		7 5 0 1 2
したいって	3	1 5 0 1
		1 5 0 1
		1 5 0 1
ひだりいって	3	5 0 1
		2 0 1
		2 0 1
みぎいって	4	5 0 1
		5 0 1
		5 0 1
		5 0 1 2
うえうえ	3	7 1
		7 1
		7 7
したした	4	1 3 1
		1 3 1
		1 3 1 1
		1 3 1 1
ひだりひだり	3	2 0 6 2
		2 6 0 6 2
		2 6 1 2
みぎみぎ	3	5 0 5 6
		5 6 0 6 5 6
		5 0 5 6
そうそう	2	1 4 1 4
		1 4 1 4
そうだよ	1	1 4 1 4 1
ちがうちがう	1	1 4 3 4
ちがうよ	2	1 4 1 4
		1 4 1 4
(不明瞭な発話)	8	2 3 2
		1 7 2 1 7 0
		6 0
		5 2
		0 2
		7 0
		3 2
		5 1 7 2 7
合計	40	

のを図4に示す。表4はQ値の平均値が高いほど影を濃く色付けをし、Q値の平均値が低いほど薄く色付けした。

表4から、上の行動はラベル7、下の行動はラベル3、左の行動はラベル2、右の行動はラベル5が他のラベルと比べて高くなっていることがわかる。右の行動に関してラベル5の値も高くなっているが、ラベル6も0.165と高くなっている。これは「みぎみぎ」の発話にラベル6が多く含まれることが原因であると考えられる。表2

表3 発話内容と分節結果（「(不明瞭な発話)」, 「そうそう」, 「そうだよ」, 「ちがうちがう」, 「ちがうよ」は除く）

発話内容	発話回数	ラベル
うえいって	3	056
		054
		054
したいって	3	354
		454
		454
ひだりいって	3	654
		654
		6654
みぎいって	4	154
		654
		154
		1542
うえうえ	3	00
		00
		06
したした	4	33
		33
		33
		33
ひだりひだり	3	66
		662
		662
みぎみぎ	3	132
		131
		12
合計	26	

表4 学習終了後の Q 値の平均 (表2のラベルを用いた場合)

ラベル \ 行動	↑	↓	←	→
0	-0.330	-0.418	-0.229	0.078
1	-0.418	0.065	-0.510	-0.326
2	-0.200	-0.328	0.365	-0.182
3	-0.178	0.338	-0.075	-0.264
4	0.046	-0.210	-0.050	0.019
5	-0.400	-0.360	-0.489	0.307
6	-0.312	-0.312	0.061	0.165
7	0.200	-0.346	-0.259	-0.208

の「した」を含む発話の分節結果に注目すると、ラベル1, ラベル3がこの順番で多く含まれている。しかし、表4で下の行動はラベル3の値がラベル1よりも高くなった。この原因として、ラベル1は「した」以外の他の発話にも含まれていることが多いが、ラベル3は「したした」の発話の分節結果に集中していることが考えら

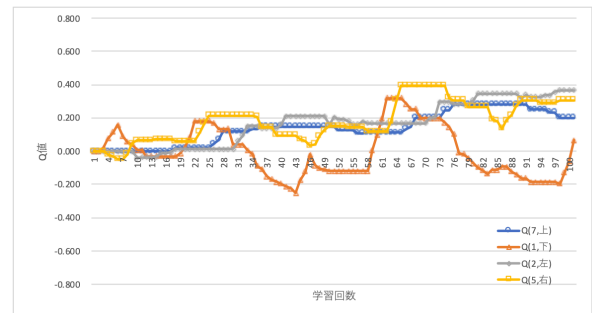


図4 Q 値 (Q(7, 上), Q(1, 下), Q(2, 左), Q(5, 右)) の学習曲線の平均値

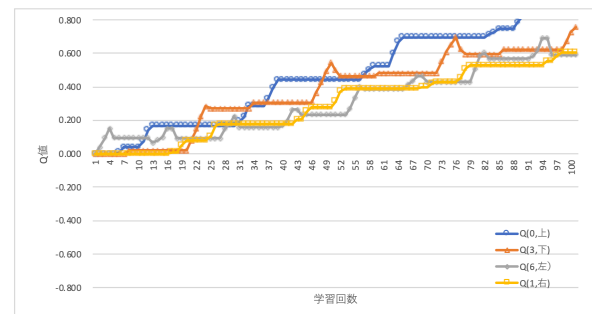


図5 Q 値 (Q(0, 上), Q(3, 下), Q(6, 左), Q(1, 右)) の学習曲線の平均値

表5 学習終了後の Q 値の平均 (表3のラベルを用いた場合)

ラベル \ 行動	↑	↓	←	→
0	0.838	-0.386	-0.444	-0.334
1	-0.405	-0.417	-0.314	0.605
2	-0.305	-0.268	-0.032	0.246
3	-0.390	0.755	-0.428	-0.203
4	-0.415	-0.146	-0.112	-0.058
5	-0.227	-0.291	-0.113	-0.015
6	-0.209	-0.599	0.589	-0.519

れる。また、図4から学習初期の段階ではQ値の学習が進まなかったことがわかる。ラベル7, ラベル2, ラベル5に関しては学習回数が増えるにつれて徐々にQ値が増加しているが、ラベル1はQ値が減少傾向にある部分がみられることから、学習がうまく進まなかったと考えられる。

次に、表3の分節結果を用いて2.4の学習方法で意味付けを試みた。1試行の学習回数を100回として10試行実験を行い、各Q値の平均を算出した。結果を表5に示す。また、Q値(Q(0,上), Q(3,下), Q(6,左), Q(1,右))の学習経過の平均値(平均値の変化)をグラフにし

たものを図5に示す。表4と同様に、Q値の平均値が高いほど影を濃く色付けをし、Q値の平均値が低いほど薄く色付けした。表5から、上の行動はラベル0、下の行動はラベル3、左の行動はラベル6、右の行動はラベル1の値が最大となっている。表3の分節結果と比較すると概ね正しく意味付けをすることができたと言える。また、ラベル2は右の行動のQ値が0.246と最大となっている。この原因として、表3の「みぎ」を含む発話の分節結果にラベル2が多く割り当てられており、学習の過程でラベル1と同時に報酬を受け取った為、ラベル2は右の行動のQ値が高くなったことが考えられる。図5から、学習回数が増えるにつれてQ値は順調に増加していることがわかる。また、図4と比べてQ値は比較的早い時期から右肩上がりになっていることがわかる。結果として、不明瞭な発話やデータ数の少ない発話を除いた場合、概ね適切に分節化することができ、正しく上下左右の意味付けをすることができた。

6 結言

6.1 まとめ

本研究では比較的自然な発話を、ノンパラメトリックベイズ二重分節解析器を用いて教師なしで単語に分節化することを検証した。実験では、迷路ゲーム内のエージェントを発話によりゴールに導くタスクを設定して実験参加者の発話データを収集し、得られた発話データを元に分節化を行った。その結果、「うえいって」、「したいって」、「ひだりいって」、「みぎいって」、「うえうえ」、「したした」、「ひだりひだり」、「みぎみぎ」という発話に限定した場合は概ね適切に単語に分節化できた。しかし、データ数の少ない発話や不明瞭な発話を含めた場合、正しく分節化できない発話が少なからずあった。次に、遅れのない報酬だけを扱う単純化したQ学習のアルゴリズムを用いて、発話（分節化された単語列）に、上、下、左、右のいずれかの方向に進めという意味付けを行うシステムを構築した。データ数の少ない発話や不明瞭な発話を含むデータに対する学習は、学習の立ち上がりが遅い傾向が見られたが、データ数の少ない発話や不明瞭な発話を除いた場合、学習がより速く進むということが観測された。

6.2 今後の展望

本研究では音響特徴量として12次元のMFCCを用いたが、近年の機械学習ではLog-Melspectrumが主要になりつつある。MFCCは音声に短時間フーリエ変換

を施しメルフィルタバンクを適用し、離散コサイン変換をすることで得られる。一方、Log-Melspectrumは音声に短時間フーリエ変換を施しメルフィルタバンクを適用し、対数変換をしたものであり、MFCCの導出過程で得られる。離散コサイン変換を行わないため必要な情報が失われにくいとされている[9]。分節精度の向上のために今後の研究ではMFCCと合わせてLog-Melspectrumを使うことを考えている。また、実験において実験参加者により自由な発話をしてもらえるように、実験前のインストラクションや迷路ゲーム画面を工夫して、自由な発話からの分節化を試みたい。

学習に関して、構築した言語獲得システムは「右」や「上」といった方向指示語に焦点を当てたものであるが、今後は入力音声「行って」やエージェントが選択した行動を評価する発話（例えば「そうそう」、「ちがうよ」など）にも意味付けを行うことを考えている。評価発話の意味が分かるようになれば、それをを用いて効率よく行動の学習を進めることができる。

参考文献

- [1] J. Saffran, R. Aslin, E. Newport: Statistical learning by 8 month old infants, *Science*, Vol. 274, No. 5294, pp. 1926–1928 (1996)
- [2] 林安紀子: 乳児における言語のリズム構造の知覚と獲得, *Journal of the Phonetic Society*, Vol. 7, No. 2, pp. 29–34 (2003)
- [3] 谷口彰, 稲邑哲也, 谷口忠大: 実ロボットを用いた自己位置と語彙の同時推定による音声言語獲得, *人工知能学会誌*, Vol. 29 (2015)
- [4] S. Thrun, W. Burgard, D. Fox: Probabilistic Robotics, *MIT Press*, (2005)
- [5] G. Neubig, M. Mimura, S. Mori, T. Kawahara: Bayesian learning of a language model from continuous speech, *IEICE*, Vol. E95-D, No. 2, pp. 614–625 (2012)
- [6] 田口亮, 岩橋直人, 船越孝太郎, 中野幹生, 能勢隆, 新田恒雄: 統計的モデル選択に基づいた連続音声からの語彙学習, *人工知能学会論文誌*, Vol. 25, No. 4, pp. 549–559 (2010)
- [7] T. Taniguchi, R. Nakashima, S. Nagasaka, Y. Tada, K. Hayashi, R. Ozaki: Nonparametric bayesian double articulation analyzer for direct language acquisition from continuous speech signals, *IEEE*, Vol. 8, Issue. 3, pp. 171–185, (2016)
- [8] M. Johnson, A. Willsky: The hierarchical dirich-

let process hidden semi-markov model, *arXiv*,
(2012)

- [9] H. Purwins, B. Li, T. Virtanen, J. Schlüter,
S. Chang, T. Sainath: Deep learning for audio
signal processing, *Journal of Selected Topics of
Signal Processing*, Vol. 13, No. 2, pp. 206–219,
(2019)