

書き起こし作業とチェック作業について

高沢 美和

田川 恭識

0. はじめに

本研究班では熊本方言談話の書き起こし作業を行った。主な作業者は、大学生、大学院生、大学院修了生である。報告者は書き起こし作業とそのデータに関するチェック作業に携わっており、作業を行う上で困難であった点や、議論として挙げたものについて、まとめることとする。

1. 機器の使用

1-1. 作業環境

まず、作業を行う前に、大容量の音声データを扱うに耐えるコンピュータ、音を十分に確認しうるヘッドフォンの貸与が必要である。

また、報告者は、音声に関する機器や、ソフトウェアの使用に慣れていなかったため、作業そのものではなく、それ以外の部分で時間を必要以上にロスしてしまったところがある。今後、このような作業を第三者に委託する場合は、あらかじめ機器の使用等についての十分な指導を行うことが望まれる。

また、書き起こしには **MultiTrans** というフリーソフトを使用した。どの作業もソフトの使用に慣れるのに時間がかかったようである。このソフトは詳しい説明があるものではなく、自分で作業をしながら慣れるしかない。また、作業を自動保存などしてくれないので、再生ボタンのすぐ横にある「閉じるボタン」を誤って押してしまった場合、長時間書き起こしたものが消えてしまうということがある。時間をロスしてしまうことが度々あった。作業に没頭している時ほど保存に注意が向かないものであるが、これは、常に保存する癖をつけるしか対処の仕様がなない。

2. 書き起こし基準に関するもの

音声データの書き起こしを行うにあたって、『熊本方言談話班 文字化の基準・規則』（以下、「基準」）を設け、作業者はそれに則って作業を行った。この基準の大枠は、大阪大学大学院社会言語学研究室における文字化規則（真田信治（2006））を参考に作成したものである。ただしこの基準は暫定的なものであり、作業を行っていく中で気付いた不備は、その都度話し合い、第 4 版まで改定した。今後も必要があれば改定される可能性がある。基準に関して議論となったものは、以下の通りである。

2-1. 音声に関する記号（タグ）の付け方

笑いやささやき声など、音声的特徴があるものにタグをつけるという規則があったが、

どのようにつけるのか、問題となった。

はじめは、笑いながら何かを発話しているものと、笑いが独立しているもので、付与方法を変えていたが、議論の結果、笑いながら何かを発話しているものも、笑い自体が独立しているものも、同様の記号を付けることにした。また、はじめは {笑い} もしくは {笑いながら} というタグであったが、上記の理由と作業の簡略化のため、{w} というタグに統一した。

また、笑いながら何かを発話している場合、どこからどこまで笑いを含んでいるのか視覚的にわかりやすくするために、該当する文節すべてにタグをつけることとした。

例) あはははっ {笑い}

{笑いながら} これ 犬が 掘ったんでしょ [↑] ?

⇒ {w} あはははっ

{w} これ {w} 犬が {w} 掘ったんでしょ [↑] ?

他に、頻発するため記号で表すことにしたタグには、言い淀み {y}、ささやき {s} がある。これらの使い方も上記の方法と同じで、該当する文節の頭につけることにした。

2-2. タグ付けの恣意性

句末や文末に付与した音調¹のタグである[↑][↓][→][↑↓][^]について、一人の作業員内でも、ある同じ音調の音声を聞いても、作業に時間がかかった場合、前後で判断にずれが生じる場合が少なくない。つまり、先ほどは上昇調と判断したが、次には上昇下降調と判断する、というような場合である。

また、作業員によってもかなり判断のずれがある。データ間でずれがあると全体として意味をなさないため、訂正する必要があるが、作業員が熟考の上つけたタグでも訂正することになるのなら、結局は最終チェック者がはじめからタグをつけたほうが早いのではないかと思われた。

2-3. 言い淀み

言い淀みに関するタグをつけるという規則があったが、「言い淀み」という語自体が曖昧なため、何をもちいて言い淀みとするのかということに、作業員間で認識のずれがあった。たとえば、繰り返して述べたものを指すのか、(例：わか わかんないなー)、ポーズを置くものを指すのか (例：わ… かんないなー) ということである。現時点ではどちらにも言い淀みタグを付けている。

¹本研究では、句末・文末のイントネーションの符号化に当たり、上村幸雄 (1989)、郡史郎 (2003)、吉沢典男 (1960) 等を参考とした。

2-4. 表記の問題

慣用的に、漢字で表記するか仮名で表記するか定かでないもの、漢字表記がいくつかあるものを、どう表記するかが問題となった。

例)

- ・たぶん/多分
- ・うさぎ/兎/兔
- ・たぬき/狸
- ・いじめ/苛め/虐め
- ・わけ/訳
- ・はたして/果たして/果して
- ・すごい/凄い
- ・やけど/火傷
- ・こども/子ども/子供
- ・とりあえず/取りあえず/取り敢えず/取敢えず
- ・(副詞)よく/良く
- ・読んで字のごとく/読んで字の如く
- ・やつ/奴

これらは、「通用として、平仮名表記が広く行われているかどうか」によって平仮名にするか漢字にするか決めては、との意見もあったが、その判断は、個人によって分かれるのではないかと考えられる。よって、データの最終チェック者が、それらの表記については目についたものから一つずつチェックして統一するという方法を取った。

2-5. 文節の区切り方

初めの基準では、学校文法の区切り方を採用することにしていた。しかし、やはりその区切り方では違和感があるという意見があり、文法化しているものや意味のまとまりの強い部分は、区切らずに書くということになった。

例えば、所謂「て形」を使った構文は区切らない。

「てください」「ています」「もいいです」「てはいけません」「てから」「てあげます」「てもらいます」「てくれます」「てやります」「てある」「ておく」「てしまう」「てみる」

(例1) 机の 上に 置いてあるのを 食べられてしまったんです

(例2) 机ん《机の》 上に 置いてあったとば《置いてあったのを》 食べられたつたい《食べられたんだよ》

他に区切らないものとしては、「～かもしれない」「～たわけじゃない」などがある。

研究目的によって「この研究における文節」とはどういうものか、というスタンスを明示する必要があるといえるが、現段階ではそれには至っていない。これから研究を行うにあたって、少なくとも問題にしたい箇所については揃えておくということが望まれる。

2-6. 句読点と時間的な間

初めの基準では、文中でポーズのある箇所に「、」、発話文末に「。」、語尾が言い淀み、文が途中で終了した場合に「…」をつける、とされていた。しかし、文が終わったのかどうかについて判断できないものも多く、その付け方が恣意的になってしまう恐れがある。そのため、句読点はすべて排除し、時間的な間があるところに「…」を付与することとした。

2-7. 複数の読み方があるもの

初めの基準で、「複数の読み方があるものを漢字で表す場合は、読み方をひらがなで‘ ’に入れて示す」（例えば、一人前…ひとりまえ、いちにんまえ、のようなもの）というものがあつたが、漢字のみでなく平仮名であっても、文字を見ただけでは何と発音しているかわからないものがある。たとえば、改行後の文頭に「は」と表記した場合、実際には前の文を受けての助詞の「は」であるという場合がある。文頭に助詞と考えられる「は」が出てきた場合、その音が「ha」ではなく「wa」であるということを示すために、「は‘わ’」と表記することとした。「へ」「を」についても同様にその音を示すことにした。

例) B3 の冒頭。

- B1: いや ちょと 待ってー[→] メイン メインの 話を やってー これ[→] ふふ
A1: ふふっふ 耳なし芳ーはー[↑↓]
B2: うん[↓]
A2: え[↑]? 全然 自信ないわー[↓]… よく 考えたら 後半しか 知らん[↓]
B3: は‘わ’ー あれでしょ? 目が 見えなくてー[↑↓]
A3: うん[↓]
B4: こー 琵琶ーか なんかを 弾いてる 人がー[↑↓]…
A4: うん[↓]
B5: いてー[↑↓]

2-8. 全平坦の音調

書き起こしを行うにあたって、熊本方言として特徴的であると考えられる、いくつかの音調について記述することを試みた。その一つが、複数の文節にまたがって高さの変化がほとんど感じられずに続く音調で、このような音調に対し「全平坦」という仮の名称を与えた。この「全平坦」の音調について、始めはその始点を「up/」 終点を「/」で区切る方法を取った。その際、共通語アクセントで平板型である語と、アクセント核を持つ語が組み合わさって文になっているものの表記について問題となった。例えば、

1) その場で 斬られるかどうかまでは 知らんけども

(「その場」→平板、「斬られる」→アクセント核を持つ)

このような文が、全平坦の音調で、(/ / は音が高いことを示す)

2) / その場で 斬られるかどうかまで / は 知らんけども

と発話されていた場合、共通語でも「その場で」は平板型であるが、音調的に一続きになっていると思われる場合は、

3) その場で up / 斬られるかどうかまで / は 知らんけども

ではなく、

4) up / その場で 斬られるかどうかまで / は 知らんけども

と記すことにした。

しかしながら、始点と終点の認定には、作業員間でばらつきが大きく一貫性が保証されないと判断されたため、最終的には文字化資料から up / / の記号を取り外すこととした。

なお、経験的には上記の音調と句末・文末の音調である[↑↓](上昇下降調、昇降調)とが共起する割合が高いようである。従って、「全平坦」の音調を見たい場合は、[↑↓]が付与された文節とそれに先行する複数の文節について調べてみることを推奨する。

2-9. その他音調に関して

その他の音調的特徴の記述として、共通語では声が下降すべきでないところで下降している音調は、dw / / で示すこととした。

例

1) dw / おおまか / だよね (大まかだよね)

2) dw / いしつ / だよね (異質だよね)

3) dw / それは / ねー

4) dw / なんでも / ない

しかし、ただでさえピッチの変化幅が抑えられることのある熊本方言話者の音声について(嵐洋子(2008))、高低 2 段階の二値的判断を行うのは時に非常な困難を伴う。結果として判断の一貫性が保証されないと判断されたため、文字化資料からは削除した。

2-10. 改行

発話者 A が続けて発話している間に、発話者 B が会話のターンを取るわけではないが相槌を打っているような場合、音声をどこで区切るのかということが、作業員によってずれがあり、修正にかなり時間がかかった。

話している途中で相手の言語的イベントが起こったとき、不自然にならないところで区切って改行することにした。

例 1

正)

A: えーと
B: うん
A: たぬきがね
B: うん
A: おばあさん 殺して 鍋に 入れるんよ
B: はー まじでー
A: そう エグいやろ
B: 最悪やね

誤) (発話を区切らなすぎた場合、次のようになる)

A: えーと たぬきがね おばあさん 殺して 鍋に 入れるんよ エグいやろー
B: うん うん はー まじで 最悪やね

例 2

正)

A: 耳なし芳一って なんか 蟹が 出てこんかった?
B: あー でてきた 気もするね
A: 出てきたよね? あれ?
B: うん

誤) (発話を区切りすぎた場合、次のようになる)

A: 耳なし芳一って
A: なんか 蟹が
A: 出てこんかった?
B: あー 出てきた 気もするね
A: 出てきたよね?
A: あれ?
B: うん

2-11. 分かち書き

初めの基準で、「この」「その」「あの」などの指示詞は分かち書きしない、というものがあつたが、指示詞なのかフィラーなのか判別しにくいものも多いため、このルールは削除した。

例) あのー 犬がねー

2-12. 吸気音、呼気音

吸気音や呼気音は、マイクが口に近いため拾っただけなのか、相手に聞こえるほどのもの（何らかの意味を持ちうるもの）なのか判別しにくい、相手に聞こえるほどの大きさであると判断されるもののみ、タグとして記すこととした。

3. おわりに

以上、簡単ではあるが、書き起こし作業を行う上で議論した点についてまとめた。未だ明確に基準が定まっていないものもあり、今後さらに検討する必要がある。

参考文献

- 嵐洋子 (2008) 「持続時間及びピッチ変動が長音の知覚に与える影響 - 東京方言話者と熊本市方言話者の比較 -」『音声言語VI』, 近畿音声言語研究会.
- 真田信治 (2006) 「奄美」『薩南諸島におけるネオ方言（中間方言）の実態調査』（課題番号：15520288）平成 15-17 年度科学研究費補助金（基盤研究 C）研究成果報告書.
- 上村幸雄 (1989) 「日本語のイントネーション」『ことばの科学 3』言語学研究会編, むぎ書房.
- 郡史郎 (2004) 「イントネーション」『朝倉日本語講座 3 音声・音韻』, 109-131, 朝倉書店.
- 吉沢典男 (1960) 「イントネーション」『話しことばの文型 (1) 一対話資料による研究一』, 秀英出版.