

ニュース音声認識のための ($n \geq 4$)-gram を併用する言語モデル

加藤直人*1 浦谷則好*2 江原暉将*1 安藤彰男*1

*1 N H K 放送技術研究所
〒 157-8510 東京都世田谷区砧 1-10-11
{katonao, eharate, ando}@str1.nhk.or.jp

*2 A T R 音声翻訳通信研究所
〒 619-0288 京都府相楽郡精華町光台 2-2
uratani@itl.atr.co.jp

あらまし 音声認識の精度向上には言語制約が強い言語モデルを構成すること必要であり、その一つの方法がタスク適応である。一方で、タスク適応しすぎると頑健性が損なわれるという問題がある。本稿では($n \geq 4$)-gram を利用することによりタスクへの適応をしつつ、2,3-gram も利用することで頑健性もそれほど損なわない言語モデルについて述べる。提案する言語モデルでは($n \geq 4$)-gram を、従来の n-gram のように宣言的知識として記憶するのではなく、単語出現位置辞書という概念を導入して手続き的知識として記憶することによりそれほどデータ量を増やすことなく利用している。本言語モデルを放送ニュースに応用し、その perplexity による評価実験を行ったところ、良好な結果を得た。

キーワード 言語モデル, 音声認識, n-gram, タスク適応, テレビニュース, perplexity

A New Language Model by using ($n \geq 4$)-gram for Broadcast News Speech Transcription

Naoto Katoh*1 Noriyoshi Uratani*2 Terumasa Ehara*1 Akio Ando*1

*1 NHK Science and Technical Research Laboratories
1-10-11 Kinuta, Setagaya-ku, Tokyo 157-8510, Japan
{katonao, eharate, ando}@str1.nhk.or.jp

*2 ATR Interpreting Telecommunications Research Laboratories
2-2 Hikaridai, Seika-cho, Soraku-gun, Kyoto 619-0288, Japan
uratani@itl.atr.co.jp

Abstract Language model adaptation is one of the important methods to construct a speech recognition system for practical use. The conventional adaptation methods adjusted n-gram estimated from various task corpora to ones from a specific task corpus. However the methods are not so effective in some tasks such as TV news, because some of TV news does not use news scripts.

This paper proposes a new language model for Broadcast news speech transcription. Our model can not only adapt to a specific task but also deal with the more tasks by dynamically using ($n \geq 4$)-gram and 2,3-gram. The proposed method can reduce amount of ($n \geq 4$)-gram data by registering it as procedural knowledge through WPD (Word Position Data). The WPD represents each position of words in a task corpus and is automatically composed of the corpus.

We conducted a series of experiments to evaluate our model and obtained a good result.

key words language model, speech recognition, n-gram, task adaptaion, Broadcast news, perplexity

1 はじめに

聴覚障害者などからテレビ番組，特に「ニュース」の音声を手書き化してほしいという要望が強い。米国では人手でリアルタイムに入力しているが，日本語の場合，漢字に変換する作業等が必要なためリアルタイムで入力するのは熟練した人でも難しい。そこで，このような作業に応用することを目的に音声認識の研究を行っている[1]。

現在，音声認識の言語モデルにはその実現の容易さから n -gram モデルが使われる。実際には記憶量の制約から $n=2,3$ (2,3-gram) が使われることが多い。さらに，認識精度を向上させるためにタスク適応が行われる[2]~[8]。タスク適応は通常，一般の言語モデルをさまざまなタスクから集めた（大量の）コーパスから構築し，あるタスク内で集められた（少量であることが多い）コーパスを用いて確率を補正することにより行われる。実際，放送ニュースの音声認識に関しても，直近の記者原稿を利用することによる効果を確認している[9]。

しかし，一方で，タスク適応は扱うタスクが非常に狭く限定されるような場合には有効であるが，扱うタスクがあまり限定できない場合にはむしろマイナスとなる。

最も望ましい言語モデルは，そのタスク内の言語表現である場合にはタスク用の言語モデルを使い，タスク内に収まらないような言語表現の場合には一般のモデルを使うことである。さらにこのようなことが自動的に行われる必要がある。

本稿では一般の言語モデルを使用しつつ，タスク適応も自動的に行う言語モデルについて述べる。提案する言語モデルでは，一般の言語モデルとし

て2,3-gramを使い，タスク適応として ($n \geq 4$)-gram を併用している。ここで，($n \geq 4$)-gram を利用するにあたって問題となるのは記憶量である。これを従来の2,3-gramのように宣言的知識として記憶していたのでは記憶量が爆発してしまう。提案する言語モデルでは ($n \geq 4$)-gram を手続き的知識として記憶することによって実現可能な記憶量となる。その際に中心的役割をするのが単語出現位置辞書である。単語出現位置辞書はタスク適応するコーパスから自動的に生成される。本言語モデルを perplexity による評価実験を行った結果についても述べる。

2 アナウンス原稿と記者原稿

ニュース番組でアナウンサーが読む原稿（アナウンス原稿）は，一部に割り込みで入ってくるものがあるものの，その多くは記者がある出来事について書いた原稿（記者原稿）を，本番直前まで手書きで校正することによって作成される（図1）。また，このような校正はあまり行われなくても少なくない。

ニュース音声認識の対象としているのはアナウンス原稿であるので，これから言語モデルを構築することが望ましい。しかし，記者原稿はワープロで書かれているので電子化されているが，アナウンス原稿は電子化されていない。そこで，言語モデルの構築には記者原稿を使用している。アナウンス原稿は記者原稿から作成されているので，記者原稿と一致している箇所も多く，極端な場合にはまったく一致する文もある。このような特徴を言語モデルに利用すれば，強い言語制約となることが期待できる。

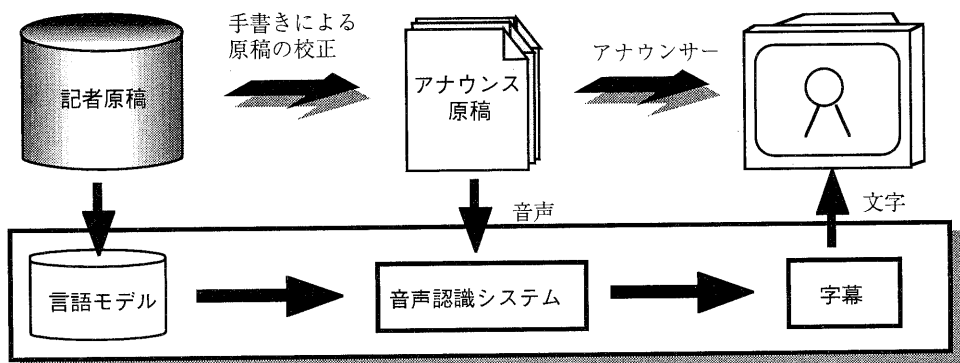


図1 記者原稿と言語モデルの関係

3 (n ≥ 4)-gram を併用する言語モデル

提案する言語モデルでは、2,3-gram とともに (n ≥ 4)-gram も用いる。2,3-gram による言語モデルの構築の際には、一般性を持たせるために過去数年分のコーパスと、ある程度のニュースの時事性(弱い適応化)を持たせるために、重みづけした(m倍した)直近のコーパスを用いている。(n ≥ 4)-gram による言語モデルでは、直近のコーパスから自動的に生成された単語出現位置辞書を利用している。

3.1 単語出現位置辞書

単語出現位置辞書とは記者原稿に出現した単語の出現位置(単語出現位置)を集めた辞書である。ここで、単語出現位置は記事の分野、その分野における記事の出現位置、その記事における単語の出現位置の3つ情報によって次のように7桁の数で定義している。

【単語出現位置の定義】

単語出現位置 XYYYYZZZ (7桁)

X (上位1桁目): 分野番号

(1: 政治, 2: 経済, 3: 社会, 4: 国際,
5: スポーツ, 6: ネット, 7: 首都圏)

YYY (2~4桁目): その分野における記事位置

ZZZ (5~7桁目): その記事における単語
出現位置

例えば、単語出現位置 7002026 は、分野が首都圏(7)、記事の出現位置が2番目(002)、単語の出現位置26番目(026)ということを表す。

単語出現位置辞書は直近(例えば前日、当日)の記者原稿から自動的に生成される(図3)。記者原稿の各単語に単語出現位置を付与し、単語ごとに集計し、単語を見出しとしてソートする。

3.2 言語モデルでの (n ≥ 4)-gram の利用

言語モデルで、単語出現位置辞書を使った (n ≥ 4)-gram の利用法について説明する。

まず、それまでの各単語で、単語出現位置辞書を引くことによりその単語単語出現位置を求める。今、i番目の単語を w^i 、それまでに連続している n-gram の長さの最大値を $N-1$ (すなわち N-gram を利用するため)、単語出現位置を c とおく。次単語 w^{i+1} ($k=1, 2, \dots, k_0$) の生起確率を p'_k を、2,3-gram による確率 p_k に基づいて式(1a), (1b)で動的に再計算することによって行う。

式(1a)は (n ≥ 4)-gram の利用する場合である。前単語の単語出現位置 C から予測される単語 w^{i+1} ($k=1$ のときとする)では、その新たな確率 p'_1 (確率値)を $N \geq 4$ のときに N に応じて計算する。ここで、定数 λ や q_0 はコーパスから求めることも可能である。式(1b)は 2,3-gram を利用する場合である。2,3-gram から予測される候補 w^{i+1}_k (それを $2 \leq k \leq k_0$ のときとする)の場合であり、2,3-gram の確率 p_k に応じて残りの $1-p'_1$ を分配する。もちろん、式(1b)の単語 w^{i+1}_k ($2 \leq k \leq k_0$) の中に w^{i+1}_1 が含まれているときはそれを除外する。

今回は、(n ≥ 4)-gram を利用するために $N \geq 4$ としているが、 $N \geq N_0$ のときのみ上記を実行するようにしておけば、(n ≥ N_0)-gram を利用する言語モデルともなる。

4 評価実験

本稿で提案する言語モデルに対して perplexity による評価実験を行った。言語モデルの構築の際には過去3年分の記者原稿と、前日と当日の直近2日間の記者原稿を用いて、以下の2つのモデルを作成し比較した。

ベースラインモデル

= 3年分の記者原稿

←一般言語モデル

+ m × 直近2日の記者原稿

←タスク適応

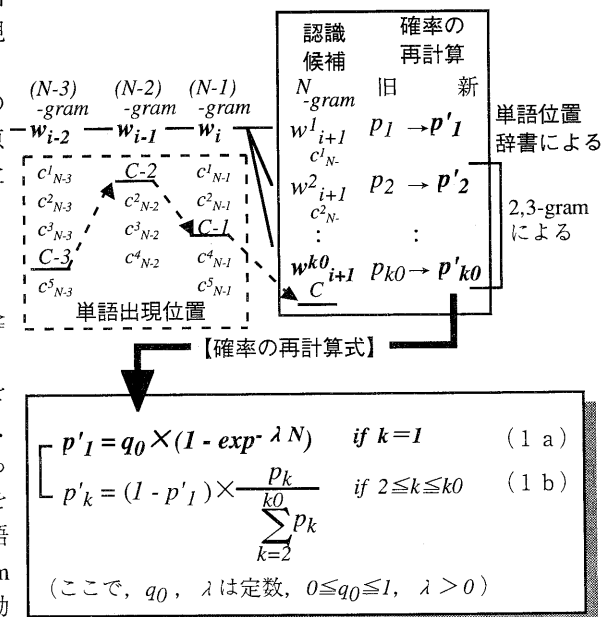


図2 単語出現位置辞書による確率の再計算

【記者原稿】

タイトル：公職選挙法違反
 作成部：首都圏
 日付：1998年09月02日

.....

・行われた町長選挙で初当選した時に加藤被告は.....
 36 37 38 39 40 41 42 43 44 45

=====

タイトル：犯人逮捕
 作成部：首都圏
 日付：1998年09月02日

.....

・埼玉県に住んでいた時に加藤さんと.....
 21 22 23 24 25 26 27 28

=====

タイトル：恐喝事件
 作成部：首都圏
 日付：1998年09月02日

.....

・警察署によりますと逮捕時に加藤容疑者は.....
 17 18 19 20 21 22 23 24

.....

記事 1

記事 2

記事 3

単語出現位置
の計算



行われた 7001036
 町長選挙 7001037
 で 7001038
 初当選 7001039
 した 7001040
 時 7001041
 に 7001042
 加藤 7001043
 被告 7001044
 は 7001045

.....

埼玉県 7002021
 に 7002022
 住んでいた 7002023
 時 7002024
 に 7002025
 加藤 7002026
 さん 7002027
 と 7002028

.....

警察署 7003017
 によりますと 7003018
 逮捕 7003019
 時 7003020
 に 7003021
 加藤 7003022
 容疑者 7003023
 は 7003024

単語ごとに集計し、
ソート



【単語出現位置辞書】

技術 1000041, 6002022, ...
 :
 放送 6002021, 7002010, ...
 :
 :

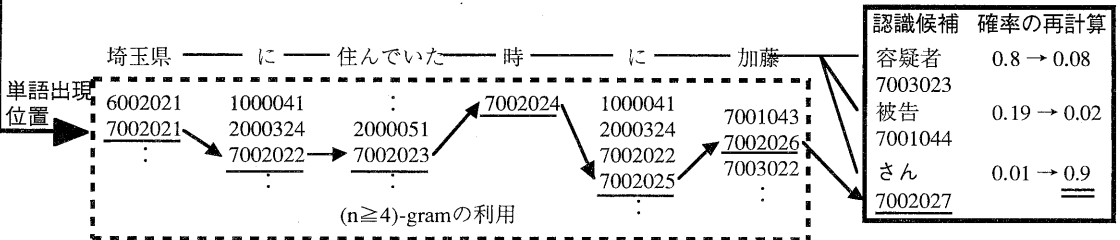
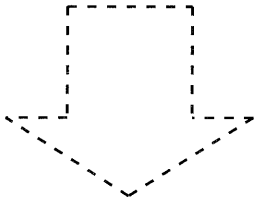


図3 単語出現位置辞書の自動生成と音声認識への応用

提案する言語モデル

- = 3年分の記者原稿 ←一般モデル
- + $m \times$ 直近2日の記者原稿 ←弱いタスク適応
- + $(n \geq 4)$ -gram ←強いタスク適応

また、パラメータである λ や q_0 は今回は $\lambda = 2, 3$, $q_0 = 0.9$ としている。

評価データには1998年9月30日のニュースのアナウンス原稿^{*}を、以下の3つの場合に分けて使用した。

評価データA：元となる記者原稿とほぼ同じ
(36記事)

評価データB：元となる記者原稿とある程度
同じ(20記事)

評価データC：元となる記者原稿なし
(14記事)

ただし、これらの評価データの分け方は、 $m=100$ におけるベースラインモデルの perplexity の値 $perplexity_0$ が

評価データA： $perplexity_0 < 20$ のとき

評価データB： $20 \leq perplexity_0 < 100$ のとき

評価データC： $100 < perplexity_0$ のとき

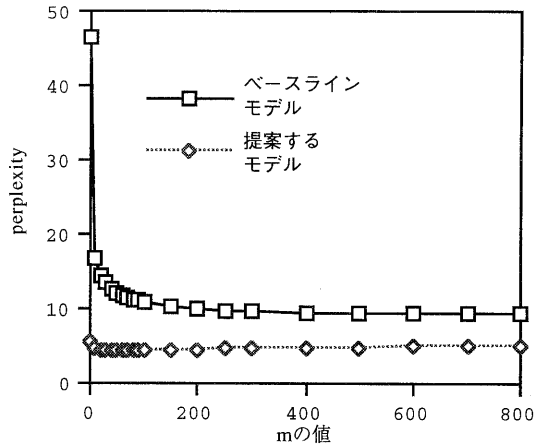
とした。この境界値は実際の perplexity の結果をいくつか見て決めた。

さまざまな m で実験した。結果を図4に示す。図4を見ると、いずれの場合も提案する言語モデルがベースラインモデルより優れていることがわかる。評価データAは現在我々が音声認識の対象として考えているものであるが、ベースラインモデルでは $m=600$ のとき、提案する言語モデルでは $m=50$ のときに最小となっている。すなわち、提案するモデルでは直近の2日の記者原稿を重みづけが小さく、 $(n \geq 4)$ -gram が有効に働いているのが

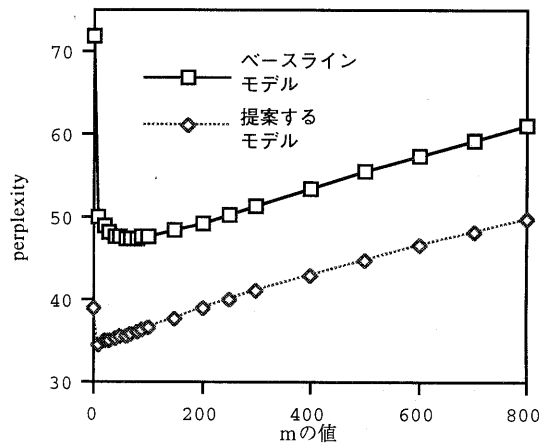
表1 評価データAでの perplexity が最小となる場合

評価データ	評価データA	評価データB	評価データC
ベースラインモデル (α)	9.6	57.4	383.5
提案するモデル (β)	4.5	35.4	249.0
削減率 (($\alpha - \beta$)/ α)	53.1%	38.3%	35.1%

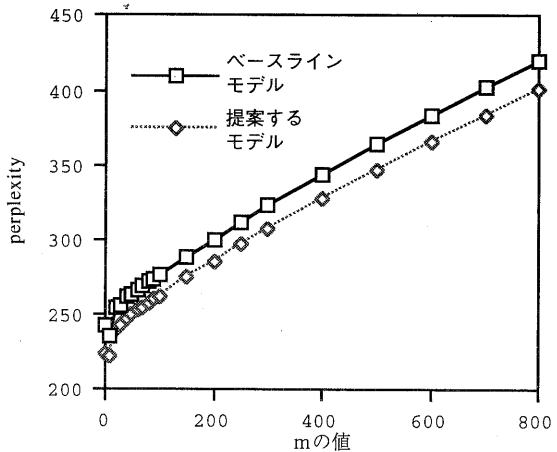
^{*}1 正確にはアナウンス原稿ではなく、アナウンサーが話したものを書き起こした原稿である。したがって、不要語等が含まれる。



(a) 評価データAに関する実験結果



(b) 評価データBに関する実験結果



(c) 評価データCに関する実験結果

図4 それぞれの評価データにおける perplexity

わかる。このときの評価データ B, C の場合の perplexity も合わせて表 1 に示す。表 1 を見ると、提案する言語モデルでは評価データ B, C でも削減率が大きいことがわかる。これはベースラインモデルが評価データ A に関してタスク適応しているために、特に評価データ C での perplexity が大きい。一般モデルはベースラインモデルの $m = 1$ の場合であり、perplexity = 242.5 である。これに対して提案するモデルでも perplexity = 249.0 であるので一般モデルとほとんど変わらなく、一般言語モデルとしての特性も失っていないことがわかる。

5 おわりに

$n \geq 4$ である n -gram を利用した言語モデルについて述べ、その有効性を評価実験により確認した。今後は実際に計算機上にインプリメントし、認識実験によって本言語モデルの検証を行う必要がある。

今回は単語出現位置が連続している場合のみを用いたが以下のような使い方や、音声認識の改良も考えられる単語出現位置辞書はまた次のような応用も考えられる。

- ・単語出現位置の連続していても値が近い単語候補を優先することにより構文構造の変化に対応。
- ・上位 1 桁目を使うことにより単語の分野制約。
- ・単語出現位置辞書により記者原稿での出現位置に句読点を挿入。
- ・実験結果のように perplexity がここまで小さくなるので、デコーダー上での第 2 パスの削除や最良優先探索の導入による音声認識の高速化。

これらの詳細について機会をあらためて報告したい。

謝辞

本研究を進めるにあたって適切な助言をいただいた当研究所音声認識・言語情報処理グループに感謝する。

【参考文献】

- [1] 今井 亨, 小林彰夫, 尾上和穂, 安藤彰男, “ニュース番組自動字幕化のための音声認識システム,” 情報処理学会研究会, HI80-11(SLP-23-11), pp.59-64, 1998.
- [2] R. Kneser and V. Steinbiss, “On the dynamic adaptation of stochastic language models,” Proc. IVASSP’93, Vol.2, pp.585-588, 1993.
- [3] P. Clarkson and A. Robinson, “Language model adaptation using mixtures and exponential decaying cache,” Proc. ICASSP’97, Vol.2, pp.799-802, 1997.
- [4] S. Matsunaga, T. Yamada, and K. Shikano, “Task adaptation in stochastic language models for continuous speech recognition,” Proc. ICASSP’92, Vol.1, pp.165-168, 1992.
- [5] 政瀧浩和, 匂坂芳典, 久木和也, 河原達也, “最大事後確率推定による N-gram 言語モデルのタスク適応,” 信学論 (D-II), Vol.J81-D-II, No.11, pp.2519-2525, 1998.
- [6] 伊藤彰則, 好田正紀, “対話音声認識のための事前タスク適応の検討,” 電子情報通信学会研究会報告, SP96-81, pp.91-98, 1996.
- [7] 赤松裕隆, 甲斐允彦, 中川聖一, “新聞・ニュース文の大語彙連続音声認識,” 情報処理学会研究会報告, SLP98-21-11, pp.97-104, 1998.
- [8] 政瀧浩和, 松永昭一, 匂坂芳典, “品詞および可変長単語列の複合 N-gram の自動生成,” 信学論 (D-II), Vol.J81-D-II, No.9, pp.1929-1936, 1998.
- [9] 小林彰夫, 今井 亨, 安藤彰男, 中林克己, “ニュース音声認識のための時期依存言語モデル,” 情報処理学会論文誌, Vol.40, No.4, pp.1421-1429, 1998.