

対話システムにおける音声認識の改善を目指した

バックオフ係数の検討

浦谷 則好, 小早川 健
NHK放送技術研究所

1. はじめに

音声認識技術は利用者が特殊な装置を利用する必要がないなどの手軽さから、いろいろな状況で利用されるようになってきている。例えば、NHKでは誰でも簡単にデジタル放送を楽しめるようにするためのTVエージェントとの対話に音声認識を利用している[1]。また、ニュース番組に対するリアルタイムの字幕付与のためにも音声認識を利用している[2]。音声認識の言語モデルとしてはNグラムモデルが使われ、データスパースネスに対処するために確率をスムージングする処理がなされる。このスムージング手法の中で主に利用されているのがKatz smoothing[3]という手法である。この手法では学習データ中に出現しない未知単語列に配分される頻度は、学習データ中に1回だけ出現した単語列の頻度と等しいとみなされ、そこから得られる重み(バックオフ係数)を使って未知単語列の出現確率が計算される。未知語列に対する配慮としては妥当に思われるが、その値の根拠は希薄である。システムを音声で制御する場合、システム側が受容できる発話の種類はかなり制限されていると思われるので、Katz smoothingで得られるバックオフ係数は過剰なものである可能性が高い[4]。また、原稿が存在するニュースや講演で、発話者が原稿にそれほど逸脱しない発話をするのが予想される場合にも、この仮定で得られるバックオフ係数は大きすぎる恐れがある。また、そこでバックオフ係数の最適化を目指して実験を行い、原稿に近い発話が予想される場合にはバックオフ係数を1/100程度にした方がよいことを確認したので、それについて報告する。

2. 実験

音声は2001年6月1日から14日までにNHKのニュースで放送されたアナウンサーおよび記者の音声である。音響モデルは不特定話者用のモデル(男声モデル1種、女声モデル1種の自動切替)である。雑音が多い、発話が不明瞭などの理由で音声認識率が極端に悪いもの(単語認識率で80%以下)や、あいさつなどの短い会話は大部分除外した。(あいさつなどは一部、テスト用に残した。)発話の数などの諸元は表1の通りである。1発話(1文章)当たり平均25単語である。

表1 対象としたニュース

ニュース項目	発話数	単語数
92	853	21, 225

2. 1. 実験1

上述のニュース音声の書き起こしから、122語の挿入(「え」とか「あの一」というフィラー)と言い直しを認定した。この122語を除いたものを推定原稿として音声認識実験を行った。また、書き起こしそのものを原稿とみなした実験も行った。バックオフ係数を1/10, 1/32, 1/100, 1/316, 1/1000

(10の-1.0, -1.5, -2.0, -2.5, -3.0乗)と変えて音声認識率の変化を求めた。なお、バックオフ係数以外は変化させていないので、累積接続確率は1.0となっていない。原稿のperplexityは書き起こしそのもの、推定原稿それぞれ7.40, 7.43であった。実験結果を表2に示す。表中で左の数字がCorrectnes(正解数/単語数)で、右の数字はAccuracy((正解数-挿入単語数)/単語数)である。表2を見ると完全にクローズドな場合(書き起こし原稿をそのまま使う場合)でもバックオフ係数を小さくするほどよいというわけではない

ことが分かる。バックオフ係数を従来の1/100程度にするのがよいと推定される。

表2 実験1の音声認識率 (%)

	推定原稿		書き起こし	
	Corr.	Acc.	Corr.	Acc.
元(1)	97.23	96.89	97.64	97.27
1/10	98.20	98.01	98.70	98.54
1/32	98.30	98.13	98.83	98.68
1/100	98.32	98.12	98.88	98.72
1/316	98.12	97.84	98.80	98.61
1/1000	97.94	97.56	98.74	98.51

2. 2. 実験2

2. 1. の実験では原稿と実際の発話の相違はわずかに0.58%しかない。「原稿に近い発話が予想される場合」という条件だとしてもその差が小さすぎる。そこで、ロバストネスを確認するために、書き起こし原稿にランダムに単語の5%を削除、置換、挿入したものを原稿として同様の実験を行った。挿入、置換に使われる単語は元の原稿に含まれていないものから選んだ。結果を表3に示す。「混合」は1.67%ずつの削除、置換、挿入を与えたものである。また、表で「排除」としているものは、元の原稿の単語から特定の単語(高頻度の機能語的なものと頻度2以下のものを除いた中頻度語から選択)をランダムに除いた実験である。それゆえ、「排除」の場合の最大認識率は95%である。なお、実験は各モード4回ずつ実施した。表中の数値は4回の平均のAccuracyとTest data perplexityである。未知語率は当然、挿入が0、排除が5%であり、削除、置換、混合はそれぞれ、

表3 実験2の音声認識率 (%)

	元(1)		1/100	
	Acc.	Per.	Acc.	Per.
削除	94.36	11.35	95.37	16.81
置換	94.27	12.13	95.76	18.09
挿入	95.81	10.69	97.45	14.12
混合	94.76	11.45	96.08	16.52
排除	88.34	9.76	89.63	9.94

0.38%, 0.37%, 0.29%であった。なお、Training data perplexityの平均は削除、置換、挿入、混合、排除、それぞれ、8.76, 9.33, 9.11, 9.03, 7.88であった。

3. むすび

Katz smoothingによるバックオフ係数の算出手法に疑問を感じて、バックオフ係数の値を変化させて音声認識実験を行った。その結果、実際の発声が原稿とあまり違わないかぎりバックオフ係数を1/100にした方が良い結果を得られることを確認した。こういうケースは対話による機器操作ばかりでなく、ニュースや講演などでも考えうるものであり、有用な知見が得られたものと考えられる。なお、まだ実験途中であるが原稿を10%変化させた場合も、バックオフ係数を1/100にした方が元のものより良い結果が得られている。

謝辞

本研究は当研究所の今井享主研をはじめとする音声認識研究チームから音響モデルと音声データの提供を受けて実施したものである。多大な協力に深謝する。

参考文献

- [1]Jun Goto et al.:A Spoken Dialogue Interface for TV Operations Basedon Data Collected by Using WOZ Method, IEICE Trans. on Information and Systems, E87-D(6):1397-1404, 2004
- [2]安藤彰男外: 音声認識を利用した放送用ニュース字幕制作システム, 電子情報通信学会論文誌, J84-DII(6):877-887, 2001
- [3]Katz, S.M.:Estimation of probabilities from sparse data for the language model component of a speech recognizer, IEEE Trans. on Acoustics, Speech, and Signal Processing, ASSP-35(3):400-401, 1987
- [4]Kneser, R. et al.:Improved Backing-off for M-gram Language Modeling, IEEE International Conf. on Acoustics, Speech, and Signal Processing, 1:181-184, 1995