

接続の方向性を考慮した 多重クラス複合N-gram言語モデル

山本博史 匂坂芳典

(株) ATR 音声翻訳通信研究所
〒619-02 京都府相楽郡精華町光台 2-2
E-mail: yama@itl.atr.co.jp

あらまし クラス2-gramにおける効率的なクラス分類を実際のコーパスから統計的に行うための手法を提案する。本手法では直前および、直後の単語への接続性を別の属性としてとらえ、各単語に対してその属性ごとに複数のクラスを割り当てる。これらのクラスは前後に接続している単語の分布に基づいて各々独立に作成されることによって、効率的でかつ信頼性の高いクラス分類となっている。さらにこの多重クラス2-gramを結合単語との多重複合N-gramに拡張することにより、千分の一以下の論理パラメーターサイズでパープレキシティ、単語認識率とも単語N-gramを上回る性能を示した。

キーワード クラスN-gram 可変長N-gram 自動クラス分類 連鎖語

MULTI CLASS COMPOSITE N-GRAM LANGUAGE MODEL BASED ON CONNECTION DIRECTION

Hirofumi Yamamoto, Yoshinori Sagisaka

ATR Interpreting Telecommunications Res.Labs.
2-2 Hikari-dai, Seika-cho, Soraku-gun, Kyoto 619-02 Japan
yama@itl.atr.co.jp

Abstract A new word-clustering technique is proposed to efficiently build statistically salient class 2-grams from language corpora. By splitting word neighboring characteristics into word-preceding and following directions, multiple (two-dimensional) word classes are assigned to each word. In each side, word classes are merged into larger clusters independently according to preceding or following word distributions. This word-clustering can provide more efficient and statistically reliable word clusters. Further, we extend it to Multi-Class Composite N-gram that unit is Multi-Class 2-gram and joined word. Multi-Class Composite N-gram showed better performance both in perplexity and recognition rates with one thousandth smaller logical parameter size than conventional word 2-grams.

key words Class N-gram, Variable Order N-gram, Automatic Clustering, Joined Word

1. はじめに

連続音声認識のための統計的言語モデルとして単語N-gramモデルが広く用いられている。単語N-gramモデルはルールベースの文法的な制約に比べて効率的かつ柔軟な場合が多い。しかしながらコンパクトな認識システムを構築する場合、必要とするメモリの多さは重大な問題点である。一方統計的な観点からは観測データの少ない(あるいは存在しない)単語ペアに対しては統計的に正確な値を予測することが難しいという問題がある。これらの問題の解決法としてはクラスN-gramあるいは可変長N-gramによる手法が提案されている。[1]-[4]

これらの手法では単語に対するクラスの定義方法はモデルの精度、サイズに直結する重要な問題点となっている。実際総単語数459, 383語、異なり単語数7, 221の訓練セットにおけるテストセットパープレキシティが単語2-gramで18.51であったのに対し、品詞情報に基づく158種類のクラスを用いたクラス2-gramはサイズを単語2-gramの0.5%に縮めることができたものの、パープレキシティは31.53にとどまった。これらのことから精度の高いクラスN-gramの構築のためには、クラスの規定方法に工夫が必要がある。

本稿では、単語に対するクラスの規定方法として単語の前方向と後方向の接続性を別の属性とみなし、各単語にその属性ごとに複数のクラスを割り当てる多重クラスを提案する。さらに、頻出する単語列を連鎖語として取り扱う多重クラス複合N-gramを提案する。実験によりこの多重クラス複合N-gramは、単語2-gramの千分の一以下の論理パラメータサイズでより高い認識性能を示すことが判明した。

2. 多重クラス

2.1 クラス2-gram

クラスN-gramにおいて直前のN-1個の単語列に続いて次単語が出現する確率は一般に次の式で与えられる。

$$P(W_N | C_N) \times P(C_N | C_{N-1}, C_{N-2}, \dots, C_1) \quad (1)$$

(ここで C_N は単語 W_N が属するクラスを表わす)

さらに $N=2$ としたクラス2-gramでは

$$P(W_2 | C_2) \times P(C_2 | C_1) \quad (2)$$

となり、これはクラス2-gramにおいて必要なパラメータ数がクラス数の2乗に比例することを表わして

おり、クラス数の増加によって急激にパラメータ数が増加することになる。

2.2 クラス2-gramにおけるクラスの問題点

ここでクラス2-gramにおけるクラスが単語間の接続性を表わすためにどのような働きをしているかをみてみることにする。一例としてある動詞を考える。動詞は活用形ごとに単語の接続性が異なるため活用形の6種類ごとに別々のクラスが割り当てられるのが普通である。しかしながら活用形が接続性に与える影響は後続の単語についてのみであり、先行の単語に対する接続性には影響を与えないと考えられる。このことは後続するクラスとして動詞が現われるかどうかを推定する場合には活用形を考慮したクラス分けが不要であることを示している。

一方動詞を自動詞と他動詞に分けた場合、これらの接続性の差は主として先行する格助詞の種類などに対して現われる。このため、先行するクラスとして動詞が現われる場合には不要なクラス分け(この場合は逆に活用形のみが重要であるため)と考えられる。

2.3 多重クラスN-gram

上記の動詞の例が示すように、先行するクラスと後続するクラスでは接続性を表わすために有効なクラスが必ずしも一致しない。このため先行、後続で用いられる場合で用いるクラスを別にする方が効率的と考えられる。この考えを(2)式に対してあてはめると次のようになる。

$$P(W_2 | C_2^f) \times P(C_2^f | C_1^f) \quad (3)$$

ここで C^f は後続する場合のクラス(以下to-クラスと呼ぶ)を表わし、 C^f は先行する場合のクラス(以下from-クラスと呼ぶ)を表わす。従来のクラスN-gramにおいては自動詞、他動詞およびその活用形を表現するためには2(自動詞、他動詞の)×6(活用形の)=12通りのクラスが必要であり、パラメータ数はクラス数の自乗であるため 12×12 となる。これに対し本提案では2通りのto-クラスと6通りのfrom-クラスで表現可能で、かつパラメータ数もto-クラス数×from-クラス数の12と、きわめて効率的なクラス分類が可能となる。(3)式は $N > 2$ の場合にも拡張可能で、クラスN-gram同様直前のN-1個の単語列に続いて次単語が出現する確率は次の式で与えられる。

$$P(W_N | C_N^f) \times P(C_N^f | C_{N-1}^f, C_{N-2}^{f-1}, \dots, C_1^{f-N+2}) \quad (4)$$

このように各単語に対して接続の方向性ごとに与えられる複数のクラスを多重クラスと呼び、この多重クラスを用いたクラスN-gramを多重クラスN-gramと呼ぶことにする。

3. クラスの自動分類

3.1 自動分類の目的

クラスN-gramにおけるクラス分類の指標としては、しばしば品詞情報が用いられる。品詞情報はコーパスに現われない単語に対しても割り当てることができるという利点がある。しかしながらN-gramにとって最も重要な単語間の統計的な接続特性を詳細に表しているとは必ずしも言い難い。また品詞情報自体は先見知識として与えられるものであるため、品詞情報自体の分類のしかたがクラスN-gramにおけるクラス分類の効率に直接影響を与えてしまうという問題点もある。従って統計的モデルの観点からは、実際のコーパスから単語間の接続性のみに着目してクラス分類を行うことが好ましい。

3.2 自動分類の方法

実際のコーパスからクラス分類を自動的に行う方法としては、[1]をはじめとするいくつかの提案がなされている。本稿においては以下に示すような手順で自動クラス分類を行った。分類の対象としたクラスは従来のクラスN-gramにおけるクラス、to-クラス、from-クラスの3つである。

1. 一単語一クラスとする。
2. 個々の単語またはクラス X に対して次のようなベクトルを考える。

$$V_i(X) = \{P_i(W_1|X), P_i(W_2|X), \dots, P_i(W_n|X)\} \quad (5)$$

$$V_f(X) = \{P_f(W_1|X), P_f(W_2|X), \dots, P_f(W_n|X)\} \quad (6)$$

ここで $P_i(W_i|X)$ 、 $P_f(W_i|X)$ は単語またはクラス X から単語 W への後向きおよび前向きの2-gramの確率値を表す。従来のクラスN-gramにおけるクラスでは前後の接続を同時に考慮するため、接続属性を表すベクトルとしては次のものを用いる。

$$V(X) = \{V_i(X), V_f(X)\} \quad (7)$$

一方to-クラス、from-クラスにおいてはそれぞれ前向き、後向きの接続のみを考慮するため次のベクトルを用いることになる。

$$V(X) = V_i(X) \quad (8)$$

$$V(X) = V_f(X) \quad (9)$$

3. マージコスト $U_{new} - U_{old}$ が最小となるようなクラスのペアを選び、統合して一つのクラスとする。ここで

$$U_{new} = \sum_w P(W) \times D(V(C_{new}(W)), V(W)) \quad (10)$$

$$U_{old} = \sum_w P(W) \times D(V(C_{old}(W)), V(W)) \quad (11)$$

であり、 C_{new} は統合後のクラス、 C_{old} は統合前のクラスを表し、 $D(V_C, V_W)$ はベクトル V_C と V_W のユークリッド距離の自乗を表すものとする。

4. 2までの手順をあらかじめ定められたクラス数になるまで繰り返す。

4. 多重クラス2-gramの性能評価

4.1 パープレキシティによる評価

前節で示した手順で求めたクラスを用いて多重クラス2-gramのパープレキシティによる評価を行った。比較対象としては、同じ10前節で求めたクラスを用いた通常のクラス2-gram、および品詞情報に基づく158のクラスを用いたクラス2-gramを用いた。訓練セットは総単語数459, 383語、異なり単語数7, 221であり、評価は訓練セットに含まれない41会話6, 828語の評価セットで行った。なお、語彙に関しては評価セットに現われる総ての単語が含まれている。また多重クラス2-gramに関してはto-クラスとfrom-クラスの数は必ずしも一致している必要はないが、比較のために同一クラス数とした。図1. に示すように品詞情報に基づくクラス2-gramと自動クラス分類によるクラス2-gramでは同一クラス数でも

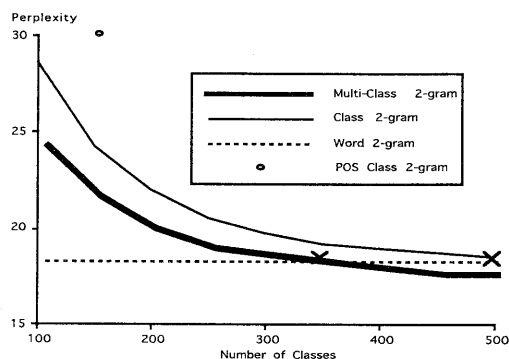


図1. 多重クラス2-gramの性能

性能に大きな差がでている。このことからクラス分類の方法としては、品詞情報に基づくよりも実際のコーパス上での接続性に基づく分類の方が効率的であることがわかる。またクラス 2-gram と多重クラス 2-gram においては、単語 2-gram と同等のパープレキシティを示す時のクラス数がクラス 2-gram で約 500、多重クラス 2-gram で約 350 (図中で×で示されている点) とクラス数で約 30%、論理パラメータ数で約半分に削減できていることがわかる。

4.2 連続単語認識による評価

パープレキシティは言語モデルの性能評価のための良い指標であるが、必ずしも単語認識率に直決しない。そこで、連続単語認識における性能評価もあわせておこなった。実験条件は次に示す通りである。

評価セット

パープレキシティの評価と同じ 41 会話
特徴量

標本化周波数 16k

フレーム周期 10msec

12 次のメルケプストラム、パワーおよびそれらの
一次回帰係数の計 26 次元

音響モデル

ML-SSS による音素環境依存 800 状態 5 混合
の HMnet [5]

男女依存モデルの自動選択

デコーディング [6]

1 バス時間同期ビタビサーチ

2 バス言語重みを変更した上でのフルサーチ

評価の対象としては単語 2-gram、クラス数 500 のクラス 2-gram、そしてクラス数 350 の多重クラス 2-gram に対しておこなった。なお、これらのモデルは前段で述べたように、いずれもほぼ同じパープレキシティを示している。評価の基準としては次に示す単語認識率で行った。

W-D-I-S

W

(W: 正解単語数、D: 脱落誤り数、

I: 挿入誤り数、S: 置換誤り数)

単語認識率は単語 2-gram で 69.05%、クラス 2-gram で 69.78%、多重クラス 2-gram で 70.29% であり、多重クラス 2-gram が最も少ないパラメータ数であるにもかかわらず最も高い性

能を示している。

5. 多重クラス複合 N-gram

5.1 単語 N-gram の導入

クラス 2-gram は訓練データが少ない場合でも少ないパラメータで頑健なモデルを構築することができるが、精度の面では単語 N-gram に劣る。これに対し、単語 N-gram では信頼性のあるモデルの構築のためには多量の訓練データを必要とする。しかしながら少量の訓練データで単語 N-gram を構築した場合でも出現回数の多いワードペアあるいはトリオ以上に関しては局所的には単語 N-gram の値に対し信頼性が確保できる。従って出現回数の少ないワードペアに対してはクラス 2-gram を用い、出現回数が多く、信頼性が確保できるワードペアあるいはトリオ以上に対しては単語 N-gram を用いるという方法が考えられる。この方法では単語ペア A, B の出現回数が十分である場合、単語 X の直後に単語列 A, B, C が出現する確率は次の式で与えられる。

$$P(A|C(A)) \times P(C(A)|C(X)) \times P(B|A) \times P(C|C(C)) \times P(C(C)|C(B)) \quad (12)$$

さらに多重クラス 2-gram を用いた場合には (12) 式は次のようになる。

$$P(A|C_i(A)) \times P(C_i(A)|C_f(X)) \times P(B|A) \times P(C|C_i(C)) \times P(C_i(C)|C_f(B)) \quad (13)$$

続いて単語ペア A, B を連鎖語として一つの単語 $A+B$ とするならば (13) 式は次のようになる。

$$P(A+B|C_i(A+B)) \times P(C_i(A+B)|C_f(X)) \times P(C|C_i(C)) \times P(C_i(C)|C_f(A+B)) \quad (14)$$

ここで $C_i(A+B)$ は連鎖語 $A+B$ の前にどのような単語が来やすいかを表すクラスであるため、連鎖語の先頭の単語 A のそれに等しいと考えられ

$$C_i(A+B) = C_i(A) \quad (15)$$

とおくことができる。同様に連鎖語の後方向の接続性は最後尾の単語に等しいと考えられ

$$C_f(A+B) = C_f(B) \quad (16)$$

とおくことができる。(15)、(16) 式を用いれば (14) 式は次のようになる。

$$P(A+B|C_i(A)) \times P(C_i(A)|C_f(X)) \times P(C|C_i(C)) \times P(C_i(C)|C_f(B)) \quad (17)$$

(17)式は単語 2-gram が連鎖語の導入により多重クラス 2-gram の形式を保ったまま、かつ新たなクラスを導入することなしに表現可能であることを示している。またこの際に増加したパラメーターは連鎖語 $A+B$ のユニグラムのみである。このことは単語 N-gram の導入に関しても同様に言え、この場合は N 連鎖語で単語 N-gram が表現されることになる。増加するパラメーターに関しても 2-gram の時と同様に N 連鎖語のユニグラムのみである。N 連鎖語により部分的に単語 N-gram を導入した多重クラス 2-gram を多重クラス複合 N-gram と呼ぶことにする。

5.2 多重クラス複合 N-gram の構築手順

多重クラス複合 N-gram は次の手順で構築することができる。

1. 初期状態として多重クラス 2-gram を与える。
2. 単語ペアのうち出現回数が一定値以上のものを連鎖語として辞書に加える。この連鎖語の to-クラスは先行単語の to-クラスと同じ、from-クラスは後続単語の from-クラスと同じとする。
3. 新たに加わった連鎖語も含めて 2. の手順を繰り返す。出現回数が一定値以上のものが存在しない場合は終了する。

5.3 品詞および可変長単語列の

複合 N-gram との比較

クラス 2-gram と単語 N-gram の短所を補い合うモデルとしては品詞および可変長単語列の複合 N-gram [2] がある。品詞および可変長単語列の複合 N-gram は品詞情報に基づくクラス 2-gram をベースとし、クラスの中から単語を分離し、独立したクラスにするという操作と分離された単語同士から連鎖語を生成し新たなクラスとするという操作をエントロピーの減少を基準にして繰り返すものである。

品詞および可変長単語列の複合 N-gram は良い性能を示すモデルであるが、次のような問題点があり、多重クラス複合 N-gram ではこの問題は解決されていると考えられる。

- ・クラスから分離された単語からでないと連鎖語を生成することができないため、すでに適切なクラス分類が行われている場合でもクラス分離を行う必要が

表 1. 品詞および可変長単語列の複合 N-gram との相違点

	品詞および可変長単語列の複合 N-gram	多重クラス複合 N-gram
分離の対象	クラスから単語	クラス 2-gram から単語 2-gram
分離の基準	エントロピーの減少	単語ペアの出現回数
N-gram の表現単位	クラス 2-gram クラス-単語 2-gram 単語-クラス 2-gram 単語 N-gram	クラス 2-gram 単語 N-gram
パラメーター数の増分	分離単語数 + 単語列数の自乗	連鎖語数

ある。

- ・クラスから分離された単語は通常十分大きなユニグラム出現数を持つが、バイグラムに対しては必ずしもそうとは言えないため新たにデータスペースの問題を引き起こすことがある。

多重クラス複合 N-gram と品詞および可変長単語列の複合 N-gram との違いをまとめると表-1 のようになる。このうち N-gram の表現単位に関しては品詞および可変長単語列の複合 N-gram の方が豊富であるが、そのうちどれが用いられるかは前後に来る単位がクラスか単語か単語列で決まってしまうため自由度という観点からは必ずしも高いとは言えない。また、分離の基準に関しては品詞および可変長単語列の複合 N-gram で出現回数を基準とすることも、多重クラス複合 N-gram でエントロピーを基準とすることも可能なため本質的な差ではない。

5.4 多重クラス複合 N-gram の性能評価

まず多重クラス複合 N-gram の性能評価としてパープレキシティによる評価を行った。多重クラス複合 N-gram の初期状態としては 4 節で用いた多重クラス 2-gram を用い、連鎖語導入のための出現回数は 20 回と設定した。4 節と同じ訓練セット、評価セットを用いた時の結果を図-2 に示す。また比較対象として同一条件における単語 2-gram、単語 3-gram、そして初期品詞クラス 158、分離クラス 1000 の品詞および可変長単語列の複合 N-gram の結果も同時に示す。

この結果からクラス数 400 においてほぼ単語 3-gram に近い性能を示していることがわかる。連鎖語導入のための出現回数が 20 回の条件で導入された連鎖語の総数は 2、

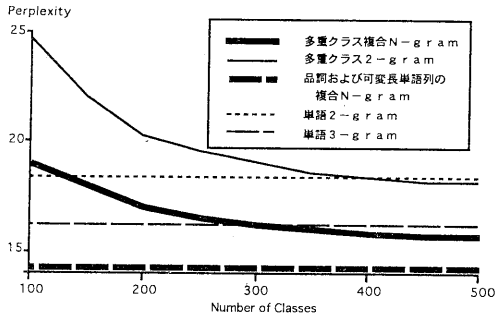


図2. 多重クラスN-gramのパープレキシティ

表2. 多重クラス複合N-gramの性能

モデルの種類	論理パラメーター数	2-gram以上のエントリ数	Perplexity	認識率
200 クラス	5.0×10^4	5.7×10^3	17.54	76.30
400 クラス	1.7×10^5	1.4×10^4	16.29	75.79
単語2-gram	5.5×10^7	6.1×10^4	18.52	68.95
単語3-gram	3.5×10^{11}	2.0×10^5	15.93	
品詞および可変長単語列の複合N-gram	1.4×10^6	5.2×10^4	14.84	75.51

212単語であり、これらの連鎖語の訓練セットにおける出現回数、すなわち単語N-gramの使用回数は116, 525で、これは全体の約20%を占める。またこのうち3単語以上の連鎖語の出現回数、すなわち単語3-gram以上の使用回数は60, 529回で、全体の約10%である。続いて同様の条件のもとでのクラス数200および400の多重クラス複合N-gramの単語認識率による評価結果を表2に示す。多重クラス複合N-gramクラス数200の場合においても単語2-gramに比べパープレキシティ、認識率ともに上回っており、その際に必要は論理パラメーター数はわずか千分の一以下となっている。またエントリサイズにおいても十分の一以下となっている。品詞および可変長単語列の複合N-gramとの比較においてもパープレキシティでは劣るものの4%以下の論理パラメーター数での同等以上の認識性能を示している。

6. まとめ

本稿ではクラスN-gramのための効率的クラス分類の方法として各単語に対して接続の方向性ごとに複数のクラスを割り当てる多重クラスを提案した。多重クラスを利用した多重クラス2-gramにおいては単語間の接続のみに着目したクラス分類を行っているため、パー

プレキシティ、単語認識率を劣化させることなく非常にコンパクトなサイズでモデルを表現することができる。

また多重クラス2-gramに対して部分的に単語N-gramを導入した多重クラス複合N-gramを提案した。多重クラス複合N-gramにおいては出現回数が多く、信頼性の高い単語N-gramが連鎖語の形で表現されており、クラス2-gramの頑健さと単語N-gramの精度を兼ね備えたモデルとなっている。また形式的にはあくまでも多重クラス2-gramの形をしているためデコーダにとっても非常に扱いやすいモデルとなっている。本モデルは一つの連鎖語(すなわち単語N-gram)の導入にあたって増加するパラメーターはそのユニグラムのみであるため、多重クラス2-gramと同等のモデルサイズでより高性能がえられる。実験結果においては単語2-gramに対し論理パラメーター数で千分の一以下のサイズでより高いパープレキシティ、単語認識率が得られ、品詞および可変長単語列の複合N-gramとの比較においても4%程度の論理パラメーター数で同等以上の単語認識率を得ることができた。

参考文献

- [1] Shuanghu Bai, Haizhou Li, Zhiwei Lin, Baosheng Yuan : "Building Class-based Language Models with Contextual Statistics, "Proc.ICASSP, vol. 1, pp. 173-176, 1998
- [2] 政瀧 浩和, 松永 昭一, 勾坂 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, vol.J81-D-II, no.9, pp.1929-1936, Sep. 1998.
- [3] P.F.Brown et al : "Class-Based n-gram Models of Natural Language, "Computational Linguistics, vol. 18, No. 4, pp. 467-479, 1992
- [4] Sabine Deligne, Frederic Bimbot : "Language modeling by variable length sequences, "Proc ICASSP, vol. 1, pp. 169-172, 1995
- [5] M.Ostendorf, H.Singer : "HMM topology design using maximum likelihood successive state splitting, " Computer Speech and Language, vol. 11, pp. 17-41, 1997
- [6] Tohru Shimizu, Hirofumi Yamamoto, Hirokazu Masataki, Shoichi Matsunaga, Yoshinori Sagisaka : "Spontaneous Dialogue Speech Recognition Using Cross-Word Context Constrained Word Graphs," Proc.ICASSP, vol. 1, pp. 145-148, 1996