

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2005-292832

(P2005-292832A)

(43) 公開日 平成17年10月20日(2005.10.20)

(51) Int.Cl.⁷

G 1 0 L 15/18

G 1 0 L 15/06

F I

G 1 0 L 3/00

5 3 7 D

G 1 0 L 3/00

5 2 1 R

テーマコード (参考)

5 D 0 1 5

審査請求 未請求 請求項の数 16 O L (全 16 頁)

(21) 出願番号 特願2005-101016 (P2005-101016)
 (22) 出願日 平成17年3月31日 (2005.3.31)
 (31) 優先権主張番号 10/814, 906
 (32) 優先日 平成16年3月31日 (2004.3.31)
 (33) 優先権主張国 米国 (US)

(71) 出願人 500046438
 マイクロソフト コーポレーション
 アメリカ合衆国 ワシントン州 9805
 2-6399 レッドモンド ワン マイ
 クロソフト ウェイ
 (74) 代理人 100077481
 弁理士 谷 義一
 (74) 代理人 100088915
 弁理士 阿部 和夫
 (72) 発明者 アレジャンドロ アチェロ
 アメリカ合衆国 98052 ワシントン
 州 レッドモンド ワン マイクロソフト
 ウェイ マイクロソフト コーポレーシ
 ョン内

最終頁に続く

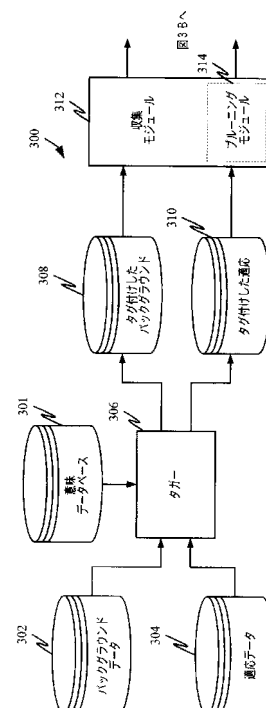
(54) 【発明の名称】 意味管理を用いた言語モデル適応

(57) 【要約】

【課題】 言語モデルを適応させるための方法および装置を提供する。

【解決手段】 言語モデルを適応させるための方法および装置が提供される。この方法および装置は、インドメイン意味情報を利用して言語モデルの管理されたクラスベースの適応を提供する。

【選択図】 図3A



【特許請求の範囲】**【請求項 1】**

N グラム言語モデルを新しいドメインに適応させる方法であって、
前記新しいドメインに向けられていない一般的テキストフレーズを示すバックグラウンドデータを受け取ることと、
前記新しいドメインで使用され、クラスに編成された意味エンティティのセットを受け取ることと、
前記バックグラウンドデータ、前記意味エンティティおよびそのクラスに基づいてバックグラウンドN グラムクラスカウントデータを生成することと、
前記バックグラウンドN グラムクラスカウントデータに基づいて言語モデルをトレーニングすることと
を備えることを特徴とする方法。

10

【請求項 2】

前記新しいドメインで使用されるテキストフレーズを示す適応データを受け取ることと、
前記適応データ、前記意味エンティティおよびそのクラスに基づいて適応N グラムクラスカウントデータを生成することと
をさらに備え、
前記言語モデルをトレーニングすることは、前記バックグラウンドN グラムクラスカウントデータおよび前記適応N グラムクラスカウントデータに基づいてトレーニングすることを備えることを特徴とする請求項 1 に記載の方法。

20

【請求項 3】

前記バックグラウンドN グラムクラスカウントデータ、前記意味エンティティおよびそのクラスに基づいてバックグラウンドN グラムワードデータを生成することと、
前記適応N グラムクラスカウントデータ、前記意味エンティティおよびそのクラスに基づいて適応N グラムワードデータを生成することと
をさらに備え、
前記バックグラウンドN グラムクラスカウントデータおよび前記適応N グラムクラスカウントデータに基づいて前記言語モデルをトレーニングすることは、バックグラウンドN グラムワードデータおよび適応N グラムワードデータを用いることを備えることを特徴とする請求項 2 に記載の方法。

30

【請求項 4】

バックグラウンドN グラムワードデータを生成することは、各データエントリが選択された数のワードを備えるマルチワード意味エンティティのバックグラウンドN グラムワードデータを生成することを備えることを特徴とする請求項 3 に記載の方法。

【請求項 5】

適応N グラムワードデータを生成することは、各データエントリが選択された数のワードを備えるマルチワード意味エンティティの適応N グラムワードデータを生成することを備えることを特徴とする請求項 4 に記載の方法。

【請求項 6】

前記バックグラウンドデータ、前記意味エンティティおよびそのクラスに基づいてバックグラウンドN グラムクラスカウントデータを生成することは、前記意味エンティティおよびそのクラスに基づいてワードレベルのバックグラウンドデータをタグ付けすることを備えることを特徴とする請求項 4 に記載の方法。

40

【請求項 7】

前記適応データ、前記意味エンティティおよびそのクラスに基づいて適応N グラムクラスカウントデータを生成することは、前記意味エンティティおよびそのクラスに基づいてワードレベルの適応データをタグ付けすることを備えることを特徴とする請求項 5 に記載の方法。

【請求項 8】

50

前記バックグラウンドデータ、前記意味エンティティおよびそのクラスに基づいてバックグラウンドNグラムクラスカウントデータを生成することは、前記タグ付けされたバックグラウンドデータの固有のクラスレベルNグラムをカウントすることを備えることを特徴とする請求項6に記載の方法。

【請求項9】

前記適応データ、前記意味エンティティおよびそのクラスに基づいて適応Nグラムクラスカウントデータを生成することは、前記タグ付けされた適応データの固有のクラスレベルNグラムをカウントすることを備えることを特徴とする請求項7に記載の方法。

【請求項10】

前記バックグラウンドデータ、前記意味エンティティおよびそのクラスに基づいてバックグラウンドNグラムクラスカウントデータを生成することは、前記タグ付けされたバックグラウンドデータからいくつかのクラスNグラムを廃棄することを備えることを特徴とする請求項8に記載の方法。

10

【請求項11】

前記適応データ、前記意味エンティティおよびそのクラスに基づいて適応Nグラムクラスカウントデータを生成することは、前記タグ付けされた適応データからいくつかのクラスNグラムを廃棄することを備えることを特徴とする請求項9に記載の方法。

【請求項12】

言語モデルを生成するステップを行うためのコンピュータ実行可能命令を有するコンピュータ可読媒体であって、前記ステップは、

20

選択されたドメインで使用され、クラスに編成された意味エンティティのセットを受け取るステップと、

前記意味エンティティのセットのクラスに相関され、一般テキストを示すバックグラウンドデータに基づくバックグラウンドNグラムクラスカウントデータを受け取るステップと、

前記意味エンティティのセットのクラスに相関され、モデル化される選択されたドメインを示す適応データに基づく適応Nグラムクラスカウントデータを受け取るステップと、

前記バックグラウンドNグラムクラスカウントデータ、前記適応Nグラムクラスカウントデータおよび前記意味エンティティのセットに基づいて言語モデルをトレーニングするステップと

30

を備えたことを特徴とするコンピュータ可読媒体。

【請求項13】

前記言語モデルをトレーニングするステップは、前記バックグラウンドNグラムクラスカウントデータおよび前記意味エンティティのセットに基づいてバックグラウンドワードカウントデータを算出するステップを備えたことを特徴とする請求項12に記載のコンピュータ可読媒体。

【請求項14】

前記言語モデルをトレーニングするステップは、前記適応Nグラムクラスカウントデータおよび前記意味エンティティのセットに基づいて適応ワードカウントデータを算出するステップを備えたことを特徴とする請求項13に記載のコンピュータ可読媒体。

40

【請求項15】

前記言語モデルをトレーニングするステップは、Nグラム相対頻度を平滑化するステップを備えたことを特徴とする請求項14に記載のコンピュータ可読媒体。

【請求項16】

平滑化するステップは、削除補間アルゴリズムを使用するステップを備えたことを特徴とする請求項15に記載のコンピュータ可読媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、言語処理で用いられる言語モデルに関する。詳細には、本発明は、所望のド

50

メインに言語モデルを適応させることに関する。

【背景技術】

【0002】

自動音声認識 (ASR: automatic speech recognition) などの言語処理システムは、しばしば、トレーニングおよびテストデータと実際のドメインデータとの間の不一致に由来するエラーによる性能劣化を扱わなければならないことが多い。よく知られているように、音声認識システムは、音響モデルおよび統計言語モデル (LM: language model) を用いて認識を提供する。音響モデルの新しいドメインへの適応は、限られた成功でしか取り組まれておらず、言語モデルの適応は、満足いく成果を達成していない。

10

【0003】

統計言語モデル (LM) は、単語系列の事前確率推定を提供する。LM は、最も可能性の高い単語系列の仮説探索 (hypothesis search) をガイドするので、ASR およびその他の形態の言語処理において重要な構成要素である。よい LM は、優れた言語処理性能のために不可欠であることが知られている。

【0004】

広く、LM はテストデータに類似していると期待される大量のトレーニングデータから集められた、平滑化された N グラム統計 (n-gram statistics) を使用する。しかしながら、類似性の定義はゆるく、対象となるある所与のドメインにどのデータソースを使用すべきかの決定は、ほとんどの場合試行錯誤により、普通、モデル作成者に委ねられる。

20

【0005】

常に、トレーニングまたはテストデータと実際のドメインまたは「インドメイン (in-domain)」データの間には不一致が存在し、それによりエラーがもたらされる。不一致の1つの源は、テストデータ中の語彙外の単語 (out-of-vocabulary words) に由来するものである。例えば、元々1つの航空会社のために設計された飛行機旅行情報システムは、問題の会社によってサービスが提供される都市名、空港名などでの不一致のために、別の会社に対してはうまく機能しないことがある。

【0006】

別の潜在的な不一致の源は、異なる言語スタイルに由来するものである。例えば、ニュースドメインでの言語スタイルは、飛行機旅行情報ドメインとは異なる。ニュースワイヤやその他の一般的テキストでトレーニングされた言語モデルは、飛行機旅行情報ドメインではあまりうまく機能しないことがある。

30

【0007】

【非特許文献1】Frederick Jelinek and Robert Mercer, "Interpolated Estimation of Markov Source Parameters from Sparse Data," In E. Gelsema and L. Kanal, editors, Pattern Recognition in Practice, pages 381-397, 1980

【発明の開示】

【発明が解決しようとする課題】

【0008】

異なる技法を用いて大量のバックグラウンドデータでトレーニングされた LM を適応させるために様々なアプローチが試されてきたが、どれも優れた結果を達成しておらず、ゆえに、LM 適応における改善が絶えず必要とされている。前述の問題の1つまたは複数に対処する方法があれば役立つであろう。

40

【課題を解決するための手段】

【0009】

言語モデルを適応させるための方法および装置が提供される。この方法および装置は、インドメイン意味情報 (in-domain semantic information) を利用して言語モデルの管理されたクラスベースの適応 (supervised class-based adaptation) を提供する。

50

【 0 0 1 0 】

一般に、適応を行うのに使用されるリソースは、一般のテキストを示すバックグラウンドデータと、選択されたドメインで使用され、クラスに編成された意味エンティティのセットから導出される。別の実施形態では、モデル化される選択されたドメインを示す適応データも使用される。

【 0 0 1 1 】

前記の別の実施形態では、そのようなデータは、意味エンティティのセットのクラスに相関され、一般のテキストを示すバックグラウンドデータに基づくバックグラウンドNグラムクラスカウントデータと、意味エンティティのセットのクラスに相関され、モデル化される選択されたドメインを示す適応データに基づく適応Nグラムクラスカウントデータとを備える。このデータから、そして意味エンティティのセットを使用して、バックグラウンドワードカウントデータおよび適応ワードカウントデータを計算し、適応データおよび意味項目のセットのドメインに言語モデルを適応させるためのベースとして使用することができる。

10

【 発明を実施するための最良の形態 】

【 0 0 1 2 】

本発明は、言語モデル適応のシステムおよび方法に関する。しかしながら、本発明をより詳細に論じる前に、本発明を使用することのできる1つの例示的な環境について論じることにする。

【 0 0 1 3 】

図1に、本発明を実施することのできる適したコンピューティング環境100の一例を示す。コンピューティングシステム環境100は、適したコンピューティング環境の一例にすぎず、本発明の使用または機能の範囲に関していかなる限定を示唆するものではない。また、コンピューティング環境100は、例示的な動作環境100に示す構成要素のいずれか1つまたはその組合せに関するいかなる依存性または要件を有するものと解釈されるべきではない。

20

【 0 0 1 4 】

本発明は、数多くの他の汎用または専用のコンピューティングシステム環境または構成で動作する。本発明と共に使用するのに適すると考えられる周知のコンピューティングシステム、環境、および/または構成の例には、それだけに限られないが、パーソナルコンピュータ、サーバコンピュータ、ハンドヘルドまたはラップトップデバイス、マルチプロセッサシステム、マイクロプロセッサベースのシステム、セットトップボックス、プログラム可能な民生用電子機器、ネットワークPC、ミニコンピュータ、メインフレームコンピュータ、上記のシステムまたはデバイスのいずれかを含む分散コンピューティング環境などが含まれる。

30

【 0 0 1 5 】

本発明は、コンピュータにより実行される、プログラムモジュールなどのコンピュータ実行可能命令の一般的コンテキストで説明することができる。一般に、プログラムモジュールには、特定のタスクを行い、または特定の抽象データ型を実装するルーチン、プログラム、オブジェクト、コンポーネント、データ構造などが含まれる。当業者は、本明細書の説明および/または図をコンピュータ実行可能命令として実施することができ、それは、以下で論じる任意の形態のコンピュータ可読媒体上で具体化することができる。

40

【 0 0 1 6 】

本発明は、通信ネットワークを介してリンクされたリモート処理装置によってタスクが行われる分散コンピューティング環境で実施することもできる。分散コンピューティング環境では、プログラムモジュールは、メモリ記憶装置を含むローカルとリモート両方のコンピュータ記憶媒体に位置することがある。

【 0 0 1 7 】

図1を参照すると、本発明を実施するための例示的なシステムは、コンピュータ110の形態の汎用コンピューティングデバイスを含む。コンピュータ110の構成要素には、

50

それだけに限られないが、処理装置 120、システムメモリ 130、およびシステムメモリを含む様々なシステム構成要素を処理装置 120 に結合するシステムバス 121 が含まれ得る。システムバス 121 は、メモリバスまたはメモリコントローラ、周辺バス、および様々なバスアーキテクチャのいずれかを使用したローカルバスを含むいくつかのタイプのバス構造のいずれでもよい。例として、限定ではなく、そのようなアーキテクチャには、ISA (Industry Standard Architecture) バス、MCA (Micro Channel Architecture) バス、拡張 ISA (EISA) バス、VESA (Video Electronics Standards Association) ローカルバス、およびメザンバスとも呼ばれる PCI (Peripheral Component Interconnect) バスが含まれる。 10

【0018】

コンピュータ 110 は、通常、様々なコンピュータ可読媒体を含む。コンピュータ可読媒体は、コンピュータ 110 によってアクセスすることができる任意の利用可能な媒体とすることができ、それには揮発性媒体と不揮発性媒体の両方、リムーバブル媒体と非リムーバブル媒体の両方が含まれる。例として、限定ではなく、コンピュータ可読媒体には、コンピュータ記憶媒体および通信媒体が含まれ得る。コンピュータ記憶媒体には、コンピュータ可読命令、データ構造、プログラムモジュールまたはその他のデータなどの情報を記憶するための任意の方法または技術で実施された、揮発性および不揮発性の両方、リムーバブルおよび非リムーバブルの両方の媒体が含まれる。コンピュータ記憶媒体には、それだけに限られないが、RAM、ROM、EEPROM、フラッシュメモリなどのメモリ 20 技術、CD-ROM、デジタル多用途ディスク (DVD) などの光ディスク記憶、磁気カセット、磁気テープ、磁気ディスク記憶などの磁気記憶装置、または所望の情報の格納に使用でき、コンピュータ 110 によってアクセスすることができる他の任意の媒体が含まれる。通信媒体は、通常、コンピュータ可読命令、データ構造、プログラムモジュールまたはその他のデータを、搬送波 WAVE または他のトランスポート機構などの変調データ信号に具体化し、任意の情報配信媒体を含む。「変調データ信号」という用語は、信号に情報を符号化するような形でその特性の 1 つまたは複数が設定または変更された信号を意味する。例として、限定ではなく、通信媒体には、有線ネットワークや直接配線接続などの有線媒体、および音響、RF、赤外線、その他の無線媒体などの無線媒体が含まれる。上記のいずれの組合せも、コンピュータ可読媒体の範囲内に含まれるべきである。 30

【0019】

システムメモリ 130 は、読取り専用メモリ (ROM) 131 やランダムアクセスメモリ (RAM) 132 などの揮発性および / または不揮発性メモリの形態のコンピュータ記憶媒体を含む。起動時などに、コンピュータ 110 内の要素間の情報転送を助ける基本ルーチンが入った基本入出力システム (BIOS) 133 は、通常、ROM 131 に格納される。RAM 132 は、通常、処理装置 120 によって直ちにアクセス可能であり、そして / または現在処理されているデータおよび / またはプログラムモジュールを収容する。例として、限定ではなく、図 1 に、オペレーティングシステム 134、アプリケーションプログラム 135、その他のプログラムモジュール 136、およびプログラムデータ 137 を示す。 40

【0020】

コンピュータ 110 は、他のリムーバブル / 非リムーバブル、揮発性 / 不揮発性コンピュータ記憶媒体を含むことができる。例にすぎないが、図 1 に、非リムーバブルの不揮発性磁気媒体との間で読取りまたは書込みを行うハードディスクドライブ 141、リムーバブルの不揮発性磁気ディスク 152 との間で読取りまたは書込みを行う磁気ディスクドライブ 151、および CD-ROM や他の光媒体など、リムーバブルの不揮発性光ディスク 156 との間で読取りまたは書込みを行う光ディスクドライブ 155 を示す。例示的な動作環境で使用することのできる他のリムーバブル / 非リムーバブル、揮発性 / 不揮発性のコンピュータ記憶媒体には、それだけに限られないが、磁気テープカセット、フラッシュメモリカード、デジタル多用途ディスク、デジタルビデオテープ、ソリッドステート RA 50

M、ソリッドステートROMなどが含まれる。ハードディスクドライブ141は、通常、インターフェース140などの非リムーバブルメモリインターフェースを介してシステムバス121に接続され、磁気ディスクドライブ151および光ディスクドライブ155は、通常、インターフェース150などのリムーバブルメモリインターフェースによってシステムバス121に接続される。

【0021】

前述し、図1に示すドライブおよび関連するコンピュータ記憶媒体は、コンピュータ可読命令、データ構造、プログラムモジュールおよびコンピュータ110のその他のデータのストレージを提供する。図1では、例えば、ハードディスクドライブ141は、オペレーティングシステム144、アプリケーションプログラム145、その他のプログラムモジュール146、およびプログラムデータ147を格納するものとして示されている。これらのコンポーネントは、オペレーティングシステム134、アプリケーションプログラム135、その他のプログラムモジュール136、およびプログラムデータ137と同じでも、異なってもよいことに留意されたい。オペレーティングシステム144、アプリケーションプログラム145、その他のプログラムモジュール146、およびプログラムデータ147には、少なくともそれらが異なるコピーであることを示すために、図では異なる番号が付与されている。

【0022】

ユーザは、キーボード162や、マイクロフォン163などの入力装置、マウス、トラックボール、タッチパッドといったポインティングデバイス161を介してコンピュータ110にコマンドおよび情報を入力することができる。他の入力装置（図示せず）には、ジョイスティック、ゲームパッド、衛星アンテナ、スキャナなどが含まれることがある。上記およびその他の入力装置は、システムバスに結合されたユーザ入力インターフェース160を介して処理装置120に接続されるが、パラレルポート、ゲームポート、ユニバーサルシリアルバス（USB）といった他のインターフェースおよびバス構造によって接続することもできる。モニタ191または他の種類の表示装置もビデオインターフェース190などのインターフェースを介してシステムバス121に接続される。モニタ191に加えて、コンピュータはスピーカ197やプリンタ196など他の周辺出力装置を含むこともでき、それらは、出力周辺インターフェース195を介して接続することができる。

【0023】

コンピュータ110は、リモートコンピュータ180など、1つまたは複数のリモートコンピュータへの論理接続を使用してネットワーク化された環境で動作することができる。リモートコンピュータ180は、パーソナルコンピュータ、ハンドヘルド装置、サーバ、ルータ、ネットワークPC、ピアデバイスおよびその他の共通ネットワークノードとすることができ、通常、コンピュータ110に関連して上述した要素の多くまたはすべてを含む。図1に示す論理接続には、ローカルエリアネットワーク（LAN）171およびワイドエリアネットワーク（WAN）173が含まれるが、他のネットワークを含むこともある。そのようなネットワーク環境は、オフィス、企業規模のコンピュータネットワーク、イントラネットおよびインターネットで一般的である。

【0024】

LANネットワーク環境で使用されるとき、コンピュータ110はネットワークインターフェースまたはアダプタ170を介してLAN171に接続される。WANネットワーク環境で使用されるとき、コンピュータ110は、通常、モデム172、またはインターネットなどのWAN173を介して通信を確立するための他の手段を含む。モデム172は、内蔵でも外付けでもよく、ユーザ入力インターフェース160または他の適当な機構を介してシステムバス121に接続することができる。ネットワーク化環境では、コンピュータ110に関連して示すプログラムモジュール、またはその部分は、リモートメモリ記憶装置に格納することもできる。例として、限定ではなく、図1に、リモートアプリケーションプログラム185をリモートコンピュータ180上にあるものとして示す。図示

10

20

30

40

50

したネットワーク接続は例示的なものであり、コンピュータ間で通信を確立する他の手段を使用できることが理解されるであろう。

【 0 0 2 5 】

本発明は、図 1 との関連で説明したようなコンピュータシステム上で実行できることに留意すべきである。しかしながら、本発明は、サーバや、メッセージ処理専用のコンピュータでも、本発明の異なる部分が分散コンピューティングシステムの異なる部分で実行される分散システムでも実行することができる。

【 0 0 2 6 】

上に示したように、本発明は、言語モデル適応のためのシステムおよび方法に関するものである。適応を行うのに使用されるリソースは、適応される必要のあるバックグラウンド LM を含む。一般に、バックグラウンド LM は、それだけに限られないが、ニュース記事など、バックグラウンドトレーニングデータの大きなコーパスから得られる。このバックグラウンドトレーニングデータを使用して、バックグラウンド言語モデルのための N グラム統計が得られる。

【 0 0 2 7 】

意味データベースまたは意味情報は、適応のための管理された情報 (supervised information) を提供する。この説明では、意味データベースは意味エンティティ (クラス) のリストを広く、模式的に表し、それらがインドメイン自然言語テキストにおいて遭遇するのと同じ形態であると想定される具現化 (realizations) のリストを各エンティティが伴う。例えば、意味データベースは、複数のクラスについて一般的に明確に定義された意味エンティティのリストの形態とすることができる。例えば、以下に一例として使用するように、予約を行う旅行情報を得るために航空会社によって使用される音声認識における言語モデルの意味項目には、その航空会社によってサービスが提供される都市および飛行先の様々な空港のリストが含まれ得る。意味エンティティおよびクラスの別の例は、企業の従業員、当月の日、当年の月のリストとすることができ、それらは、おそらく、スケジューリング用途のインドメインに含まれるであろう。

【 0 0 2 8 】

意味クラスは、開いたクラスと閉じたクラスに分類することができる。開いたクラス中のクラスメンバは、ドメイン間で変化するが、閉じたクラス中のクラスメンバは変化しない。例えば、飛行機旅行用途での意味クラスは以下のものとすることができるであろう。

【 0 0 2 9 】

開いたクラス : { AIRLINE (航空会社)、AIRPORT (空港)、CITY (都市)、STATE (州) }

閉じたクラス : { DAYS (日)、MONTH (月)、INTERVAL (間隔)、CLASS OF SERVICE (サービスクラス)、ZONE (ゾーン)、FOOD SERVICE (フードサービス)、GROUND SERVICE (地上サービス) }

用途ごとに、意味クラスの数およびタイプは大きく変化する可能性がある。しかしながら、飛行機旅行用途などいくつかの用途では、その意味クラスが特定されると、その意味エンティティまたは単語 (word) レベルの具現化だけで、別の航空会社による使用のためにその言語モデルを本質的に適応させるために変更する必要があるすべてとすることができる。

【 0 0 3 0 】

言語モデル適応で使用されるオプションの第 3 のリソースは適応データである。適応データは、インドメインのアプリケーションでのクラスの使用例となりうるセンテンス (sentences)、フレーズ (phrases)、テキストセグメント (text segments) などの形態の実際のまたはインドメインのデータを備える。バックグラウンドデータに比べて、適応データは、普通、バックグラウンドデータより何桁も少ない。一実施形態では、インドメインデータは、適応開発データ (adaptation development data) と適応トレーニングデータ (adaptation

10

20

30

40

50

n training data) にサブ分割される。適応トレーニングデータは、バックグラウンドトレーニングセットと組み合わせられてより大きなトレーニングセットになり、両セットからの N グラムカウントは等しい重みで混合される（ただし、他の混合スキームも可能であり、N グラムカウントは、MAP 適応などのように、異なる重みで混合することもできる）。適応開発データは、厳密に、バックグラウンド言語モデルおよび適応言語モデルの両方を平滑化するために使用される。開発セットからの N グラムは、バックグラウンド / 適応言語モデルには含まれない。

【0031】

例示的な実施形態では、すべてのデータセットは単語レベルの自然言語テキストである。

10

【0032】

（クラスベースの適応）

管理された意味情報 (supervised semantic information) がクラスベースの言語モデルの使用を通じて言語モデルに組み込まれる。簡潔には、単一の意味クラス c_3 に属する新しい単語 w_3 の確率推定は以下のように行うことができる。

【0033】

$$Pr(w_3 | w_2 w_1) = Pr(w_3 | c_3) \cdot Pr(c_3 | w_2 w_1)$$

(1)

ここで、 $Pr(w_3 | c_3 w_2 w_1) = Pr(w_3 | c_3)$ というモデルリングを前提とする。

20

【0034】

例えば、 $Pr(city\ name | fly\ to)$ (都市名 | へ飛行) は以下を用いて推定される。

【0035】

$$Pr(city\ name | fly\ to) = Pr(city\ name | CITY (都市)) \cdot Pr(CITY | fly\ to)$$

ここで、 $Pr(CITY | fly\ to)$ は、意味クラスでタグ付けされたトレーニングデータを用いて推定され、 $Pr(city\ name | CITY)$ は、インドメイン意味データベースを用いて適応される。以前のインドメイン知識が利用可能である場合、よくある都市名は、まれな都市名より高い確率で割り当てることができ、そうでなければ、都市名の一樣分布が想定される。クラスベースの適応アプローチを用いることの利点は以下の通りである。

30

【0036】

単語コンテキストを仮定した意味クラスの確率はうまく推定することができる。上記の例では、 $Pr(city\ name | fly\ to)$ は、トレーニングデータおよび適応データで非常に類似していると考えられる。

【0037】

インドメイン意味データを用いて $Pr(w_3 | c_3)$ を適応させることにより、高速 LM 適応を行なうことができる。適応された確率 $Pr(w_3 | c_3)$ は、ドメイン特有の言語モデルを再トレーニングする新しいトレーニングテキストを収集することなく、カウント「 $w_1 w_2 w_3$ 」と組み合わせられる。

40

【0038】

語句 (word phrases) が意味クラスにカプセル化されるので、より広範な単語コンテキストで確率推定を実現することができる。例えば、5 グラム「los angeles to new york」は、トライグラム (trigrams) 「los angeles to」、「angeles to new」および「to new york」のシーケンスとしてモデル化されるよりも直感的に満足の行くクラストライグラム「CITY to CITY」としてモデル化される。

【0039】

50

(適応手順)

図 2 に例示的な適応手順 200 を示す。図 3 A および 3 B に、手順 200 を行うための例示的システム 300 を示す。上に示したように、適応データの使用はオプションであるが、本発明の別の実施形態である。両方を用いた実施形態について以下に説明するが、これは必須または限定とみなすべきではない。また、説明に進む前に、手順 200 およびシステム 300 が、一般に同時にバックグラウンドデータおよびオプションの適応データ上で動作するものとして記述されていることに留意されたい。しかしながら、これは、理解を簡単にするためであり、必要または限定とみなすべきではない。

【 0040 】

ステップ 202 は、一般に、バックグラウンドデータおよび適応データの両方のためにタグ付けデータを得ることを表している。例示した実施形態では、これは、202 に示すように単語レベルのデータをタグ付けすることを含む。特に、トレーニング (バックグラウンドおよび適応) データが、図 2 のステップ 202 で意味クラスと先ずタグ付けされる。当然ながら、タグ付けされたデータが存在する場合、このステップは不要である。図 3 A では、意味データベースが 301 で示され、トレーニングデータはコーパス 302 および 304 にあり、この場合、タグ付けはタガー 306 によって行われる。

10

【 0041 】

タガー 306 は、コーパス 302 および 304 によって提供された単語レベルのテキストを変更し、そこで認識された意味エンティティのクラスを示すタグを付加する。例えば、「fly from san francisco to」が与えられ、「san francisco」が意味クラス「CITY」に属すると知っている場合、タガー 306 からの出力は、「fly from CITY to」になることになる。意味エンティティのいくつかが対応する意味クラスで置き換えられた単語レベルのトレーニングデータは 308 および 310 で示されている。

20

【 0042 】

一実施形態では、タグ付けにヒューリスティックス (heuristics) を適用することができる。そのようなヒューリスティックスは、タグ付けのための単純な文字列マッチングアプローチを含んでもよい。タガー 306 は、所与のデータベースエントリをテキスト中の単語のシーケンスと合致させ、そのようにして特定された最長のフレーズにクラスラベルを割り当てる。別の実施形態では、単語のあいまい性が異なるクラス間で発生した場合、その語句はタグなしのままとされる。別の実施形態では、各意味クラス候補に確率を割り当てることによって、ソフトなタグ付けを行うことができるであろう。

30

【 0043 】

ステップ 202 でタグ付けが行われた後、タグ付けデータが、別途、提供されない場合、手順はステップ 204 に進んですべてのトレーニングテキストからクラス N グラムカウントを収集し、そうでない場合は、タグ付けデータに含まれる固有の N グラムをカウントする。図 3 A で、このステップは収集モジュール 312 によって行われる。

【 0044 】

クラス N グラムカウントのブルーニングを備えるオプションのステップ 206 を必要に応じて行うこともできる。クラスベースの適応では、クラス N グラムが単語 N グラムに拡張されると、言語モデルのサイズは、各意味クラス中の要素数によって強く影響される。例えば、クラストライグラム「PERSON joins COMPANY (「人」が「会社」に入社する)」(この場合、「PERSON」および「COMPANY」が意味クラスを含む) は、「PERSON」と「COMPANY」がそれぞれ何千ものクラス要素を含むとき何百万もの単語トライグラムを生じる。それゆえ、言語モデルブルーニングが、言語モデルのサイズを扱いやすくするために必要になることがある。一実施形態では、複数の意味クラスを含む N グラムは廃棄される。計算リソースが利用可能であれば、それらを保持することもできるであろう。加えて、単語 N グラムに拡張する前に、クラス N グラムのカウントカットオフブルーニングを用いることができる。図 3 A には、収集モジュール 312 がブルーニングモジュール 314 を使用することによってこの機能を行うもの

40

50

してと示されている。収集モジュール 3 1 2 からの出力は、図 3 B に示すバックグラウンド N グラムカウントデータ 3 1 6 および適応 N グラムカウントデータ 3 1 8 を備える。

【 0 0 4 5 】

ステップ 2 0 8 で、クラス N グラムが意味データベース 3 0 1 を用いて単語 N グラムに拡張される。図 3 B では、このステップはワード N グラムジェネレータ 3 2 0 によって行われる。一実施形態では、ワード N グラムジェネレータ 3 2 0 は、以下の拡張アルゴリズムを実施し、バックグラウンド N グラムワードカウントデータ 3 2 2 および適応 N グラムワードカウントデータ 3 2 4 を生成することができる。

【 0 0 4 6 】

(a) クラス N グラムを仮定し、クラスタグをそのクラス要素のそれぞれによって置き換える。 10

【 0 0 4 7 】

例えば、クラストライグラム「analyst for COMPANY」は、単語 4 グラム「analyst for x . y .」を作成することができ、ここで、「x . y .」は意味データベース中の会社名 (Verizon Wireless など) である。

【 0 0 4 8 】

(b) クラス N グラムカウントから単語 N グラムカウントを算出する。

【 0 0 4 9 】

単語 N グラムカウントは、 $Pr(word | class)$ に応じて、その対応するクラス N グラムカウントの一部として算出される。 20

【 0 0 5 0 】

意味クラス「COMPANY」の確率が、

$$Pr(microsoft | COMPANY) = 0 . 5$$

$$Pr(oracle | COMPANY) = 0 . 25$$

$$Pr(verizon wireless | COMPANY) = 0 . 25$$

であり、

N グラム「analyst for COMPANY」が 5 カウントであったと想定すると、

単語レベルの N グラムカウントデータは、 30

$$「analyst for microsoft」= 2 . 5$$

$$「analyst for oracle」= 1 . 25$$

$$「analyst for verizon wireless」= 1 . 25$$

になるであろう。

【 0 0 5 1 】

上記の例では、生成された単語 4 グラム「analyst for x . y .」のカウントは、

$$\#(「analyst for COMPANY」) \cdot Pr(「x . y .」| COMPANY)$$

に等しい。 40

【 0 0 5 2 】

(c) しかしながら、クラスベースの N グラムは、マルチワード (multi - word) 意味エントリのために特定の N グラムのトレーニングと動作しない単語レベルの N グラムを生成することができることに留意されたい。例えば、3 ワードの N グラム言語モデルが望まれていると想定すると、「analyst for verizon wireless」は正しい形態のものではない。この状況では、スライディングウィンドウを用いてより低次の単語 N グラムが生成される。上記の例では、「analyst for verizon」も 1 . 25 のカウントを有するであろうし、「for verizon wireless」も 1 . 25 のカウントを有するであろう。

【 0 0 5 3 】

しかしながら、クラスがNグラム中の他の場所、すなわち、右端の位置以外に現れた場合、マルチワード意味項目拡張 (multi-word semantic item expansion) について二重カウンティングを避けるのに以下のステップを行うことができる。先の例と同様に、拡張に関するステップ (a) と計算に関するステップ (b) が同じ方式で行われる。しかし、ステップ (c) は行われず、むしろ、Nグラムのコンテキストは、拡張後に所望の数の右端の単語だけを取ることににより短縮される。

【0054】

例として、カウント5を有する「COMPANY analyst said」のクラストライグラムを、

$$\Pr(\text{microsoft} | \text{COMPANY}) = 0.5$$

$$\Pr(\text{oracle} | \text{COMPANY}) = 0.25$$

$$\Pr(\text{verizon wireless} | \text{COMPANY}) = 0.25$$

の意味クラス「COMPANY」の同じ確率で想定すると、その単語レベルのNグラムデータは、

$$\text{「microsoft analyst said」} = 2.5$$

$$\text{「oracle analyst said」} = 1.25$$

$$\text{「wireless analyst said」} = 1.25$$

になり、ここで、「wireless analyst said」は、トライグラムに右端の3ワードだけを取ることににより実現されたものである。

【0055】

意味データベース301がタガー306およびワードNグラムジェネレータ320と共に動作可能である場合が示されているが、データベース301のインスタンスのそれぞれでの内容は、多くの用途で異なることがあり、それによりこの方法がより役立つことを理解されたい。

【0056】

ステップ210で、言語モデル326が、バックグラウンドデータおよびオプションの適応データの生成された単語Nグラムカウントを用いてトレーニングされ、ここでは、トレーニングモジュール328によって行われる。必要ならば、単語Nグラムに関してカウントカットオフブルーニングを行って言語モデルのサイズをさらに低減することもできる。

【0057】

トレーニングは、Nグラム相対頻度推定 (n-gram relative frequency estimates) を平滑化することを含むことができる。例えば、参照により本明細書に組み込まれる非特許文献1に記載されている削除補間法 (deleted-interpolation method) を、Nグラム相対頻度推定を平滑化するために使用することができる。簡潔には、この再帰的削除補間式 (recursive deleted-interpolation formula) は以下のように定義される。

【0058】

【数1】

$$\Pr_I(w|w_1^{n-1}) = (1 - \lambda_{w_1^{n-1}}) \cdot f(w|w_1^{n-1}) + \lambda_{w_1^{n-1}} \cdot \Pr_I(w|w_2^{n-1})$$

$$\Pr_I(w) = (1 - \lambda) \cdot f(w) + \lambda \cdot \frac{1}{V}$$

【0059】

ここで、

【0060】

10

20

30

40

【数 2】

$$f(w|w_{1_k}^{n-1})$$

【0 0 6 1】

は単語 N グラムの相対頻度を表し、

【0 0 6 2】

【数 3】

$$w_{1_k}^{n-1}$$

【0 0 6 3】

は前の $n - 1$ 語に及ぶ単語履歴である。異なるコンテキスト順で均一な単語分布 $1 / V$ の N グラムモデルが線形に補間される。補間重み

【0 0 6 4】

【数 4】

$$\lambda_{w_1^{n-1}}$$

【0 0 6 5】

は、周知の最尤法を用いて推定することができる。データのまばらさのために、補間重みは、普通、単語コンテキストをクラスにグループ化することにより推定パラメータ数を低減するために結び付けられる。1つの可能な方法は、ある所与の単語コンテキストの出現数に基づいてパラメータをバケット化することである。

【0 0 6 6】

ステップ 2 1 0 により、管理された言語モデル適応を完了し、この例では、削除補間された言語モデルが提供される。言語処理システムにおける削除補間された言語モデルの実装は、標準 ARPA 形式のバックオフ言語モデルへの変換を含んでもよい。2004 年 3 月 26 日に出願された「REPRESENTATION OF A DELETED INTERPOLATION N - GRAM LANGUAGE MODEL IN ARPA STANDARD FORMAT」というタイトルの同時継続中の米国特許出願に、ARPA 形式への変換の一例について記載されており、それをワンパスシステムに用いることができる。

【0 0 6 7】

本発明を特定の実施形態を参照して説明してきたが、本発明の精神および範囲を逸脱することなく、形態および詳細に変更を加えることができることを当業者は理解するであろう。

【図面の簡単な説明】

【0 0 6 8】

【図 1】本発明を実施することのできる一般的なコンピューティング環境を示すブロック図である。

【図 2】言語モデルを適応させるための流れ図である。

【図 3 A】言語モデルを適応させるためのシステムを示すブロック図である。

【図 3 B】言語モデルを適応させるためのシステムを示すブロック図である。

【符号の説明】

【0 0 6 9】

- 1 0 0 コンピューティング環境
- 1 1 0 コンピュータ
- 1 2 0 処理装置
- 1 3 0 システムメモリ
- 1 3 1 ROM
- 1 3 2 RAM
- 1 3 3 BIOS

10

20

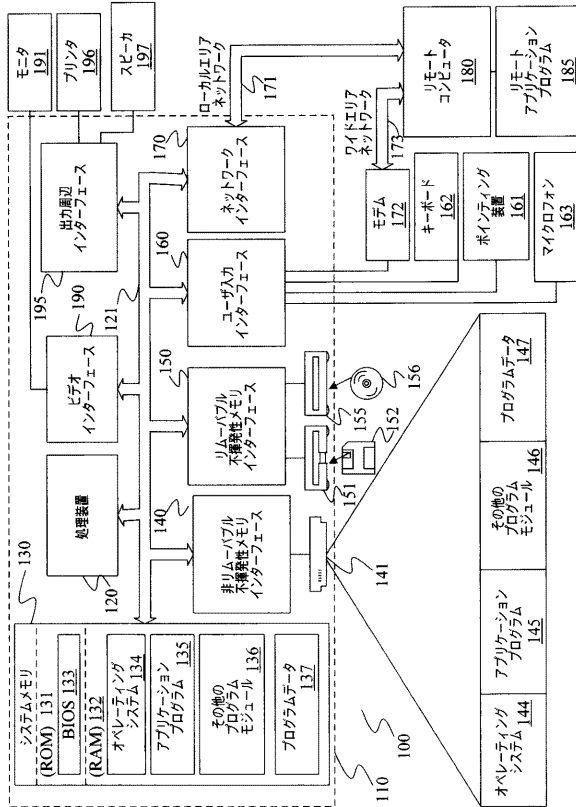
30

40

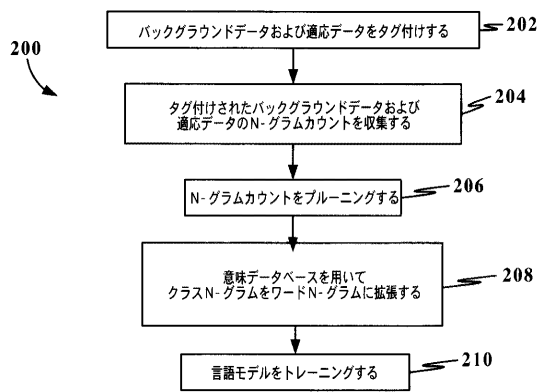
50

1 3 4	オペレーティングシステム	
1 3 5	アプリケーションプログラム	
1 3 6	その他のプログラムモジュール	
1 3 7	プログラムデータ	
1 4 0	非リムーバブル不揮発性メモリインターフェース	
1 4 1	ハードディスクドライブ	
1 4 4	オペレーティングシステム	
1 4 5	アプリケーションプログラム	
1 4 6	その他のプログラムモジュール	
1 4 7	プログラムデータ	10
1 5 0	リムーバブル不揮発性メモリインターフェース	
1 5 1	磁気ディスクドライブ	
1 5 2	リムーバブルの不揮発性磁気ディスク	
1 5 5	光ディスクドライブ	
1 5 6	リムーバブルの不揮発性光ディスク	
1 6 0	ユーザ入力インターフェース	
1 6 1	ポインティング装置	
1 6 2	キーボード	
1 6 3	マイクロフォン	
1 7 0	ネットワークインターフェース	20
1 7 1	ローカルエリアネットワーク	
1 7 2	モデム	
1 7 3	ワイドエリアネットワーク	
1 8 0	リモートコンピュータ	
1 8 5	リモートアプリケーションプログラム	
1 9 0	ビデオインターフェース	
1 9 1	モニタ	
1 9 5	出力周辺インターフェース	
1 9 6	プリンタ	
1 9 7	スピーカ	30
3 0 0	例示的システム	
3 0 1	意味データベース	
3 0 2	バックグラウンドデータ	
3 0 4	適応データ	
3 0 6	タグ	
3 0 8	タグ付けしたバックグラウンド	
3 1 0	タグ付けした適応	
3 1 2	収集モジュール	
3 1 4	ブルーニングモジュール	
3 1 6	バックグラウンドN - グラムクラスカウントデータ	40
3 1 8	適応N - グラムクラスカウントデータ	
3 2 0	ワードN - グラムジェネレータ	
3 2 2	バックグラウンドN - グラムワードカウントデータ	
3 2 4	適応N - グラムワードカウントデータ	
3 2 6	適応された言語モデル	
3 2 8	トレーニングモジュール	

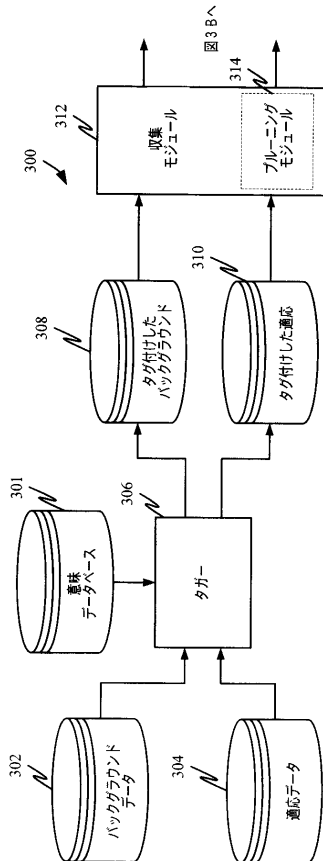
【図 1】



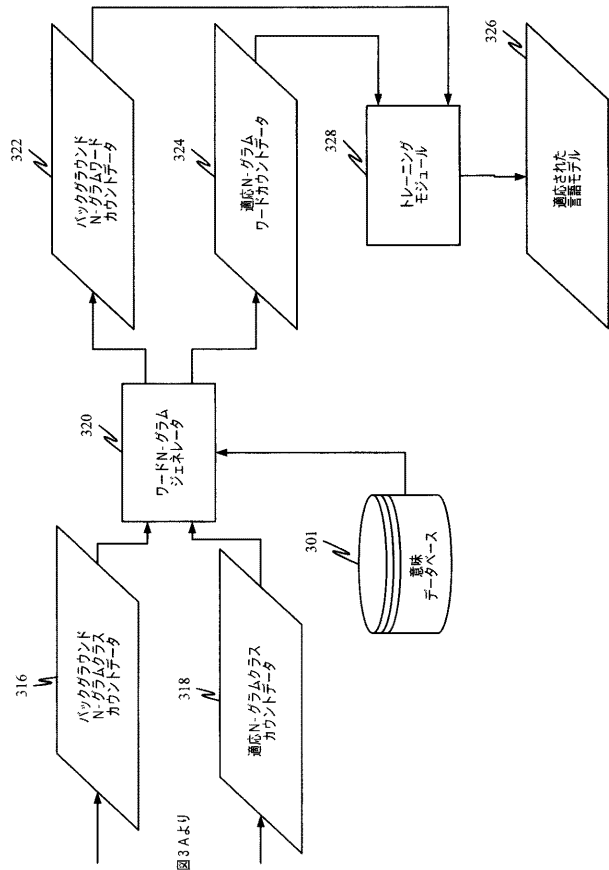
【図 2】



【図 3 A】



【図 3 B】



フロントページの続き

- (72)発明者 シプリアン アイ . ケルバ
アメリカ合衆国 9 8 0 5 2 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マ
イクロソフト コーポレーション内
- (72)発明者 ミラインド マハジャン
アメリカ合衆国 9 8 0 5 2 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マ
イクロソフト コーポレーション内
- (72)発明者 イーチュン タム
アメリカ合衆国 9 8 0 5 2 ワシントン州 レッドモンド ワン マイクロソフト ウェイ マ
イクロソフト コーポレーション内

F ターム(参考) 5D015 HH23