

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第6044996号  
(P6044996)

(45) 発行日 平成28年12月14日(2016.12.14)

(24) 登録日 平成28年11月25日(2016.11.25)

(51) Int.Cl. F 1  
**G 0 6 F 1 7 / 2 8 ( 2 0 0 6 . 0 1 )** G 0 6 F 1 7 / 2 8 6 2 7

請求項の数 5 (全 26 頁)

(21) 出願番号	特願2013-149869 (P2013-149869)	(73) 特許権者	000004226 日本電信電話株式会社 東京都千代田区大手町一丁目5番1号
(22) 出願日	平成25年7月18日(2013.7.18)	(73) 特許権者	504132272 国立大学法人京都大学 京都府京都市左京区吉田本町36番地1
(65) 公開番号	特開2015-22508 (P2015-22508A)	(74) 代理人	110001519 特許業務法人太陽国際特許事務所
(43) 公開日	平成27年2月2日(2015.2.2)	(72) 発明者	須藤 克仁 東京都千代田区大手町二丁目3番1号 日 本電信電話株式会社内
審査請求日	平成27年8月27日(2015.8.27)	(72) 発明者	永田 昌明 東京都千代田区大手町二丁目3番1号 日 本電信電話株式会社内

最終頁に続く

(54) 【発明の名称】 文字列対応付け装置、方法、及びプログラム

(57) 【特許請求の範囲】

【請求項1】

異なる第1の言語及び第2の言語にそれぞれ属する同じ意味の文字列の組み合わせである文字列組について、前記第1の言語の文字列と、前記第2の言語の文字列との間で文字の対応付けを行う文字列対応付け装置であって、

前記文字列組を複数組記憶した文字列組データベースに記憶された前記文字列組の各々に対して、前記文字列組の各文字列を、前記文字列の先頭から順番に、他方の言語の部分文字列と翻字関係にない0文字以上の部分文字列を示す前置非翻字セグメントと、前記他方の言語の部分文字列と翻字関係にある0文字以上の部分文字列を示す翻字セグメントと、前記他方の言語の部分文字列と翻字関係にない0文字以上の部分文字列を示す後置非翻字セグメントとで構成したときに、前記第1の言語の部分文字列が、前記第2の言語の部分文字列と翻字関係にない非翻字部分である確率を表す非翻字モデル選択確率と、前記第2の言語の部分文字列が、前記第1の言語の部分文字列と翻字関係にない非翻字部分である確率を表す非翻字モデル選択確率と、前記第1の言語の部分文字列が、前記第2の言語の部分文字列と翻字関係にある翻字部分であり、かつ前記第2の言語の部分文字列が、前記第1の言語の部分文字列と翻字関係にある翻字部分である確率を表す翻字モデル選択確率と、前記第1の言語の文字列のうちの前記前置非翻字セグメントの部分文字列、及び前記後置非翻字セグメントの部分文字列の各々に対する前記第1の言語における生成確率を表す非翻字モデル生成確率と、前記第2の言語の文字列のうちの前記前置非翻字セグメントの部分文字列、及び前記後置非翻字セグメントの部分文字列の各々に対する前記第2の

10

20

言語における生成確率を表す非翻字モデル生成確率と、前記第 1 の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第 2 の言語の文字列のうちの前記翻字セグメントの部分文字列との間の部分文字列の各ペアに対する同時生成確率を表す翻字モデル生成確率と、に基づいて尤もらしくなるように、前記文字列組の各文字列を前記前置非翻字セグメント、前記翻字セグメント、及び後置非翻字セグメントで構成し、かつ、前記第 1 の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第 2 の言語の文字列のうちの前記翻字セグメントの部分文字列との間で文字の対応付けを行う対応付け計算部

を含み、

前記対応付け計算部は、

前記第 1 の言語の前記非翻字モデル選択確率と、前記第 2 の言語の前記非翻字モデル選択確率と、前記第 2 の言語の各部分文字列に対する前記翻字モデル選択確率と、前記第 1 の言語の各部分文字列に対する前記非翻字モデル生成確率と、前記第 2 の言語の各部分文字列に対する前記非翻字モデル生成確率と、前記第 1 の言語の部分文字列と前記第 2 の言語の部分文字列との間の部分文字列の各ペアに対する前記翻字モデル生成確率と、に対して初期値を各々設定する初期値設定部と、

前記初期値設定部によって設定され、又は前回更新された、前記非翻字モデル選択確率、前記翻字モデル選択確率、前記非翻字モデル生成確率、及び前記翻字モデル生成確率に基づいて、前記文字列組の各々に対して、前記第 1 の言語の文字列のうちの部分文字列と、前記第 2 の言語の文字列のうちの部分文字列との間の部分文字列の各ペアについて、前記ペアが翻訳関係にある期待値を計算し、前記第 1 の言語の文字列のうちの各部分文字列について、前記部分文字列が非翻字部分である期待値を計算し、前記第 2 の言語の文字列のうちの各部分文字列について、前記部分文字列が非翻字部分である期待値を計算する期待値計算部と、

前記文字列組の各々に対して前記期待値計算部によって計算された各ペアに対する前記翻訳関係にある期待値、前記第 1 の言語の各部分文字列についての前記非翻字部分である期待値、及び前記第 2 の言語の各部分文字列についての前記非翻字部分である期待値に基づいて、前記非翻字モデル選択確率、前記翻字モデル選択確率、前記非翻字モデル生成確率、及び前記翻字モデル生成確率を更新するパラメータ更新部と、

予め定められた停止条件が満たされたか否かを判定し、前記停止条件が満たされるまで、前記期待値計算部による計算、及び前記パラメータ更新部による更新を繰り返す停止条件判定部と、を含み

前記文字列組の各々に対して、前記非翻字モデル選択確率、前記翻字モデル選択確率、前記非翻字モデル生成確率、及び前記翻字モデル生成確率の各々に基づいて、前記文字列組の各文字列を前記前置非翻字セグメント、前記翻字セグメント、及び後置非翻字セグメントで構成し、かつ、前記第 1 の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第 2 の言語の文字列のうちの前記翻字セグメントの部分文字列との間で文字の対応付けを行う文字列対応付け装置。

#### 【請求項 2】

異なる第 1 の言語及び第 2 の言語にそれぞれ属する同じ意味の文字列の組み合わせである文字列組について、前記第 1 の言語の文字列と、前記第 2 の言語の文字列との間で文字の対応付けを行う文字列対応付け装置であって、

前記文字列組を複数組記憶した文字列組データベースに記憶された前記文字列組の各々に対して、前記文字列組の各文字列を、前記文字列の先頭から順番に、他方の言語の部分文字列と翻字関係にない 0 文字以上の部分文字列を示す前置非翻字セグメントと、前記他方の言語の部分文字列と翻字関係にある 0 文字以上の部分文字列を示す翻字セグメントと、前記他方の言語の部分文字列と翻字関係にない 0 文字以上の部分文字列を示す後置非翻字セグメントとで構成したときに、前記第 1 の言語の部分文字列が、前記第 2 の言語の部分文字列と翻字関係にない非翻字部分である確率を表す非翻字モデル選択確率と、前記第 2 の言語の部分文字列が、前記第 1 の言語の部分文字列と翻字関係にない非翻字部分である確率を表す非翻字モデル選択確率と、前記第 1 の言語の部分文字列が、前記第 2 の言語

10

20

30

40

50

の部分文字列と翻字関係にある翻字部分であり、かつ前記第 2 の言語の部分文字列が、前記第 1 の言語の部分文字列と翻字関係にある翻字部分である確率を表す翻字モデル選択確率と、前記第 1 の言語の文字列のうちの前記前置非翻字セグメントの部分文字列、及び前記後置非翻字セグメントの部分文字列の各々に対する前記第 1 の言語における生成確率を表す非翻字モデル生成確率と、前記第 2 の言語の文字列のうちの前記前置非翻字セグメントの部分文字列、及び前記後置非翻字セグメントの部分文字列の各々に対する前記第 2 の言語における生成確率を表す非翻字モデル生成確率と、前記第 1 の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第 2 の言語の文字列のうちの前記翻字セグメントの部分文字列との間の部分文字列の各ペアに対する同時生成確率を表す翻字モデル生成確率と、に基づいて尤もらしくなるように、前記文字列組の各文字列を前記前置非翻字セグメント、前記翻字セグメント、及び後置非翻字セグメントで構成し、かつ、前記第 1 の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第 2 の言語の文字列のうちの前記翻字セグメントの部分文字列との間で文字の対応付けを行う対応付け計算部

10

を含み、

前記対応付け計算部は、

前記文字列組の各々に対して、前記文字列組の各文字列を前記前置非翻字セグメント、前記翻字セグメント、及び後置非翻字セグメントで構成し、かつ、前記第 1 の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第 2 の言語の文字列のうちの前記翻字セグメントの部分文字列との間で文字の対応付けを行って、対応付けの初期設定を行う初期対応設定部と、

20

前記初期対応設定部によって設定され、又は前回更新された、前記複数組の文字列組のうちの処理対象の文字列組以外の文字列組の各々についての前記対応付けに基づいて、前記第 1 の言語の前記非翻字モデル選択確率と、前記第 2 の言語の前記非翻字モデル選択確率と、前記翻字モデル選択確率と、前記処理対象の文字列組の前記第 1 の言語の文字列のうちの各部分文字列に対する前記非翻字モデル生成確率と、前記処理対象の文字列組の前記第 2 の言語の文字列のうちの各部分文字列に対する前記非翻字モデル生成確率と、前記処理対象の文字列組の前記第 1 の言語の文字列のうちの部分文字列と、前記第 2 の言語の文字列のうちの部分文字列との間の部分文字列の各ペアに対する前記翻字モデル生成確率と、を計算するフィルタリング部と、

前記フィルタリング部によって計算された前記非翻字モデル選択確率、前記翻字モデル選択確率、前記非翻字モデル生成確率、及び前記翻字モデル生成確率に基づいて、前記処理対象の文字列組の各文字列を前記前置非翻字セグメント、前記翻字セグメント、及び後置非翻字セグメントで構成し、かつ、前記第 1 の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第 2 の言語の文字列のうちの前記翻字セグメントの部分文字列との間で文字の対応付けを行って、前記処理対象の文字列組に対する前記対応付けを更新するサンプリング部と、

30

予め定められた停止条件が満たされたか否かを判定し、前記停止条件が満たされるまで、各文字列組を処理対象とした前記フィルタリング部による計算及び前記サンプリング部による更新を繰り返す停止条件判定部と、を含む文字列対応付け装置。

### 【請求項 3】

40

対応付け計算部を含み、異なる第 1 の言語及び第 2 の言語にそれぞれ属する同じ意味の文字列の組み合わせである文字列組について、前記第 1 の言語の文字列と、前記第 2 の言語の文字列との間で文字の対応付けを行う文字列対応付け装置における文字列対応付け方法であって、

対応付け計算部によって、前記文字列組を複数組記憶した文字列組データベースに記憶された前記文字列組の各々に対して、前記文字列組の各文字列を、前記文字列の先頭から順番に、他方の言語の部分文字列と翻字関係にない 0 文字以上の部分文字列を示す前置非翻字セグメントと、前記他方の言語の部分文字列と翻字関係にある 0 文字以上の部分文字列を示す翻字セグメントと、前記他方の言語の部分文字列と翻字関係にない 0 文字以上の部分文字列を示す後置非翻字セグメントとで構成したときに、前記第 1 の言語の部分文字

50

列が、前記第2の言語の部分文字列と翻字関係にない非翻字部分である確率を表す非翻字モデル選択確率と、前記第2の言語の部分文字列が、前記第1の言語の部分文字列と翻字関係にない非翻字部分である確率を表す非翻字モデル選択確率と、前記第1の言語の部分文字列が、前記第2の言語の部分文字列と翻字関係にある翻字部分であり、かつ前記第2の言語の部分文字列が、前記第1の言語の部分文字列と翻字関係にある翻字部分である確率を表す翻字モデル選択確率と、前記第1の言語の文字列のうちの前記前置非翻字セグメントの部分文字列、及び前記後置非翻字セグメントの部分文字列の各々に対する前記第1の言語における生成確率を表す非翻字モデル生成確率と、前記第2の言語の文字列のうちの前記前置非翻字セグメントの部分文字列、及び前記後置非翻字セグメントの部分文字列の各々に対する前記第2の言語における生成確率を表す非翻字モデル生成確率と、前記第1の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第2の言語の文字列のうちの前記翻字セグメントの部分文字列との間の部分文字列の各ペアに対する同時生成確率を表す翻字モデル生成確率と、に基づいて尤もらしくなるように、前記文字列組の各文字列を前記前置非翻字セグメント、前記翻字セグメント、及び後置非翻字セグメントで構成し、かつ、前記第1の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第2の言語の文字列のうちの前記翻字セグメントの部分文字列との間で文字の対応付けを行い、

10

前記対応付け計算部は、初期値設定部と、期待値計算部と、パラメータ更新部と、停止条件判定部と、を含み

前記初期値設定部によって、前記第1の言語の前記非翻字モデル選択確率と、前記第2の言語の前記非翻字モデル選択確率と、前記第2の言語の各部分文字列に対する前記翻字モデル選択確率と、前記第1の言語の各部分文字列に対する前記非翻字モデル生成確率と、前記第2の言語の各部分文字列に対する前記非翻字モデル生成確率と、前記第1の言語の部分文字列と前記第2の言語の部分文字列との間の部分文字列の各ペアに対する前記翻字モデル生成確率と、に対して初期値を各々設定するステップと、

20

前記期待値計算部によって、前記初期値設定部によって設定され、又は前回更新された、前記非翻字モデル選択確率、前記翻字モデル選択確率、前記非翻字モデル生成確率、及び前記翻字モデル生成確率に基づいて、前記文字列組の各々に対して、前記第1の言語の文字列のうちの部分文字列と、前記第2の言語の文字列のうちの部分文字列との間の部分文字列の各ペアについて、前記ペアが翻字関係にある期待値を計算し、前記第1の言語の文字列のうちの各部分文字列について、前記部分文字列が非翻字部分である期待値を計算し、前記第2の言語の文字列のうちの各部分文字列について、前記部分文字列が非翻字部分である期待値を計算するステップと、

30

前記パラメータ更新部によって、前記文字列組の各々に対して前記期待値計算部によって計算された各ペアに対する前記翻字関係にある期待値、前記第1の言語の各部分文字列についての前記非翻字部分である期待値、及び前記第2の言語の各部分文字列についての前記非翻字部分である期待値に基づいて、前記非翻字モデル選択確率、前記翻字モデル選択確率、前記非翻字モデル生成確率、及び前記翻字モデル生成確率を更新するステップと、

停止条件判定部によって、予め定められた停止条件が満たされたか否かを判定し、前記停止条件が満たされるまで、前記期待値計算部による計算、及び前記パラメータ更新部による更新を繰り返すステップと、を含み

40

前記文字列組の各々に対して、前記非翻字モデル選択確率、前記翻字モデル選択確率、前記非翻字モデル生成確率、及び前記翻字モデル生成確率の各々に基づいて、前記文字列組の各文字列を前記前置非翻字セグメント、前記翻字セグメント、及び後置非翻字セグメントで構成し、かつ、前記第1の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第2の言語の文字列のうちの前記翻字セグメントの部分文字列との間で文字の対応付けを行う文字列対応付け方法。

#### 【請求項4】

対応付け計算部を含み、異なる第1の言語及び第2の言語にそれぞれ属する同じ意味の

50

文字列の組み合わせである文字列組について、前記第1の言語の文字列と、前記第2の言語の文字列との間で文字の対応付けを行う文字列対応付け装置における文字列対応付け方法であって、

対応付け計算部によって、前記文字列組を複数組記憶した文字列組データベースに記憶された前記文字列組の各々に対して、前記文字列組の各文字列を、前記文字列の先頭から順番に、他方の言語の部分文字列と翻字関係にない0文字以上の部分文字列を示す前置非翻字セグメントと、前記他方の言語の部分文字列と翻字関係にある0文字以上の部分文字列を示す翻字セグメントと、前記他方の言語の部分文字列と翻字関係にない0文字以上の部分文字列を示す後置非翻字セグメントとで構成したときに、前記第1の言語の部分文字列が、前記第2の言語の部分文字列と翻字関係にない非翻字部分である確率を表す非翻字モデル選択確率と、前記第2の言語の部分文字列が、前記第1の言語の部分文字列と翻字関係にない非翻字部分である確率を表す非翻字モデル選択確率と、前記第1の言語の部分文字列が、前記第2の言語の部分文字列と翻字関係にある翻字部分であり、かつ前記第2の言語の部分文字列が、前記第1の言語の部分文字列と翻字関係にある翻字部分である確率を表す翻字モデル選択確率と、前記第1の言語の文字列のうちの前記前置非翻字セグメントの部分文字列、及び前記後置非翻字セグメントの部分文字列の各々に対する前記第1の言語における生成確率を表す非翻字モデル生成確率と、前記第2の言語の文字列のうちの前記前置非翻字セグメントの部分文字列、及び前記後置非翻字セグメントの部分文字列の各々に対する前記第2の言語における生成確率を表す非翻字モデル生成確率と、前記第1の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第2の言語の文字列のうちの前記翻字セグメントの部分文字列との間の部分文字列の各ペアに対する同時生成確率を表す翻字モデル生成確率と、に基づいて尤もらしくなるように、前記文字列組の各文字列を前記前置非翻字セグメント、前記翻字セグメント、及び後置非翻字セグメントで構成し、かつ、前記第1の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第2の言語の文字列のうちの前記翻字セグメントの部分文字列との間で文字の対応付けを行い、

前記対応付け計算部は、初期対応設定部と、フィルタリング部と、サンプリング部と、停止条件判定部と、を含み、

前記初期対応設定部によって、前記文字列組の各々に対して、前記文字列組の各文字列を前記前置非翻字セグメント、前記翻字セグメント、及び後置非翻字セグメントで構成し、かつ、前記第1の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第2の言語の文字列のうちの前記翻字セグメントの部分文字列との間で文字の対応付けを行って、対応付けの初期設定を行うステップと、

前記フィルタリング部によって、前記初期対応設定部によって設定され、又は前回更新された、前記複数組の文字列組のうちの処理対象の文字列組以外の文字列組の各々についての前記対応付けに基づいて、前記第1の言語の前記非翻字モデル選択確率と、前記第2の言語の前記非翻字モデル選択確率と、前記翻字モデル選択確率と、前記処理対象の文字列組の前記第1の言語の文字列のうちの各部分文字列に対する前記非翻字モデル生成確率と、前記処理対象の文字列組の前記第2の言語の文字列のうちの各部分文字列に対する前記非翻字モデル生成確率と、前記処理対象の文字列組の前記第1の言語の文字列のうちの部分文字列と、前記第2の言語の文字列のうちの部分文字列との間の部分文字列の各ペアに対する前記翻字モデル生成確率と、を計算するステップと、

前記サンプリング部によって、前記フィルタリング部によって計算された前記非翻字モデル選択確率、前記翻字モデル選択確率、前記非翻字モデル生成確率、及び前記翻字モデル生成確率に基づいて、前記処理対象の文字列組の各文字列を前記前置非翻字セグメント、前記翻字セグメント、及び後置非翻字セグメントで構成し、かつ、前記第1の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第2の言語の文字列のうちの前記翻字セグメントの部分文字列との間で文字の対応付けを行って、前記処理対象の文字列組に対する前記対応付けを更新するステップと、

前記停止条件判定部によって、予め定められた停止条件が満たされたか否かを判定し、

10

20

30

40

50

前記停止条件が満たされるまで、各文字列組を処理対象とした前記フィルタリング部による計算及び前記サンプリング部による更新を繰り返すステップと、を含む文字列対応付け方法。

【請求項 5】

コンピュータを、請求項 1 又は請求項 2 に記載の文字列対応付け装置の各部として機能させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、文字列対応付け装置、方法、及びプログラムに係り、特に、異なる言語の文字列の組における文字の対応付けを行う文字列対応付け装置、方法、及びプログラムに関する。

10

【背景技術】

【0002】

ある言語の音韻体系で表記された語句を別の言語の音韻体系での表記に変換する機械翻字を、統計モデルとして表現するために、互いが対応する単語の組を統計モデルの学習のためのデータとして利用して、単語を構成する文字同士の対応関係を推定することが広く行われている。例えば、非特許文献 1 では、英語の音韻表現と日本語におけるカタカナ語のローマ字化された表記との間での音韻記号-ローマ字間の 1 対多の対応付け方法について記している。さらに、非特許文献 2 では、英語の文字と音韻表記との多対多の対応付け

20

【0003】

一方で、統計モデルの学習に利用する単語の組を大量に収集しようとする、ある程度の誤りの混入は避けられない。いわゆる「カタカナ語」と英語の対応で言えば、日英対訳辞書の項目において日本語側がカタカナで表記されているものでも、「コンピュータ」と "computer" のように翻字関係となっているものもあれば、「カブトムシ」と "beetle" のように、カタカナで表記されるが翻字関係とはなっていないものもある。こうした誤った単語対応を統計モデルの学習に利用することでノイズが混入し、統計モデルの質を低下することは避けるべきである。この問題に対して、翻字関係となっている文字間対応の統計モデルと、翻字関係となっておらず 2 言語間で独立な文字列の統計モデルを利用した翻字対応付け方法が提案されており、有効に働くことが示されている（例えば、非特許文献 4）。対訳辞書の存在を仮定しない非特許文献 5 のような「統計的機械翻訳」と呼ばれる技術分野においては、対訳文中の共起関係等を用いて自動的に単語対応を得ている。この自動的に得られた単語対応から翻字対応関係を得ようとする、単語対応に誤りが含まれる可能性も高くなるが、非特許文献 4 の方法により、1 対 1 の単語対応組を、翻字となっている単語組と翻字となっていない単語組を自動的に分類し、翻字となっている単語組からのみ翻字対応の統計モデルを学習することが可能となる。また、非特許文献 6 には、上記非特許文献 3 の技術を、翻字でない文字列に対応させた場合について記載されている。

30

【先行技術文献】

40

【非特許文献】

【0004】

【非特許文献 1】 Kevin Knight and Jonathan Graehl, 「Machine Transliteration」、Computational Linguistics、1998、Volume 24、Number 4、p.599-612

【非特許文献 2】 Sittichai Jiampojarn 他 2 名, 「Applying Many-to-Many Alignments and Hidden Markov Models to Letter-to-Phoneme Conversion」、Proceedings of NAAC LHLT、2007、p.372-379

【非特許文献 3】 Andrew Finch and Eiichiro Sumita, 「A Bayesian Model of Bilingual Segmentation for Transliteration」、Proceedings of International Workshop on Spoken Language Translation、2010

50

【非特許文献4】Hassan Sajjad他2名、「A Statistical Model for Unsupervised and Semisupervised Transliteration Mining」、Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics、2012、p.469-477

【非特許文献5】Philipp Koehn他2名、「Statistical Phrase-Based Translation」、Proceedings of HLT-NAACL、2003、p.48-54

【非特許文献6】Ohnmar Htun他3名、「Improving Transliteration Mining by Integrating Expert Knowledge with Statistical Approaches」、International Journal of Computer Applications、2012、58(17)、p.12-22

【発明の概要】

【発明が解決しようとする課題】

10

【0005】

非特許文献4の方法は単語対応が「翻字である」か「翻字でない」かを識別するため、複合語等、複数の単語の間の対応関係において「部分的には翻字であるが、その他の部分は翻字になっていない」ような場合を識別することができない。非特許文献5のような句に基づく統計的機械翻訳においては、2言語間の句と句の対訳関係を自動的に推定するため、部分的に訳語でないものが含まれることがあり、非特許文献4の方法では十分な識別を行うことが期待できず、そこから得られる翻字対応および翻字モデルの正確性に問題が生ずる。例えば、「コンピュータ」に対して”the computer”という句が対応しているという状況においては、“the”に対応する文字列がカタカナ語側にはないため、「翻字である」か「翻字でない」かの2値分類は適さない。

20

【0006】

本発明は、上記の事情を鑑みてなされたもので、異なる言語の文字列の組における文字の対応付けを精度よく行うことができる文字列対応付け装置、方法、及びプログラムを提供することを目的とする。

【課題を解決するための手段】

【0007】

上記の目的を達成するために本発明に係る文字列対応付け装置は、異なる第1の言語及び第2の言語にそれぞれ属する同じ意味の文字列の組み合わせである文字列組について、前記第1の言語の文字列と、前記第2の言語の文字列との間で文字の対応付けを行う文字列対応付け装置であって、前記文字列組を複数組記憶した文字列組データベースに記憶された前記文字列組の各々に対して、前記文字列組の各文字列を、前記文字列の先頭から順番に、他方の言語の部分文字列と翻字関係にない0文字以上の部分文字列を示す前置非翻字セグメントと、前記他方の言語の部分文字列と翻字関係にある0文字以上の部分文字列を示す翻字セグメントと、前記他方の言語の部分文字列と翻字関係にない0文字以上の部分文字列を示す後置非翻字セグメントとで構成したときに、前記第1の言語の部分文字列が、前記第2の言語の部分文字列と翻字関係にない非翻字部分である確率を表す非翻字モデル選択確率と、前記第2の言語の部分文字列が、前記第1の言語の部分文字列と翻字関係にない非翻字部分である確率を表す非翻字モデル選択確率と、前記第1の言語の部分文字列が、前記第2の言語の部分文字列と翻字関係にある翻字部分であり、かつ前記第2の言語の部分文字列が、前記第1の言語の部分文字列と翻字関係にある翻字部分である確率を表す翻字モデル選択確率と、前記第1の言語の文字列のうちの前記前置非翻字セグメントの部分文字列、及び前記後置非翻字セグメントの部分文字列の各々に対する前記第1の言語における生成確率を表す非翻字モデル生成確率と、前記第2の言語の文字列のうちの前記前置非翻字セグメントの部分文字列、及び前記後置非翻字セグメントの部分文字列の各々に対する前記第2の言語における生成確率を表す非翻字モデル生成確率と、前記第1の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第2の言語の文字列のうちの前記翻字セグメントの部分文字列との間の部分文字列の各ペアに対する同時生成確率を表す翻字モデル生成確率と、に基づいて尤もらしくなるように、前記文字列組の各文字列を前記前置非翻字セグメント、前記翻字セグメント、及び後置非翻字セグメントで構成し、かつ、前記第1の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記

30

40

50

第2の言語の文字列のうちの前記翻字セグメントの部分文字列との間で文字の対応付けを行う対応付け計算部を含んで構成されている。

【0008】

本発明に係る文字列対応付け方法は、対応付け計算部を含み、異なる第1の言語及び第2の言語にそれぞれ属する同じ意味の文字列の組み合わせである文字列組について、前記第1の言語の文字列と、前記第2の言語の文字列との間で文字の対応付けを行う文字列対応付け装置における文字列対応付け方法であって、対応付け計算部によって、前記文字列組を複数組記憶した文字列組データベースに記憶された前記文字列組の各々に対して、前記文字列組の各文字列を、前記文字列の先頭から順番に、他方の言語の部分文字列と翻字関係にない0文字以上の部分文字列を示す前置非翻字セグメントと、前記他方の言語の部分文字列と翻字関係にある0文字以上の部分文字列を示す翻字セグメントと、前記他方の言語の部分文字列と翻字関係にない0文字以上の部分文字列を示す後置非翻字セグメントとで構成したときに、前記第1の言語の部分文字列が、前記第2の言語の部分文字列と翻字関係にない非翻字部分である確率を表す非翻字モデル選択確率と、前記第2の言語の部分文字列が、前記第1の言語の部分文字列と翻字関係にない非翻字部分である確率を表す非翻字モデル選択確率と、前記第1の言語の部分文字列が、前記第2の言語の部分文字列と翻字関係にある翻字部分であり、かつ前記第2の言語の部分文字列が、前記第1の言語の部分文字列と翻字関係にある翻字部分である確率を表す翻字モデル選択確率と、前記第1の言語の文字列のうちの前記前置非翻字セグメントの部分文字列、及び前記後置非翻字セグメントの部分文字列の各々に対する前記第1の言語における生成確率を表す非翻字モデル生成確率と、前記第2の言語の文字列のうちの前記前置非翻字セグメントの部分文字列、及び前記後置非翻字セグメントの部分文字列の各々に対する前記第2の言語における生成確率を表す非翻字モデル生成確率と、前記第1の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第2の言語の文字列のうちの前記翻字セグメントの部分文字列との間の部分文字列の各ペアに対する同時生成確率を表す翻字モデル生成確率と、に基づいて尤もらしくなるように、前記文字列組の各文字列を前記前置非翻字セグメント、前記翻字セグメント、及び後置非翻字セグメントで構成し、かつ、前記第1の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第2の言語の文字列のうちの前記翻字セグメントの部分文字列との間で文字の対応付けを行うステップを含む。

【0009】

また、本発明に係る前記対応付け計算部は、前記第1の言語の前記非翻字モデル選択確率と、前記第2の言語の前記非翻字モデル選択確率と、前記第2の言語の各部分文字列に対する前記翻字モデル選択確率と、前記第1の言語の各部分文字列に対する前記非翻字モデル生成確率と、前記第2の言語の各部分文字列に対する前記非翻字モデル生成確率と、前記第1の言語の部分文字列と前記第2の言語の部分文字列との間の部分文字列の各ペアに対する前記翻字モデル生成確率と、に対して初期値を各々設定する初期値設定部と、前記初期値設定部によって設定され、又は前回更新された、前記非翻字モデル選択確率、前記翻字モデル選択確率、前記非翻字モデル生成確率、及び前記翻字モデル生成確率に基づいて、前記文字列組の各々に対して、前記第1の言語の文字列のうちの部分文字列と、前記第2の言語の文字列のうちの部分文字列との間の部分文字列の各ペアについて、前記ペアが翻字関係にある期待値を計算し、前記第1の言語の文字列のうちの各部分文字列について、前記部分文字列が非翻字部分である期待値を計算し、前記第2の言語の文字列のうちの各部分文字列について、前記部分文字列が非翻字部分である期待値を計算する期待値計算部と、前記文字列組の各々に対して前記期待値計算部によって計算された各ペアに対する前記翻字関係にある期待値、前記第1の言語の各部分文字列についての前記非翻字部分である期待値、及び前記第2の言語の各部分文字列についての前記非翻字部分である期待値に基づいて、前記非翻字モデル選択確率、前記翻字モデル選択確率、前記非翻字モデル生成確率、及び前記翻字モデル生成確率を更新するパラメータ更新部と、予め定められた停止条件が満たされたか否かを判定し、前記停止条件が満たされるまで、前記期待値計算部による計算、及び前記パラメータ更新部による更新を繰り返す停止条件判定部と、を

10

20

30

40

50



含み前記文字列組の各々に対して、前記非翻字モデル選択確率、前記翻字モデル選択確率、前記非翻字モデル生成確率、及び前記翻字モデル生成確率の各々に基づいて、前記文字列組の各文字列を前記前置非翻字セグメント、前記翻字セグメント、及び後置非翻字セグメントで構成し、かつ、前記第1の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第2の言語の文字列のうちの前記翻字セグメントの部分文字列との間で文字の対応付けを行うようにすることができる。

【0010】

また、本発明に係る前記対応付け計算部は、前記文字列組の各々に対して、前記文字列組の各文字列を前記前置非翻字セグメント、前記翻字セグメント、及び後置非翻字セグメントで構成し、かつ、前記第1の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第2の言語の文字列のうちの前記翻字セグメントの部分文字列との間で文字の対応付けを行って、対応付けの初期設定を行う初期対応設定部と、前記初期対応設定部によって設定され、又は前回更新された、前記複数組の文字列組のうちの処理対象の文字列組以外の文字列組の各々についての前記対応付けに基づいて、前記第1の言語の前記非翻字モデル選択確率と、前記第2の言語の前記非翻字モデル選択確率と、前記翻字モデル選択確率と、前記処理対象の文字列組の前記第1の言語の文字列のうちの各部分文字列に対する前記非翻字モデル生成確率と、前記処理対象の文字列組の前記第2の言語の文字列のうちの各部分文字列に対する前記非翻字モデル生成確率と、前記処理対象の文字列組の前記第1の言語の文字列のうちの部分文字列と、前記第2の言語の文字列のうちの部分文字列との間の部分文字列の各ペアに対する前記翻字モデル生成確率と、を計算するフィルタリング部と、前記フィルタリング部によって計算された前記非翻字モデル選択確率、前記翻字モデル選択確率、前記非翻字モデル生成確率、及び前記翻字モデル生成確率に基づいて、前記処理対象の文字列組の各文字列を前記前置非翻字セグメント、前記翻字セグメント、及び後置非翻字セグメントで構成し、かつ、前記第1の言語の文字列のうちの前記翻字セグメントの部分文字列と、前記第2の言語の文字列のうちの前記翻字セグメントの部分文字列との間で文字の対応付けを行って、前記処理対象の文字列組に対する前記対応付けを更新するサンプリング部と、予め定められた停止条件が満たされたか否かを判定し、前記停止条件が満たされるまで、各文字列組を処理対象とした前記フィルタリング部による計算及び前記サンプリング部による更新を繰り返す停止条件判定部と、を含むようにすることができる。

【0011】

本発明に係るプログラムは、コンピュータを、本発明に係る文字列対応付け装置の各部として機能させるためのプログラムである。

【発明の効果】

【0012】

以上説明したように、本発明の文字列対応付け装置、方法、及びプログラムによれば、異なる第1の言語及び第2の言語にそれぞれ属する同じ意味の文字列の組み合わせである文字列組について、文字列組の各文字列を、文字列の先頭から順番に、前置非翻字セグメントと、翻字セグメントと、後置非翻字セグメントとで構成したときに、翻字モデル選択確率と、非翻字モデル選択確率と、非翻字モデル生成確率と、翻字モデル生成確率と、に基づいて尤もらしくなるように、文字列組の各文字列を前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントで構成し、かつ、第1言語の文字列のうちの翻字セグメントの部分文字列と、第2言語の文字列のうちの翻字セグメントの部分文字列との間の文字を対応付けることにより、異なる言語の文字列組における文字の対応付けを精度よく行うことができる、という効果が得られる。

【図面の簡単な説明】

【0013】

【図1】本発明の第1の実施の形態に係る文字列対応付け装置の構成を示す概略図である。

【図2】本発明の第1の実施の形態に係る文字列対応付け装置における文字列対応付け処

理ルーチンの内容を示すフローチャートである。

【図3】本発明の第2の実施の形態に係る文字列対応付け装置の構成を示す概略図である。

【図4】本発明の第2の実施の形態に係る文字列対応付け装置における文字列対応付け処理ルーチンの内容を示すフローチャートである。

【図5】対応付けの一例として、本実施の形態に係る文字列対応付け装置に入力された英語とカタカナとを示す図である。

【図6】英語とカタカナの対応付け結果の例を示す図である。

【発明を実施するための形態】

【0014】

以下、図面を参照して本発明の実施の形態を詳細に説明する。

【0015】

<発明の概要>

本発明の実施の形態は、対応付けの対象となる文字列組について、「翻字である」か「翻字でない」かの区別を部分文字列の単位で識別する。また、単純に部分文字列単位での識別を行うと文字列組の中で「翻字である」部分と「翻字でない」部分が頻繁に入れ替わってしまうことがあり文字列対応方法として適さないため、「翻字である」部分は文字列組において高々1箇所であるという制約を加える。具体的には、上記非特許文献1記載の多対多の文字列対応付け方法において、文字列の生成モデルとして「翻字モデル」「非翻字モデル」の2つを同時に学習することで、文字列組の翻字となっている部分を識別し、非翻字部分が混在したデータにおける適切な文字列間の対応付けを実現する。

【0016】

本発明の実施の形態では、上記非特許文献3や上記非特許文献4と同様の「翻字である」文字列の確率モデル（以下、「翻字モデル」と称する）と、「翻字でない」文字列の確率モデル（以下、「非翻字モデル」と称する）の2種類のモデルの存在を仮定する。翻字モデルは原言語と目的言語の部分文字列の同時確率  $P(s, t)$  ( $s$  は原言語の部分文字列、 $t$  は目的言語の部分文字列) のモデルであり、非翻字モデルは「原言語の部分文字列の確率  $P(s)$  のモデル」と「目的言語の部分文字列の確率  $P(t)$  のモデル」との独立した2つの確率モデルを含んで構成される。

【0017】

本発明の実施の形態では、上記非特許文献4や、上記非特許文献6のように、学習に用いられる原言語と目的言語の文字列組が、文字列全体として翻字であるか、翻字でないかを区別するのではなく、文字列のどの部分が翻字であり、どの部分が翻字でないかを区別するように確率モデルを学習する。また、翻字である部分は高々1箇所であると仮定する。つまり、原言語と目的言語の文字列の各々は、

【0018】

- ・「0文字以上の翻字でない部分」（以下、前置非翻字セグメントと称する。）
- ・「0文字以上の翻字である部分」（以下、翻字セグメントと称する。）
- ・「0文字以上の翻字でない部分」（以下、後置非翻字セグメントと称する。）

【0019】

の順で構成されると仮定する。すなわち、文字列組の各文字列を、文字列の先頭から順番に、他方の言語の部分文字列と翻字関係のない0文字以上の部分文字列を示す前置非翻字セグメントと、他方の言語の部分文字列と翻字関係にある0文字以上の部分文字列を示す翻字セグメントと、他方の言語の部分文字列と翻字関係のない0文字以上の部分文字列を示す後置非翻字セグメントとで構成した場合を想定する。

なお、文字列全体が翻字でない場合は、すべてが前置非翻字セグメントに属するものとする。本実施の形態における文字列の対応付けは、原言語と目的言語の文字列組を構成する文字列を上記セグメントに分割し、非翻字セグメントにおいては原言語と目的言語で独立な非翻字モデルに基づいて部分文字列が生成され、翻字セグメントにおいては原言語と目的言語との翻字モデルに基づいて部分文字列の組が生成される場合の尤度が最大となる

10

20

30

40

50

ような対応付けを求める過程である。翻字モデルや非翻字モデル自体の構成及びモデル最適化アルゴリズムは特に規定しないが、翻字モデルは、原言語の 0 文字以上の文字列と目的言語の 0 文字以上の文字列の組の同時生成確率（ただし双方とも 0 文字となる場合は除く）モデルであり、非翻字モデルは原言語と目的言語で独立な、1 文字以上の文字列の生成確率モデルであるとする。

【 0 0 2 0 】

また、ある文字列が翻字モデルと非翻字モデルとのどちらから生成されるかを表す確率変数（以下、「翻字モデル選択確率」、及び「非翻字モデル選択確率」と称する）も同時に考慮する。従って、上記の前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントで文字列が構成されることを前提とすると、原言語の非翻字モデル選択確率は、原言語の部分文字列が、目的言語の部分文字列と翻字関係にない非翻字部分である確率を表しており、目的言語の非翻字モデル選択確率は、目的言語の部分文字列が、原言語の部分文字列と翻字関係にない非翻字部分である確率を表している。

10

また、翻字モデル選択確率は、原言語の部分文字列が、目的言語の部分文字列と翻字関係にある翻字部分であり、かつ目的言語の部分文字列が、原言語の部分文字列と翻字関係にある翻字部分である確率を表している。

【 0 0 2 1 】

モデルの最適化においては、上記非特許文献 4 で用いられている前向き後向き（forward-backward）アルゴリズムに基づく期待値最大化（EM）アルゴリズムや、その変種といえる、上記特許文献 1 で用いられている、ギブス（Gibbs）サンプリングに基づく前向きフィルタリング・後ろ向きサンプリングアルゴリズムにおいて、上記の 3 セグメントでの構成を考慮した上で、非翻字セグメントでは非翻字モデル生成確率、翻字セグメントでは翻字モデル生成確率を利用して期待値を計算し、その結果を利用してモデルの更新を行う過程を繰り返し行えばよい。

20

ここで、原言語の部分文字列の非翻字モデル生成確率は、原言語の前置非翻字セグメント又は後置非翻字セグメントにおける部分文字列の、原言語における生成確率を表し、目的言語の部分文字列の非翻字モデル生成確率は、目的言語の前置非翻字セグメント又は後置非翻字セグメントにおける部分文字列の、目的言語における生成確率を表す。

また、翻字モデル生成確率は、原言語の文字列のうちの翻字セグメントの部分文字列と、目的言語の文字列のうちの翻字セグメントの部分文字列との間の部分文字列の各ペアに対する同時生成確率を表している。

30

【 0 0 2 2 】

〔第 1 の実施の形態〕

<システム構成>

本発明の第 1 の実施の形態に係る文字列対応付け装置 1 0 0 は、原言語（第 1 の言語）の文字列（単語）と目的言語（第 2 の言語）の文字列（単語）との対訳である複数組の文字列組を入力とし、文字列組の各々について、文字列の対応付けを行う。この文字列対応付け装置 1 0 0 は、CPU と、RAM と、後述する文字列対応付け処理ルーチンを実行するためのプログラムを記憶した ROM とを備えたコンピュータで構成され、機能的には次に示すように構成されている。図 1 に示すように、文字列対応付け装置 1 0 0 は、入力部 1 0 と、演算部 2 0 と、出力部 3 0 とを備えている。

40

【 0 0 2 3 】

入力部 1 0 は、文字列の対応付けを行う対象である複数組の文字列組を受け付ける。具体的には、入力部 1 0 は、翻字又は対訳になっていることが期待され、かつ空白文字等も含んだ文字列組を、入力装置、記憶媒体もしくはネットワークを通じて複数組読み込む。

【 0 0 2 4 】

演算部 2 0 は、文字列組データベース 2 1、文字列組データ読み込み部 2 2、及び対応付け計算部 2 3 を備えている。

【 0 0 2 5 】

文字列組データベース 2 1 には、入力部 1 0 により受け付けた複数組の文字列組が格納

50

される。

【 0 0 2 6 】

文字列組データ読み込み部 2 2 は、文字列組データベース 2 1 から全ての文字列組を読み込む。

【 0 0 2 7 】

対応付け計算部 2 3 は、文字列組毎に、原言語の文字列と、目的言語の文字列との間で部分文字列同士の対応付けを行う。

【 0 0 2 8 】

対応付け計算部 2 3 は、初期値設定部 2 3 1、期待値計算部 2 3 2、パラメータ更新部 2 3 3、停止判定部 2 3 4、及び文字対応付け処理部 2 3 5 を備えている。

10

【 0 0 2 9 】

初期値設定部 2 3 1 は、原言語の部分文字列に対する非翻字モデル選択確率と、目的言語の部分文字列に対する非翻字モデル選択確率と、原言語の部分文字列と目的言語の部分文字列とのペアに対する翻字モデル選択確率と、原言語の文字列のうちの各部分文字列に対する非翻字モデル生成確率と、目的言語の文字列のうちの各部分文字列に対する非翻字モデル生成確率と、原言語の文字列のうちの部分文字列と、目的言語の文字列のうちの部分文字列との間の部分文字列の各ペアに対する翻字モデル生成確率と、に対して初期値を設定する。ここで、初期値は、出現頻度や共起頻度の比を用いて設定することが広く行われているが、一様分布としてもよい。

なお、対応付けの単位として上記非特許文献 4 のように 1 文字単位のものしか考慮しない場合は、非翻字モデル生成確率は単純な文字ユニグラム確率となるため、平滑化を考慮しなければ単純な文字の出現頻度に基づく出現確率分布に従って求められ、以後の処理において確率を更新せず、固定するようによい。

20

また、翻字モデル選択確率・非翻字モデル選択確率については、上記非特許文献 4 ではそれぞれ 0 . 5 を初期値としているが、0 より大きく、和が 1 となるような任意の初期値を設定してもよい。

【 0 0 3 0 】

期待値計算部 2 3 2 は、初期値設定部 2 3 1 によって設定され、又は前回更新された、非翻字モデル選択確率、翻字モデル選択確率、非翻字モデル生成確率、及び翻字モデル生成確率に基づいて、文字列組の各々に対して、原言語の文字列のうちの部分文字列と、目的言語の文字列のうちの部分文字列との間の部分文字列の各ペアについて、当該ペアが翻訳関係にある期待値を計算し、原言語の文字列のうちの各部分文字列について、当該部分文字列が非翻字部分である期待値を計算し、目的言語の文字列のうちの各部分文字列について、当該部分文字列が非翻字部分である期待値を計算する。期待値の計算方法は上記非特許文献 4 に記載の方法に類似するが、本実施の形態では、文字列に、前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントの順で、各セグメントが割り当てられて構成されるため、翻字セグメントのみが存在することを仮定している上記非特許文献 4 の方法を拡張した方法を用いる。

30

【 0 0 3 1 】

上記非特許文献 4 に記載の方法で用いられている前向き後ろ向きアルゴリズムでは、動的計画法を用いて、文字列組の先頭から対応付けの確率を記録する前向き確率表と、逆に文字列組の末尾から対応付けの確率を記録する後向き確率表とに、計算した確率を順次記録するが、本実施の形態では、前置非翻字セグメント、翻字セグメント、後置非翻字セグメントのそれぞれに対応する前向き確率表と後向き確率表の、合計 6 個の確率表を利用する。前向き確率及び後向き確率の計算と記録とに際しては、前置非翻字セグメント、翻字セグメント、後置非翻字セグメントの順を考慮する。前置非翻字セグメントの前向き確率表に記録される確率の計算にあたっては、そこより前の位置における前置非翻字セグメントの前向き確率のみを参照し、翻字セグメントの前向き確率表に記録される確率の計算にあたっては、そこより前の位置における前置非翻字セグメントの前向き確率及び翻字セグメントの前向き確率を参照し、後置非翻字セグメントの前向き確率表に記録される確率の

40

50

計算にあたっては、そこより前の位置における翻字セグメントの前向き確率及び後置非翻字セグメントの前向き確率を参照する。

【 0 0 3 2 】

後向き確率表についてはこれと逆の順序となるため、後置非翻字セグメントの後向き確率表に記録される確率の計算にあたっては、そこより後の位置における後置非翻字セグメントの後向き確率のみを参照し、翻字セグメントの後向き確率表に記録される確率の計算にあたっては、そこより後の位置における後置非翻字セグメントの後向き確率及び翻字セグメントの後向き確率を参照し、前置非翻字セグメントの後向き確率表に記録される確率の計算にあたっては、そこより後の位置における翻字セグメントの後向き確率及び前置非翻字セグメントの後向き確率を参照する。また、確率表の初期値については、前置非翻字セグメントの前向き確率表の先頭と、後置非翻字セグメントの後向き確率表の末尾を1とし、それ以外を0とする。

10

【 0 0 3 3 】

期待値計算部 2 3 2 は、初期値設定部 2 3 1 によって設定され、又は前回更新された、非翻字モデル選択確率、翻字モデル選択確率、非翻字モデル生成確率、及び翻字モデル生成確率に基づいて、上記のように、前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントのそれぞれに対応する前向き確率表の各前向き確率を計算して記録すると共に、前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントのそれぞれに対応する後ろ向き確率表の各後ろ向き確率を計算して記録する。なお、前向き及び後ろ向き確率表は、文字列組ごとに計算される。

20

期待値計算部 2 3 2 は、これらの前向き確率表と後ろ向き確率表を用いて、文字列組中の任意の文字対応付けについて翻訳関係にある期待値を計算する。原言語側の文字位置  $j$  から  $j' - 1$  までの文字列と目的言語側の文字位置  $i$  から  $i' - 1$  までの文字列が翻字となる期待値  $E_t(j, j', i, i')$  は、以下のように計算できる（なお、文字列位置は 0 から始まるものとする）。

【 0 0 3 4 】

【数 1】

$$E_t(s_{j \rightarrow j'}, t_{i \rightarrow i'}) = \frac{(P_{f1}(j, i) + P_{f2}(j, i)) \times p_t \times P_t(s_{j \rightarrow j'}, t_{i \rightarrow i'}) \times (P_{b2}(j', i') + P_{b3}(j', i'))}{P_{f1}(J, I) + P_{f2}(J, I) + P_{f3}(J, I)} \dots (1)$$

30

【 0 0 3 5 】

ここで、 $P_t(s_{j \rightarrow j'}, t_{i \rightarrow i'})$  は原言語側の文字位置  $j$  から  $j' - 1$  までの文字列と目的言語側の文字位置  $i$  から  $i' - 1$  までの文字列の同時生成確率（翻字モデルの確率）、 $p_t$  は翻字モデル選択確率、 $P_{f1}(j, i)$ 、 $P_{f2}(j, i)$ 、 $P_{f3}(j, i)$  はそれぞれ原言語側の文字位置  $j$ 、目的言語側の文字位置  $i$  までの前置非翻字セグメント、翻字セグメント、後置非翻字セグメントに対する前向き確率、 $P_{b1}(j, i)$ 、 $P_{b2}(j, i)$ 、 $P_{b3}(j, i)$  はそれぞれ原言語側の文字位置  $j$ 、目的言語側の文字位置  $i$  より後の前置非翻字セグメント、翻字セグメント、後置非翻字セグメントに対する後ろ向き確率、 $J$  は原言語側の文字列長、 $I$  は目的言語側の文字列長である。

40

期待値計算部 2 3 2 は、全ての文字列組  $d$  毎に、文字列組  $d$  の部分文字列の各ペアについて、上記 (1) 式の期待値を計算する。

【 0 0 3 6 】

期待値計算部 2 3 2 は、前向き確率表と後ろ向き確率表を用いて、非翻字部分の期待値も

50

同様に計算する。具体的には、原言語側の文字位置  $j$  から  $j' - 1$  までの部分文字列が、目的言語側の文字位置  $i$  の後で目的言語側から（翻字ではなく）独立に生成される期待値  $E_{src}(s_{j \rightarrow j'}, i)$  を、以下の式のように前置非翻字セグメントで現れる場合と後置非翻字セグメントで現れる場合の期待値の和として計算する。期待値計算部 232 は、前向き確率表と後向き確率表を用いて、目的言語側についても同様に、目的言語側の文字位置  $i$  から  $i' - 1$  までの部分文字列が、原言語側の文字位置  $j$  の後で原言語側から（翻字ではなく）独立に生成される期待値  $E_{trg}(s_{i \rightarrow i'}, j)$  を計算する。

【0037】

【数2】

$$E_{src}(s_{j \rightarrow j'}, i) = \frac{P_{f1}(j, i) \times p_{src} \times P_{src}(s_{j \rightarrow j'}) \times (P_{b1}(j', i') + P_{b2}(j', i'))}{P_{f1}(J, I) + P_{f2}(J, I) + P_{f3}(J, I)} + \frac{P_{f3}(j, i) \times p_{src} \times P_{src}(s_{j \rightarrow j'}) \times P_{b3}(j', i')}{P_{f1}(J, I) + P_{f2}(J, I) + P_{f3}(J, I)} \quad \dots (2)$$

10

20

【0038】

ここで、 $P_{src}(s_{j \rightarrow j'})$  は原言語側の文字位置  $j$  から  $j' - 1$  までの部分文字列の生成確率（原言語側の非翻字モデルの確率）、 $p_{src}$  は原言語側の非翻字モデル選択確率である。

【0039】

期待値計算部 232 は、全ての文字列組  $d$  毎に、文字列組  $d$  の原言語側の部分文字列の各々について、上記(2)式の期待値を計算する。また、期待値計算部 232 は、全ての文字列組  $d$  毎に、文字列組  $d$  の目的言語側の部分文字列の各々について、上記(2)式と同様の期待値を計算する。

期待値計算部 232 は、全ての文字列組  $d$  に対する計算結果に基づいて、原言語及び目的言語の部分文字列のペア  $(s, t)$  の各々について、当該ペア  $(s, t)$  に関する期待値の計算結果を集計して当該ペア  $(s, t)$  に対する翻字モデルからの同時生成の期待値  $E_t(s, t)$  を計算する。期待値計算部 232 は、原言語の各部分文字列  $s$  について、当該部分文字列  $s$  に関する期待値の計算結果を集計して当該部分文字列  $s$  に対する非翻字モデルからの生成の期待値  $E_{src}(s)$  を計算する。また、期待値計算部 232 は、目的言語の各部分文字列  $t$  について、当該部分文字列  $t$  に関する期待値の計算結果を集計して当該部分文字列  $t$  の非翻字モデルからの生成の期待値  $E_{trg}(t)$  を計算する。

30

また、上記の確率表および期待値の計算において、前置非翻字セグメント、後置非翻字セグメント内では、必ず原言語側が先に生成されるものとし、同じ文字列の生成に際して複数回の数え上げが起こらないようにする。

40

【0040】

パラメータ更新部 233 は、文字列組の各々に対して、期待値計算部 232 によって計算された各ペアに対する翻訳関係にある期待値、原言語の各部分文字列についての非翻字部分である期待値、及び目的言語の各部分文字列についての非翻字部分である期待値に基づいて、非翻字モデル選択確率、翻字モデル選択確率、非翻字モデル生成確率、及び翻字モデル生成確率を更新する。具体的には、パラメータ更新部 233 は、期待値計算部 232 において計算された期待値をもとに、翻字モデルの同時生成確率、非翻字モデルの生成確率、及び翻字モデル選択確率・非翻字モデル選択確率を更新する。翻字モデルの同時生成確率及び非翻字モデルの生成確率については、期待値が 0 より大きい文字列組もしくは文字列に対して、期待値の総和との比をとって更新後の同時生成確率又は生成確率とする

50

。つまり、翻字モデルの同時生成確率については、原言語の部分文字列  $s$  と目的言語の部分文字列  $t$  の同時生成確率  $P_t(s, t)$  を以下のように更新する。

【0041】

【数3】

$$P_t(\bar{s}, \bar{t}) \leftarrow \frac{E_t(\bar{s}, \bar{t})}{\sum_{s,t} E_t(s, t)} \quad \dots (3)$$

10

【0042】

非翻字モデルの生成確率については、原言語側の部分文字列  $s$  の生成確率  $P_{src}(s)$  を以下のように更新する。目的言語側の部分文字列  $t$  の生成確率  $P_{trg}(t)$  も同様に更新する。

【0043】

【数4】

$$P_{src}(\bar{s}) \leftarrow \frac{E_{src}(\bar{s})}{\sum_s E_{src}(s)} \quad \dots (4)$$

20

【0044】

翻字モデル選択確率・非翻字モデル選択確率については、原言語と目的言語でそれぞれ翻字セグメントに属する文字数・非翻字セグメントに属する文字数の割合に基づくと考え、原言語の非翻字モデル選択確率は非翻字セグメントで現れる原言語の文字数の割合、目的言語の非翻字モデル選択確率は非翻字セグメントで現れる目的言語の文字数の割合とする。つまり、部分文字列の期待値を利用して計算すると、原言語の非翻字モデル選択確率  $p_{src}$  は以下のように更新される。 $|s|$  は部分文字列  $s$  の長さ（文字数）を表す。目的言語の非翻字モデル選択確率  $p_{trg}$  についても同様に更新される。

30

【0045】

【数5】

$$p_{src} \leftarrow \frac{\sum_s |s| + E_{src}(s)}{\sum_s |s| \times E_{src}(s) + \sum_{s,t} |s| \times E_t(s, t)} \quad \dots (5)$$

40

【0046】

そして、翻字モデル選択確率  $p_t$  は、以下の式に示すように、 $(1 - \text{原言語の非翻字モデル選択確率}) \times (1 - \text{目的言語の非翻字モデル選択確率})$  で更新される。

【0047】

【数6】

$$p_t \leftarrow (1 - p_{src}) \times (1 - p_{trg}) \quad \dots (6)$$

50

## 【 0 0 4 8 】

停止判定部 2 3 4 は、予め定められた停止条件が満たされたか否かを判定し、当該停止条件が満たされるまで、期待値計算部 2 3 2 による計算、及びパラメータ更新部 2 3 3 による更新を繰り返す。停止条件としては、文字列組データの尤度が一定以上の数値になった、尤度の変動幅が一定以下の数値になった、パラメータ更新の繰り返し回数が一定の回数を超えた、などが考えられる。本実施の形態では、文字列組データの尤度が一定以上の数値になることを、停止条件とする。

## 【 0 0 4 9 】

文字対応付け処理部 2 3 5 は、文字列組の各々に対して、パラメータ更新部 2 3 3 により最終的に更新された非翻字モデル選択確率、翻字モデル選択確率、非翻字モデル生成確率、及び翻字モデル生成確率の各々に基づいて、当該文字列組の各文字列を、前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントで構成し、かつ、原言語の文字列のうちの翻字セグメントの部分文字列と、目的言語の文字列のうちの翻字セグメントの部分文字列との間で文字の対応付けを行う。具体的には、文字対応付け処理部 2 3 5 は、得られた最終的な非翻字モデル選択確率、翻字モデル選択確率、非翻字モデル生成確率、及び翻字モデル生成確率を用い、当該非翻字モデル選択確率、翻字モデル選択確率、非翻字モデル生成確率、及び翻字モデル生成確率の下で最尤な文字対応付けをビタビアルゴリズムによって求める。ビタビアルゴリズムは当該分野で広く知られた、動的計画法によって最尤な系列を求めるアルゴリズムである。具体的には、期待値計算部 2 3 2 における前向き確率表の確率計算とは異なり、前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントの各々の前向き確率表に対して、最も確率の高い経路の確率およびその経路の情報を保持する形で末尾まで計算を行い、その結果得られた確率表を末尾から先頭に向かって辿っていくことで最尤経路を得ることができる。経路が翻字セグメントの前向き確率表内にあるものは該当する部分が翻字の関係にあり、非翻字セグメントの前向き確率表内にあるものは翻字の関係にないことを表す。

## 【 0 0 5 0 】

出力部 3 0 は、文字対応付け処理部 2 3 5 で対応付けられた文字列組の各々の対応付けの情報を出力する。具体的には、出力部 3 0 は、対応付けの情報を付与した文字列組を端末に表示、もしくは記憶媒体やネットワークを通じて書き出す。

## 【 0 0 5 1 】

< 文字列対応付け装置の作用 >

次に、第 1 の実施の形態に係る文字列対応付け装置 1 0 0 の作用について説明する。まず、対訳となっている第 1 の言語体系の文字列及び第 2 の言語体系の文字列の組である文字列組が、文字列対応付け装置 1 0 0 に複数入力されると、文字列対応付け装置 1 0 0 によって、入力された複数の文字列組が、文字列組データベース 2 1 に格納される。そして、文字列対応付け装置 1 0 0 によって、図 2 に示す文字列対応付け処理ルーチンが実行される。

## 【 0 0 5 2 】

まず、ステップ S 1 0 0 において、文字列組データ読み込み部 2 2 によって、文字列組データベース 2 1 から、全ての文字列組を読み込む。

## 【 0 0 5 3 】

ステップ S 1 0 2 において、初期値設定部 2 3 1 によって、原言語及び目的言語の各々の非翻字モデル選択確率と、翻字モデル選択確率と、原言語及び目的言語の各々の各部分文字列に対する非翻字モデル生成確率と、原言語及び目的言語の部分文字列の各ペアに対する翻字モデル生成確率と、に対して初期値を設定する。

## 【 0 0 5 4 】

次に、ステップ S 1 0 4 において、期待値計算部 2 3 2 によって、全ての文字列組の各々について、周知の前向き後ろ向きアルゴリズムを用いて、上記ステップ S 1 0 2 で設定され、又は後述するステップ S 1 1 0 で前回更新された、非翻字モデル選択確率、翻字モ

10

20

30

40

50



デル選択確率、非翻字モデル生成確率、及び翻字モデル生成確率に基づいて、前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントのそれぞれに対応する前向き確率表の各前向き確率を計算して前向き確率表に記録すると共に、前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントのそれぞれに対応する後ろ向き確率表の各後ろ向き確率を計算して後ろ向き確率表に記録する。

【 0 0 5 5 】

ステップ S 1 0 6 において、期待値計算部 2 3 2 によって、全ての文字列組の各々について、上記ステップ S 1 0 4 で計算された当該文字列組の前向き確率表と後ろ向き確率表を用いて、当該文字列組の部分文字列の各ペアに対し、翻字関係にある期待値を計算する。

10

【 0 0 5 6 】

ステップ S 1 0 8 において、期待値計算部 2 3 2 によって、全ての文字列組の各々について、上記ステップ S 1 0 4 で計算された当該文字列組の前向き確率表と後ろ向き確率表を用いて、当該文字列組の原言語側の各部分文字列に対し、非翻字部分となる期待値を計算すると共に、当該文字列組の目的言語側の各部分文字列に対し、非翻字部分となる期待値を計算する。

【 0 0 5 7 】

ステップ S 1 1 0 において、パラメータ更新部 2 3 3 は、上記ステップ S 1 0 6 で計算された翻訳関係にある期待値、上記ステップ S 1 0 8 で計算された原言語の各部分文字列についての非翻字部分である期待値、及び目的言語の各部分文字列についての非翻字部分である期待値に基づいて、原言語及び目的言語の各々の非翻字モデル選択確率、翻字モデル選択確率、原言語の各部分文字列に対する非翻字モデル生成確率、目的言語の各部分文字列に対する非翻字モデル生成確率、及び原言語及び目的言語の部分文字列の各ペアに対する翻字モデル生成確率を更新する。

20

【 0 0 5 8 】

ステップ S 1 1 4 において、停止判定部 2 3 4 によって、予め定められた停止条件が満たされたか否かを判定する。そして、停止条件が満たされていない場合には、上記ステップ S 1 0 4 へ戻り、上記ステップ S 1 0 4 ~ ステップ S 1 1 0 の処理を実行する。一方、停止条件が満たされている場合には、ステップ S 1 1 6 へ進む。

【 0 0 5 9 】

ステップ S 1 1 6 において、文字対応付け処理部 2 3 5 によって、文字列組の各々に対して、上記ステップ S 1 1 0 で最終的に更新された非翻字モデル選択確率、翻字モデル選択確率、非翻字モデル生成確率、及び翻字モデル生成確率の各々に基づいて、当該文字列組の各文字列を、前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントで構成し、かつ、原言語の文字列のうちの翻字セグメントの部分文字列と、目的言語の文字列のうちの翻字セグメントの部分文字列との間で文字の対応付けを行う。

30

【 0 0 6 0 】

ステップ S 1 1 8 において、出力部 3 0 によって、上記ステップ S 1 1 6 で対応付けた結果を出力して、文字列対応付け処理ルーチンを終了する。

【 0 0 6 1 】

以上説明したように、本発明の第 1 の実施の形態に係る文字列対応付け装置によれば、原言語及び目的言語において対訳となる文字列の組み合わせである文字列組について、文字列組の各文字列を、文字列の先頭から順番に、前置非翻字セグメントと、翻字セグメントと、後置非翻字セグメントとで構成したときに、翻字モデル選択確率と、非翻字モデル選択確率と、非翻字モデル生成確率と、翻字モデル生成確率と、に基づいて尤もらしくなるように、文字列組の各文字列を前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントで構成し、かつ、原言語の文字列のうちの翻字セグメントの部分文字列と、原言語の文字列のうちの翻字セグメントの部分文字列との間で文字の対応付けを行うことにより、異なる言語の文字列組における文字の対応付けを精度よく行うことができる。

40

【 0 0 6 2 】

50

また、文字列の生成モデルとして「翻字モデル」「非翻字モデル」の2つを同時に学習することで、文字列の組の翻字となっている部分を識別し、非翻字部分が混在したデータにおける適切な文字列間の対応付けを実現することができる。

【0063】

また、統計的機械翻訳で用いられる自動的に抽出された対訳語句対のような、必ずしも翻字の組となっておらず、部分的に翻字となっているような文字列組データに対しても、翻字となっている文字列を選択的に対応づけることが可能となる。

【0064】

〔第2の実施の形態〕

<システム構成>

次に、第2の実施の形態について説明する。なお、第1の実施の形態と同様の構成となる部分については、同一符号を付して説明を省略する。

【0065】

第2の実施の形態では、ギブスサンプリングを用いて、文字列組における文字の対応付けを求めている点が、第1の実施の形態と異なっている。

【0066】

本実施の形態で用いるギブスサンプリングは、上記非特許文献3や上記非特許文献6で用いられており、当該分野では広く知られる方法である。

上記非特許文献3及び非特許文献6に記載の方法のように、参考文献7 (Daichi Mochi hashi他2名、「Bayesian Unsupervised Word Segmentation with Nested Pitman-Yor Language Modeling」、Proc. of ACL-IJCNLP、2009)に記載の方法に類似したアルゴリズムによって対応付けを繰り返し更新する。

【0067】

この方法では、文字列組の集合Dから1つの文字列組dを選び、dを除く文字列組の集合D-dにおける対応付けから推定される事後確率分布に基づいてd上の対応付け結果をサンプリングする、という過程を、dを入れ替えながら繰り返し行う。サンプリングの結果はdに対する一意の対応付け結果であるため、EMアルゴリズムを利用する第1の実施の形態の構成とは異なり、停止判定部によって学習が終了したと判定された時点ですべてのdに対する対応付け結果が得られるため、文字対応付け処理部235は必要ない。ただし、停止判定後に1度フィルタリングステップ・サンプリングステップを繰り返すことによって対応付けを再度更新してもよい。

【0068】

図3に示すように、第2の実施の形態に係る文字列対応付け装置200の演算部220は、文字列組データベース21、文字列組データ読み込み部22、及び対応付け計算部24を備えている。

【0069】

対応付け計算部24は、ギブスサンプリングを用いて、複数の文字列組の各々に対して対応付けを行う。対応付け計算部24は、初期値対応部241、フィルタリング部242、サンプリング部243、及び停止判定部244を備えている。

【0070】

初期値対応部241は、文字列組の各々に対して、文字列組の各文字列を前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントで構成し、かつ、原言語の文字列のうちの翻字セグメントの部分文字列と、目的言語の文字列のうちの翻字セグメントの部分文字列との間で文字の対応付けを行って、対応付けの初期設定を行う。ここで、ギブスサンプリングを用いる場合は、確率分布の初期値でなく、適当な初期対応付けを与える必要がある。従って、初期対応付の決定方法としては、何かの規則により対応付けを行う、ランダムに対応付けを行う、別途簡便なモデルでの対応付け結果を利用するなどの方法が考えられる。本実施の形態では、ランダムに対応付けを行う。

【0071】

フィルタリング部242は、フィルタリング部242によって前回計算された、原言語

10

20

30

40

50

の非翻字モデル選択確率と、目的言語の非翻字モデル選択確率と、翻字モデル選択確率と、処理対象の文字列組の原言語の文字列のうちの各部分文字列に対する非翻字モデル生成確率と、処理対象の文字列組の目的言語の文字列のうちの各部分文字列に対する非翻字モデル生成確率と、処理対象の文字列組の原言語の文字列のうちの部分文字列と、目的言語の文字列のうちの部分文字列との間の部分文字列の各ペアに対する翻字モデル生成確率とに基づいて、上記第1の実施の形態における前向き確率の計算の場合と同様の方法で、処理対象の文字列組dに対して、文字列組の先頭から前向き確率を計算し、前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントのそれぞれに対応する前向き確率表に、計算した各前向き確率を記録する。ここで、確率が非常に小さい経路を大量に保持すると計算量が大きくなる場合があるため、他の対立する対応付け経路と比較して非常に確率が小さいような場合はその経路を無視してもよい。

10

また、フィルタリング部242は、初期対応設定部241によって設定され、又は後述するサンプリング部243によって前回更新された、複数組の文字列組のうちの処理対象の文字列組以外の文字列組の各々についての対応付けに基づいて、原言語の非翻字モデル選択確率と、目的言語の非翻字モデル選択確率と、翻字モデル選択確率と、処理対象の文字列組の原言語の文字列のうちの各部分文字列に対する非翻字モデル生成確率と、処理対象の文字列組の目的言語の文字列のうちの各部分文字列に対する非翻字モデル生成確率と、処理対象の文字列組の原言語の文字列のうちの部分文字列と、目的言語の文字列のうちの部分文字列との間の部分文字列の各ペアに対する翻字モデル生成確率と、を計算する。翻字モデル生成確率 $P_{gt}$ (EMアルゴリズムによる構成の場合の $P_t$ に相当)は、上記

20

【0072】

【数7】

$$P_{gt}(s_{j \rightarrow j'}, t_{i \rightarrow i'}) = p_t \times \frac{c(s_{j \rightarrow j'}, t_{i \rightarrow i'}) + \alpha BM(s_{j \rightarrow j'}, t_{i \rightarrow i'})}{C + \alpha} \quad \dots (7)$$

【0073】

30

ここで、 $c(s_{j \rightarrow j'}, t_{i \rightarrow i'})$ は、原言語側の文字列 $s_{j \rightarrow j'}$ と目的言語側の文字列 $t_{i \rightarrow i'}$ がdを除く文字組列の集合で対応付けられている回数、 $BM(s_{j \rightarrow j'}, t_{i \rightarrow i'})$ は基底測度、 $\alpha$ はハイパーパラメータ、 $C$ はdを除く文字列組データでの翻字となっているすべての対応付けの数である。非翻字モデルについても同様であり、原言語・目的言語の非翻字モデル生成確率も、同様に以下のように定義できる(EMアルゴリズムによる構成の場合の $P_{src}$ 、 $P_{trg}$ に相当)。

【0074】

【数8】

$$P_{gsrc}(s_{j \rightarrow j'}) = p_{src} \times \frac{c_{src}(s_{j \rightarrow j'}) + \alpha BM(s_{j \rightarrow j'})}{C_{src} + \alpha} \quad \dots (8)$$

40

$$P_{gtrg}(t_{i \rightarrow i'}) = p_{trg} \times \frac{c_{trg}(t_{i \rightarrow i'}) + \alpha BM(t_{i \rightarrow i'})}{C_{trg} + \alpha} \quad \dots (9)$$

【0075】

翻字モデル選択確率 $p_t$ 、非翻字モデル選択確率 $p_{src}$ 、 $p_{trg}$ は前記EMアルゴ

50

リズムの場合と同様、翻字になっている文字数と非翻字になっている文字数の割合で与えることができる。

【0076】

サンプリング部243は、フィルタリング部242によって計算された非翻字モデル選択確率、翻字モデル選択確率、非翻字モデル生成確率、及び翻字モデル生成確率に基づいて、処理対象の文字列組の各文字列を前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントで構成し、かつ、原言語の文字列のうちの翻字セグメントの部分文字列と、目的言語の文字列のうちの翻字セグメントの部分文字列との間で文字の対応付けを行って、処理対象の文字列組に対する対応付けを更新する。具体的には、サンプリング部243は、フィルタリング部242で計算された前向き確率表を末尾から順に辿っていき、前向き確率に基づいて対応付けを決定していき、dに対する対応付けを更新する。より具体的には、文字列組データ中のある位置 $(j, i)$ に対して、その位置に至る1つ前の位置 $(j', i')$ を選択する際に、 $(j', i')$ から $(j, i)$ に至る確率に基づく重み付きサンプリングを行う(確率が高い経路ほど選ばれやすいようにする)。 $(j', i')$ から $(j, i)$ に至る確率は、 $(j', i')$ における前向き確率と $(j', i')$ から $(j, i)$ への経路の確率 $(s_{j', j}, t_{i', i})$ の翻字あるいは非翻字確率)である。こうして決定された対応付けをdの新しい対応付けとする。

10

【0077】

全ての文字列組の各々を、処理対象の文字列組として、フィルタリング部242及びサンプリング部243による処理を繰り返し行う。

20

停止判定部244は、予め定められた停止条件が満たされたか否かを判定し、停止条件が満たされるまで、各文字列組を処理対象とした、フィルタリング部242による計算及びサンプリング部243による更新を繰り返す。本実施の形態では、ギブスサンプリングを利用しており、文字列組dを一つ選択するごとにフィルタリング部242によるフィルタリング処理と、サンプリング部243によるサンプリング処理とを行うので、文字列組の集合Dの全てに対してフィルタリング処理とサンプリング処理とを適用した後に停止判定を行う。ただし、途中で停止してもよい。停止条件はEMアルゴリズムを利用する第1の実施の形態の場合と同様に設定が可能である。

【0078】

<文字列対応付け装置の作用>

30

次に、第2の実施の形態に係る文字列対応付け装置200の作用について説明する。なお、第1の実施の形態と同様の処理については、同一符号を付して説明を省略する。

【0079】

まず、対訳となっている原言語の文字列及び目的言語の文字列の組である文字列組が、文字列対応付け装置200に複数入力されると、文字列対応付け装置200によって、入力された複数の文字列組が、文字列組データベース21に格納される。そして、文字列対応付け装置200によって、図4に示す文字列対応付け処理ルーチンが実行される。

【0080】

ステップS100において、文字列組データベース21から、全ての文字列組を読み込む。

40

【0081】

そして、ステップS202において、初期値対応部241によって、文字列組の各々に対して、文字列組の各文字列を前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントで構成し、かつ、原言語の文字列のうちの翻字セグメントの部分文字列と、目的言語の文字列のうちの翻字セグメントの部分文字列との間で文字の対応付けを行って、対応付けの初期設定を行う。

【0082】

次にステップS204において、フィルタリング部242によって、処理対象として、1つの文字組のデータdを設定する。

【0083】

50

ステップS 2 0 6において、フィルタリング部 2 4 2によって、上記ステップS 2 0 4で設定された処理対象の文字列組のデータdについて、本ステップで前回計算された、原言語の非翻字モデル選択確率と、目的言語の非翻字モデル選択確率と、翻字モデル選択確率と、処理対象の文字列組の原言語の文字列のうちの各部分文字列に対する非翻字モデル生成確率と、処理対象の文字列組の目的言語の文字列のうちの各部分文字列に対する非翻字モデル生成確率と、処理対象の文字列組の原言語の文字列のうちの部分文字列と、目的言語の文字列のうちの部分文字列との間の部分文字列の各ペアに対する翻字モデル生成確率とに基づいて、上記第1の実施の形態における前向き確率の計算の場合と同様の方法で、処理対象の文字列組dに対して、文字列組の先頭から前向き確率を計算し、前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントのそれぞれに対応する前向き確率表に、計算した各前向き確率を記録する。

10

## 【 0 0 8 4 】

ステップS 2 0 8において、フィルタリング部 2 4 2によって、上記ステップS 2 0 2で設定され、又は後述するステップS 2 1 0で前回更新された、複数組の文字列組のうちの処理対象の文字列組以外の文字列組の各々についての対応付けに基づいて、原言語及び目的言語の各々の非翻字モデル選択確率と、翻字モデル選択確率と、処理対象の文字列組の原言語及び目的言語の各々の、文字列のうちの各部分文字列に対する非翻字モデル生成確率と、処理対象の文字列組の目的言語の文字列のうちの各部分文字列に対する非翻字モデル生成確率と、翻字モデル生成確率と、を計算する。

## 【 0 0 8 5 】

20

ステップS 2 1 0において、サンプリング部 2 4 3によって、上記ステップS 2 0 8で計算された非翻字モデル選択確率、翻字モデル選択確率、非翻字モデル生成確率、及び翻字モデル生成確率に基づいて、処理対象の文字列組の各文字列を前置非翻字セグメント、翻字セグメント、及び後置非翻字セグメントで構成し、かつ、原言語の文字列のうちの翻字セグメントの部分文字列と、目的言語の文字列のうちの翻字セグメントの部分文字列との間で文字の対応付けを行って、処理対象の文字列組に対する対応付けを更新する。

## 【 0 0 8 6 】

ステップS 2 1 2において、文字列組の集合Dのうちの全ての文字列組dについて、上記ステップS 2 0 4～ステップS 2 1 0の処理を実行したか否かを判定する。そして、全ての文字列組dについて、上記ステップS 2 0 4～ステップS 2 1 0の処理を実行していない場合には、上記ステップS 2 0 4へ戻り、新たな文字列組dを処理対象として設定する。一方、全ての文字列組dについて、上記ステップS 2 0 4～ステップS 2 1 0の処理を実行した場合には、ステップS 2 1 4へ進む。

30

## 【 0 0 8 7 】

ステップS 2 1 4において、停止判定部 2 4 4によって、予め定められた停止条件が満たされたか否かを判定する。そして、停止条件が満たされていない場合には、上記ステップS 2 0 4へ戻り、再び、全ての文字列組dについて、上記ステップS 2 0 4～ステップS 2 1 2の処理を実行する。一方、停止条件が満たされている場合には、ステップS 2 1 6へ進む。

## 【 0 0 8 8 】

40

ステップS 2 1 6において、出力部 3 0によって、上記ステップS 2 1 0で最終的に対応付けた結果を出力して、文字列対応付け処理ルーチンを終了する。

## 【 0 0 8 9 】

なお、第2の実施の形態に係る文字列対応付け装置 2 0 0の他の構成及び作用については、第1の実施の形態と同様であるため、説明を省略する。

## 【 0 0 9 0 】

以上説明したように、第2の実施の形態に係る文字列対応付け装置によれば、ギブサンプリングにより、非翻字モデル選択確率と、翻字モデル選択確率と、非翻字モデル生成確率と、翻字モデル生成確率とを計算し、処理対象の文字列組に対する対応付けを更新することにより、異なる言語の文字列組における文字の対応付けを精度よく行うことができ

50

る。

【0091】

<実施例>

次に本発明に係る実施の形態を実施した例について示す。実施例では原言語として日本語（カタカナ語）、目的言語として英語を利用した。なお、本実施例ではギブスサンプリングを利用した第2の実施の形態によって、対応付けの計算を行う。

【0092】

図5はそれぞれ本実施例で用いた日本語と英語との文字列組データ（およそ5万語対）を抜粋したものである。同じ行に記された文字列組が対応していることを示す。各記号（カタカナ及びアルファベット）は1文字ずつ区切り文字（空白文字）によって分割されている。空白文字自体を対応付けの対象文字として扱いたい場合には、区切り文字として別の文字（例えば”：”）を利用するか、空白文字を別の記号（例えば”[s p]”）に置換すればよい。

10

【0093】

第2の実施の形態によって、上記図5に示した文字列組データに対して対応付けを行った。まず、上記文字列組データを記録した電子ファイルを文字列組データ読み込み部22により読み込んだ。また、カタカナとアルファベットの任意長さの文字列間の対応を考慮すると計算量が大きくなるため、本実施例においては、カタカナ（日本語）側は最大2文字、アルファベット（英語）側は最大3文字までの対応に限定することとした。

【0094】

本実施例における初期値対応部241では、ランダムな対応付けを与えるために、フィルタリング部242によるフィルタリング処理における前向き確率表の前向き確率の値をすべて1と仮定した上で（ただし、カタカナ側3文字以上、アルファベット側4文字以上となる対応付けは許さないようにする）、サンプリング部243によるサンプリング処理で利用されるサンプリングアルゴリズムを適用した。これにより、ランダムな初期対応付けを得ることができる。

20

【0095】

続いて、文字列組データの各文字列組に対して、フィルタリング処理とサンプリング処理の順に繰り返し計算を行う。各文字列組に対するフィルタリング処理・サンプリング処理を文字列組データのデータ数分行う処理を1ラウンドとして、本実施例では30ラウンドの処理を行った。また、英語の記号” ”は通常カタカナに対応するものではないため、常に日本語の長さ0の部分文字列に対応するものとして扱った。なお、ギブスサンプリングにおいてはデータを一個ずつ処理して対応付けとモデルの確率が更新されるため、データの処理順序が学習に与える影響が大きいことが広く知られている。したがって、文字列組データ中のデータの処理順序は毎ラウンド開始時にランダムに入れ替えることとした。

30

【0096】

30ラウンドの繰り返しの後、ビタビアルゴリズムによって再度文字列対応付けを行った。その結果を図6に示す。図6は日本語と英語の文字列がタブ文字で区切られて表記されており、対応付けられる部分文字列が記号”：”で区切られている。つまり、”：”の数は日本語側と英語側で同一となっており、1番目の要素同士が対応する部分文字列であることを示す。また、”<noise>”となっているのは、対応する部分文字列がない（ノイズである）ことを示している。

40

【0097】

図6の結果から、”コンピュータ”と”computers”の組では、英語側末尾の”s”がノイズであること、”バーン”と”burn in”の組では、英語側の”in”がノイズであること、”シンメトリー”と”asymmetry”との組では、英語側先頭の”a”がノイズであること、”サイド”と”side effect”の組では、英語側の”effect”がノイズであること、などが見て取れる。また、その他の文字列組についても、妥当な文字列対応付け結果が得られた。

50

## 【 0 0 9 8 】

なお、同じデータを上記非特許文献 4 の方法を実現するコンピュータプログラムによって対応付けした場合、及び上記非特許文献 6 の方法を実現するコンピュータプログラムによって対応付けした場合は、本実施例において部分的にノイズであると判定された文字列組は、すべてノイズであると判定されたため、本実施例は部分的にノイズが含まれるような文字列組の対応付けに適した方式であることが認められた。

## 【 0 0 9 9 】

なお、本発明は、上述した実施形態に限定されるものではなく、この発明の要旨を逸脱しない範囲内で様々な変形や応用が可能である。

## 【 0 1 0 0 】

例えば、文字列組データベース 2 1 は、外部に設けられ、文字列対応付け装置とネットワークで接続されていてもよい。

## 【 0 1 0 1 】

上述の文字列対応付け装置は、内部にコンピュータシステムを有しているが、「コンピュータシステム」は、WWWシステムを利用している場合であれば、ホームページ提供環境（あるいは表示環境）も含むものとする。

## 【 0 1 0 2 】

また、本願明細書中において、プログラムが予めインストールされている実施形態として説明したが、当該プログラムを、コンピュータ読み取り可能な記録媒体に格納して提供することも可能である。

## 【 符号の説明 】

## 【 0 1 0 3 】

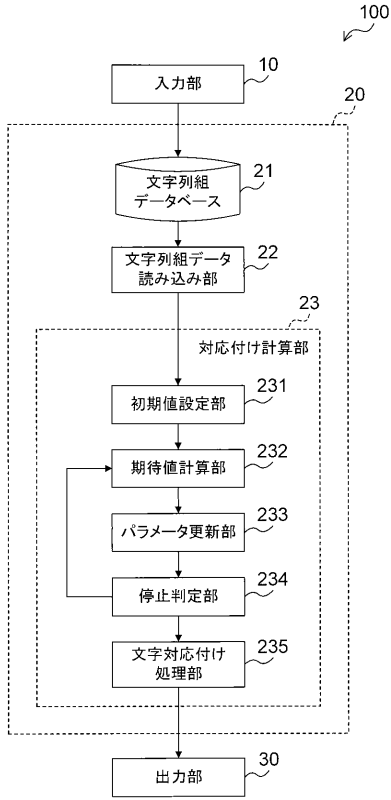
- 1 0 入力部
- 2 0、2 2 0 演算部
- 2 1 文字列組データベース
- 2 2 文字列組データ読み込み部
- 2 3、2 4 対応付け計算部
- 3 0 出力部
- 1 0 0、2 0 0 文字列対応付け装置
- 2 3 1 初期値設定部
- 2 3 2 期待値計算部
- 2 3 3 パラメータ更新部
- 2 3 4 停止判定部
- 2 3 5 文字対応付け処理部
- 2 4 1 初期値対応部
- 2 4 2 フィルタリング部
- 2 4 3 サンプリング部
- 2 4 4 停止判定部

10

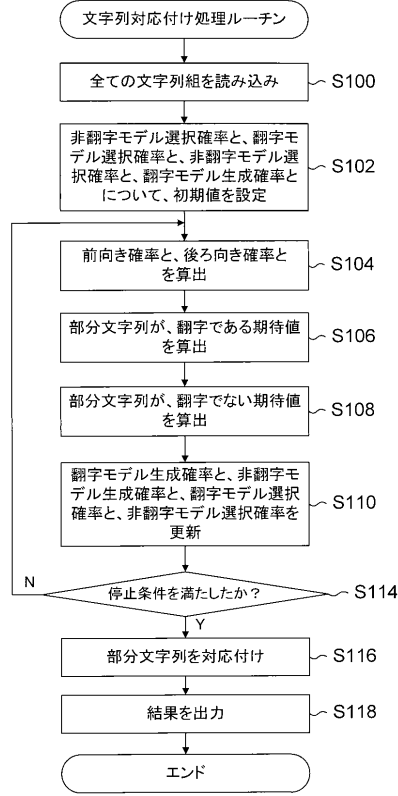
20

30

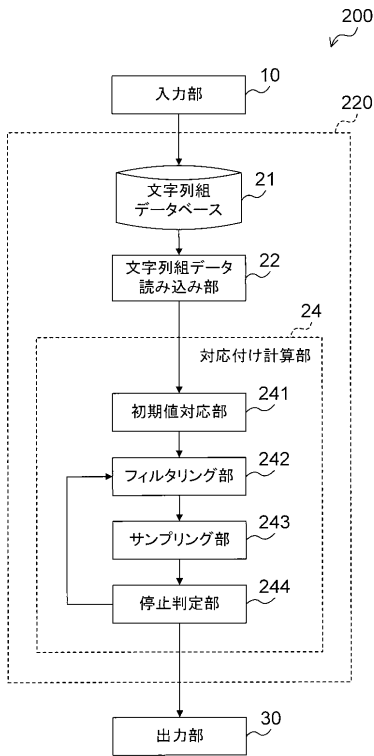
【図1】



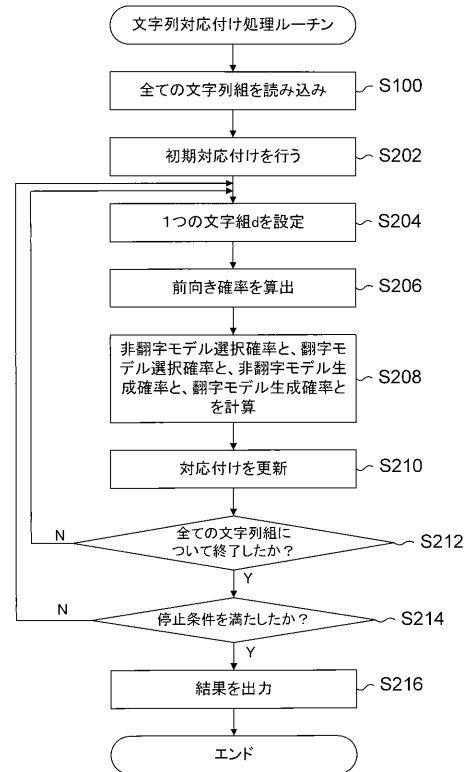
【図2】



【図3】



【図4】





【 図 5 】

(A) コンピュータ  
 バーン  
 シンメトリー  
 サイド  
 アーゴノミックス  
 アーゴノメトリックス  
 ダム  
 アーチ  
 リング  
 グループ

(B) computers  
 burn-in  
 asymmetry  
 side-effect  
 ergonomics  
 ergonometrics  
 dam  
 arch  
 ring  
 group

【 図 6 】

コン:ピ:ユ:ー:タ:<noise> c o:m p:u:t e:r:s  
 バ:ー:ン:<noise> b u:r n:i:-i:n  
 <noise>シ:ン:メ:ト:リ:ー a:s y:m:e:t r:y s i:d e:i:-e:f f:i:e c:t  
 サ:イ:ド:<noise><noise><noise><noise>  
 ア:ー:ゴ:ノ:ミ:ク:ス e:r:g o:n o:m i:c:s  
 ア:ー:ゴ:ノ:メ:ト:リ:ッ:ク:ス e:r:g o:n o:m e:t:r i:c:s  
 ダ:ム d a:m  
 ア:ー:チ a r:c h  
 リ:ン:グ r i:n:g  
 グ:ル:ー:ブ g:r:o u:p

---

フロントページの続き

(72)発明者 森 信介

京都府京都市左京区吉田本町3番地1 国立大学法人京都大学内

審査官 成瀬 博之

(56)参考文献 特開2001-142877(JP,A)

特開2007-156545(JP,A)

特開2012-185679(JP,A)

特開2003-263432(JP,A)

特開2005-092682(JP,A)

(58)調査した分野(Int.Cl., DB名)

G06F 17/27 - 17/28