

【公報種別】特許法第 17 条の 2 の規定による補正の掲載

【部門区分】第 6 部門第 3 区分

【発行日】平成23年6月30日(2011.6.30)

【公開番号】特開2006-48683(P2006-48683A)

【公開日】平成18年2月16日(2006.2.16)

【年通号数】公開・登録公報2006-007

【出願番号】特願2005-216527(P2005-216527)

【国際特許分類】

G 0 6 F 17/30 (2006.01)

【F I】

G 0 6 F 17/30 3 4 0 A

G 0 6 F 17/30 3 5 0 C

G 0 6 F 17/30 2 1 0 D

【誤訳訂正書】

【提出日】平成23年5月9日(2011.5.9)

【誤訳訂正 1】

【訂正対象書類名】明細書

【訂正対象項目名】0 0 0 9

【訂正方法】変更

【訂正の内容】

【0 0 0 9】

本システムは、更に、フレーズが文書内の他のフレーズの存在を予測する能力に基づいて、相互に関連するフレーズを識別するようになっている。より詳細には、2 フレーズの実際の共出現率を、その 2 フレーズの期待される共出現率と関連させる予測指標を用いる。実際の共出現率と期待される共出現率との比としての情報ゲインは、このような予測指標の 1 つである。予測指標が所定の閾値を超える場合に 2 フレーズを関連付ける。その場合、第 2 フレーズは、第 1 フレーズに対して大きな情報ゲインをもつ。意味的に、関連するフレーズは、「President of the United States」と「White House」等の所与の主題または概念を検討または説明するために普通に用いるものである。所与のフレーズに対して、それぞれの予測指標に基づく関連性または重要性に従って、関連フレーズをランク付けできる。

【誤訳訂正 2】

【訂正対象書類名】明細書

【訂正対象項目名】0 0 6 6

【訂正方法】変更

【訂正の内容】

【0 0 6 6】

文書コレクションを最終的に通過させることにより、見込フレーズリストは、大規模コーパスのフレーズ使用の期待される分布により、比較的短くなる。従って、例えば 10 回の通過（例えば、10, 000, 000 文書）の場合、フレーズは、まさに最初の回だけに出現し、その時点で良好フレーズとなる可能性はほとんどない。そのフレーズは使用されるようになったばかりの新規フレーズかもしれず、従って、後続のクロールをしている間に次第に共通となる。この場合、それぞれのカウントは増大し、最終的に良好フレーズとなる閾値を満たすことになる。

【誤訳訂正 3】

【訂正対象書類名】明細書

【訂正対象項目名】0 0 6 8

【訂正方法】変更

【訂正の内容】

【0068】

上記のように、共出現マトリックス212は、良好フレーズと関係付けられる格納データの $m \cdot m$ マトリックスである。マトリックスの各行 j は、良好フレーズ g_j を表し、各列 k は良好フレーズ g_k を表す。各良好フレーズ g_j について期待値 $E(g_j)$ が計算される。期待値 E は、 g_j を含むと期待されるコレクション内の文書の割合である。例えば、これは、 g_j を含む文書の数と、クロールされているコレクション内の文書の総数 T との比、 $P(j)/T$ として計算される。

【誤訳訂正4】

【訂正対象書類名】明細書

【訂正対象項目名】0069

【訂正方法】変更

【訂正の内容】

【0069】

上記のように、 g_j を含む文書の数は、 g_j が1文書に出現する毎に更新される。 $E(g_j)$ の値は、 g_j のカウントが1つ増加する毎に、またはこの第3段階の間に更新できる。

【誤訳訂正5】

【訂正対象書類名】明細書

【訂正対象項目名】0071

【訂正方法】変更

【訂正の内容】

【0071】

i) 期待値 $E(g_k)$ を計算する。 g_j および g_k の期待される共出現率 $E(j, k)$ は、両者が関連性のないフレーズの場合、 $E(g_j) * E(g_k)$ である；

【誤訳訂正6】

【訂正対象書類名】明細書

【訂正対象項目名】0072

【訂正方法】変更

【訂正の内容】

【0072】

ii) g_j および g_k の実際の共出現率 $A(j, k)$ を計算する。これは文書総数 T で除した生の共出現カウント $R(j, k)$ である；

【誤訳訂正7】

【訂正対象書類名】明細書

【訂正対象項目名】0073

【訂正方法】変更

【訂正の内容】

【0073】

iii) 実際の共出現率 $A(j, k)$ が或る閾値量だけ前記期待される共出現率 $E(j, k)$ を超える場合、 g_j は g_k を予測するといえる。

【誤訳訂正8】

【訂正対象書類名】明細書

【訂正対象項目名】0090

【訂正方法】変更

【訂正の内容】

【0090】

最初に、思い出すべきは、共出現マトリックス212が良好フレーズ g_j を含み、それぞれは、情報ゲイン閾値より大きな情報ゲインをもつ少なくとも1つの他の良好フレーズ g_k を予測するということである。次いで、関連フレーズを識別する(400)のために、良好フレーズの各対(g_j 、 g_k)毎に、その情報ゲインを関連フレーズ閾値、例えば1

00、と比較する。すなわち、 g_j および g_k は、次の場合、関連フレーズである：

【誤訳訂正9】

【訂正対象書類名】明細書

【訂正対象項目名】0092

【訂正方法】変更

【訂正の内容】

【0092】

この高い閾値を用いて、統計的期待率を十分超える良好フレーズの共出現を識別する。統計的に、これはフレーズ g_j および g_k が、期待される共出現率より100倍多く共出現する、ということの意味する。例えば、文書や「Monica Lewinsky (モニカ・ルインスキー)」というフレーズを考えると、フレーズ「Bill Clinton (ビル・クリントン)」は同一文書に100倍以上出現する可能性があり、ひいては、フレーズ「Bill Clinton」は、任意のランダムに選択される文書にも出現する可能性がある。言いかえると、出現率が100:1なので、予測精度は99.999%である。

【誤訳訂正10】

【訂正対象書類名】明細書

【訂正対象項目名】0100

【訂正方法】変更

【訂正の内容】

【0100】

一実施の形態では、クラスタ番号は、フレーズ間の直交関係をも示すクラスタービットベクトルにより決定される。クラスタービットベクトルは、良好フレーズリスト208内の良好フレーズ数 n の長さのビットのシーケンスである。所与の良好フレーズ g_j について、ビット位置はソートした g_j の関連フレーズ R に対応する。ビットは、 R の関連フレーズ g_k がフレーズ g_j と同一のクラスタにある場合に設定される。より一般的には、これは、 g_j と g_k との間のいずれかの方向の情報ゲインがある場合に、クラスタービットベクトルの対応するビットが設定される、ということの意味している。

【誤訳訂正11】

【訂正対象書類名】明細書

【訂正対象項目名】0101

【訂正方法】変更

【訂正の内容】

【0101】

次いで、クラスタ番号は、結果的に得られるビット列の値となる。この実装は、多くの方向、または一方向の情報ゲインをもつ関連フレーズが同一クラスタに出現するという特性をもつ。

【誤訳訂正12】

【訂正対象書類名】明細書

【訂正対象項目名】0223

【訂正方法】変更

【訂正の内容】

【0223】

「blue merle」::「Australian Shepherd」, 「red merle」, 「tricolor」, 「aussie」;

【誤訳訂正13】

【訂正対象書類名】明細書

【訂正対象項目名】0224

【訂正方法】変更

【訂正の内容】

【0224】

「agility training」:「weave poles(編んだ柱)」, 「teeter(シーソー)」, 「tunnel(トンネル)」, 「obstacle(障害)」, 「border collie(ボーダーコリー)」。