

コルモゴロフ複雑性に基づく IDFの単語N-gramへの適用

白川 真澄, 原 隆浩, 西尾 章治郎

大阪大学 大学院情報科学研究科

本研究の貢献

① IDFに新たな理論的説明を付与

② ①に基づきIDFを単語N-gramに適用した手法を提案

③ ②を計算する手法を実装

接尾辞木とウェーブレット木を用いた
計算方法を実装・公開した

<http://github.com/iwnsew/ngweight>

(計算済みのデータ, デモシステムを含む)

今日は話しませんので詳しくは論文で

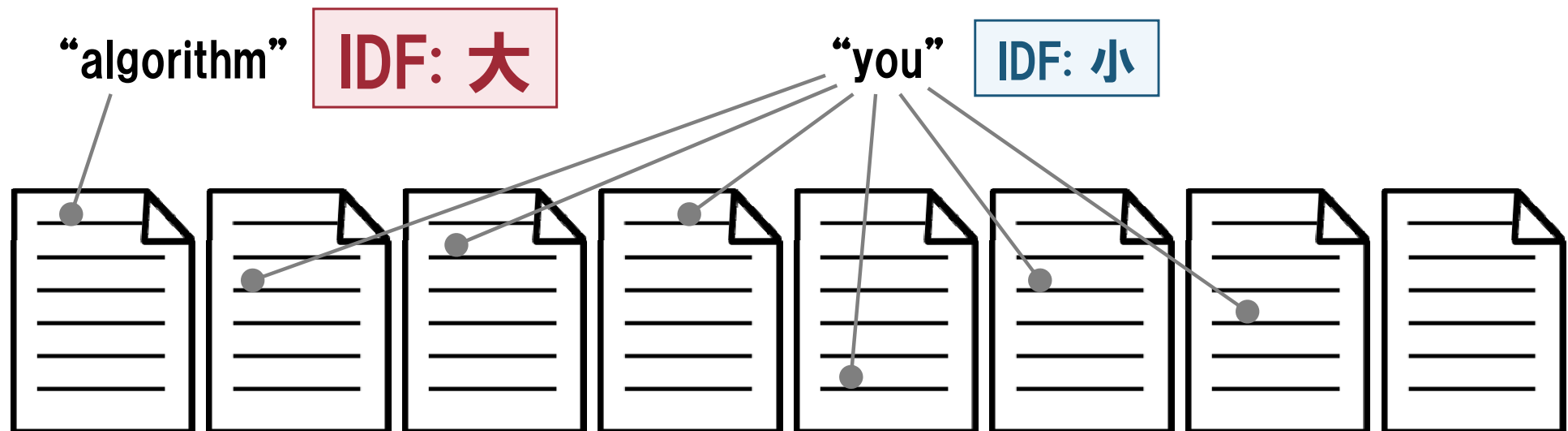
④ ②の頑健性を評価実験にて確認

IDF (Inverse Document Frequency, 逆文書頻度)

より少数の文書にしか出現しない語に大きい重みを付与

$$IDF(t) = \log \frac{|D|}{df(t)}$$

t : 語
 $df(t)$: 語 t の文書頻度
 $|D|$: 文書集合 D の文書数



IDFの欠点

単語N-gram(複合語)ではうまく機能しない

何故か？

文書頻度が小さい単語N-gram → 特徴的！

つながりが不自然な単語N-gram → 文書頻度が小さい！

つながりが不自然な単語N-gram → 特徴的！！！！！！！！！！

Google検索による文書頻度の推定値

$$df(\text{"Osaka University"}) = 2,490,000$$
$$df(\text{"Osaka be"}) = 7,790$$
$$df(\text{"Osaka be to"}) = 4$$

単語間結合度の測定

「語の重み付け」とは別に研究されてきた
Multiword Expression (MWE) とか言われる

代表的な手法

- ✓ PMI (Pointwise Mutual Information)
- ✓ SCP [Silva+, MOL99]
- ✓ EMI [Zhang+, Expert Systems with Applications, 2009]
- ✓ MED [Bu+, COLING10]
- ✓ などなど…

本研究は, IDFに単語間結合度の測定を
「**理論的**」に組み込むことを目指す

コルモゴロフ複雑性と 情報距離

コルモゴロフ複雑性 (Kolmogorov Complexity) (1/2)

[Kolmogorov, Sankhya, 63]

有限長のデータ列(文字列など)の複雑さを表す指標

$K(x)$: データ列 x の複雑性

Q1: x_1 と x_2 はどちらが複雑でしょうか？

$x_1 = \text{"010101010101010"} \rightarrow \text{"01"} \times 7 + \text{"0"}$

$x_2 = \text{"011101100010110"} \rightarrow \text{"011101100010110"}$

A1: たぶん $K(x_1) < K(x_2)$

コルモゴロフ複雑性 (Kolmogorov Complexity) (2/2)

[Kolmogorov, Sankhya, 63]

有限長のデータ列(文字列など)の複雑さを表す指標

$K(x|y)$: データ列 y が与えられた条件での x の複雑性

Q2: x_1 と x_2 はどちらが複雑でしょうか?(ただし y を用いてもよい)

$y = \text{"01110110001011"}$

$x_1 = \text{"010101010101010"}$ \rightarrow $\text{"01"} \times 7 + \text{"0"}$

$x_2 = \text{"011101100010110"}$ \rightarrow $y + \text{"0"}$

A2: たぶん $K(x_1|y) > K(x_2|y)$

$K(x, y) = K(x|y) + K(y) = K(y|x) + K(x)$ が成り立つ
(一方の文字列を, 他方を表すのに再利用できる)

情報距離 (Information Distance)

[Bennett+, TOIT, 98]

コルモゴロフ複雑性に基づき定義される距離

$E(x, y)$: データ列 x と y の情報距離

$$E(x, y) = \max(K(x|y), K(y|x)) = K(x, y) - \min(K(y), K(x))$$

情報距離とは、一方を他方に変換するためのエネルギー

1ビット変換するのにかかるエネルギー = $1kT \ln(2)$

[Landauer, IBM Journal of Research and Development, 61]

↑
↑
絶対温度
ボルツマン定数

MED: Multiword Expression Distance

[Bu+, COLING10]

情報距離に基づく普遍的な単語間結合度の測定手法

$$MED(g) = \log \frac{df(w_1, \dots, w_N)}{df(g)} \quad g: w_1 \dots w_N \text{ で表される } N\text{-gramの語}$$

語のコンテキストとセマンティックの間情報距離

g のコンテキスト: g を含む文書集合

g のセマンティック: ユニグラム w_1, \dots, w_N を全て含む文書集合

Inspired by [Cilibrasi+, TKDE, 2007] (Normalized Google Distanceの論文)

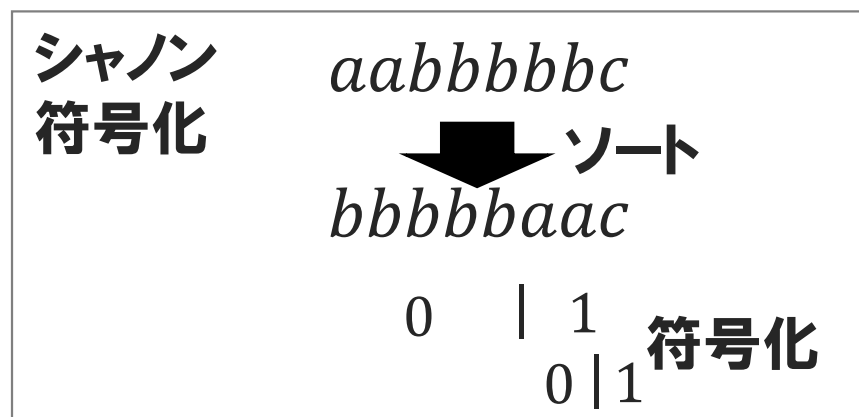
MEDの導出 (1/2)

コンテキスト(文書集合) x の確率を,
文書頻度 $|x|$ に比例すると考える

$$P(x) = \frac{|x|}{\sum_{x_i \in X} |x_i|} \quad X: \text{全てのコンテキストの集合}$$

すると, シャンン符号化*1を用いて, 以下のように
コルモゴロフ複雑性を近似できる [Li+2008]

$$K(x) \approx -\log P(x)$$
$$K(x, y) \approx -\log P(x, y)$$



MEDの導出 (2/2)

$\phi(g)$ と $\mu(g)$ の情報距離は、
 g のコンテキスト g のセマンティック

$$\begin{aligned} & E(\phi(g), \mu(g)) \\ &= K(\phi(g), \mu(g)) - \min(K(\phi(g)), K(\mu(g))) \\ &= -\log P(\phi(g), \mu(g)) + \min(\log P(\phi(g)), \log P(\mu(g))) \\ &= -\log|\phi(g) \cap \mu(g)| + \max(\log|\phi(g)|, \log|\mu(g)|) \\ &= -\log|\underline{\phi(g)}| + \log|\underline{\mu(g)}| \end{aligned}$$

シャノン符号化

確率は文書頻度に比例

$\phi(g) \subseteq \mu(g)$ なので

g を含む文書の頻度

w_1, \dots, w_N を全て含む文書の頻度

すなわち, $MED(g) = \log \frac{df(w_1, \dots, w_N)}{df(g)}$

IDFと情報距離の関係

IDF

$$IDF(g) = \log \frac{|D|}{df(g)}$$

g のコンテキストと 空文字 ε のコンテキスト との情報距離

ε はすべての文書に含まれる

$$E(\phi(g), \phi(\varepsilon))$$

$$= -\log|\phi(g) \cap \phi(\varepsilon)| + \max(\log|\phi(g)|, \log|\phi(\varepsilon)|)$$

$$= -\log|\phi(g)| + \log|\phi(\varepsilon)| \leftarrow \phi(g) \subseteq \phi(\varepsilon) \text{ なので}$$

$$= \log \frac{df(\varepsilon)}{df(g)} = \log \frac{|D|}{df(g)} = IDF(g)$$

IDFと一致

シャノン符号化 &
確率は文書頻度に
比例

本研究の貢献

① IDFに新たな理論的説明を付与

語を、それが出現する文書の集合で表現し、シャノン符号化により文書集合のコルモゴロフ複雑性を近似したとき、語と空文字との情報距離がIDFとなる

② ①に基づきIDFを単語N-gramに適用した手法を提案

③ ②を計算する手法を実装

接尾辞木とウェーブレット木を用いた計算方法を実装・公開した
<http://github.com/iwnsew/ngweight>
(計算済みのデータ, デモシステムを含む)

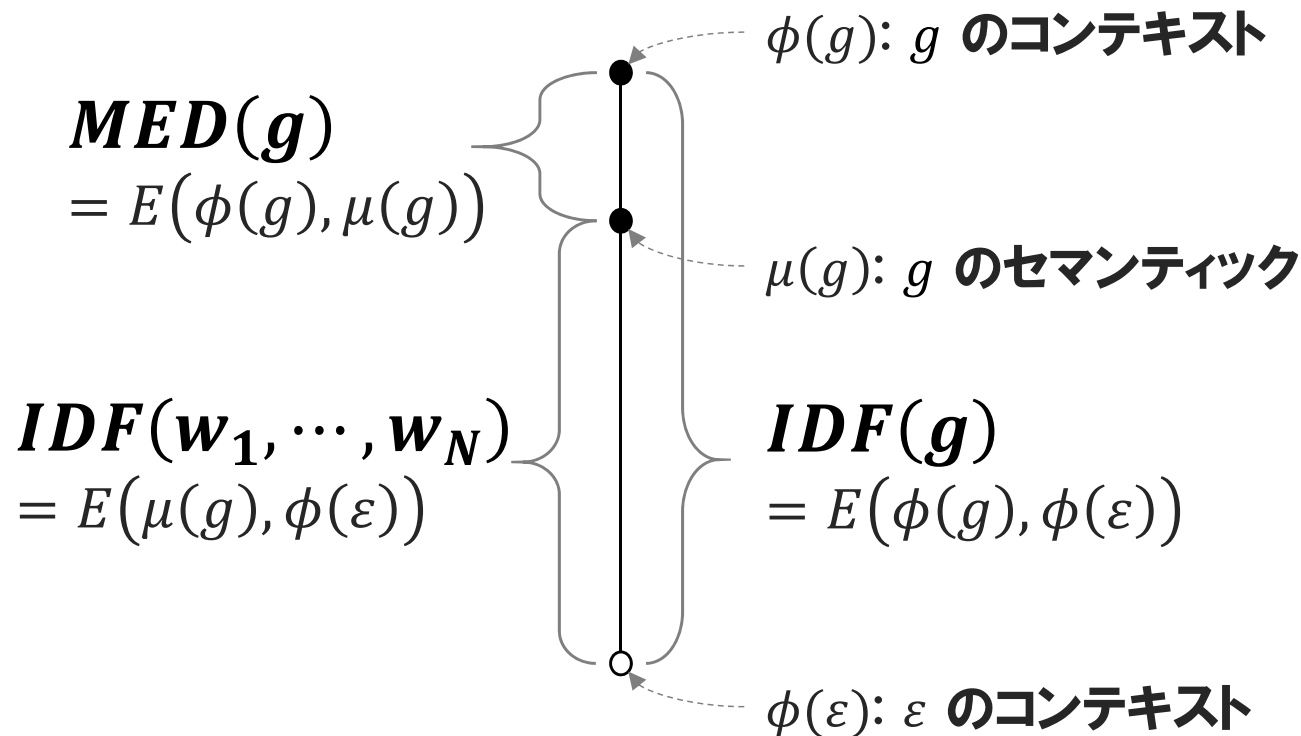
今日は話しませんので詳しくは論文で

④ ②の頑健性を評価実験にて確認

情報距離空間におけるIDFとMED

単語N-gram g のIDFは, g を単語に分解したときのIDFと g のMEDの和である

でもMEDはむしろ小さいほうが, 単語間の結合が自然なのは?

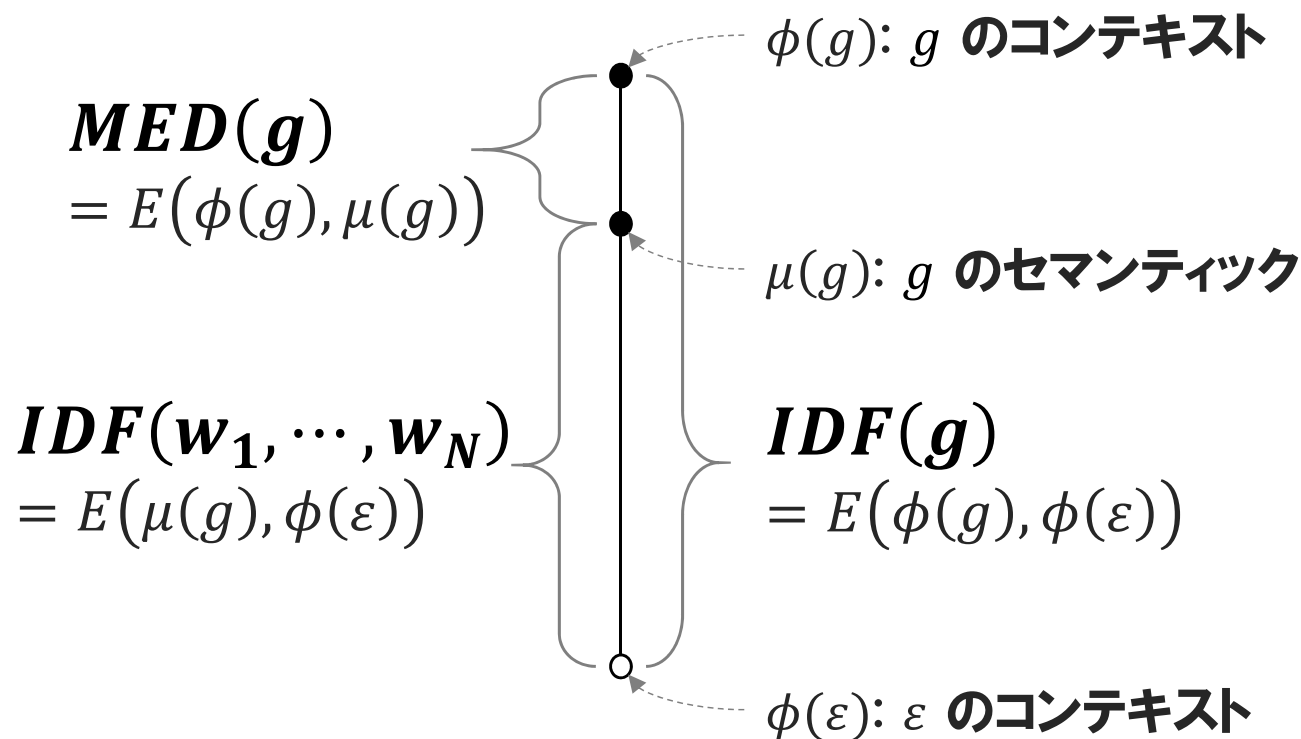


提案手法: IDF_{N-gram}

情報距離に基づいてIDFを再設計

IDFは大きく、且つMEDが小さいときに大きい重みを付与

$$IDF_{N-gram}(g) = IDF(w_1, \dots, w_N) - MED(g) = \log \frac{|D| \cdot df(g)}{df(w_1, \dots, w_N)^2}$$



重みのみを用いた特徴語抽出

入力: "Alice's Adventures in Wonderland - Kindle edition by Lewis Carroll"

単語N-gram	IDF_{N-gram}	単語N-gram	IDF_{N-gram}
kindle edition	12.043	adventures	7.101
kindle	11.653	kindle edition by	6.739
alice s adventures in wonderland	11.496	lewis	6.192
adventures in wonderland	10.906	edition	4.836
s adventures in wonderland	10.804	adventures in	4.280
wonderland	9.670	s adventures	3.586
lewis carroll	9.498	alice s	3.507
alice s adventures	9.385	s adventures in	2.255
alice s adventures in	9.348	by lewis	1.768
in wonderland	8.762	s	1.030
carroll	8.152	by	0.820
by lewis carroll	7.461	in	0.154
alice	7.234	edition by	-0.875

本研究の貢献

① IDFに新たな理論的説明を付与

語を、それが出現する文書の集合で表現し、シャノン符号化により文書集合のコルモゴロフ複雑性を近似したとき、語と空文字との情報距離がIDFとなる

② ①に基づきIDFを単語N-gramに適用した手法を提案

単語N-gramに対応した重み付け手法を提案し、さらに、重みのみを用いて任意の長さの特徴語をテキスト中から抽出できることを発見した

③ ②を計算する手法を実装

接尾辞木とウェーブレット木を用いた計算方法を実装・公開した

<http://github.com/iwnsew/ngweight>

(計算済みのデータ, デモシステムを含む)

今日は話しませんので詳しくは論文で

④ ②の頑健性を評価実験にて確認

評価1: 特徴語抽出

Wikipediaの第一段落1,678記事分をデータセットとして利用

アンカーテキストと太字の語を正解として適合率(上位R個, Rは正解の数)を算出

品詞タグ
不要!



品詞タグ
必要

手法	適合率
<i>TF-IDF_{N-gram}</i>	0.377
名詞句抽出 + 各単語のTF-IDFの和 [Hasan+, COLING10]	0.386
名詞句抽出(全部小文字) + 各単語のTF-IDFの和 [Hasan+, COLING10]	0.369
名詞句抽出 + 各単語のTF-IDFの平均	0.367
名詞句抽出(全部小文字) + 各単語のTF-IDFの平均	0.352
名詞句抽出 + 単語N-gram自体のTF-IDF	0.369
名詞句抽出(全部小文字) + 単語N-gram自体のTF-IDF	0.355
TF-IDF (単語のみ抽出)	0.229

評価2: クエリ分割

情報検索ベースのクエリ分割データセット [Roy+, SIGIR12] を利用
分割した語句をWeb検索の条件とし, 得られた検索結果を評価

nDCG: 正解が上位にあるほど高いスコア

外部情報
不要!

外部情報
必要

手法	nDCG (上位5位)	nDCG (上位10位)
IDF_{N-gram}	0.730	0.742
Mishra手法(クエリログ使用) [Mishra+, WWW11]	0.706	0.737
Mishra手法+Wikipediaタイトル [Roy+, SIGIR12]	0.725	0.750
PMI(クエリログ使用) [Roy+, SIGIR12]	0.716	0.736
PMI(Webコーパス使用) [Roy+, SIGIR12]	0.670	0.707
分割なし	0.655	0.689

本研究の貢献

① IDFに新たな理論的説明を付与

語を、それが出現する文書の集合で表現し、シャノン符号化により文書集合のコルモゴロフ複雑性を近似したとき、語と空文字との情報距離がIDFとなる

② ①に基づきIDFを単語N-gramに適用した手法を提案

単語N-gramに対応した重み付け手法を提案し、さらに、重みのみを用いて任意の長さの特徴語をテキスト中から抽出できることを発見した

③ ②を計算する手法を実装

接尾辞木とウェーブレット木を用いた計算方法を実装・公開した

<http://github.com/iwnsew/ngweight>

(計算済みのデータ, デモシステムを含む)

今日は話しませんので詳しくは論文で

④ ②の頑健性を評価実験にて確認

特徴語抽出およびクエリ分割において、追加の情報を必要とする既存手法と同等の性能を、IDFの重みのみを用いて達成した

今後の課題

IDF_{N-gram} の計算にかかる処理速度の向上

日本語などの単語間の区切りがない言語への対応

本研究のまとめ

① IDFに新たな理論的説明を付与

語を、それが出現する文書の集合で表現し、シャノン符号化により文書集合のコルモゴロフ複雑性を近似したとき、語と空文字との情報距離がIDFとなる

② ①に基づきIDFを単語N-gramに適用した手法を提案

単語N-gramに対応した重み付け手法を提案し、さらに、重みのみを用いて任意の長さの特徴語をテキスト中から抽出できることを発見した

③ ②を計算する手法を実装

接尾辞木とウェーブレット木を用いた計算方法を実装・公開した

<http://github.com/iwnsew/ngweight>

(計算済みのデータ, デモシステムを含む)

今日は話しませんので詳しくは論文で

④ ②の頑健性を評価実験にて確認

特徴語抽出およびクエリ分割において、追加の情報を必要とする既存手法と同等の性能を、IDFの重みのみを用いて達成した