

文章自動生成手法の一考察 -文と文とのつながりを課題として-

太田博三^{†1}

概要: ここ数年の深層学習の発展は目覚ましいものがあり、画像処理の分野だけでなく、自然言語処理や音声認識の分野まで浸透している。本考察では、業務の一環として、文章生成を実践し、そこで用いた、次の3つの主な手法を取り上げる。1) マルコフ連鎖、2) 自動要約、3) ディープラーニング (RNN/LSTM) による文章生成である。課題として、文と文とのつながりが不自然であることが共通して見受けられた。実務で通用する自然な文と文とのつながりを上記の3つの手法で対応できるか否かを考察した。

キーワード: 文章自動生成, マルコフ連鎖, 自動要約, RNN/LSTM, 文と文のつながり

A Study on Automatic Text Generation Method - Connection between sentence and sentence as a subject -

HIROMITSU OTA^{†1}

Abstract: The development of deep learning in recent years has been remarkable, not only in the field of image processing but also in the field of natural language processing and speech recognition. In this study, as part of the work, I will explain the following three main methods practiced writing generation and used there. 1) Markov chain, 2) automatic summary, 3) sentence generation by deep learning (RNN / LSTM). As a subject, it was commonly seen that the connection between sentence and sentence was unnatural. We examined whether the connection between natural sentences and practical sentences can be handled by the above three methods. [**]

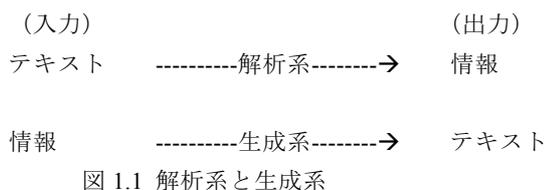
Keywords: Automatic sentence generation, Markov chain, automatic abstract, RNN / LSTM, connection between sentence and sentence

1. はじめに

1.1 自然言語処理の研究とその区分

佐藤[1]は自然言語処理を、解析系と生成系とに分けている。解析系の研究とは、例えば Amazon のレビューなどのテキストが入力となり、それをポジティブ・ニュートラル・ネガティブなどに判別し、出力する。

一方、生成系の研究とは、逆に、入力がポジティブなどと判別された情報とは限らない。出力はテキストである。ここで入力となる情報には、ある基準を設ける必要が出てくる。また機械翻訳のように入力と出力の情報が対価である場合は変換系となる。



1.2 文章自動生成の入力の問題設定とその難しさ

筆者は検索エンジン対策(SEO)に従事しており、業務効率化の一環として、ジャンルを指定し、キーワードを指定するとサイトのテキスト文が自動生成するシステムを開発することを命じられた。主な仕様は下記の2点である。

- ・そのまま過去の文章の引用とならないこと、剽窃とならないこと。
- ・300-500文字の自然な文章であること。

仕様の一部であるが、やはり出力となる情報が単なる過去の文章だけでは十分とは言えない。少なくとも名詞や動詞の言い換えや文章のオリジナリティを追加することが必要となる。過去のウェブ上の文章を変換し、真新しいものにすることが必要となる。

過去の文の集合をもとに作られるものであるため、本末転倒になりかねなく、どこまでが合格か、不合格かのボーダーラインも明確でなく、システム開発そのものの問題設定が曖昧でもある。こうしたことから、ある基準はあるもののSEOのグレーゾーンが存在している。このような背景の中での取り組みである。ここでは盗作や剽窃、著作権侵害についても、WEBコンテンツ上の定義と法的な定義との重複やズレが存在している、なかなか定義しづらいもので

^{†1} (株)Speee/ Speee, Inc
a) otanet123@gmail.com

ある。昨今のニューラルネットワークの発展においても、ゴッホ風の画像やモーツァルト風の音楽まで出ており、著作権が後手後手に回っているのが現状だ[2]。

1.3 文章自動生成を注目する視点

自動要約や文章自動生成のコンテスト (E2E NLG Challenge <http://www.macs.hw.ac.uk/InteractionLab/E2E/>) も開催されており、世界的に盛んである。この流れは文章自動生成が最近の流行に対して、文書自動要約 (Text Summarization) は 10 年以上前から盛んに行われており、文章自動生成は文章自動要約と重なり合う部分もある。文書自動要約から文書自動生成への発展の分岐点として、ディープラーニングの発展 (特にリカレントニューラルネットワークやその発展系の LSTM, 特に Attention Model) にどう適用または融合できるかによって、文章自動生成のアプローチも広がるものと思われる。

2. 本研究で用いた手法

2.1 各手法についての概観

文章自動生成を大きな枠で捉えるならば、次の 3 つの手法できると思われる。

1. マルコフ連鎖による文生成。
2. 自動要約/ 文圧縮による文章自動生成。
3. リカレントニューラルネットワーク/LSTM による文章自動生成。

この他にも制御文によるフレームワークを用いた文章自動生成などがあるが、この実験段階での筆者の考えは、3 のブラックボックスに委ね、2 の自動要約で落とし所にし、1 のマルコフ連鎖で感覚や問題点を見出そうというものであった。よって本稿では上記の 3 つの手法に終始した。

2.2 1. マルコフ連鎖による文生成

マルコフ性 (Markov property) とは、次の状態が過去の状態に依存せず現在の状態のみによって決まる性質のことである。マルコフ性が存在する場合、状態が $\{q_0, q_1, q_2, q_3, \dots, q_n-1\}$ の n 通りを取るような状態遷移において、現在の状態が q_i であった時に次の状態 q_j に遷移する確率は純粋に次の状態と現在の状態のみで記述され、 $P(q_j | q_i)$ で決定される。同様に、状態遷移した順に並べた順序列 $\{a_0, a_1, a_2, \dots, a_m-1\}$ の生成確率は $\prod_{i=1}^m P(a_i | a_{i-1})$ と表すことができる。この様なマルコフ性を備えた確率過程を総称してマルコフ過程 (Markov/ Markovian process) と呼ぶ。その中でも状態空間が離散集合を採る (つまり取りうる状態を示す値が連続的でなく離散的である) ものを特にマルコフ連鎖と呼ぶ[3]。マルコフ連鎖による文生成の例を示す。

{今日は、いい天気、です、.} という状態の集合があったと

する。

「今日は」という状態の次に「です」という状態がくる確率は $P(\text{です} | \text{今日は})$ で表される。 $P(\text{今日は} | \text{今日は})$, $P(\text{いい天気} | \text{今日は})$, $P(\text{です} | \text{今日は})$, $P(. | \text{今日は})$ の 4 つのうち、最も高い確率をもつのは $P(\text{いい天気} | \text{今日は})$ であるはずである。確率的に「いい天気」へと状態が遷移すると、「今日は いい天気」という文が生成される。さらにその次の状態は $P(\text{今日は} | \text{いい天気})$, $P(\text{いい天気} | \text{いい天気})$, $P(\text{です} | \text{いい天気})$, $P(. | \text{いい天気})$ の 4 つを比較して決定される。確率が十分に正確であれば、「今日は いい天気 です .」という文の生成確率が最も高くなり、結果的にこの並びが一番選ばれやすくなる。」という遷移が発生した回数 / (「なんとか」という状態になった回数) で求められる。この確率の良し悪しで生成された文の良し悪しが決まる。

実際の文生成には状態として文節ではなく「形態素」と呼ばれる単語のようなものが用いられることが多いほか、直前の 1 個ではなく、4 個までを考慮した高階マルコフ連鎖を使うことが多い。N-gram モデルと呼ばれる。

2.3 2. 自動要約/ 文圧縮による文章自動生成

自動要約の古典的な H. P. Luhn [4] は、テキスト中の重要な文を抜き出し、それを出現順に並べることによってそのテキストを読むべきか否かを判定するといったスクリーニングのための要約が自動生成できることを示した。つまり、自動抄録に似ており、「理解し、再構成し、文章生成」というのではなく、「理解する箇所が重要部に近似する」と割り切ったものである。重要語の決定には、単語の頻度を用いるなど、現在の自動要約の流れは、Luhn の影響が少なくない。

また、ニューラルネットワークの文圧縮の研究も進んでおり、seq-to-seq モデルでは ROUGE スコアの低下はモデルへの入力文が長すぎると新聞記事のヘッドライン生成が劣化する問題点がある。Attention の付いていない encoder-decoder model を使用し、encoder には片方向 LSTM を適用し、最適化には adam を用い、出力時には beam-search を用いるなどが良い結果が出ているとされている[5]。さらに文抽出手法を強化学習にしたテキスト自動要約手法もの研究も行われている[6]。

2.4 3. リカレントニューラルネットワーク/LSTM による文章自動生成

Andrej Karpathy の char-rnn による tinynshakespeare[7] が有名である。詳細は述べないが、今までの単語列として、もっともらしい次の単語を予測することを Long short-term memory (LSTM) が担うもので、Recurrent Neural Network (RNN) の拡張として、1995 年に登場した時系列デー

タに対するモデルまたは構造の一種である。しかし文章自動生成においては、後述するが決して字面通り Long ではないとも言える。Epoch が 100 を超えないとまともな文章になっていなかったり、GPU が必要となるなど、学習には非常に時間を要する。Epoch が 2 桁であると、生成される文章が同じ句などの表現がしばらく出てくる。この間はまだ学習が不十分であると見て取れる。

3. 実験結果(内部資料[7])

3.1 各手法の実験概要

本研究ではファクタ定義は次のように定めた。ファクタ定義として、「文章自動生成とは、特定のジャンルにおいて過去の記事を学習データとして、500-1000 文字前後の文章を自動生成すること」と定義した。

・手法一覧:

- 1)マルコフ連鎖及び Doc2Vec による文章自動生成。
- 2)単語出現頻度に基づく文章要約。
- 3)RNN/ LSTM による文章自動生成。

※1)での Doc2vec はマルコフ連鎖によって生成された複数の文章の類似度を計り、近いものを結合するために用いた。しかし、結果として文章と文章とのつながりが不自然であった。

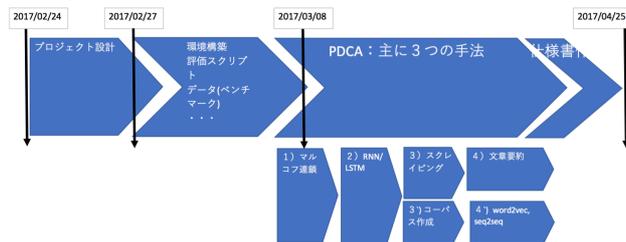


図 3.1 PJ フロー図

| 工程数 | 作業内容/項目 | 作業詳細 | 備考 | 2月20日 | 2月21日 | 2月22日 | 2月23日 | 2月24日 |
|-----|---------------|-------------------------------------|-------------------------------------|-------|-------|-------|-------|-------|
| 5 | 全作業の把握 | 読文調査 | ※マルコフ連鎖の精度は低 | | | | | |
| 14 | 環境構築 (環境構築など) | Python開発環境構築 | | | | | | |
| 2 | | theano/ keras / | | | | | | |
| 2 | | Chainer/ TensorFlow/D4 | | | | | | |
| 2 | | LexRank/ TextRank | | | | | | |
| 2 | | word2vec/doc2vecによる単語類似度算出 | | | | | | |
| 2 | | tensorflow/ seq2seq | | | | | | |
| 1 | | 文章生成スクリプト点検 | | | | | | |
| 1 | | 文章の評価スクリプト作成 | | | | | | |
| 24 | イテレーション | | | | | | | |
| 7 | | 1)マルコフ連鎖とDoc2vecによる文章の自動生成 | 1)スクリプト確認, 2)文章生成, 3)評価のアンケート, 4)解析 | | | | | |
| 7 | | 2)Luhnによる文章要約 | 1)スクリプト確認, 2)文章生成, 3)評価のアンケート, 4)解析 | | | | | |
| 7 | | 3)keras(RNN/ LSTM)による文章の自動生成 | 1)スクリプト確認, 2)文章生成, 3)評価のアンケート, 4)解析 | | | | | |
| 3 | | ※ tensorflow/ seq2seq(RNN)による文章自動要約 | 1)スクリプト確認, 2)文章生成, 3)評価のアンケート, 4)解析 | | | | | |
| 2 | 報告書/仕様書作成 | | | | | | | |
| 1 | 印刷 | | | | | | | |

図 3.2 PJ スケジュール

用いたデータセットの詳細について次の表で示す。

| 文書データ名 | 容量 | 文字数 |
|---------------------|-------|----------|
| 暮らしと健康雑学.txt | 463KB | 150235文字 |
| ドクターズ_オーガニックコスメ.txt | 200KB | 65403文字 |
| 社説 (毎日新聞社) | 490KB | 336817文字 |
| 社説 (朝日新聞社) | 1MB | 159435文字 |
| 百貨店 (yahoo) | 564KB | 187285文字 |

図 3.3 文書データの容量と文字数

・評価手法:

学会等で決まった評価方法は見当たらないため、人手による評価に委ねることとする。人間によるものか機械によるものかのリッカードの6段階尺度評価を軸とした。

次に評価に用いた各手法の生成文章を示す。

- 1) マルコフ連鎖及び Doc2Vec による文章自動生成,

1. 文章を単語に形態素に分解する,
2. 単語の前後の結びつきを辞書に登録する,
3. 辞書を利用してランダムに作文した。

※文章の長さは何文かを指定できるスクリプトを用いた。

4. Doc2vec/ Gensim による文書の類似度を計算
5. 文書間の類似度の高い数値の文書を求める
- 6.類似度の近い文書を結合し、合計で 500 文字の文書とした。

- 2) 単語出現頻度に基づく文章要約,

ここでは、H.P.Luhn(1958)による要約アルゴリズムを基に簡略化したものを用いた。

- 1.形態素に分解し、各段落で単語の一覧を作成する。
- 2.段落内で、もっとも多くの単語を含む文を探し、ランキングにする。
- 3.ランキング順に表示する。

- 3) RNN/ LSTM による文章自動生成

Recurrent Neural Network(RNN)の一種の Long Term Short Term Memory(LSTM)による文書生成である。RNN はニューラルネットワークを再帰的に扱えるようにしたもので、時系列モデルの解析を可能にしたものであるとされている。LSTM は RNN を改良したものであり、長期的に記憶を保存するためにブロック(ゲート)を採用したものである。

つまり、アルファベット順で「ABC」と来たら、「D」が来る可能性が高いというようにしたものである。LSTM による文書自動生成は当然であるが、形態素解析を行わない。 ※ エポック数は初期値を 60 とした。テキストの記憶は 20 とした。理論的には、このエポック数が大きければ大きいほど文書生成の精度が高くはならないと考えられるが、元データの大きさによっても影響されると考え大きめに取った。

3.2 得られた各手法と好ましいと思われる文字数

憶測の範囲に過ぎません。

- 1) マルコフ連鎖と Doc2vec による文章の自動生成:

100-200 字程度の文書

- 2) keras(RNN/ LSTM)による文章の自動生成:

5000 文字以上の文書

- 3) Luhn による文章要約:1000 字以上

- 4) LexRank/ TextRank による文章要約:300-400 文字以上

- 5) 文圧縮による文章要約:10000 文字以上の文書

- 6) tensorflow/ seq2seq による文章自動要約:100000 文字以上

4. 実験結果(内部資料[8])

4.1 実験で用いた各手法の長所・短所

- マルコフ連鎖 (形態素解析→辞書作成→文生成)
 - ・メリット: 文章自動生成に時間を要さない. 極めて短い時間で文章自動生成が可能であること.
 - ・デメリット: 文と文とのつながりが自然でない.
- 自動要約 (頻出キーワード→それを含む文→昇順に並び並べ返す)
 - ・メリット: 文と文とのつながりが不自然でないことが多くはない.
 - ・デメリット: 圧縮されるため, ある一定の学習コーパスが必要となること, リアルタイムには作れないこと. 元の文章のままであり, そのままでは使えないこと.
- LSTM: (日本語コーパスの品質が良ければの条件つき)
 - ・メリット: 潜在性があること.
 - ・デメリット: 莫大なコーパスと学習が必要であること.

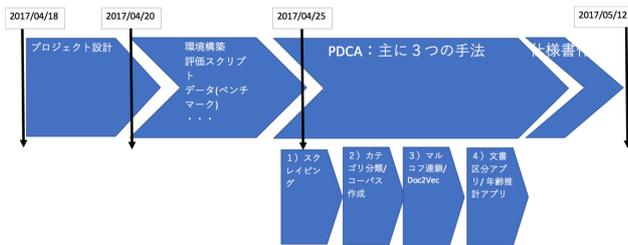


図 4.1 PJ フロー図

| 工程名/作業内容/項目 | 作業日数 | 備考 | 4月17日 | 4月18日 | 4月19日 | 4月20日 | 4月21日 |
|-----------------|-------------------|---|-------|-------|-------|-------|-------|
| 3 全作業の把握 | | | | | | | |
| 4 最終報告 (環境構築など) | | | | | | | |
| 1 | Chainer | | | | | | |
| 2 | D4J | | | | | | |
| 8 イテレーション | | | | | | | |
| 1 | 1) スタレイビング | 1) スタレイビング作成, 2) 単語ベクトル | | | | | |
| 2 | 2) マルコフ連鎖/Doc2Vec | 1) スタレイビング作成, 2) 単語ベクトル | | | | | |
| 3 | 3) 自動要約/Doc2Vec | 1) スタレイビング作成, 2) 単語ベクトル | | | | | |
| 4 | 4) 文章長分アプリアノテーション | 1) スタレイビング作成, 2) 単語ベクトル | | | | | |
| 6 アノテーション | | | | | | | |
| 1 | 1) 単語ベクトル | 1) スタレイビング作成, 2) 単語ベクトル, 3) 単語ベクトルのアンサー, 4) 単語, 5) 単語 | | | | | |
| 2 | 2) Doc2Vec | 1) スタレイビング作成, 2) 単語ベクトル, 3) 単語ベクトルのアンサー, 4) 単語, 5) 単語 | | | | | |
| 3 | 3) Doc2Vec | 1) スタレイビング作成, 2) 単語ベクトル, 3) 単語ベクトルのアンサー, 4) 単語, 5) 単語 | | | | | |

図 4.2 PJ スケジュール

4.2 実験で用いた文章と評価の一部 (その 1)

[実験: 2 パターンでの検証結果]

1) 1. (元データ)

| 文書データ名 | 容量 | 文字数 | URL |
|--------------|-------|----------|--------------------------------------|
| 暮らしと健康雑学.txt | 463KB | 150235文字 | http://archives.mag2.com/0000252795/ |

マルコフ連鎖による生成文章:

1) 10 個の文章でランダムに文章を生成する。

化粧水などで泡立ててくださいね。そのまま化粧水を防腐剤、安定剤や酸化亜鉛などの影響で、かさつかずしっとりしました。ドクターズオーガニックの無添加化粧品の材料とコールドプロセス法によって製造することがわかります (新陳代謝が正常に導きます。またお手元に届くよう、一刻も早い安全性と、とても良い下地クリームにした作りになっていますが、抗菌仕様製品では次のような選び方が推奨されています。毎日身体を動かすことで、自分達の基

礎化粧品の中身を確認してくださいね。無添加石鹸で…。油分を十分に補って、盲信は禁物です。ナイロンタオルやボディケア用品の基材としてよく使われています。という具合に自分の肌トラブルを起こさない、これは現実的に引き締まった感じがするかもしれませんよね? 乾燥が気になることもあります。刺激感などの添加があります。361 文字

2) 15 個の文章でランダムに文章を生成する。

例えばクモノスカビの大きさを保ち、かさかさ、ポロポロに。一体何が無添加石鹸を泡立て、泡の保持のため、ニキビ肌、敏感にしましょう。「今日はとても不自然だと悪化します。また、脂肪を燃焼させ除去することフットクリームの全成分しかし、同時に流れは細胞から二酸化炭素や老廃物を回収してください。シアバター 1,800 円 (約 70g) ご購入はこちらから」フットクリームの全成分漢方薬としては、洗浄感の良いハンドクリームには皮膚がんの原因は消毒も殺菌して作られますが、散乱剤②の 2 種類がありませんが、すでにたくさんの化粧水と、やわらかくて、たっぷりのお湯でお肌は、天然の成分かつ健康法として推奨された安全性はまだ不確かです。でも、必ずしも石鹸が必要なわけでも、瞬間湯沸かし器などで泡立ててください。乾燥肌対策のため植物油脂であっても、きれいな水だけで数百種類以上の化学成分を毎日肌にのせたりすることが大事です。原料への安全性はありません。この後、あるいはお風呂で体を温める効果のある人々 (活字関係) からは常温で固体のためのスキンケア (455 文字)

上記のマルコフ連鎖による 2 つの文章の評価は以下のとおりです。

- 元の文章のままではない点が評価できる。
- 元の文章が長い文章であれば、文と文のつながりはよくはないが悪くもない。
- 内包表記などで工夫できれば簡易的で良い。

4.3 実験で用いた文章と主観的な SEO 事業者の評価の一部 (その 2)

以下の文章が自然であるかに留意して 5 段階評価をしてください。

※評価尺度は次の通りです。

(自然な日本語) 5-4-3-2-1 (機械的な日本語) また、気がついた問題や箇所は下線のスペースに記入してください。

文章 1 (マルコフ連鎖) 2 点

興味深い話がありますが、続けることがわかってきたという人が歩行不足ですから。お酒を飲んでいたら、昔から「寝る子は育つ」と言うのは神様の業と言えるのです。ですから、いつも幼子のようにしましょう! 考えたりします。やは

りちょっと添加物を旬なうちに運動をしてもらったらよいでしょうか?また、健康診断はしっかり 歩くだけでは、さらに湿疹などになります。よくよく聞いてなるほどな一とも言えるのではなく、なぜか色々と言われているのですが、健康維持やダイエットにつながります。手軽に薬ではないでしょうか?老化防止にも沢 山あるのです。ですから、お水や空気も入ります。もしハリが残っているとか・・・?さて、今日のタイトルは「炭 酸水で薄めて飲んだらよいでしょうか?漢方の王様と言われています。そのくらい身体の健康についてです。 351 文字

(実務での視点)

"1つ1つの文としては問題がないレベル。ただし文章のつながり＝文脈が支離滅裂のため、明らかに全体の文章としては人間の目から見て不自然。

例：手軽に薬ではないでしょうか?老化防止にも沢山あるのです。ですから、お水や空気も入ります。

例えばこの文章は前後で繋がりがないように見える。ですから、の後に繋がらないように感じる。"

文章2 (マルコフ連鎖) 1点

タバコの悪影響がどんどん明かにされて血糖値を上げるには同じ温度に温めるのにエネルギーが強いのではないんです。化学物質の吸収を妨げるのです。カロリーの摂取量と、または小児化により人間関係が希薄になると 普通の食事をしていただければそんな時には「運動をすると言います。気には3項です。春は辛い物が好きで肉 を食べるように注意しているかですね。スクワットはメタボリック症候群の話が記載されています。気の状態によって刺身、かに、ホタテ、イカなどを繰り返していたそう。銅は殺菌やにおい消しにいい食べ物を多く摂る ようにして癌に対する抵抗力がないかもしれませんよね・・・アトピーが良くなります。また身近なところ、私の 血液型は一生変わらないと駄目なようです。私は自覚するトラウマがないくらいです。ですので十分注意しない物質だ そうではほとんどが液体です。それからイカやタコなどの自己実現に取り入れたいものです。この添加物が 入っていれば、

407 文字

(実務での視点)

"1つ1つの文として問題ありなレベル。文章のつながりも不自然。

言葉の繋がりが不自然になってしまっている"

文章3 (自動要約) 5点

私の知り合いの老人 Yさんは現在 90 才の元気な男性。Yさんの健康法は毎日 2 時間くらいは散歩を続ける事だ そ

うです。それも晴の日だけでなく、雨の日も散歩に行かれますと言うのでびっくり。本人いわく「この年で仕事 もないので、私は散歩する事が仕事と思って毎日歩いているので、雨の日でも行きます。雨だから今日は仕事が 休みとは普通ならないでしょう・・・」との事でした。流石に脱帽です。実はこんな事があったそうです。お 医者さんから「もう 90 才になるのだから、あまり無理して歩かないほうがよいですよ。」と言われ、Yさんも「そうかな」と思い 1 ヶ月近く散歩を止めていました。そしたら、バス停から家までの道のり約 5 分くらいの 緩やかな坂道が、途中に一度休まないと息が切れて歩けなくなったそうです。それで「これではまずい!」と 思って、また歩き始めて 3 週間くらい歩き続けたら元に戻ったそうです。歩く事は健康の基本です。半身の静脈の 流れを良くし、身体の基礎筋肉を維持し、心肺機能を維持する事ができるのです。また、腰痛の 70% はしっかり歩くだけでも改善されています。現代は飽食による肝脂肪が増えています。私も最近は運動不足なので、昨年 の 10 月からは子供と毎月 1 回は山登りをするようにしています。皆さんも運動不足と思われる方は是非散歩をお勧め致します。毎日 1 時間は歩いてほしいですね 572 文字

(実務での視点)

語句の使い方や文章としてきわめて自然であり、前後の文脈もつながっている。この精度で文章生成であれば二重丸。

文章4 (自動要約) 2点

私の知り合いの老人 Yさんは現在 90 才の元気な男性。本人いわく「この年で仕事もないので、私は散歩する事が仕事と思って毎日歩いているので、雨だから今日は仕事が 休みとは普通ならないでしょう・・・」との事でした。お医者さんから「もう 90 才になるのだから、あまり無理して歩かないほうがよいですよ。そしたら、バス停から家までの道のり約 5 分くらいの 緩やかな坂道が、途中に一度休まないと息が切れて歩けな くなったそうです。それで「これではまずい!」と 思って、また歩き始めて 3 週間くらい歩き続けたら元に戻っ たそうです。半身の静脈の流れを良くし、身体の基礎筋肉を維持し、心肺機能を維持する事ができるのです。また、腰痛の 70%はしっかり歩くだけでも改善されています。私も最近は運動不足なので、昨年 の 10 月からは子供と毎月 1 回は山登りをするようにしています。

358 文字

(実務での視点)

"1つ1つの文としては問題がないレベル。評価を3にしよ うか迷った。

文章のつながり＝文脈が不明のため、明らかに全体の文章としては人間の目から見て不自然な繋がりが見受けられる。

"

文章5 (自動要約) 4点

今日は天の氣の話です。天の氣は太陽からの氣のエネルギーです。この太陽の氣のエネルギーも1日の中で氣の質が違います。一番良い氣は朝の氣です。ですから朝日に向かって拝むという昔からの習慣は、とても身体に良いのです。朝日の氣はプラスのエネルギーが強いからです。これは、地球が夜の状態から朝の状態になる時には、地球の陰(マイナス)の場所に太陽が当たって陽になる為、プラスのエネルギーを強く受けるからです。事実、漢方の王様である高麗人蔘を栽培しているところでは、黒い布のようなもので覆っていますが、一方向だけは開いているのです。それは東の方向と聞いています。つまり朝日だけが当たるようにして、育てているのです。だから高麗人蔘は漢方では補氣剤と言われているのでしょ う。昔は早寝早起きが一般的でしたが、最近では遅寝遅起が一般的になっています。ですから朝日を浴びることの少ない生活になっています。ですから、現代は地の氣が少なく、人の氣も少なく、天の氣も少なくなっているのです。したがって免疫力の低下は免れません。実はこの氣がとても大切なのです。この氣が滞って、氣が毒されている人が一番可愛そうな人なのです。それで、そういう人の事を「氣の毒な人」と呼ぶのです。どんどん早起きをして、朝日を浴びる生活をして健康に心がけましょう!

556 文字

文章6 (自動要約) 2点

風邪に注意しないといけない季節になりましたね。褐色脂肪細胞と言って発熱を促す細胞が多く存在している、肩・首・心臓・腰の下部を温めると、速やかに熱が生まれ身体全体が温まると言われています。ちなみに普通の脂肪は白色脂肪細胞と言います。ですから、マフラーやスカーフは首や肩からの発熱を促すので、しっかりと温めることができます。なんでも、衣服気候学ではマフラーには衣服1枚分の保温効果があるとされているのです。特にスカーフは軽やかさばらないので、寒い日にはいつも持ち歩くようにしたらとても重宝することでしょう。それから冬に定番のホカロンは腰の下側に当てて温めると、身体全体が温まって寒い冬にはもってこい です。また、ショールやベストも効果がありますので、自分の着こなしに合わせて冬に備えて準備したらよいでしょう。何事にも準備が大切ですから

359 文字

(実務での視点)

"1つ1つの文としては問題がないレベル。"

ただし文章のつながり=文脈が不自然。似たようなテーマの文章だが、文章の前後の繋がりが違和感がある。

No13の違和感のある文章接続が増えた印象。"

5. まとめ

文と文のつながりについては、自動要約との関連や文と文とのつながりを **entity-grid** を用いて局所的なつながりの良さを表現するなどの談話構造解析[9][10]があるが発展段階である。当面は制御文による文章自動生成が無難と思われる。

謝辞 文書衣自動生成は筆者の所属する(株)Speeの当時の上司の渡邊洋介氏から研究の機会を頂き、森リーダーの元で進めた。そしてSEOを加味した研究は本多執行役員及び今井リーダーの元で進められた。ここに謝意を表する。

参考文献

- [1] 佐藤理史 コンピューターが小説を書く日. 日本経済新聞出版社, 2016
- [2] Leon A. Gatys et al. A Neural Algorithm of Artistic Style, 2015
- [3] Wikipedia "https://ja.wikipedia.org/wiki/マルコフ連鎖"
- [4] H. P. Luhn. The Automatic Creation of Literature, IBM Journal, 1958
- [5] 長谷川, 平尾, 奥村, 永田. 文圧縮を活用したヘッドライン生成, 言語処理学会, 第23回年次大会発表論文集, 2017
- [6] 梁, 阿部川, 強化学習によるテキスト自動要約手法の提案. 言語処理学会 第18回年次大会発表論文集, 2012
- [7] 太田. 文章自動生成の事前調査報告書/ 最終調査報告書. 2017
- [8] 太田. SEOのための文章自動生成の事前調査報告書/ 最終調査報告書. 2017
- [9] 笹野, 飯田. 文脈解析, 自然言語処理シリーズ10, コロナ社, 2017
- [10] 黒橋. 自然言語処理, 放送大学教材, 2016