

学生の提出レポート解析に 文脈ベクトルを使う

吉 武 春 光

1. まえがき

筆者は、吉武(2016)において潜在意味解析の中の1つの手法である Word2Vec を使うための環境整備を行った。実際に Word2Vec を解析に用いた研究は、2015年から2016年にかけて、筆者が指導していた博士後期課程の学生と修士課程の学生が行った。今回、筆者は、Word2Vec の延長線上にある Doc2Vec という研究手法を用いて複数文の間の類似度を計算する研究を行うことにした。

Word2Vec は、語と語の間の関係をベクトル化するものであったが、文と文の間の関係をベクトル化する手法として開発されたのが Doc2Vec である。筆者は、吉武(2007, 2009)において、学生のレポート文章の解析について論じたが、意味記述の煩雑さのために、研究が進んでいなかった。これを打破するために、今回、Doc2Vec に注目した。

2. Doc2Vec による文脈ベクトルの生成

本章では、Doc2Vec の概要を述べた上で、Doc2Vec を使用するために必要となる、前処理、文脈ベクトル生成処理について述べる。

2.1 Doc2Vec の使い方

Doc2Vec (別名 paragraph2vec または sentence embeddings) は Word2Vec のアルゴリズムを『文、段落、文章などの連続する表現』に拡張したものである(<http://radimrehurek.com/gensim/models/doc2vec.html>)。

Doc2Vec の C言語版は存在せず、Python Gensim に実装されているもののみである。Gensim 上の Doc2Vec のクラスは基本的に Word2Vec のクラスを継承しているので、Python における Doc2Vec の使用法は、基本的に Word2Vec と同じである。

なお、Word2Vec に組み込まれているアルゴリズムは “continuous bag of words” (cbow) と “skip-gram” (sg) の2つであったが、Doc2Vec では “distributed memory” (dm) と “distributed bag of words” (dbow) という2つのアルゴリズムが組み込まれており、dbow が default になっている。

2.2 Doc2Vec の入力

Doc2Vec の入力はプログラミング言語 Python の LabeledSentence objects の iterator(イテレータ)となっている。Python の iterator とは、「複数の要素を持ったデータの要素を順番に取り出すとき、より汎用的な手法を提供する方法」のことで、具体的には、Python のタプルやリストや辞書のデータ型のどれでも良いということである。例えば、Doc2vec tutorial のページに挙げてある例では、英文の「some words here」は、クラス LabeledSentence を用いて、下記のように 'some' と 'words' と 'here' を変数 words にリスト構造として入れている。

```
sentence = LabeledSentence(words=[u'some', u'words', u'here'], labels=[u'SENT_1'])
```

上記では、英文の「some words here」に対して 'SENT_1' という文ラベルを付加している。

なお、各々の文に対して、文ラベルを考えるのが煩雑なので、各文に文ラベルを自動的に付与する TaggedLineDocument というクラスも用意してある。TaggedLineDocument の入力は、半角空白で区切られた語が連なったテキストファイルとなっており、文の識別のためのラベルは、先頭文が 0 である文番号を自動的に割り当てるようになっている。入力が複数行ある場合は、第 1 入力行の最後で改行し、次の行に第 2 文目を入れる。

解析対象の単位は『改行までの1文』なので、もし解析対象の文が文の途中で改行されている場合は、改行コードを取り除く必要がある。

2.3 ベクトル生成

Python プログラムに読み込ませた語の連なりに対して、クラス `Doc2Vec` を用いてベクトル生成をさせる。数多くの引数があるが、主要なものを次に挙げる。

本研究では Ubuntu 16.04 という OS 上で Python3.5 を使用したが、`Doc2Vec` のソースコードは `/usr/local/lib/python3.5/dist-packages/gensim/models/doc2vec.py` に入っている。

表1 Doc2Vec クラスの主要引数

引数名	説明
<code>documents</code>	入力語の連なり。
<code>dm</code>	default は 1 で “distributed memory” (dm) を使う。 0 を指定した場合は “distributed bag of words” (dbow) モデルを使う。
<code>size</code>	文脈ベクトルの次元サイズ (default は 300)。
<code>window</code>	文脈ベクトルを計算する際に、対象語から調べる前後の語の数。 大きな値を指定すると、より広範囲の語まで処理対象とする (default は 8)。
<code>hs</code>	default は 1 で hierarchical sampling を使う。
<code>min_count</code>	出現回数が <code>min_count</code> より少ない語は無視される (default は 5)。
<code>workers</code>	並列計算を行う際のスレッド数 (default は 1)。
<code>alpha</code>	initial learning rate (徐々に <code>min_alpha</code> に近づく)
<code>min_alpha</code>	最終的な learning rate

2.4 繰り返し学習

`Doc2Vec` では、10回～20回ほどの繰り返し学習が必要となっている。まず、クラス `Doc2Vec` の中のメソッド `build_vocab` を用いて内部にボキャブラリーを作っておいてから、クラス `Doc2Vec` の中のメソッド `train` を10回

～20回ほど実行させることが必要になる。

2.5 文の類似度の検索

学習が終わったデータに対しては、`docvecs` プロパティに対して、クラス `Doc2Vec` の中のメソッド `most_similar` を使って、類似度が高い文番号と、その類似度の値をタプルにしたものの、任意の上位 n 個(`topn` 個)のリストを得ることが出来る。

また、`docvecs` プロパティに対して、クラス `Doc2Vec` の中のメソッド `similarity` を使って、2つの文番号を指定すると、2つの文番号の間の類似度を得ることが出来る。

2.6 予備実験

吉武(2015) で用いた日本経済新聞社の本紙 (2014年1年分) を使って、予備実験を行った。記事を `MeCab` (工藤2002) を使って分かち書きし、語幹を半角スペース区切りにしたものである。`MeCab` を使って分かち書きを行うプログラムは、プログラミング言語 `Python` を使って自作した (`text2mecab.py` 付録1)。なお、`Doc2Vec` の意味の単位は、「入力された最初から改行コードまで」の文字列である。そこで、日経本紙では1つの記事が複数の文から構成されている場合は、文末のピリオド (句点) の位置で改行せずに、次の文をつなげ、記事の最後でのみ改行コードの挿入を行った。下記の、日経本紙2014年の第2記事 (記事番号1) では記事が3つの文から構成されていることが分かる。

但し、動作確認を目的とするために、第1記事 (記事番号0) をファイルの最後にコピー追加した (記事番号180497)。更に、第2記事を複製し、第2記事 (記事番号1) と第3記事 (記事番号2) を同じものとした。

静岡県は十日、ヒロインターナショナル (浜松市、谷口一博社長) が運営する飲食店「*」で米国産牛肉や豪州産牛肉などを松阪牛と表示しているたとして、景品表示法などに基づく改善指導を行うた。

図1 日経本紙2014年の第1文 (記事番号0) (MeCab 処理済)

勤務先から売上金を盗むだことを隠すため、強盗の被害に遭うたと虚偽の通報をするたとして、警視庁東大和署は二日、東京都武蔵村山市、アルバイト、A容疑者（48）を窃盗と軽犯罪法違反（虚偽申告）の容疑で逮捕するた。調べによると、A容疑者は先月三十一日午前一時ごろ、勤務先の同市内のガソリンスタンドから売上金約八十二万円入りのバッグを盗む、その後「ナイフを持った二人組の男にバッグを奪われるた」と一〇番通報をするた疑い。A容疑者がスタンド内に隠したバッグを別の店員が発見。同署が追及するたところ「借金で困るておる、生活費が欲しいてやるた」と認めるた。

図2 日経本紙2014年の第2文（記事番号1）（MeCab処理済）

次の図に、python コマンドを使った実行の様子を示す。

```
yositate@dual:/zpool/lz4/hy-svr.hdd/doc2vec/NKM_2004_all_test$ python3
Python 3.5.2 (default, Nov 17 2016, 17:05:23)
[GCC 5.4.0 20160609] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>> import gensim
>>> docsrc =
gensim.models.doc2vec.TaggedLineDocument('NKM_2004_all_test.txt')
>>> m = gensim.models.doc2vec.Doc2Vec(docsrc, size=300, window=10,
hs=1, min_count=2, workers=8, sample=1e-5, alpha=.025, min_alpha=.025)
ここで時間がかかる。
>>> m.docvecs.most_similar(0,topn=3)
[(180497, 0.8200676441192627), (5130, 0.4211444556713104), (46957,
0.41085657477378845)]
記事番号 0 に最も近いのは 記事番号 180497 と表示されており、近似度
は 0.82 であった。
>>> m.docvecs.most_similar(1,topn=3)
[(2, 0.8915579915046692), (7216, 0.4085385203361511), (15521,
0.4077302813529968)]
記事番号 1 に最も近いのは 記事番号 2 と表示されており、近似度は 0.89
であった。
そこで、繰り返し学習を 10 回、行わせた。
>>> m.train(docsrc,total_examples=m.corpus_count,epochs=10)
ここで時間がかかる。
132822934
>>> m.docvecs.most_similar(0,topn=3)
[(180497, 0.9427038431167603), (151616, 0.36943313479423523),
(165145, 0.36827799677848816)]
記事番号 0 に最も近いのは 記事番号 180497 と表示されており、近似度
は 0.94 に上がった。
>>> m.docvecs.most_similar(1,topn=3)
[(2, 0.9340255856513977), (172317, 0.34578341245651245), (114136,
0.34412682056427)]
記事番号 1 に最も近いのは 記事番号 2 と表示されており、近似度は 0.93
に上がった。
>>>
```

図3 日経データを使った予備実験

近似度は 0 から 1 の間の値を取り、1 に近い値ほど似ていることを示している。この結果を見ると、異なる記事の近似度は大きくても 0.34 程度であるが、同一記事（記事番号0 と記事番号180497）の近似度は 0.94 という 1 に近い値になったことが分かった。また、別の同一記事（記事番号1 と記事番号2）の近似度は 0.93 になったことも分かった。

以上の予備実験により、似ている記事を確かに検出できるという確信を得た。

3. 解析対象データと解析実験

さて、今回の実験対象は e-Learning システム Moodle 上の吉武の担当科目「情報ネットワーク論」に提出された課題である。[2012年度の課題：経路制御](133答案) と、吉武(2007)で使用したデータ [2005年度の課題：日本の文字コード](78答案)を使用した。

3.1 Moodle からの課題の抽出と前処理

Moodle の課題は、[オンラインテキスト] と [ファイル提出] という 2 種類の形式がある。

3.1.1 [オンラインテキスト] 形式の場合

[オンラインテキスト] 形式は、ダウンロードしたら HTML 形式になっている。その中から HTML タグと呼ばれる制御コードを取り除く必要がある。更に、2012年度から運用している Moodle バージョン2 以降では内部漢字コードが UTF8 になっているが、2011年度まで運用していた Moodle バージョン1.9 では内部漢字コードが EUC になっているために、計算機処理するためには、漢字コードを EUC から UTF8 に変換する必要がある。[2005年度の課題：日本の文字コード] の全78答案に対して HTML タグを取り除き、更に UTF8 に変換を行い、更に MeCab を使って分かち書きし語幹に変換を行った。以下は処理済みの一例である。

日本の文字コードは、16ビットを使用するASCIIコードがある、ASCIIコードにヨーロッパの各国用の拡張を施したISO 8859*1、ASCIIコードに拡張を施したJIS*0201の8単位符号がある。またJIS*0201の8単位符号にはISO 8859*1の拡張部分と重なるコードがある。半角カタカナコードの問題点は電子メールで文字化けすることである、半角カタカナに所属する記号も文字化けするので注意が必要である。

図4 [2005年度の課題：日本の文字コード]の一例 (MeCab 処理済)

3.1.2 [ファイル提出]形式の場合

[ファイル提出]形式は、アップロードしてあるファイル形式のままダウンロードされるが、Word 2007以降で標準のdocx形式以外に、Word 2003までで標準のdoc形式もある。

docx形式のWordファイル中の文字を抜き出すためには、2つの方法がある。

1) docx形式のWordファイル中の文字を抜き出すサービスを行っているサイトを利用する方法。

処理対象ファイルを1つずつ指定するため、処理に時間がかかる。

2) プログラミング言語Perlで書かれたdocx2txt.plを使う方法。

複数のファイルを次々に処理できる。

doc形式のWordファイル中の文字を抜き出すためには、プログラミング言語で書かれた処理プログラムも存在するのだが、試してみたら問題が生じたために1)のサイトのみを使用した。

docx形式のWordファイルから文字を抜き出すためには1)で全てのファイルを処理させ、失敗したもののみを2)で処理した。

[2012年度の課題：経路制御]の全141答案(docx形式とdoc形式が混在)に対して上記の文字を抜き出す処理を行い、更にMeCabを使って分かち書きし語幹に変換を行った。以下は処理済みの一例である。

経路制御とは、インターネット上で IP パケットを特定の目的地に転送するための、パケットの通り道（経路）についての情報を管理する、複数ある経路のうちから、最適だ経路を選択する仕組みのことをいう。まず、経路制御の必要性だがルータのネットワークインタフェースには、IP アドレスとネットマスクが管理者により設定される。その場合、該当インタフェースに直接接続されているネットワークについては、経路表に自動的に掲載される。しかし、ルータのインタフェースに直接繋がっていない離れたネットワークの場合、その離れたネットワークがどこに（どのインタフェースの方向に）存在するのか、直接的な判断は不可能である。

図5 [2012年度の課題：経路制御]の一例 (MeCab 処理済)

3.2 実験1 [2012年度の課題：経路制御]

まず、MeCab 処理済の133答案の全てに対して、各々、意味ベクトルの生成を行ってから、繰り返し学習を行い、その結果に対して、各答案に対して最も近似度が高い答案を調べる Python プログラムを作成した(d2v_out.py 付録2)。繰り返し学習の回数は20回とした。

3.2.1 実験結果

得られた近似度を降順にソートした結果は次の通りである。

答案番号	答案番号	近似度
37	35	0.897456169
35	37	0.897456169
65	10	0.796049595
10	65	0.796049595
74	62	0.665772617
62	74	0.665772617
...
18	0	0.284812748
50	39	0.280477494
130	87	0.277597368
17	14	0.262090504

答案番号 37 に最も近似度が高いのは 答案番号 35 だと判り、その時の近似

度は 0.897456169 となっている。また、答案番号 65 と答案番号 10 との近似値は 0.796049595 となっている。答案番号37、答案番号35、答案番号65、答案番号10 の全記事は次の通りである

経路制御は根幹をなす技術である、大変重要な役割を持っている。経路制御の必要性を挙げていく。まずは一つ目はインターネットワークは、網の目状のようになるている。よって、目的の場所へ行くために、可能だ経路はいくつも存在することである。いくつもの経路が存在するということは、すなわちどの経路を通るかということを決める必要がある。二つ目はもしも経路に何らかの問題が発生する、不通になっている場合に、どのような対処を取るかを考える必要がある、ルーティングは経路の情報をあらかじめネットワーク機器に設定しておくスタティックルーティングと、経路情報を常に更新するダイナミックルーティングとに分かれる。

簡単に原理を説明するていくと、TCP/IP ネットワークでパケットを送るうとするとき、経路制御表（ルーティング・テーブル）を参照するてパケットの通り道（経路）についての情報を管理する、複数ある経路のうちから、最適だ経路を選択する。ある端末から他の端末へとパケットを送るうとする場合、目的の端末が自ネットワーク内にない場合、端末内にある経路制御表を参照する、パケットを中継するせる端末を決定する。経路制御表はインターネットの状態を示すものだ、刻々と姿を変えている。そのため問題が存在するても、また他の経路から通じるということ。

経路制御表（ルーティング・テーブル）とは、宛先アドレスにたどり着くためには次にどのルータに行くば良いのか、を書いた表のことである、駅の乗り換え時刻表の・ようなもの。経路制御（ルーティング）にはタイプがある、*ホスト・ルーティング*ネットワーク・ルーティング*デフォルト・ルーティングが存在する。他ネットワーク間での接続にはルータが用いられるが、IP アドレスを見て相手のホストを探す接続することを*ホスト・ルーティングという。しかしその IP アドレスを持つホストを含むネットワークのルータを探して、ルータ同士を接続する方法*ネットワーク・ルーティング）を取る必要がある。単純には、すべてのネットワーク*アドレスとそのルータの*アドレスの一覧表を持つばよいのですが、世界中には膨大なネットワークがあるので、そのようなことは不可能だためルータによく接続する相手、最近接続する相手については経路制御表（ルーティングテーブル）を作成するておく、それに合致するたものはそれを利用する、そのリストにない場合には上位のルータに接続する*デフォルトルートという）。経路制御表において幾つかの経路の候補がある場合に、どの経路を選択するか決める手順をルーティング・アルゴリズムという。手順には静的（スタティック）と動的（ダイナミック）が存在する。静的（スタティック）だルーティングはルータなどが、管理者が予め設定するた固定だ経路表（ルーティングテーブル）に基づいて経路選択を行うことである。要は人がや管理することだ、手軽だは会うが維持管理が大変だある。それに対して、動的（ダイナミック）だルーティングはルータなどが経路情報を交換する合う、自動的に生成・更新する続ける経路表*ルーティングテーブル*に基づいて経路選択を行なうこと。

図 6 [2012年度の課題：経路制御]の記事番号 37 の全記事 (MeCab 処理済)

情報ネットワーク論（月2）

AG ○○○

経路制御は根幹をなす技術である、大変重要な役割」を持っている。経路制御の必要性を説明する。まず一つ目にインターネットは、網の目上になっているということである。よって、目的場所に行くために可能な経路はいくつも存在するという事である。いくつもの経路が存在するという事は、すなわちどの経路を通るのかということを決める必要がある。二つ目はもしも経路に何らかの問題が発生する*不通になっている場合に、どのような対処を取るかを考える必要がある、ルーティングは経路の情報をあらかじめネットワーク機器に設定しておく*と、経路情報を常に更新する動的ルーティングとに分かれる。

簡単に原理を説明していくと、TCP/IP ネットワークでパケットを送るうとするとき、経路制御表（ルーティング・テーブル）を参照してパケットの道路（経路）についての情報を管理する、複数ある経路のうちから、最適な経路を選択する。ある端末からほかの端末へパケットを送るうとする場合、端末内にある経路制御表を参照する、パケットを中断させる端末を決定する。経路制御表は「インターネットの状態を示すもの、姿を次々と変えている。そのため問題が存在するても、またほかの経路から通じるということである。

経路制御表（*）とは、宛先アドレスにたどり着くためには次にどのルータに行くか良いのかを書いた表のことである、駅の乗り換え時刻表のようなもの。経路制御（*）にはタイプがある、*ホスト・ルーティング*ネットワークルーティング*デフォルトルーティングが存在する。他ネットワーク間での接続ではルータが用いられるが、IPアドレスを見て相手のホストを探す接続することを*ホスト・ルーティングという。しかしそのIPアドレスを持つホストを含むネットワークのルータを探す、ルータ同士を接続させる方法*ネットワーク・ルーティング）を取る必要がある。単純には、全てのネットワークIPアドレスとそのルータのIPアドレスの一覧表を持つ良いのですが、世界中には膨大なネットワークがあるため、そのようなことは不可能だためルータによく接続する相手、最近接続した相手については経路制御表（ルーティング・テーブル）を作成しておく、それに合致するものはそれを利用する、そのリストにない場合には上位のルータに接続する*デフォルトルートという）。

経路制御表においていくつかの経路の候補があるばあいには、どの経路を選択するか決める手順をルーティング・アルゴリズムという。手順には静的（スタティック）と動的（ダイナミック）が存在する。静的（スタティック）なルーティングはルータなどが、管理者があらかじめ設定した固定した経路表（ルーティング・テーブル）に基づく経路選択を行うことである。要は人が管理することだ、手軽にはあるが維持管理が大変だのが欠点である。それに対して、動的（ダイナミック）なルーティングはルータなどが経路情報を交換する合、自動的に生成・更新する続ける経路表（ルーティング・テーブル）に基づく経路選択を行うことが大事だ。

図7 [2012年度の課題：経路制御]の記事番号35の全記事（MeCab 処理済）但し、答案中に氏名が入っていたので、氏名を○に変更してある。

経路制御の必要性・原理・ルーティングアルゴリズムについて

AF ○○ ○○○

経路制御の必要性についてだが、まずインターネットは網の目状になっているため、自分がつながりたいと思う相手のもとに行くためには様々な網の目の中からそして様々な経路の中から探し出す必要になる。そのために経路制御を利用することだ、どの経路を通るべしだのかを判断する必要がある。また、もしもその経路が不通になっている場合にはどのようにして対処するべしのかも考えるておく必要がある。そこに動的な経路制御アルゴリズムを採用することによって、動的に制御することができる。これが必要性としてあげられる。

次に経路制御の原理についてだが、経路制御は電車の乗換駅に例えることができる。インターネットにつながるコンピュータは駅に該当する、駅には乗換のために時刻表が置かれている。これらの前提を踏まえて、インターネットのネットワークでは、パケットを送るときにあらかじめ路線や時刻を調べるということを行わない。まず最初にパケットをインターネットに送り出す。そうすると、送り出されるパケットは、駅にて時刻表に相当する経路制御表を調べだし、次にどの経路に乗るべしかを判断する。その経路の終点にたどり着くと、そこで再び経路制御表を自ら調べだし、次に乗るべし経路を設定する。この作業を繰り返す、目的の地にとり着くのである。この、経路制御表はインターネットに路線の状態を示すものがある、刻々と変化している。もし事故（不通）のために途中で路線が通れないという状況になっても、この経路制御表が無事に通れる経路を常に示してくれるので、それに従って迂回することができるのである。

それではこの経路制御表の構成についてだが、この経路制御表は、宛先アドレスにとり着くためには次にどのルーターに行くべしのかを書く表である。種類としては、デフォルトルーティング・*が種類としてあげられる。実際には、このデフォルトルーティング・*を組み合わせて使用する。

最後にルーティングアルゴリズムについてだが、ルーティングアルゴリズムとは経路制御表において経路の中にいくつかの候補がある場合に、どの経路を選ぶべしだのかを決める手順のことである。この手順には二つのルーティングが存在している、それが静的ルーティングと動的ルーティングである。静的ルーティングは人間が行先・経由・セグメント生きのルーター経由などを指示する。静的ルーティングには周りのネットワーク環境の変化に応じてそれに追従して設定を行う必要がある場合がある、矛盾のないルーティングを行うように気を配る必要がある。それは、比較的気軽に利用できるという利点があるが、維持管理が大変であるという欠点もある。

次に動的ルーティングであるが、この動的ルーティングは、自ら宛先への距離情報を隣合うルーター同士で交換する、これにより経路制御表を完成させる、票の中で最も距離が短いものを選択する機能を持つ。この「距離」として定義されるものの違いでその種類ごとに距離情報交換用のプロトコルが存在する。

図8 [2012年度の課題：経路制御]の記事番号65の全記事 (MeCab 処理済) 但し、答案中に氏名が入っていたので、氏名を○に変更してある。

経路制御の必要性・原理・経路制御表・ルーティングアルゴリズムについて

AF ○○ ○

経路制御の必要性についてであるが、まずインターネットというものは網の目状になるため、ある相手の場所に行くためには様々な網の目の中からいくつもある経路から探し出す必要がない。そのために経路制御を利用することだ、どの経路を通るかを判断する必要がある。また、もしもその経路が不通になる場合にはどのようにして対処をするべしだのかを考えるておくべしであるが、そこに動的な経路制御アルゴリズムを採用することだ、動的に制御を行うことができる。これが第一に必要性としてあげられる。

次に経路制御の原理についてであるが、経路制御は鉄道の乗換駅に例えられることとなる。さらにインターネットにつながるコンピュータは駅に該当する、駅には乗換のために時刻表がおくである。これら前提を踏まえて、インターネットワークでは、パケットを送る時に、あらかじめ路線や時刻を調べるたりするというを行うない。まず最初にパケットをインターネットに送り出す。そうすると、送り出されるパケットは、駅にて時刻表に相当する経路制御表を調べる出す、次にどの経路に乗るべしかを判断する。その路線の終点にたどり着くと、そこで再び経路制御表を自ずから調べだし、次に乗るべし路線を決定する。この作業を繰り返す、目的地にたどり着くのである。この、経路制御表はインターネットに路線の状態を示すものがある、刻々と変化している。もし事故（不通）のために途中で路線が通れないという状況になるても、この経路制御表が無事に通れる経路を常に示してくれるので、それに従って迂回することが未然にできるのだある。

それではこの経路制御表の構成についてであるのだが、この経路制御表は、宛先アドレスにたどり着くためには次にどのルーティングに行くべしのかを書いた表である。種類としては、ホストルーティング・*が種類としてあげられる。実際には、このホストルーティング・*を組み合わせて利用する。

最後に、ルーティングアルゴリズムについてであるが、ルーティングアルゴリズムとは経路制御表においていくつかの経路の候補がある場合に、どの経路を選択するかを決める手順の事である。これにはふたつのルーティングが存在する、静的ルーティングと動的ルーティングがある。静的なルーティングは人間が行く先、経由、セグメント行きのルータ経由などを指示する。静的なルーティングには、回りのネットワーク環境の変化に応じてそれに追従して設定を行う必要がある場合がある、矛盾のないルーティングを行うように気を配る必要があるが、比較的気軽に利用できる。しかし、維持管理が大変であるという問題点もある。

次に動的なルーティングであるが、この動的ルーティングは、自ずから宛先への距離情報を隣合うルータ同士で交換する、これにより経路制御表を完成させる、表の中で最も距離が短いものを選択する機能を持つ。この「距離」として定義されるものの違いでその種類ごとに距離情報交換ようのプロトコルが存在する。

図9 [2012年度の課題：経路制御]の記事番号10の全記事（MeCab処理済）但し、答案中に氏名が入っていたので、氏名を○に変更してある。

3.2.2 実験結果の考察

原文を見ると分かるように、答案番号35と37は酷似した内容であった。また、答案番号10と65も酷似した内容であった。[2012年度の課題：経路制御]は、講義資料を元にまとめるだけで答案が作成出来るものであったので、ある程度、答案が似ても仕方ないと思われる。しかし、これだけ似ていると答案を写したのかと疑いたくなってしまふ。

3.3 実験2 [2005年度の課題：日本の文字コード]

次に、吉武(2007)で使用したデータ [2005年度の課題：日本の文字コー

ドJ(78答案)に対して、MeCab 処理を行い、実験1で使用した Python プログラム (d2v_out.py)を使って、意味ベクトルの生成を行ってから、繰り返し学習を行った。

3.3.1 実験結果

次に、繰り返し学習を20回と指定して、近似度を計算させた。得られた近似度を降順にソートした結果は次の通りである。

答案番号	答案番号	近似度
42	1	0.979680717
1	42	0.979680717
44	3	0.978571534
3	44	0.978571534
69	47	0.975920677
47	69	0.975920677
...
12	25	0.759830475
28	64	0.735853076
38	32	0.709882736
10	38	0.616832972

答案番号 42 に最も近似度が高いのは 答案番号 1 だと判り、その時の近似度は 0.979680717 となっている。また、答案番号 44 と答案番号 3 との近似値は 0.978571534 となっている。答案番号 42、1、44、3 の全記事は次の通りである。

日本の文字コードは、シフトJIS、EUC-JP、ISO*2022*JPがあります。ASCIIコードにヨーロッパ各国の言語を拡張するた、ISO 8859*1コードがある、日本にはASCIIコードにカタカナ用の拡張を施したJIS*0201があります。しかし、それぞれの拡張部分には重なるところがある、ISO 8859*1とJIS*0201の8単位符号では同時に使用できないという問題が生じます。JIS*0201では「*」と表示されるものが、ISO 8859*1では「*」で表示されるので、電子メールなどを使用するとJIS*0201の重なるているコードが文字化けをおこすことがあります。このようだ問題によりインターネットの世界では半角カタカナ、重なりコードを使用することは禁止されるています。

図10 [2005年度の課題：日本の文字コード]の記事番号42の全記事
(MeCab 処理済)

日本の文字コードは、16ビットを使用するたASCIIコードがある、ASCIIコードにヨーロッパの各国用の拡張を施したISO 8859*1、ASCIIコードに拡張を施したJIS*0201の8単位符号がある。またJIS*0201の8単位符号にはISO 8859*1の拡張部分と重なるコードがある。半角カタカナコードの問題点は電子メールで文字化けすることだある、半角カタカナに所属する記号も文字化けするので注意が必要だある。

図11 [2005年度の課題：日本の文字コード]の記事番号1の全記事
(MeCab 処理済)

日本の文字コードの種類は、*シフトJIS*EUC-JP（UNIX用*ISO*2022*JP（電子メール用）半角カタカナと呼ばれるコードの問題点は、電子メールに入れるた半角カタカナは文字化けを起こすものが多い。インターネットの世界では使用が禁止されるている。

図12 [2005年度の課題：日本の文字コード]の記事番号44の全記事
(MeCab 処理済)

日本語の文字コードJIS X 0201*ASCIIコードにカタカナ用の文字コードを施したもの*JIS X 0208（漢字コード）、シフトJIS、EUC-JP、ISO-2022-JP
問題点 半角カタカナはJIS X 0201の*単位符号の右半分に表記されるている文字（つまり半角カタカナと「」、。などの記号）を入力すると、それがインターネットの世界で表示されるた際に文字化けをするため、使用禁止になるている。

図13 [2005年度の課題：日本の文字コード]の記事番号3の全記事
(MeCab 処理済)

3.3.2 実験結果に対する考察

残念ながら、答案番号 42 と答案番号 1 の近似度が高いというのは、納得できないものであった。更に、答案番号 44 と答案番号 3 の近似度が高いというのも、納得できないものであった。

その原因であるが、[2005年度の課題：日本の文字コード]は、課題がサンプルであったので、どの答案も似通っていて文字数が少ないものが多かった。一方、実験 1 [2012 年度の課題：経路制御] では答えるべき内容が多かったので、文字数が多い答案が多かった。そのため、実験 1 [2012 年度の課題：経路制御] があれだけ高精度の近似度を出していることから判断して、文字数が少ない答案は深層学習の多次元ベクトル計算には向いていないと思われる。

4. あとがき

本研究では、前講で述べた Word2Vec を発展させた Doc2Vec を使い、実際に学生が提出したレポートの文面の近似度を調べた。確かに、近似度が 0.9 ぐらいと判断されたレポートは、目視でも文面を僅かだけ修正したものではないかと推測された。長らく進展していなかった近似度の研究が、新しい手法で大きく進展した意義は大きい。

実験 1 と実験 2 で使用した Doc2Vec のパラメータは、試行して良い結果が得られた値を使用した。パラメータ決定に関しては、検討が必要かもしれない。

なお、Deep Learning は、計算させる度に結果の僅かな変動が生じるものらしい。実際の実験結果でも、小数点以下 5 桁程度は変動しているようである。なので、小数点以下 3 桁程度を有効値と見なす方が良さそうである。これは、気持ち悪いものであるが、仕組み上、仕方ないそうである。

本研究では Doc2Vec の有効性を確認できたが、Deep Learning の手法は、他にも幾つもある。今後は、他の手法へも挑戦したい。また、Deep Learning の手法を画像処理などにも適用してみたい。

参考論文

吉武春光(2007): "SDRTによる談話の意味記述", 西南学院大学商学論集, Vol. 53, No.3 & 4, pp. 211-238, 2007.2月

吉武春光(2009): "Prologを使ったSDRT表記の実現", 西南学院大学商学論集, Vol. 55, No. 4, pp. 257-304, 2009.3月

吉武春光(2016): "日経記事の解析に文脈ベクトルを使うための環境整備", 西南学院大学商学論集, Vol. 62, No.3&4, pp. 263-284, 2016.3月

工藤 拓(2002): MeCab: Yet Another Part-of-Speech and Morphological Analyzer

<http://mecab.googlecode.com/svn/trunk/mecab/doc/index.html>

Genism <https://radimrehurek.com/gensim/>

RadimRehurek: Gensim のDoc2Vec サイト

<http://radimrehurek.com/gensim/models/doc2vec.html>

RadimRehurek: Doc2vec tutorial

<https://rare-technologies.com/doc2vec-tutorial/>

付録

1) text2mecab.py

```
#!/usr/bin/env python3
import MeCab
import csv
import sys
args = sys.argv
t = MeCab.Tagger("")

with open(args[1]) as csvfile:
    dataReader = csv.reader(csvfile, delimiter=',')
    for row in dataReader:
        if row[0] != "articleid":
            result = t.parse(row[6])
            for col in result.split('\n'):
                if col != "EOS" and col != "":
                    parts = col.split(',')
                    print(parts[6], end=' ')
            print()
```


2) d2v_out.py

```
#!/usr/bin/env python3
import gensim
import sys
import csv

docsrc = gensim.models.doc2vec.TaggedLineDocument('./2014-all.mecab')
m = gensim.models.doc2vec.Doc2Vec(docsrc, size=300, window=10, hs=1,
min_count=3, workers=4, sample=1e-5)

print('¥n 繰り返し学習の開始')
for epoch in range(20):
    print('第 {} 回'.format(epoch + 1))
    m.train(docsrc, total_examples=m.corpus_count, epochs=20)
    m.alpha -= (0.025 - 0.0001) / 19
    m.min_alpha = m.alpha

length = len(list(m.docvecs))
print("max line number=", length)
f = open('out.csv', 'w')
writer = csv.writer(f, lineterminator='¥n')
allout = []

for x in range(length):
    each = m.docvecs.most_similar(x, topn=1)
    out = (x, each[0][0], each[0][1])
    allout.append(out)
    writer.writerow(out)
f.close()

from operator import itemgetter
allout.sort(key=itemgetter(2,0), reverse=True)

f = open('out_sorted.csv', 'w')
writer = csv.writer(f, lineterminator='¥n')
```