

二つの時系列データを対象とした特定着目点の動向についての記述文生成

中野由加子¹ 小林一郎¹

¹ お茶の水女子大学

{g1920532,koba}@is.ocha.ac.jp

概要

近年、様々な分野において数値データの収集が容易になり、表や時系列チャートなど多様な表現形式での数値データについて自然言語で記述する Data-to-Text の研究が注目されている。Data-to-Text の研究のうちデータの分析的な意味について説明する研究においては、一つの時系列データの動向について説明する研究は行われている。その一方で時系列データ全体の動向を捉えることが可能でかつ時系列データ同士を時系列で比べて関係性を捉える研究は、著者らが知る限りほぼ提案されていない。本研究は、二つの時系列データにおいて特定着目点における動向を説明する自然言語文生成手法を提案する。

1 はじめに

技術の発展に伴い様々な分野における数値データの収集が容易になりつつあるが、大量のデータから重要な部分を見つけ出し、データを理解することが困難な場合は少なくない。データから重要な情報を抽出して自然言語文で説明することができれば、データをより深く理解することが容易になる。そのため近年、表や時系列チャートなど多様な形式で表現されている数値データについて自然言語で記述する Data-to-Text の研究が盛んに行われている [1, 2, 3]。多くの Data-to-text の研究が、特定のドメインにおけるデータを説明する自然言語と同様な自然言語文を生成することを対象にしている [4, 5] が、数値データの分析を通じて特徴を捉え、それについて説明することもデータの内容を容易に理解することに大いに役に立つ [6, 7]。実際、説明対象である時系列データについての研究は異常検知 [8] や株価変動の傾向予測 [9] など、数値情報の分析に関する研究が多い。これらは分析することに特化し

ており、自然言語で説明をするということは対象にしていない。これに対し、観測データにおいて異常な動向が見られた際に異常について自然言語文で説明を提示することができれば、障害に対して迅速に対応することが可能になる。また、時系列データにおける急激な増加や減少などの動向について、原因や背景情報とともにその理由に対する説明文を生成することができれば、より容易にデータの挙動について理解することができる。Data-to-Text の分野において、時系列データの動向について説明をする研究も行われているが、これらは一つの時系列データを対象としている [6, 7]。これに対して本研究では、とくに二つの時系列データにおける関係性を捉え特定の着目点からそれらの動向を説明する自然言語文生成を行う。また、提案モデルの訓練のために二つのデータセットを提案し、作成した。

2 時系列データにおける説明対象

動向の関係性 時系列データにおいて人が顕著に感じる動きとして、「増加 (increase)」、「減少 (decrease)」、「ピーク (peak)」、「ディップ (dip)」の4つが挙げられる (cf. [6])。本研究では、二つの時系列データの振る舞いをこれらの数値的变化量がある4種類の動向の組み合わせで構成する「協調」または「反対」の関係及び上記の4種類の動向と数値的变化をほとんど取らない「平坦 (flat)」の組み合わせである「片方平坦」の関係を説明対象とする。

動向の時間的關係 二つの時系列データにおける動向の数値的变化量を取る範囲の比較結果を説明文に反映させた。片方の時系列データの動向が数値的变化を取り終わった地点が他方の動向の数値的变化の開始点よりも前の場合に「before」や「after」を用いて説明を行う。説明文における接続詞に、二つの時系列データにおける動向の時間的關係が反映される。

表 1 接続詞の条件

関係	end1 < start2	end2 < start1	左二つ以外
協調	before	after	and
反対	after	before	while
片方平坦	while	while	while

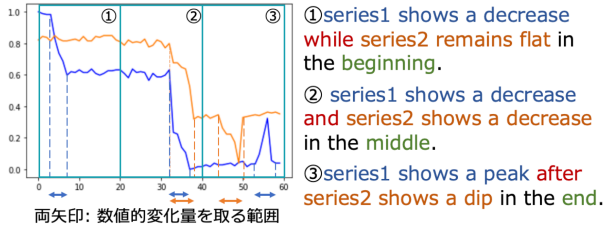


図 1 単語を置き換えた具体例

3 データセット

特定の時期の動向についての説明のために、2 種類のデータセットを作成した。一つ目のデータセットは時系列データと正解文章で構成され、二つ目のデータセットは時系列データと時期の指示文と正解文で構成される。正解文章はどちらも英語で記述されている。最初にそれぞれのデータセットの共通部分について説明をした後にそれぞれのデータセットについて説明をする。

時系列データ 時系列データ長を 60 とし、0-19 を「初期」、20-39 を「中期」、40-59 を「末期」とする。それぞれの時期（3 種類）において 16 種類の動向の組み合わせ（協調 4 種類、反対 4 種類、片方平坦 8 種類）のうちのひとつをとり、時系列データが全体として取る動向の種類は $16^3 = 4,096$ 通りとなる。数値的变化量を取る上記の 4 種類の動向を示す範囲以外では大きな数値的变化量を取らない。

正解文 正解文はテンプレートを用いて、動向/時期/接続詞の単語を置き換えることで生成した。表 1 に接続詞の条件を示す¹⁾。一つの時期の動向における説明文テンプレートを以下に示す。

- 協調/反対の動向の場合
「Series1 shows a/an (動向 1) (接続詞) series2 shows a/an (動向 2) in the (時期).」
- 片方平坦の場合
「Series(1/2) shows a/an (動向 1) (接続詞) series(1/2) remains flat in the (時期).」

図 1 にテンプレートにおいて単語を置き換えた具体例を示した。初期においては減少/平坦の動向を

1) (1/2) は 1 または 2 を指す。series(1/2) の動向における数値的变化の開始点、end(1/2) は series(1/2) における数値的变化の終了点を表す。

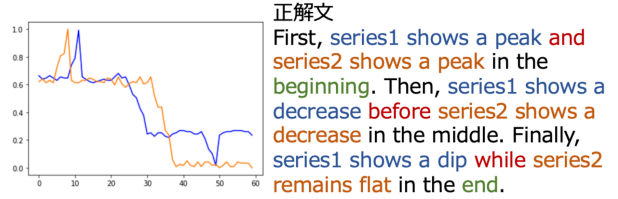


図 2 データセット 1

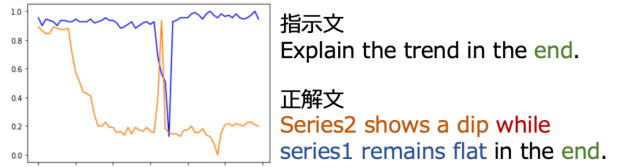


図 3 データセット 2

取り、接続詞は while をとる。中期においては共に減少の動向を取り、数値的变化を取る時期が被っているため接続詞は and となる。末期においてピーク/ディップを取り、series2 のディップが数値的变化を取り終わった後に series1 のピークが数値的变化を取り始めるため接続詞は after となる。

3.1 データセット 1

時系列データと正解文章のペア (図 2) で構成する。正解文章は初期/中期/末期の全ての動向について説明をする。図 2 の例のように動向についての説明を、First, Then, Finally で繋いでいる。

3.2 データセット 2

時系列データ、指示文、正解文のペア (図 3) で構成する。指示文は説明対象の時期を指示する。こちらもテンプレートを用いて作成されており、「Explain the trend in the (時期).」の (時期) の単語を置き換えることで生成する。正解文は指示文で指定された時期における動向についての説明を行う。こちらはデータセットの共通部分で説明をした正解文となる。図 3 においては、末期における動向について説明するよう指示する文を受けて、末期の平坦/ディップについての説明文を正解文とする。

4 提案手法

4.1 提案モデル

モデルは Transformer[10] を拡張して構築されており (図 4)、入力である時系列データと指示文から特定の時期の動向についての説明文を生成する。二つの動きの比較を可能にするために positional

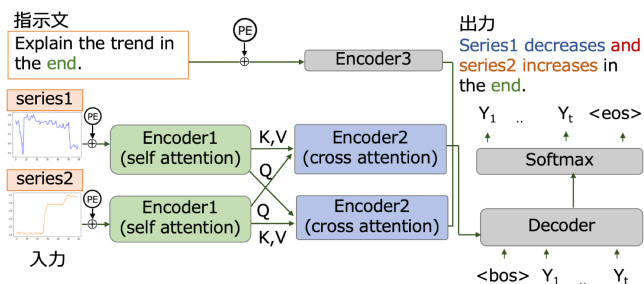


図4 提案モデルの概要

embedding をそれぞれの時系列データに追加することにより時刻を共通化している。エンコーダ部分では、まず Encoder1 でそれぞれの時系列データの特徴を捉えた後に、Encoder2 で二つの時系列データの関係性を捉える。こちらの二つのエンコーダ層の違いは、Encoder1 は transformer の self-attention メカニズムで生成されるベクトルである key, query, value は全て入力と同じ、つまり自分自身とする。一方で、Encoder2 は key と value は先ほどと同じである一方で query については series1 と series2 を入れ替えることで、二つの時系列データの関係性を捉えられるようにしている。Encoder3 は、指示文の解析を行う。Encoder2 と Encoder3 の出力を concat し、デコーダに入力する。デコーダは時系列データの解析結果と指示文の解析結果から特定の時期の動向についての説明文を生成する。

入力形式 時系列データの振る舞いを捉えるために、「スライディングウィンドウ (sliding window)」と呼ばれる、ある程度の範囲を時間方向に少しずつシフトさせた特徴量によって時系列データを表現することが多い。本研究においてもそのことを踏襲し、入力形式として時系列データにスライディングウィンドウを適用したものを採用する。ウィンドウサイズを K とし、時刻 t における値を $W_t = \{x_{t-K+1}, x_t, \dots, x_t\}$ とする。このとき、入力は、 w_t のリストとなり、長さ T の時系列データの場合、 $W = \{W_1, W_2, \dots, W_T\}$ と表される。この形式をとることで、変化量がわかりやすくなり、時系列データの動向が捉えやすくなる。本研究では、この値に続けてそれぞれの値の差を追加することでモデルに与える情報を増やし、精度を高めた。

4.2 訓練手法

本研究では2つのステップで訓練を行う。1段階目に、すべての時期における動向について説明をしている情報量の多いデータセット1での学習を行

表2 実験設定

	step	epoch
訓練1	2000	50
訓練1&2	1000	100
訓練2	1000	60

表3 訓練設定

Embedding	128
隠れ層	512
損失関数	cross entropy
勾配法	Adam
学習率	0.0001
ドロップアウト	0.1

表4 実験結果

	BLEU	ppl.	T	P	C	T/P	T/C	P/C	T/P/C
訓練1	96.0	1.53	0.938	-	0.865	-	0.938	-	-
訓練1&2	78.3	1.63	0.731	0.795	0.70	0.704	0.609	0.620	0.582
訓練2	77.7	1.65	0.645	1.0	0.657	0.645	0.500	0.657	0.473

う。この訓練を通してモデルはそれぞれの単語における数値的な意味を学習し、特定の時期の動向説明のための単語のベクトル表現を獲得する。獲得したベクトルは、Encoder2 とデコーダで共有する。2段階目に、指示文で言及された時期における動向について説明可能とする訓練を行う。こちらの訓練で、指定された時期における動向を捉え記述可能とする。本研究では、2段階で訓練を行なった場合と2段階目のみで訓練を行なった場合の比較を行う²⁾。

5 実験

5.1 実験設定

それぞれのデータセットは172,032 データで構成されており、それぞれが train, test, valid で 17,0112 : 1,152 : 768 ずつとなる。データセット2の正解文で言及されている傾向(16種類)と時期(3種類)の種類の比率が均一になるように設計した。訓練では特定の step ごとに評価データを用いて文生成を行い、BLEU 値が最大の時のモデルを実験に使用した(表2)。入力のウィンドウサイズを8、バッチサイズを8とする。訓練設定は表3に示す。

5.2 評価手法

自動評価手法として BLEU[11] と perplexity(ppl.) を使用した。また、正解文と生成文における動向(trend)、時期(period)、接続詞(conjunction)の単語の正解率についても評価対象とした。

5.3 実験結果及び考察

訓練1においては高い精度を示しており、それぞれの単語における数値的な意味を学習し、特定着目点の動向についての説明のための単語のベクトル

2) 1段階目の訓練をした場合を訓練1、2段階で訓練した場合を訓練1&2、2段階目のみで訓練を行なった場合を訓練2と表す。

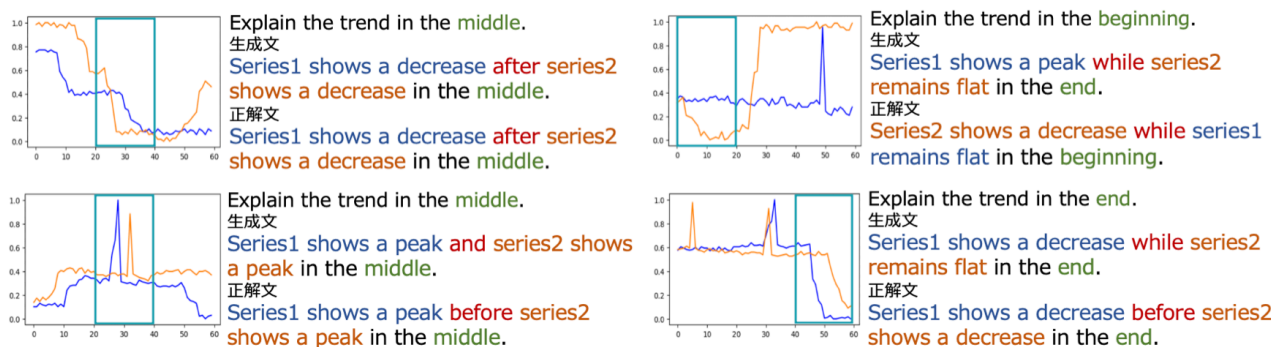


図5 生成結果例; 上部が訓練 1&2, 下部が訓練 2 で訓練したモデルの生成結果例

表5 接続詞/動向の正解確率

	before	after	and	while	協調	反対	片方平坦
訓練 1	0.878	0.663	0.678	0.946	0.937	0.912	0.952
訓練 1&2	0.542	0.636	0.529	0.790	0.795	0.774	0.693
訓練 2	0.193	0.446	0.265	0.914	0.580	0.566	0.719

表現を獲得していることが確認できる(表4)。訓練 1&2 は、訓練 1 で獲得したベクトル表現を用いているが、指示文を説明文に反映させるタスクが加わったため全体として訓練 1 よりも精度が下がった。訓練 1&2 と訓練 2 の BLEU スコアは、ほぼ変わらないが時系列データを捉えることで正しく説明ができる動向 (T) と接続詞 (C) における精度は訓練 1&2 の方が高い。訓練 2 の時期 (P) の正解率の高さは、指示文の単語が反映が要因と考えられる。時期の高い精度にも関わらず、T/P/C の精度は訓練 1&2 を行った場合の方が高いことが確認できた。

5.4 分析

図5に生成結果例を提示する。左上の場合は、指示文通り中期における減少の動向について説明できしており、series1 よりも前に series2 が減少し終わることも接続詞に反映されている。図5右上の場合は、指示とは異なる時期の動向について説明をしている。こちらは時系列データの動向を捉えられていたが説明文に指示文を反映させることができなかったと考えられる。図5左下の図においては、中期の動向のピークについては正しく説明ができていますが、時間的関係の説明(接続詞)が不正解である。訓練 2 の接続詞の正解率は while 以外訓練 1&2 の場合を下回る(表5)。片方が平坦の場合の接続詞が一律 while となることから時間的関係を捉えずに時間的關係(接続詞)の説明が可能となる。一方で、while (片方平坦) 以外の場合は二つの時系列データの時間的關係を捉える必要があるため、時間的關係の説明がより難しくなることが要因として考えられる。

図5右下の場合は、末期の series2 の減少を平坦と説明してしまっている。訓練 1&2 の協調/反対の動向の正解率は訓練 2 の正解率より 20%以上高いが、片方平坦の場合は訓練 2 が約 2%上回った(表5)。片方が平坦の場合、両方の時系列データに大きな数値的变化がある場合に比べて動向を捉えやすいたことが要因として考えられる。

訓練 2 の場合は、特定の時期における動向の説明のための単語の数値的な意味の学習と説明文に指示文を反映させるという二つのタスクを同時に行う。そのため訓練 1&2 よりも全体的に精度は下がるが、学習がより容易な動向や接続詞についての精度はより高い。一方で、訓練 1&2 の場合、指示文が反映できていない場合も見られたが、訓練 1 で獲得した単語のベクトル表現を用いてより正確に動向と時間関係を捉えられることが確認できた。時系列データがより複雑になるほど動向や時間的關係を捉えることは難しくなるため、訓練 2 と訓練 1&2 の精度の差はより大きくなると考えられる。

6 おわりに

本論文では二つの時系列データの特定制目点の動向について説明を行う文生成モデルを提案する。モデルの訓練のために 2 つの dataset を作成し、動向と時期と時間的關係の学習、指示された時期の特定制目点の動向を捉えて説明をする学習を行なった。生成結果から、単語における数値的な意味を学習してベクトル表現を獲得してから二つの時系列データの特定制目点の動向について説明を行う訓練を行うことで、時系列データの関係性を捉えて十分な精度で正しい文生成が行えることが確認された。

今後の研究では、生成文の精度を上げるとともに説明文の情報量を増やすことで表現の幅を広げていきたい。

謝辞

本研究は JSPS 科研費 18H05521 の助成を受けたものです。

参考文献

- [1] Shankar Kantharaj, Rixie Tiffany Ko Leong, Xiang Lin, Ahmed Masry, Megh Thakkar, Enamul Hoque, and Shafiq Joty. Chart-to-text: A large-scale benchmark for chart summarization, 2022.
- [2] Miao Chen, Xinjiang Lu, Tong Xu, Yanyan Li, Jingbo Zhou, Dejing Dou, and Hui Xiong. Towards table-to-text generation with pretrained language model: A table structure understanding and text deliberating approach, 2023.
- [3] Hadi Banaee, Mobyen Uddin Ahmed, and Amy Loutfi. Towards NLG for physiological data monitoring with body area networks. In **Proceedings of the 14th European Workshop on Natural Language Generation**, pp. 193–197, Sofia, Bulgaria, August 2013. Association for Computational Linguistics.
- [4] Li Gong, Josep Crego, and Jean Senellart. Enhanced transformer model for data-to-text generation. In **Proceedings of the 3rd Workshop on Neural Generation and Translation**, pp. 148–156, November 2019.
- [5] Jason Obeid and Enamul Hoque. Chart-to-text: Generating natural language descriptions for charts by adapting the transformer model, 2020.
- [6] Harsh Jhamtani and Taylor Berg-Kirkpatrick. Truth-conditional captions for time series data. In **Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing**, pp. 719–733, November 2021.
- [7] Soichiro Murakami, Sora Tanaka, Masatsugu Hangyo, Hidetaka Kamigaito, Kotaro Funakoshi, Hiroya Takamura, and Manabu Okumura. Generating weather comments from meteorological simulations. In **Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume**, pp. 1462–1473, April 2021.
- [8] Shreshth Tuli, Giuliano Casale, and Nicholas R. Jennings. Tranad: Deep transformer networks for anomaly detection in multivariate time series data. **CoRR**, Vol. abs/2201.07284, , 2022.
- [9] Harsimrat Kaeley, Ye Qiao, and Nader Bagherzadeh. Support for stock trend prediction using transformers and sentiment analysis, 2023.
- [10] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, L ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, **Advances in Neural Information Processing Systems**, Vol. 30. Curran Associates, Inc., 2017.
- [11] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: A method for automatic evaluation of machine translation. In **Proceedings of the 40th Annual Meeting on Association for Computational Linguis-**

tics, ACL '02, p. 311–318, USA, 2002. Association for Computational Linguistics.

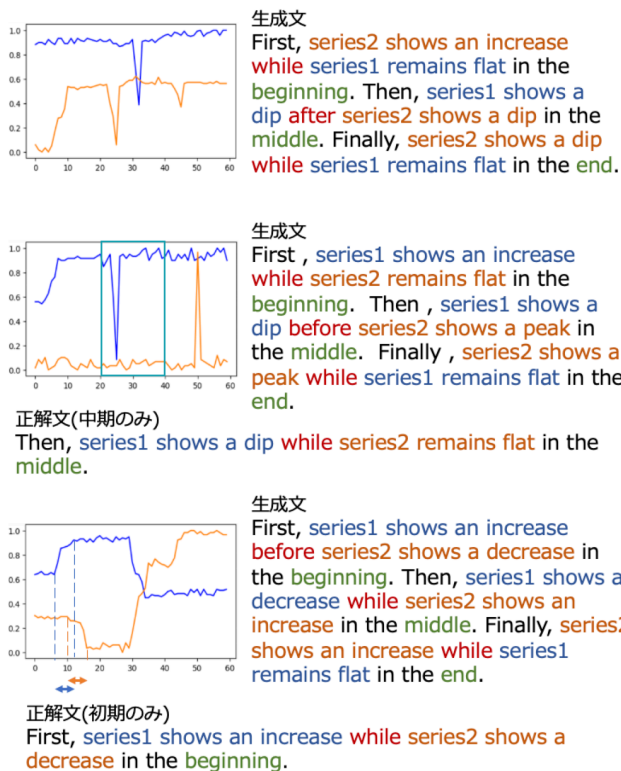


図6 訓練1における生成結果例

A 分析結果

A.1 訓練1分析結果

訓練1の生成例を図6に示す。訓練1の場合は多くの場合で正しい文生成が行えているため、生成文における不正解の時期のみ正解文を提示する。図6の上の場合は、初期/中期/末期の動向と時間的關係について正しく説明ができています。図6の中央の場合、中期におけるディップ/平坦の動向をディップ/ピークとして説明をしている。訓練1における生成文で動向についての説明が間違っている場合の多くは、協調/反対の動向を片方平坦として説明をする場合や片方平坦の動向を協調/反対として説明をする場合であった。接続詞については、動向が片方平坦の場合は while が用いられるため時間的關係が説明文に反映されなかった。図6の上の場合は、初期における動向について正しく説明をしているが、時間關係について正しく説明ができていない。この場合は動向が数値的变化を取る範囲が被っており、反対の動向を取るため接続詞は while となる。この場合のように接続詞のみが正しくない場合も見られた。



図7 訓練1&2における生成結果例

図8 訓練2における生成結果例

A.2 訓練1と訓練1&2の比較

訓練1&2は、訓練1で獲得したベクトル表現を用いているが、指示された時期のみについて説明をするというタスクが加わったため全体として訓練1よりも精度が下がった。指示されていない時期についての説明文生成を行ってしまう場合があるため、訓練1で他の単語よりも精度が高い単語が訓練1&2において他の単語よりも精度が低くなる場合も見られた(表5)。

A.3 訓練1と訓練1&2の生成結果例

図7に訓練1&2の生成結果例、図8に訓練2の生成結果例を提示する。上段が正解例で下段が不正解例となる。図7の下段においては、指示された末期の動向ではなく初期の動向について説明をしている。図8の下段では、初期の増加をピークと説明しており、正しく動向を捉えられていない。