# Universal Transformer を使用した対話破綻検出 Dialogue Breakdown Detection using Universal Transformer

桑原 健太<sup>\*1</sup> 大村 英史<sup>\*1</sup> 桂田 浩一<sup>\*1</sup> Kenta Kuwahara Hidefumi Ohmura Kouichi Katsurada

\*1 東京理科大学理工学研究科情報科学専攻

Department of Information Sciences, Graduate School of Science and Technology, Tokyo University of Science

In recent years, a lot of dialogue systems have been developed at many institutions. However, most of these systems sometimes cause the problem of dialogue breakdown because of the systems' inappropriate utterances, which makes it impossible to continue the conversation between a user and the system. In this paper, we propose a dialogue breakdown detection system using a universal transformer. Since our system can model the detailed structure of input sentence, it is expected to predict the dialogue breakdown accurately than the previous systems. To show its performance, we conducted an experiment in which our system is compared with a previous study using an RNN. The experimental result shows that our system can obtain a higher F-measure metrics than the previous method in the task of dialogue breakdown challenge.

### 1. はじめに

近年, 非タスク指向の対話システムが注目を集めている. 非タスク指向の対話システムは日常的に利用できる利便性 の高いシステムであるが,現状の技術水準ではしばしばシ ステムの不適切な発話によってユーザとの対話が続行不可 能になる状態(対話破綻)をしばしば引き起こすという問 題がある.そこでこの問題に対処するために対話破綻を検 出・回避する技術やリカバリする技術が検討されてきた.

このうち対話破綻の検出に関しては、対話破綻を引き起 こすかどうかを予測させる「対話破綻検出チャレンジ」が 近年開催されている[1]. このチャレンジにおいて小林らは Recurrent Neural Network(RNN)で構成された Neural Convers ational Model(NCM)を用いて文脈をモデル化し、対話破綻 を検出する方法を提案した[2]. さらに久保らは小林らのモ デルの問題点であった「対話破綻検出チャレンジ」の配布 データ以外の対話データが必要であるという点を、学習時 のデータに対話破綻状態か否かの情報を加えることによっ て改善し、高い性能を示した[3]. しかし自身の検出器の更 なる性能向上には、文構造に関する特徴量を入力に与える ことが必要であると示唆していた.

我々は久保らのモデルをベースに、Universal Transformer [5]を用いた対話破綻検出器を提案する. Universal Transformerは、自動翻訳の分野で広く用いられ始めた Transformer [4]というモデルを翻訳以外の自然言語タスク に適応可能にしたモデルである. これらのモデルでは注意 機構(Attention)という、2つの系列データの要素間の関係性 を学習する機能を持つ. これによって、入力文の文構造を 学習することができるため、久保らが述べた問題を解決で きると考えられる.

本稿では2章で対話破綻検出チャレンジの概要と小林らの手法,久保らの手法について述べ,3章で提案手法を説明する.その後,4章で評価実験について解説し,5章で本研究のまとめを述べる.

### 2. 対話破綻検出チャレンジの概要と関連研究

#### 2.1 対話破綻検出チャレンジ

対話破綻検出チャレンジは、ユーザとシステムの対話ロ グ中の各システム発話において対話破綻を引き起こしてい るか否かを判定するタスクを行うコンペティションである. データとして対話システムとユーザの行った対話が提供さ れており、システム発話には複数人のアノテータによって 「O」、「T」、「X」のラベルがつけられている.それぞ れ、「破綻を引き起こしていない」、「破綻を引き起こし ていないが違和感がある」、「破綻を引き起こしている」と いう意味を持つ.チャレンジの参加者は、このラベルの確 率分布を予測することである.以下に破綻の例とアノテー ションの例を示す.

ユーザ:猫に付いて行ったらなにか発見できますか。 システム:猫は元気です O:0 T:9 X: 21

各ラベルの隣の数字は、そのラベルをつけた人数を表している.この例ではシステムがユーザの質問を無視した発話をしているため、多くのアノテータはシステムの発話が 破綻を引き起こしたと判断している.

### 2.2 対話モデルを使用した対話破綻検出

対話破綻を検出するには、まず各発話を破綻検出に必要 な特徴量を含むベクトル表現に変換し、このベクトル表現 を用いて各ラベルの確率分布を予測するモデルを構築すれ ばよい.小林らは Vinyals らが提案した Neural Conversational Model (NCM)という対話モデルを用いる方法 を提案した[2].NCM は RNN を使用した Encoder-Decoder 型の対話モデルである.Encoder はユーザ発話を入力とし て受け取り,発話の意味を表す潜在表現へと圧縮する. Decoder は Encoder で得られた潜在表現を素性としてシステ ム発話を生成させる.図1上に小林らが用いたモデルを示 す.小林らは事前に学習した NCM にユーザ発話とシステ ム発話のペアを与え、Encoder と Decoder の最終時刻の内部 状態を特徴量として検出器を構築した.しかし、対話破綻

連絡先:桑原 健太,東京理科大学大学院理工学研究科情報 科学専攻,千葉県野田市山崎 2641,04-7122-9373, 6315045@ed.tus.ac.jp



The 33rd Annual Conference of the Japanese Society for Artificial Intelligence, 2019



検出チャレンジで配布されているデータには破綻している データが含まれている.このデータを用いて NCM を学習 させると,破綻している対話も学習することになる.する と NCM の内部状態は,破綻しているか否かを区別しなく なるため,得られる特徴量を使った破綻の検出は困難とな る.そこで小林らは別途正常な対話データを用いて学習を 行っていた.

これに対し久保らは NCM の学習を行う際にシステムの 発話文の末尾に「O」,「T」,「X」をアノテーション分 布に従って確率的に挿入するという手法を提案した[3].図 1 下に久保らが用いたモデルを示す.破綻してない対話と 破綻している対話が末尾のラベルによって区別できるため, この手法によって破綻している対話かどうかを区別しなが ら学習することが可能になった.その結果,第二回対話破 綻検出チャレンジにおいて配布データのみを使用しながら も最も高い F値を示した.しかしながら,自身のモデルに 対し,ユーザ発話とシステム発話に同じ単語が出現した場 合には,破綻であっても破綻していないと判断してしまう 場合があるとの問題点を挙げていた.この問題を解決する ためには文構造や体裁をチェックする特徴量が必要である と述べている.

### 3. 提案モデル

本稿で提案するモデルは、久保らのモデルをベースに、 NCM を Universal Transformer で学習したものである.以下 ではまず Universal Transformer の概要を説明した後に、本 手法の詳細について述べる.

### 3.1 Universal Transformer

Universal Transformer のベースとなる Transformer [4]は注 意機構(Attention)の導入によって,系列データの各要素の 処理において別の系列の特定の要素に着目した処理を行え るようにしたニューラルネットワークの一種である.注意 機構によって構文や意味構造を学習できることから, Transformer は機械翻訳タスクにおいて RNN モデルより優 れた性能を示している.しかし, Transformer は再帰的構造 を持たないため RNN より帰納バイアスが弱く,学習デー タに含まれない長さのデータには対応できなかった.この 問題に対し, Transformer に再帰的構造を付与したモデルが Universal Transformer [5]である.

Universal Transformer は Attention とフィードフォワードネ ットワーク(FFN)のみで構成された Encoder-Decoder モデル のニューラルネットワークである. 系列 Xの処理における 系列Yの各要素の重要度を表現する Attention は、学習され る重みQ, K, Vを用いて以下の式によって定式化される.

## QX = Q \* X, KY = K \* Y, VY = V \* YAttention $(X, Y) = \text{softmax} (QX * KY^T / \sqrt{d}) * VY$

ここでdはスケーリング因子のパラメータである. softmax 関数によってXの処理に対する、Yの各要素の影響を 求めている. つまり、AttentionによってXのある単語とYの ある単語との意味的つながりや文法的つながりを学習可能 になる. X = Yとする Attention は Self-Attention と呼ばれ、 Xを文とした場合、単一の文内の構文や意味構造を表現す る. Universal Transformer の Encoder 部は、Self-Attention と 単語ごとに処理される FFN の 2つのネットワークを1つの 層としている. 層への入力は最初に self-attention で処理さ れ、次に単語ごとに FFN によって処理される. この層を用 いて多層にしたネットワークが Encoder となる. 再帰的構





図 3:提案手法(Universal Transformer を用いた破綻検出器)

造を付与するため, Universal Transformer では Encoder の出 力を再び Encoder の入力として与えている. この処理は複 数回行われる. 最終的な Encoder の出力は最終時刻の最終 層の出力となる. Decoder は Encoder とほとんど同じ構成で あるが, Self-Attention と単語ごとの FFN の間に, Encoder の出力との Attentionの処理が加わる. この Attentionによっ て Encoder の入力との関係をモデル化することができる.

#### 3.2 Universal Transformer を用いた破綻検出法

提案手法では久保らが提案した RNN を用いた手法と同 様に最初に学習データのシステム発話の末尾に破綻ラベル をアノテーション分布に従って挿入する.次にユーザ発話 とシステム発話のペアを用いて Universal Transformer を学 習させる. 破綻を検出する際には、まずユーザ発話とシス テム発話のペアを学習した Universal Transformer に与え, Encoderの最終時刻の出力 $H_E$ と、Decoderの最終時刻の出力  $H_{\rm D}$ を得る.この2つの特徴量には self-attention によって, 発話の構文や意味構造などの情報が含まれていると見なす ことができる. また, Hpにはユーザ発話とシステム発話の 関係性もモデル化されていると考えられる.従って $H_{E}, H_{D}$ は発話の文法的な表現や意味的表現、ユーザ発話とシステ ム発話の関係を表していると期待される. Universal Transformer によって得られた $H_E$ ,  $H_D$ は, 破綻検出のための 特徴量として判別器の Support Vector Machine (SVM)に与え られる.以上に述べた提案手法の対話モデルを図2に,破 綻検出システム全体の構成を図3に示す.

### 4. 評価実験

### 4.1 実験概要

RNNを使用したモデルと提案手法の比較を行った.モデル構築に使用したデータセットは第一回対話破綻検出チャ

レンジの学習用データ(rest1046)と開発用データ (DBDC1dev),および第二回対話破綻検出チャレンジの学 習用データ(DBDC2dev)である.rest1046,DBDC1devは同 じ対話システム DCM を用いて収集されており,それぞれ 1046対話,20対話が収録されている.DBDC2devに関して は DCM に加え DIT, IRS という対話システムを加えた3つ の対話システムを用いて得られた対話から構成されており, それぞれ50対話ずつが収録されている.対話モデルの学習 にはこれらすべてを用いた.破綻検出器である SVM の学 習には対話モデルの学習に用いたデータから rest1046 を除 いたデータを用いた.rest1046 は破綻ラベルの偏りが大き いため破綻検出器の学習データには適さないと判断したた めである.

評価データには第二回対話破綻検出チャレンジの評価用 データ(DBDC2test)を用いた.DBDC2testもDBDC2devと同 様に DCM, DIT, IRS を用いて 50 対話ずつ収集されてい る.第二回対話破綻検出チャレンジでは各対話システムの データごとに評価しているため、本実験においてもデータ ごとに評価する.以下ではそれぞれの評価データを DCMtest, DITtest, IRStestと表す.評価には対話破綻検出 チャレンジで用いられた以下の5つの評価尺度を用いる.

- Accuracy: 全ラベルの一致率
- F-measure(X):破綻ラベル Xの検出の精度を示す F値. Xの適合率と再現率の調和平均によって表される.
- F-measure(T+X): ラベル Tを X と同一のラベルとして 扱った際の F 値
- JS Divergence(O,T,X): Jensen-Shannon Divergence によるアノテーション分布との分布間距離
- Mean Squared Error(O,T,X): アノテーション分布との 分布間の平均二乗誤差

データ	モデル	Accuracy	F-measure	F-measure	JSD	MSE
		_	(X)	(T+X)	(O,T,X)	(O,T,X)
DCMtest	RNN	0.480	0.511	0.682	0.134	0.070
	Universal Transformer	0.475	0.549	0.811	0.131	0.067
DITtest	RNN	0.596	0.705	0.882	0.087	0.045
	Universal Transformer	0.605	0.707	0.907	0.081	0.041
IRStest	RNN	0.525	0.609	0.763	0.139	0.074
	Universal Transformer	0.500	0.616	0.813	0.143	0.076

#### 表1:破綻検出の結果

表 2:各モデルの Precision と Recall

データ	モデル	Precision	Recall	Precision	Recall
		(X)	(X)	(T+X)	(T+X)
DCMtest	RNN	0.418	0.657	0.779	0.607
	Universal Transformer	0.395	0.904	0.762	0.866
DITtest	RNN	0.562	0.947	0.849	0.917
	Universal Transformer	0.555	0.973	0.857	0.964
IRStest	RNN	0.491	0.801	0.743	0.784
	Universal Transformer	0.457	0.944	0.711	0.950

対話モデルの学習にはミニバッチ学習を行い, バッチサ イズは 10とした.中間層は 256次元とし,中間層は 1層と した.最適化手法は学習率 0.1 の確率的勾配降下法を用い, 目的関数には予測単語ごとの交差エントロピー誤差を用い た.前処理として Mecabによる形態素解析を行った.辞書 には mecab-ipadic を用いた. Universal Transformer に関して は,再帰回数は 6回に設定した.

### 4.2 実験結果と考察

破綻検出の結果を表1に示す. Universal Transformer を用いた手法が RNN を用いた手法を概ね上回っている事が分かる.以下に RNN で破綻検出に失敗して Universal Transformer では成功したシステム発話とその直前のユーザ 発話の例を示す.

ユーザ:明日は猛暑らしいですから システム:猛暑は欲しいですね

この例はどちらにも「猛暑」という同じ単語が含まれて いるため、2.2節で述べた通り RNN では破綻と判定するこ とができなかったと考えられる.しかし、このような例に 対しても Universal Transformer では.破綻を検出できるよ うになった.

Universal Transformer では Accuracy に関して RNN を上回 ることができていない. この原因として, O のラベルが付 与された発話の多くを X と誤識別してしまったことが挙げ られる. 両者のモデルの特徴を調べるため, 破綻のラベル X が付与されたデータの Precision (適合率) と Recall (再 現率), および T を X と同一のラベルとした際の Precision と Recall の値を表 2 に示す. DITtest の Precision(T+X)を除 いて Precision は RNN の方が高く, Recall は Universal Transformer の方が高いことが分かる. この結果から Universal Transformer は O や T の発話を X と出力したケー スが多いことを示しており, Accuracy が低いのもこの誤認 識の多さが原因であると考える. 今後は Precision を向上す るための方法を検討する必要がある.

### 5. おわりに

本研究では Universal Transformer で学習させた対話モデ ルを用いた対話破綻検出法を提案した. RNNを用いた破綻 検出器と比較した結果,破綻検出精度の F 値に関して Universal Transformer が上回るという結果が得られた. これ は Self-Attention によって文構造や意味構造が学習され, RNNを用いたモデルにおける問題点を解決できたためと考 えられる. しかしながら対話破綻検出の適合率は RNN を 用いたモデルを下回る結果となったため,今後は適合率を 向上するための特徴量について検討したい.

#### 参考文献

- 東中竜一郎,船越孝太郎,稲葉通将,荒瀬由紀,角森 唯子,"対話破綻検出チャレンジ 2",第78回言語・音 声理解と対話処理研究会(第7回対話システムシンポジ ウム),人工知能学会研究会資料 SIG-SLUD-B505-19, pp.64-69,2016.
- [2] 小林颯介,海野裕也,福田昌昭,"再帰型ニューラルネットワークを用いた対話破綻検出と言語モデルのマルチタスク学習",第75回言語・音声理解と対話処理研究会(第6回対話システムシンポジウム),人工知能学会研究会資料 SIG-SLUD-075-09, pp41-46, 2015
- [3] 久保隆宏,中山光樹, "Neural Conversational Model を 用いた対話と破綻の同時学習",第78回言語・音声理 解と対話処理研究会(第7回対話システムシンポジウム), 人工知能学会研究会資料 SIG-SLUD-B505-26, pp94-97, 2016.
- [4] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Us zkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Il lia Polosukhin, "Attention is all you need", arXiv:1706.037 62, 2017
- [5] Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, Lukasz Kaiser, "Universal Transformers", arxive:1807.03819, 2018