実世界における直観に反した運動の予測と言語生成

2247010



研究代表者

お茶の水女子大学 大学院 人間文化創成科学研究科

博士後期課程

黒 田 彗 莉

[研究の目的]

過去10年の自然言語処理研究により、ヒト の活動や生活をサポートするような対話型ロ ボット(ヒトの発話を受け取り、内容に即して 回答するロボット)の開発が進んだ。また、大 規模言語モデル (LLM) の開発に伴い、今後 の対話型ロボットはより自然に返答できるよう になると期待される。しかし現在の LLM は. 答えのない問いへの回答や、ヒトの会話のよう に曖昧な表現を含んだ返答が苦手である。また. ヒト同士の会話と、ヒトとロボットのコミュニ ケーションは決定的に異なる。会話においてヒ トはただ相手の話を聞くだけでなく、相手の振 る舞いから次の発言を予測しながら発話する。 一方でロボットは、あくまでヒトからの一方的 な指示に答えるに過ぎず、相手(ヒト)の言動 を予測しているわけではない。ヒトからロボッ トへの一方的な命令という関係を超えて、今後 ヒトと機械が共存していくには、ロボットが自 律的に状況を判断できるようになる必要がある。 さらに、判断した状況から様々な将来を予測す ることで、ヒトの言動を先回りしてサポートで きるロボットの内部モデル構築を目指す。

たとえば、緑の円柱が青い立方体に近づいている状況があったとする。これらの物体は、この後どうなるだろうか。直観的には「緑の円柱が青い立方体にぶつかる」と想像できる。一方で、直観に反して「もし緑の円柱が止まった」とすると、「青い立方体と緑の円柱はぶつから

ない」と予測できる。このように直観に反した 状況でも、ヒトは柔軟に推論を変更することで、 その場面を画像として想像し、言語として発話 できる。しかし予測研究の多くは、環境にある 物体を視覚的(画像)もしくは物理的(物体の 速さなど)にとらえ、直観に即した予測をする にとどまっている。また、直観に反した予測を 扱う研究は、ほとんどの場合、物理シミュレー ターで実世界を表し、物理情報を直接修正する ことで仮想環境を生成している。そのため、ヒ トのように実世界からの情報をもとに柔軟に推 論を修正し、起こりえない予測を生成できるモ デルはまだない。

これまでの研究(例: Engelcke et al., 2020; Burgess et al., 2019)は、入力の系列情報に対する決まった(直観に即した)予測のみを扱ってきた。それに対し本研究では、直観に即した従来の予測だけでなく、直観に反する条件が与えられたときに、これまでの推論を柔軟に変更することで、状況に即したもっともらしい環境を予測できる推論モデルを構築する。与えられた条件に適合するように推論を修正して、新し

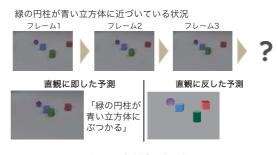


図1 本研究の概要図

い予測を生成するようなモデルはまだ提案されていない。

また、予測推論モデルから画像を生成する手法として、シーングラフから画像を生成する研究(Johnson et al., 2018)がある。しかし、環境にある物体の速度や移動方向といった物理的な特徴や物体同士の位置関係を表したグラフ構造から画像を生成した先行研究はいまだ存在しない。

このように自身の研究に適した形で従来の研究を改善する方法を独自に考え,新たな予測モデルの構築や画像生成手法を提案することで課題を解決する。

[研究の内容,成果]

1. 提案手法

1.1 物理特性を含んだ訓練データセット作成

本研究ではデータセットとして CLEVRER (Yi et al., 2019) をメインで用いた。CLEVRER は「画像内で3種類の物体が動き、物理的事象は衝突のみ」というシンプルなデータセットである。また、提案モデルの入力情報は画像だけでなく、画像から得られた物理情報(位置や速度など)の2種を想定している。そのため、元となるデータセットから物理情報のデータセットも新たに作成する。作成の手順を図2に示す。

(1) 画像内の環境の獲得

CLEVRER に写っている物体の種類や位置情報を考慮したグラフベースの訓練データを作成するために、YOLACT を用いた。YOLACT を用いて物体検知をし、画像内の2次元位置情報と物体の形状・色を取得する。検知可能な種類数は物体の色8種類・形状3種類・素材2種類の組み合わせ48種類で行った。

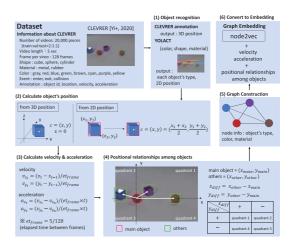


図2 訓練データセット作成の概要図

(2) 物体の座標・環境を表現したグラフの構築

獲得したバウンディングボックスの座標から物体の2次元位置情報を決定し、グラフを作成した。本研究ではグラフを構築した後、それらのグラフを埋め込みベクトル作成においては word2vec の Skip-gram から着想を得た node2vec と doc2vec の PV-DBOW から着想を得た graph2vec の2種類の手法を用いた。node2vec を用いた埋め込みでは1フレーム内に写っている物体一つずつをノードとみなし、構築したグラフのノードをもとに埋め込みベクトルを作成した。一方で graph2vec では、観測環境 (CLEVRER 1フレーム) を表したグラフ全体から埋め込みベクトルを作成した。

(3) 環境内の物体の速度と加速度の算出

環境に写っている各物体についての物理特性を捉えるために、物体の速度および加速度の算出を行った。時刻 t_k における物体の位置を x_{tk} , y_{tk} とし、次の時刻 t_{k+1} における物体の位置を x_{tk+1} , y_{tk+1} とする。また時刻 t_k から次の時刻 t_{k+1} の経過時間は、5秒の動画を 128 フレームに分割しているので、5/128 秒となる。加速度算出における初速 v_0 は、環境における各物体は停止しているところから始まるため v_0 =0と設定した。

(4) 物体間の位置関係の表現

環境内の物体に関して、物体同士の位置関係も重要な情報である。ヒトを例に上げると、自分を中心にして他の物体がどの位置関係にあるかを瞬時に捉えている。ここでは全ての物体同士の位置関係について、各物体が中心となったときの他物体の位置方向を算出した。

1.2 予測モデル構築

本研究では、PredNet に対して画像特徴量の変化から環境の移り変わりのタイミングを捉える研究である Variational Temporal Abstraction (VTA) の機能を組み合わせ、環境中の視覚情報と物理特性の両方の変化から環境の移り変わりを予測する予測モデルを構築した。VTA は系列情報から階層的な抽象度を見つける状態空間モデルであり、入力情報の変化点を抽出することができる。

モデルの構造は、物体の動きなどの物理特性を予測する機構と、環境の画像情報を予測する機構の二つの階層構造を並列にし、各モデルのそれぞれの機構に VTA の機構である変化点判別フラグ m を取り入れて構築した。

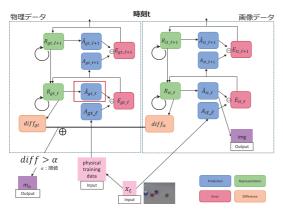


図3 提案する予測モデルの概要図

1.3 言語データのレンプレート作成

変化点予測モデルを用いて予測した埋め込みベクトルから言語を生成するために、言語情報を新たに学習する必要がある。そのため言語データセットとして、物理特性を表すグラフ表現の埋め込みベクトルとそのときの状態を説明した文章のペアデータを作成した。グラフ表現

の埋め込みベクトルは、CLEVRERのアノテーションデータから作成した。またペアとなる文章は、3(衝突前・衝突・衝突後)×3(文章の種類)の9種類のテンプレートに当てはまるように作成した。テンプレートの詳細を以下に示す。衝突した2つの物体はA・Bと表し、AとBはそれぞれ「"灰、赤、青、緑、茶、水、紫、黄"色の"球、円柱、立方体"」の情報をもつ。また衝突のデータだけでなく、物体が衝突の前に近づくと衝突後に離れるときをデータセットとして作成した。近づくときは衝突の5フレーム前、離れるときは衝突の5フレーム後を対象とした。作成したペアデータの例は図4に示す。

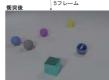
文章テンプレート例:衝突する物体(青色の球・灰色の球)



「青色の球と灰色の球が近づく」 「青色の球が灰色の球に近づく」 「灰色の球が青色の球に近づく」



「青色の球と灰色の球がぶつかる」 「青色の球が灰色の球にはじかれる」 「灰色の球が青色の球にはじかれる」



「青色の球と灰色の球が離れる」 「青色の球から灰色の球が離れる」 「灰色の球から青色の球が離れる」

図4 言語データのテンプレート作成の例

1.4 言語生成モデル

言語生成モデルは、Transformer の Decoder のみを用いた。図 5 に Decoder モデルを示す。従来の Transformer は Encoder-Decoder モデルで構築されているが、本研究では提案モデルの変化点予測モデルにおけるグラフの埋め込みベクトルの予測結果を Encoder における出力結果とみなし、この結果を Encoder から Decoder への入力とした。Decoder の学習は図4で作成したペアデータを用いた。学習データは 219303 個、テストデータは 10965 個とした。

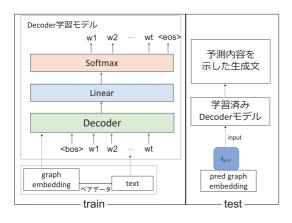


図5 言語モデル

2. 実験

実験では、予測して物理情報の変化点を言語 情報として生成し、予測内容について解釈可能 にすることを目的とする。

2.1 結果と考察

例1では、緑色の球と赤色の円柱がぶつかる 状況に対して, 赤色の円柱が反対方向に動いた ときの動きを検証した。言語生成モデルを用い たときは「緑色の円柱が赤色の円柱からはなれ る」という文章になり、予測画像は図6に示す ものになった。これは想定する正解文に近い文

条件 宝画像 赤色の円柱が反対に動く 正解文 「緑色の球と<mark>赤色</mark>の**円柱**はぶつからない」 「緑色の球が**赤色**の円柱からはなれる」 「赤色の円柱が緑色の球からはなれる」 予測画像 緑色の円柱が赤色の円柱 からはなれる 物体の色○, 形状×, 状況○,画像× 例ii



冬件

青色の円柱が反対に動く

正解文

「灰色の球と青色の円柱はぶつからない」 「灰色の球が青色の円柱からはなれる」 「青色の円柱が灰色の球からはなれる」

灰色の球と青色の立方体が はなれる

> 物体の色○, 形状×, 状況○, 画像×

図6 実験結果

章が生成できたといえる。一方で、予測画像に ついては物体の物理特性を表すグラフ構造から 画像を作成したが、ノイズがかかり、また動き が大きく変わった画像を生成することはむずか しいことがわかった。

例2では、灰色の球と青色の円柱がぶつかる 状況に対して、青色の円柱が反対に動いたとき の動きを検証した。言語生成モデルを用いたと きは「灰色の球と青色の立方体がはなれる」と いう文章になり、予測画像は図6に示すものに なった。例1と同様に、生成文は正解文に近い 文章が生成できたが、予測画像については想定 していた画像とは異なるものが生成された。

2.2 文章生成の評価

次に各モデルで生成した文章を自動評価指標 を用いて精度を算出した。生成した文を自動評 価指標で評価した結果を表1に示す。自動評価 指標には、言語生成タスクで用いられている BLEU. キャプション生成タスクで用いられる METEOR, CIDEr を使用した。生成文1文に 対し正解データは3文用意したため、各スコア の平均を生成文のスコアとした。予測推論内容 から生成した生成文と正解文のスコアは全体を 通して約75程度であり、観測環境の状況を高 い精度で生成することができたといえる。

表1 文章生成の評価結果

モデル	BLEU@2	BLEU@3	BLEU@4	METEOR	CIDEr
PredNet ベース	79.3	75.1	72.3	70.4	73.1

3. まとめと今後の課題

本研究では観測した環境を視覚的・物理的に 予測できるモデルを構築し、状況を説明できる モデルを提案した。また環境の物理的な特徴を 予測した結果を言語として生成し、モデルによ る予測内容を解釈可能にした。現実には起こり 得ない状況についての文章は高い精度で生成で きたが、そのときの場面を鮮明に描写すること には課題が残った。

また実験に使用したデータセットは、我々人

間が目にする実環境よりも簡単なデータになっているため、予測モデルとしても言語生成としてもまだ改善の余地はある。今後の課題としてより実世界に近いデータセットを利用した実験を考えている。

[成果の発表, 論文など]

Eri Kuroda & Ichiro Kobayashi: Predictive Inference Models for Real-world Physical Environments, Journal of Advanced Computational Intelligence and Intelligent Informatics, Vol. 29, No. 3, pp. 456-468 (2025).

- Eri Kuroda, & Ichiro Kobayashi. Verbal Representation of Object Collision Prediction Based on Physical CommonSense Knowledge. 2025 17th International Conference on Machine Learning and Computing (ICMLC2025), Long Paper, Oral, Guangzhou, China, Feb 14th-17th, 2025.
- 3. Eri Kuroda, Yuki Taya & Ichiro Kobayashi. Verbal Description Focusing on Physical Properties of Real-World Environments. 2024 Joint 13th International Conference on Soft Computing and Intelligent Systems and 25th International Symposium on Advanced Intelligent Systems (SCIS&ISIS2024), Long Paper, Oral, Hyogo, Japan, Nov 9th-12th, 2024.