

# 深層学習を用いた動画像からの危険認知手法のための 基礎的研究

7 月 7 日 (火)

小松 大起

## 1 本実験の目的

危険認知手法実現のために動画像を用いた深層学習によって次フレームの画像を用いた深層学習によって次フレームの画像を予測し生成を行い, その画像と実際の画像を比べる.

## 2 PredNet

PredNet は Deep Recurrent Convolutional Neural Network の 1 種で 神経科学の概念である Predictive Coding を組み込んで作られたモデルである. 2016 年に William Lotter, Gabriel Kreiman, David Cox の 3 氏によって公開された.

### 2.1 PredNet の層構造

PredNet の各層には 4 つの素子が存在しておりそれぞれ、Target, Representation, Prediction, Error と呼ぶ. Target は下層からの出力である誤差信号をエンコード, 符号化する. Representation は Recurrent unit で, 上層からの出力, 側方からの誤差信号, 1 ステップ前の自分の出力を受け取る. Representation unit は Target の予測をする Prediction unit に投射し, 入力の予測が出力される. Error は Prediction と Target の誤差であり, Error は上層に送られる. この Error が小さくなるように学習を進めていく. また, 層構造の例を図 2 に示す.

更新式は (1) から (4) のように表される.

## 3 使用した動画像

### 3.1 KITTI

ドイツの都市環境を運転している車の屋根に取り付けられたカメラによってキャプチャされたデータセッ

トである.City, Residential, Road のカテゴリに分かれており, それぞれ都市街, 住宅地, 高速道路というようなシーン分けがされている. それぞれのカテゴリには City には 28 シーン, Residential には 21 シーン, Road には 12 シーン存在している. 合計 61 のシーンがあるが, それぞれのカテゴリから 10 フレームがテストデータとしてサンプリングされ予測画像生成に用いられ, 57 シーンが訓練に用いられ, 4 シーンが検証に用いられる. また, 用いられる画像は中央でトリミングされて,  $128 \times 160$  ピクセルの画像となっている. 1 フレームは 0.1 秒である. 最大 5 フレーム先の画像を生成することが可能.

## 4 オートエンコーダ

オートエンコーダの目的は次元削減を行うこと.

## 5 先週までの作業

- PredNet について
- 授業課題
- PredNet の出力を得られた 図に示す. 1 フレーム先の画像を予測するように訓練されている.

## 6 今週の作業

- prednet 数式理解
- 人間の視覚情報による平均的反応時間は 0.18 秒から 0.20 秒ぐらいであるが, 0.6 秒先の画像生成で危険予測ができるのかを考える.

## 7 来週以降の作業

- オートエンコーダを prednet に組み込む

$$A_t^i = \begin{cases} x_t & \text{if } i = 0 \\ \text{MAXPOOL}(\text{ReLU}(\text{CONV}(E_{t-1}^i))) & i > 0 \end{cases} \quad (1)$$

$$\hat{A}_t^i = \text{ReLU}(\text{CONV}(R_t^i)) \quad (2)$$

$$E_t^i = [\text{ReLU}(A_t^i - \hat{A}_t^i), \text{ReLU}(\hat{A}_t^i - A_t^i)] \quad (3)$$

$$R_t^i = \text{CONV LSTM}(E_t^{i-1}, R_t^{i-1}, R_{t+1}^i) \quad (4)$$

Fig. 2: PredNet の更新式

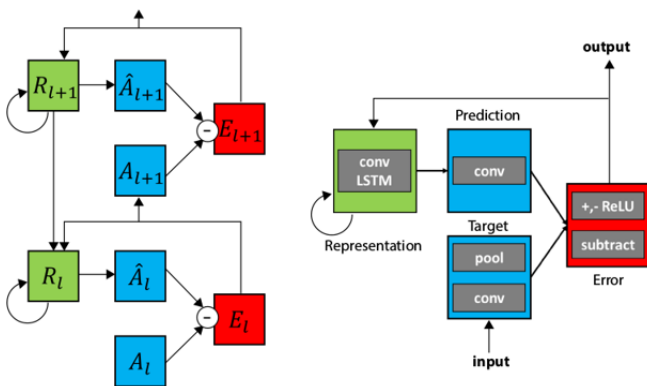


Fig. 1: PredNet の層構造の例



Fig. 3: PredNet の出力結果