

“Football Action Recognition using Hierarchical LSTM”の 実装報告

福本 将司^{1, a)}

概要：筆者は2017年にTsunodaらが発表した論文“Football Action Recognition using Hierarchical LSTM”の実装を行っている。この論文の目的は、映像における複数の個人を検出し、集団行動を推定することである。本論文では主にサッカーでの行動（パス、ドリブル、シュート）の推定を行う。提案手法では、映像に映る人数が可変な場合に対応しており、高い認識精度を実現している。

1. はじめに

本論文[1]ではチームのスポーツでの行動を認識するための方法を提案している。2ストリームのLSTM(Long Short Term Memory)に基づく既存手法[2]の改善を目指し、行動カテゴリーの識別可能性と、映像に映るプレイヤーの数の変動に対する堅牢性の向上を検証する。主に本論文ではサッカー映像に対する行動認識を行っており、映像を入力としてパス、ドリブル、シュートの3つの行動を認識する。

2. 実装する論文の概要

本論文では、CNN(Convolutional Neural Network)と2層のLSTMで構成されている。すなわち、CNN+L_k、LSTM1、LSTM2で構成されている。ここで、CNN+L_kはプレイヤーに関する特徴である。これはCNNによって出力される特徴と選手の位置、ボールの位置などのメタデータの連結された特徴である。図1に、提案手法のモデルを示す。

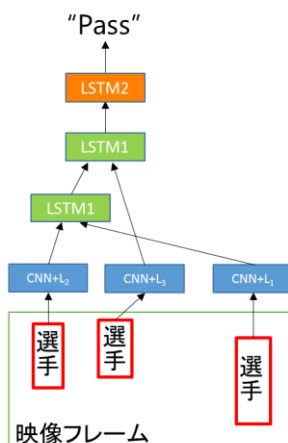


図1 提案手法のモデル

チームスポーツの行動認識のターゲットは常に人間およびボールであり、年齢、性別、服装など認識ターゲットの外観は類似している。したがって、CNNで選手の部位などの高レベルの特徴は使えない。一方で形状などの低レベルの特徴は効果的であると考えられる。そのため本論文では

中間層の途中で特徴を抽出し、それらを連結してCNNの特徴として利用する。サッカーのプレーを認識するには、メタデータが重要であると考えられる。例えば、ゴールエリアにボールに近い選手がいる場合、シュートの可能性が高くなる。そこで、選手とボールの3次元位置と選手のチームID、選手とボール間の距離、および選手とカメラ間の距離をメタデータとしてCNN特徴と連結する。最後に、この連結特徴をCNN+L_kとして使用する。

LSTMはニューラルネットワークの1つで、以前に計算された情報を記憶して後の計算に利用することができるため、時系列データを扱う際に大きな効力を発揮する。選手は常に移動し続けるため、カメラで撮影された選手の総数は毎フレーム変化する。したがって、本論文ではLSTMを用いて再帰的な方法でCNN+L_kを結合する。第1のLSTM(LSTM1)は可変数のCNN+L_kを再帰的に結合し、第2のLSTM(LSTM2)は結合された複数の選手の特徴を時系列に結合する。また、サッカーの行動予測を行う

3. 実装の進捗状況

実装環境としてPythonを用いており、TensorflowとKerasを機械学習のライブラリとして用いている。

現在は本論文の既存手法であるLRCNの実装を行っている[2]。本論文で使われているデータセットが公開されていないため、UCF101という101の行動ラベル付きの動画データセットを用いて評価をしている。現在は既存手法の実装を目標にし、余裕があれば本論文の実装に入りたいと考えている。

参考文献

- [1] T.Tsunoda, Y.Komori, M.Matsugu, T.Harada. Football Action Recognition using Hierarchical LSTM. 2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops
- [2] J. Donahue, L. A. Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell. Long-term recurrent convolutional networks for visual recognition and description. arXiv preprint arXiv:1411.4389, 2014.

¹ 静岡大学工学部数理システム工学科
静岡県浜松市中区城北 3-5-1
^{a)} fukumoto.masashi.17@shizuoka.ac.jp