
Classifying Humans Activities Using Human Poses Key-Points

Jérémy Trudel
Computer Science
Université de Montréal
Montreal, Quebec

Caroline Dakouré
Computer Science
Université de Montréal
Montreal, Quebec

Camille Felx-Leduc
Computer Science
Université de Montréal
Montreal, Quebec

Nicolas, Lemieux
Computer Science
École de Technologie Supérieure
Montreal, Quebec

Helgi Tomas Gislason
Computer Science
Université de Montréal
Montreal, Quebec

Abstract

fdsfjdsfjdsf The abstract paragraph should be indented 1/2 inch (3 picas) on both the left- and right-hand margins. Use 10 point type, with a vertical spacing (leading) of 11 points. The word **Abstract** must be centered, bold, and in point size 12. Two line spaces precede the abstract. The abstract must be limited to one paragraph.

1 Introduction

The goal of this project was to characterise how much the body position of one person could be correlated to the activity he or she was doing. Spatial points associated with body markers such as an ankle, knee, elbow, head and so on were used to complete this task. Three datasets were gathered, which allowed us to investigate the problem from different perspectives.

The first dataset is the MPII Human Pose Annotations Dataset. In this dataset, there is a total of 30,000 annotated picture examples (body-pose) labeled in 420 classes, each corresponding to a given activity.

The second dataset is named Stanford 40 Actions. It has 40 different activities with around 200 - 300 picture examples each. Here we do not have the pose data, but they will be obtained using a pose estimation algorithm namely OpenPose.

Finally the last dataset is UTD-Multimodal Human Action Dataset. This dataset is formed by 8 subjects performing 27 labeled activities about 4 times each and provides pose information in a time series corresponding to the video frames during which the activity was conducted on the videos provided.

The three classifiers decided to look into were: K-NN, Random Forest and AdaBoost.

2 Data

2.1 UTD-MHAD

2.2 MPII

3 Pre-processing

3.1 UTD-MHAD

The UTD-MHAD database is composed by videos from 8 subjects performing 27 different activities (classes) about 4 times each. Pre-processing for this database was done by computing the Covariance of 3D joint descriptor, a technique first described by (référence). The idea is to compute the covariance matrices of the joints through a sequence. First, we align the 20 joints 3D coordinates from each frame in a 60×1 vector that is denoted by S . Then the covariance matrix for that sequence is computed by $C(S) = \frac{1}{T} \sum_{t=1}^T (S - \bar{S})(S - \bar{S})^T$, where \bar{S} is the sample mean of S , and the T is the transpose operator. As the covariance matrix is symmetric, we then take all the elements of the upper matrix to make the 3D joint descriptor that is in this case an $60(60+1)/2 = 1830$ elements vector, on which the classification task can be performed. Although as noted by the authors :

The 3D cov descriptors capture the dependence of locations of different joints on one another during the performance of an action. However, it does not capture the order of motion in time. Therefore, if the frames of a given sequence are randomly shuffled, the covariance matrix will not change. This could be problematic, for example, when two activities are the reverse temporal order of one another, e.g. “push” and “pull”.

3.2 MPII