

Unstructured Human Activity Detection from RGBD Images

Authors : Jaeyong Sung, Colin Ponce, Bart Selman and Ashutosh Saxena

Presented By:

Ankit Goyal and Avi Singh



Introduction

Detecting and recognizing human activities is of great interest in several fields, with one being personal robotics. There are numerous challenges in activity recognition. Most of the activities occur in cluttered environments, in an uncontrolled manner. Different people have different speeds and mannerisms while performing these activities. This paper presents an approach to detect activities under such conditions. Some of the activities detected and recognized are cooking, drinking water, brushing teeth.

Previous work has focused mainly on RGB videos, or on the usage of RFID sensors. RGB videos lead to poor accuracy even in case of uncluttered environments, while the RFID methods are too intrusive as they require the placement of RFID tags on people and objects.

The model used in this paper is a two-layered maximum entropy Markov model. This model exploits the inherent hierarchical nature of human activities. For example, brushing a teeth involves several sub-activities such as picking up the toothbrush, squeezing the tooth-paste, actual brushing etc. The graphical model is not fixed, and an on-the-fly graph structure selection techniques is described.

Related Work

A lot of work has been done in the field of human activity recognition. Some of the common approaches and their limitations have been listed below:

1. One approach is to use space-time features to model points of interest in the video. Some authors have suggested methods to add more information to these features. However, this approach is only capable of classifying, rather than detecting activities.
2. Other approaches include filtering techniques and sampling of video patches.
3. Hierarchical techniques for activity recognition have been used as well, but these typically focus on neurologically-inspired visual cortex-type models. There is a blind adhere to models of the visual cortex which may not always be correct.
4. Other approaches includes the use of Hidden Markov Models (HMMs). However it has been argued in literature that CRFs and MEMMs overcome limitations posed by HMMs. CRFs and MEMMs enables longer term interaction among observations which HMMs don't.

Proposed Model

The model should incorporate different nuances in the human activity. An activity comprises of a series of sub-activities done in some particular order. In order to incorporate the hierarchical nature of human activity, a maximum entropy Markov model is proposed [Fig 1.] and its salient features are explained as below:

- x^t denote the features extracted from the articulated skeleton model at time frame t .
- Every frame is connected to high-level activities through the mid-level sub-activities. High-level activities do not change every frame, we do not index them by time. Rather, we simply write z^i to denote the i^{th} high-level activity. Activity i occurs from time t_{i+1} to time t_i .
- Every frame is connected to a sub activity. y^t represents the sub activity connected to frame at time t . The sub activities are intern

connected in sets to an activity. Thus, $\{y^{t_{i-1}+1}, \dots, y^{t_i}\}$ is the set of sub-activities connected to activity z^i

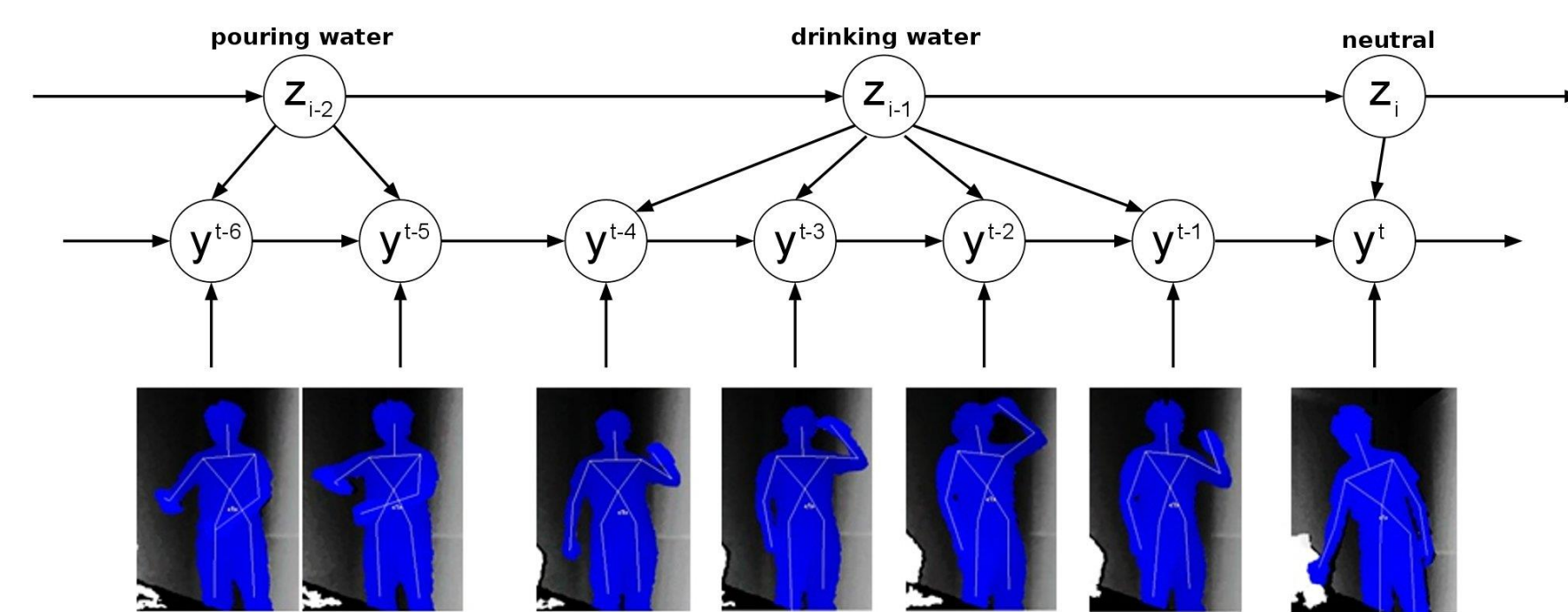


Figure 1: Proposed Model

Feature Extraction

Skeletal Features

PrimeSense provides a system for skeleton tracking from RGBD data. It gives us the three-dimensional Euclidean coordinates for fifteen joints, and rotation matrices for their orientations. All measurements are with respect to the frame of the sensor. Features computed from this data are as follows:

1. **Body pose:** Since features with respect to the sensor are not useful, the rotation matrices for 10 joints are calculated with respect to the torso, and represented using quaternions. Whether a person is standing or sitting, or leaning over is computed using the position each foot with respect to the torso, and by computing head and hip joint angles with respect to the vertical.
2. **Hand Position:** The position of the two hands are computed with respect to the head and the torso. The hand positions are also observed for the last six frames, and the maximum and minimum hand positions are extracted from them.
3. **Motion Information** Nine frames are selected from the last three seconds, and the joint rotations that have occurred are computed with respect to each these nine frames.

HOG Features

Histogram of Oriented Gradient (HOG) features are computed, which give us a count of how frequently a particular gradient is seen in the ROI of an image. Using the skeletal data from PrimeSense, a bounding box is drawn around the head, torso, left arm, right arm. HOG features are then computed inside each of these bounding boxes.

Learning and Inference

Learning

A Gaussian Mixture Model is used to cluster the original data, and each individual cluster is treated as a separate sub-activity. Clusters are also generated from some negative examples (no activity happening), so that the system is not prone to errors on observing random activities. We need to evaluate the values for the following terms from the training data in order to perform inference.

1. $P(y^t|x^t)$: This term models the dependence of the sub-activity label y^t on the input features x^t . The GMM used in the previous step is used here too.

2. $P(y^{t_i-m}|y^{t_i-m-1}, z_i)$: For all activities except neutral, this table is built from observing the transition of posterior probability from the soft cluster of Gaussian Mixture model at each frame. For neutral activities, $P(y^{t_i-m}|y^{t_i-m-1}, z_i=N) \propto 1 - \sigma_{z_i \neq N} P(y^{t_i-m}|y^{t_i-m-1}, z_i=N)$
3. $P(z_i|z_{i-1})$ Set manually.

Inference

The joint probability is computed as follows:

$$\begin{aligned} P(z_i, y^{t_{i-1}+1} \dots y^{t_i} | O_i, z_{i-1}) \\ = P(z_i | O_i, z_{i-1}) P(y^{t_{i-1}+1} \dots y^{t_i} | z_i, O_i, z_{i-1}) \\ = P(z_i | z_{i-1}) \cdot \prod_{t=t_{i-1}+2}^{t_i} P(y^t | y^{t-1}, z_i, x^t) \\ \cdot \sum_{y^{t_{i-1}}} P(y^{t_{i-1}+1} | y^{t_{i-1}}, z_i, x^{t_{i-1}+1}) P(y^{t_{i-1}}) \end{aligned}$$

The unknown terms in the equation above are derived as :

$$P(y^t | y^{t-1}, z_i, x^t) = \frac{P(y^{t-1}, z_i, x^t | y^t) P(y^t)}{P(y^{t-1}, z_i, x^t)}$$

From the model, the following conditional independence assumption are made: y^{t_1} and z^i are independent from x^t given y^t is made and under this we get

$$P(y^t | y^{t-1}, z_i, x^t) = \frac{P(y^t | y^{t-1}, z_i) P(y^t | x^t)}{P(y^t)}$$

Finally the entire formula can be written as :

$$\begin{aligned} P(z_i, y^{t_{i-1}+1} \dots y^{t_i} | O_i, z_{i-1}) \\ = P(z_i | z_{i-1}) \\ \cdot \sum_{y^{t_{i-1}}} \frac{P(y^{t_{i-1}+1} | y^{t_{i-1}}, z_i) P(y^{t_{i-1}+1} | x^{t_{i-1}+1})}{P(y^{t_{i-1}+1})} P(y^{t_{i-1}}) \\ \cdot \prod_{t=t_{i-1}+2}^{t_i} \frac{P(y^t | y^{t-1}, z_i) P(y^t | x^t)}{P(y^t)} \end{aligned}$$

This formula can be factorized as follows :

$$P(z_i, y^{t_{i-1}+1} \dots y^{t_i} | O_i, z_{i-1}) = \mathcal{A} \cdot \prod_{t=t_{i-1}+2}^{t_i} \mathcal{B}(y^{t-1}, y^t)$$

To maximize the probability the individual terms are maximized as:

$$\begin{aligned} \max P(z_i, y^{t_{i-1}+1} \dots y^{t_i} | O_i, z_{i-1}) = \max_{y^{t_{i-1}+1}} \mathcal{A} \\ \cdot \max_{y^{t_{i-1}+2}} \mathcal{B}(y^{t_{i-1}+1}, y^{t_{i-1}+2}) \dots \max_{y^{t_i}} \mathcal{B}(y^{t_i-1}, y^{t_i}) \end{aligned}$$

Graph Selection

Using the above results we can find the set of y^t 's that maximize the joint probability $P(z_i, y^{t_{i-1}+1} \dots y^{t_i} | O_i, z_{i-1})$, the probability of an activity z^i being associated with the i^{th} substructure and the previous activity. Our task is to use that to compute the probability of z^i given all observations up to this point. Simply trying all the possibilities would be intractable and so we use a dynamic programming approach which is explained as follows :

- We are at some time t ; we wish to select the optimal graph structure given everything we have seen so far. We will define the graph structure inductively based on graph structures that were chosen at previous points in time. Let $G_{t'}$ represent the graph structure that was chosen at some time $t' < t$. As a base case, G_0 is always the empty graph.

- For every $t' < t$, define a candidate graph structure $\tilde{G}_t^{t'}$ consisting of $G_{t'}$, followed by a single substructure from time $t' + 1$ to time t connected to a single high-level node z^i .
- Given the set of candidate structures $\{\tilde{G}_t^{t'} | 1 \leq t' < t\}$, the plan is to find the graph structure and high-level activity to maximize the likelihood given the set of observations so far $z_i \in Z$.
- Let O be the set of all observations so far. Then $P(z_i | O; \tilde{G}_t^{t'})$ is given by the following equation :

$$\begin{aligned} P(z_i | O; \tilde{G}_t^{t'}) &= \sum_{z_{i-1}} P(z_i, z_{i-1} | O; \tilde{G}_t^{t'}) \\ &= \sum_{z_{i-1}} P(z_{i-1} | O; \tilde{G}_t^{t'}) P(z_i | O, z_{i-1}; \tilde{G}_t^{t'}) \\ &= \sum_{z_{i-1}} P(z_{i-1} | O; G_{t'}) P(z_i | O_i, z_{i-1}) \end{aligned}$$

- The first factor inside the summation is calculated through dynamic programming and the second factor is calculated from $P(z_i, y^{t_{i-1}+1} \dots y^{t_i} | O_i, z_{i-1})$ as described earlier.

- The optimal probability of having node i be a specific activity z^i is computed as follows and is stored for the purpose of dynamic programming.

$$P(z_i | O; G_t) = \max_{t' < t} P(z_i | O; \tilde{G}_t^{t'})$$

- Thus the prediction at time t is made by the following :

$$\text{activity}_t = \arg \max_{z_i} P(z_i | O) = \arg \max_{z_i} \max_{t' < t} P(z_i | O; \tilde{G}_t^{t'})$$

The selected graph is shown to be optimal and the time complexity for the entire calculation is $O(n \cdot m \cdot T^2 \cdot t)$.

Results

The results show an average precision/recall of 84.7%/83.2% in detecting the correct activity when the person was seen before in the training set and 67.9%/155.5% when the person was not seen before. The results are summarized in the table below.

| Location | Activity | "New Person" | | | | | | "Have Seen" | | | | | |
|-----------------|------------------------|------------------|----------------|---------|------------|---------------|------------|------------------|----------------|------------|---------------|------------|------|
| | | Naive Classifier | One-layer MEMM | RGB HOG | Full Model | Skel-Skel HOG | Full Model | Naive Classifier | One-layer MEMM | Full Model | Skel-Skel HOG | Full Model | |
| bathroom | ringing nose | 77.9 | 82.3 | 82.9 | 81.3 | 81.1 | 81.4 | 71.3 | 80.7 | 81.1 | 81.1 | 81.1 | 81.1 |
| | brushing teeth | 64.5 | 20.5 | 83.3 | 57.7 | 50.7 | 30.8 | 73.4 | 16.6 | 88.5 | 55.3 | 81.5 | 75.6 |
| | waiting contact lens | 82.0 | 80.2 | 81.5 | 80.7 | 44.2 | 40.4 | 23.2 | 85.2 | 78.6 | 88.3 | 87.8 | 71.9 |
| | Average | 74.7 | 53.1 | 78.9 | 70.2 | 45.7 | 48.2 | 58.3 | 57.8 | 72.7 | 65.0 | 80.9 | 66.9 |
| | Average | 82.0 | 52.8 | 82.0 | 25.6 | 0.0 | 0.0 | 15.6 | 8.8 | 63.2 | 48.3 | 70.2 | 87.2 |
| bedroom | drinking water | 19.2 | 12.1 | 19.1 | 12.1 | 0.0 | 0.0 | 3.0 | 0.1 | 70.0 | 71.7 | 64.1 | 31.6 |
| | talking on the phone | 95.6 | 65.9 | 95.6 | 65.9 | 60.6 | 54.8 | 33.8 | 36.5 | 95.0 | 57.4 | 48.7 | 52.3 |
| | opening pill container | 65.0 | 36.9 | 65.0 | 36.9 | 20.2 | 11.6 | 17.4 | 12.2 | 76.1 | 39.4 | 61.0 | 30.3 |
| | Average | 65.0 | 36.9 | 65.0 | 36.9 | 20.2 | 11.6 | 17.4 | 12.2 | 76.1 | 39.4 | 61.0 | 30.3 |
| | Average | 33.3 | 56.9 | 33.2 | 57.4 | 56.1 | 90.0 | 50.9 | 74.2 | 45.6 | 43.3 | 78.9 | 29.0 |
| kitchen | cooking (chopping) | 44.2 | 29.3 | 45.6 | 31.4 | 58.0 | 4.0 | 94.5 | 11.1 | 24.8 | 17.7 | 44.6 | 45.8 |
| | drinking water | 72.5 | 21.2 | 0.0 | 0.0 | 0.0 | 0.0 | 91.8 | 23.9 | 95.4 | 75.3 | 52.2 | 51.5 |
| | opening pill container | 76.9 | 6.2 | 25.8 | 6.2 | 83.6 | 33.5 | 24.1 | 30.0 | 91.9 | 55.2 | 17.9 | 62.4 |
| | Average | 56.8 | 28.4 | 26.6 | 29.7 | 89.4 | 31.9 | 75.1 | 36.1 | 64.4 | 47.9 | 48.4 | 47.2 |
| | Average | 69.7 | 0.9 | 83.3 | 25.0 | 0.0 | 0.0 | 21.0 | 11.8 | 91.5 | 48.5 | 34.1 | 69.7 |
| living room | drinking water | 57.1 | 53.1 | 52.8 | 55.8 | 0.0 | 0.0 | 1.2 | 0.0 | 54.3 | 69.3 | 80.2 | 48.7 |
| | talking on couch | 71.5 | 35.4 | 27.4 | 91.3 | 42.7 | 99.4 | 23.2 | 62.2 | 73.2 | 43.7 | 91.4 | 50.7 |
| | reclining on couch | 97.2 | 76.4 | 95.8 | 78.6 | 0.0 | 0.0 | 100.0 | 21.5 | 31.3 | 21.1 | 95.7 | 96.5 |
| | Average | 73.9 | 41.2 | 72.3 | 62.7 | 10.0 | 10.0 | 46.4 | 24.1 | 52.6 | 45.9 | 72.4 | 55.9 |
| | Average | 60.5 | 31.0 | 60.6 | 31.5 | 17.5 | 6.7 | 2.7 | 0.6 | 69.4 | 48.2 | 80.4 | 52.3 |
| office | talking on the phone | 47.1 | 73.3 | 48.2 | 74.1 | 41.2 | 25.1 | 64.0 | 97.0 | 75.5 | 81.3 | 42.5 | 59.3 |
| | writing on whiteboard | 41.1 | 12.4 | 51.2 | 23.2 | 0.0 | 0.0 | 0.0 | 0.0 | 67.1 | 68.8 | 53.4 | 36.7 |
| | drinking water | 93.7 | 76.3 | 92.5 | 76.8 | 100.0 | 11.9 | 100.0 | 29.0 | 83.4 | 40.0 | 89.2 | 69.3 |
| | working on computer | 69.5 | 48.2 | 62.6 | 31.2 | 39.2 | 10.9 | 49.2 | 21.7 | 59.2 | 38.2 | 69.2 | 69.3 |
| | Average | 66.3 | 41.7 | 67.2 | 58.2 | 33.1 | 23.5 | 49.3 | 33.0 | 67.9 | 55.5 | 66.4 | 56.0 |
| Overall Average | 66.3 | 41.7 | 67.2 | 58.2 | 33.1 | 23.5 | 49.3 | 33.0 | 67.9 | 55.5 | 66.4 | 56.0 | |

Figure 2: Comparison of different models