

Identifying Human Interactions in Indoor Environments

Anand Panangadan Maja Mataric' Gaurav Sukhatme

Computer Science Department
University of Southern California

942 West 37th Place

Los Angeles, CA 90089

{anandvp,maja,gaurav}@robotics.usc.edu

Abstract

A laser range-finder-based system for detecting interactions between people in indoor environments is presented. An entropy-based measure is used to recursively segment tracks of individuals into distinct activities. Probability distributions corresponding to co-occurring activities are compared and similar distributions are flagged as representing possible interactions. Experimental results from two environments are described.

1. Introduction

The ability to detect and model interactions between people has many applications, ranging from surveillance to designing environments that facilitate interaction. In outdoor environments, most interacting people are close to each other to offset significant background noise. Thus, distance between people can be used to detect interactions outdoors. However, this is not always true in indoor environments. For instance, people following each other along a corridor may not be interacting, while ping-pong players across a table are. Most work on detecting interactions therefore relies on pre-defined models [2, 3]. In contrast, we represent a person's activity as a probability distribution over a discrete space of displacements. We use an entropy based metric (Kullback-Leibler divergence) to detect interactions by measuring the similarity between the corresponding probability distributions. Thus, a model for each type of interaction is not necessary.

2. Approach

We use multiple laser range-finders, placed along the edges of the environment, for tracking people. Unlike vision systems, the identities of people cannot be detected by the laser range-finders. Each detected object is tracked by

instantiating a particle filter. Laser range-finders have been used elsewhere to track people for activity modeling [4], but not for detecting interactions. The output of our laser tracker is a sequence of positions (in global coordinates) for each tracked object. This sequence represents a series of consecutive activities performed by a single person. For instance, a person may sit at a desk for a while, walk to another desk and then sit at that desk. We use a recursive entropy-based algorithm to split the sequence into distinct activity sub-sequences. This requires that every activity is represented as a probability distribution.

Activity as a probability distribution: Every sequence of positions obtained from the tracker is converted into a sequence of displacements. Let $(r_i, \theta_i), i = 0, 1, \dots, n$ be such a sequence. The displacements are discretized into one of nine canonical displacements. If $r_i \geq 0.2m$, then it is discretized into one of eight sectors depending on θ_i , else it is discretized as a separate displacement (0.2m is the minimum distance moved by a mobile person in a time-step). Thus, the continuous 2-dimensional space of displacements is reduced into a discrete space, \mathcal{X} , with 9 values.

Given a sequence $D = \{d_1, d_2, \dots, d_N\}, d_i \in \mathcal{X}$ of discrete displacements, the corresponding Maximum Likelihood (ML) probability distribution p over the discrete space \mathcal{X} is obtained by counting each type of displacement and dividing by the total number of displacements in the sequence. We define an activity to be a set of canonical displacements drawn from a fixed distribution. We assume that different activities give rise to different ML probability distributions of displacements. Thus, segmentation must divide each activity sequence into consecutive subsequences such that these are distinct in a probabilistic sense. We compute the difference between subsequences using an entropic measure known as the Jensen-Shannon (JS) divergence. Let p_1, p_2 be two probability distributions defined over the discrete space of canonical displacements, \mathcal{X} . The JS divergence between p_1, p_2 is defined as

$$JS(p_1, p_2) = H(p_{12}) - (w_1 H(p_1) + w_2 H(p_2))$$

where, $w_1, w_2 \geq 0$, $w_1 + w_2 = 1$, $p_{12}(x) = w_1 p_1(x) + w_2 p_2(x) \forall x \in \mathcal{X}$ is the weighted average distribution, and $H(p)$ is the entropy of the probability distribution p . We utilize the lengths of the sequences giving rise to p_1, p_2 to weight the divergence.

We obtain all activity subsequences from a track of displacements by recursively splitting the track at the point that gives rise to maximum divergence of the two component tracks [1]. The recursive segmenting procedure is terminated when the confidence of seeing the JS divergence between the left and right subsequences due to a real difference in their underlying probability distributions falls below a certain threshold.

Detecting interactions: We assume that two interacting people generate displacements from similar distributions. The similarity between two distributions p_1, p_2 is compared using the symmetric version of the Kullback-Leibler (KL) divergence between probability distributions, (Jeffrey’s divergence): $KL(p_1, p_2) = D(p_1, p_2) + D(p_2, p_1)$ where

$$D(p_1, p_2) = \sum_{x \in \mathcal{X}} p_1(x) \log \frac{p_1(x)}{p_2(x)}$$

If the KL distance between the distributions of a pair of activity segments falls below a certain threshold, the segments are assumed to come from interacting persons.

3. Results

We monitored people in two indoor environments: a laboratory measuring approximately 8m × 8m, and a long, narrow corridor measuring approximately 27m × 1.5m. Typical activities included sitting at desks, moving between desks, playing ping-pong, walking along the corridor, and waiting by a door. We conducted three motion capture sessions in each environment, each lasting approximately three hours. Normal activity proceeded throughout the recording sessions.

Quality of segmentation: The ability to detect interactions is related to the quality of segmentation. We stop the recursive segmentation algorithm when the confidence in the JS divergence falls below 90%. Segmentation works well when a stationary person starts to move or vice-versa. In these cases, there is 100% agreement with manual segmentation. During ping-pong, segmentation should ideally occur only when the ball goes out of play (when one of the players moves away to pick it). We considered a playing segment of approximately 600 seconds and compared the activity subsequences that were segmented out with manual segmentation using video (Table 1).

Interactions: The system detects two people playing ping-pong as an interaction. This interaction is particularly interesting because the two players are separated by a significant distance, are constantly moving, and movement can

	Hand segmented	True positive	False positive	False negative
Player 1	13	9	13	4
Player 2	13	7	12	6

Table 1. Comparison of JS and manual segmentation for a ping-pong sequence

occur in many directions. Thus measures like distance between people, or direction of movement, will fail to pick up this interaction. Though the positions of the players are influenced by each other, the statistical correlation between the sequences is not significant since the players do not move in lock-step. Our entropy-based measure detects the interaction when the players are hitting the ball.

Other detected interactions are two people following each other along the corridor and standing/sitting.

4. Conclusion

We presented a system that detects interactions between people in an indoor environment without the use of a pre-defined model of interaction. Detected interactions included people playing ping-pong, which is difficult using only inter-personal distance measures. The recursive segmentation method splits tracks of individuals into sub-tracks each representing a distinct activity; the resulting segments agree well with the results of hand segmentation.

5. Acknowledgement

Helen Yan developed the laser tracking system. This work was supported in part by the DoE RIM grant DE-FG03-01ER45905 and by the ONR MURI grant SA3319.

References

- [1] P. Bernaola-Galvan, R. Roman-Roldan, and J. Oliver. Compositional segmentation and long-range fractal correlations in dna sequences. *Physical Review E*, 53(5):5181–5189, 1996.
- [2] S. Hongeng and R. Nevatia. Multi-agent event recognition. In *IEEE International Conference on Computer Vision (ICCV)*, volume 2, pages 84–91, Vancouver, Canada, July 7-14 2001.
- [3] Y. Ivanov, C. Stauffer, A. Bobick, and W. Grimson. Video surveillance of interactions. In *1999 Computer Vision and Pattern Recognition (CVPR) Workshop on Visual Surveillance, June 23-25, 1999, Fort Collins, Colorado, USA, 1999*.
- [4] H. Yan and M. J. Matarić. General spatial features for analysis of multi-robot and human activities from raw position data. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 2770–2775, EPFL, Switzerland, Oct. 2002.