Anticipating Activity from Multimodal Signals

Tiziana Rotondo¹, Giovanni Maria Farinella¹, Davide Giacalone², Sebastiano Mauro Strano², Valeria Tomaselli², Sebastiano Battiato¹ ¹Department of Mathematics and Computer Science, University of Catania, Italy ² STMicroelectronics, Catania

tiziana.rotondo@unict.it, {davide.giacalone, mauro.strano, valeria.tomaselli} @st.com,

{gfarinella, battiato}@dmi.unict.it

Abstract

Images, videos, audio signals, sensor data, can be easily collected in huge quantity by different devices and processed in order to emulate the human capability of elaborating a variety of different stimuli. Are multimodal signals useful to understand and anticipate human actions if acquired from the user viewpoint? This paper proposes to build an embedding space where inputs of different nature, but semantically correlated, are projected in a new representation space and properly exploited to anticipate the future user activity. To this purpose, we built a new multimodal dataset comprising video, audio, tri-axial acceleration, angular velocity, tri-axial magnetic field, pressure and temperature. To benchmark the proposed multimodal anticipation challenge, we consider classic classifiers on top of deep learning methods used to build the embedding space representing multimodal signals. The achieved results show that the exploitation of different modalities is useful to improve the anticipation of the future activity.

Problem Definition

Problem: Let $\mathbf{y}_t = (\mathbf{v}_t, a_t, \mathbf{s}_t)^T$ be the input vector at time *t* where \mathbf{v}_t is the video signal, \mathbf{a}_t is the audio signal and \mathbf{s}_t is related to other sensor data, we define

ST Multimodal Dataset

The dataset comprises the following modalities: video (collected with a smartphone camera), audio, tri-axial acceleration, angular velocity, tri-axial magnetic





FUTURE Features Temporal

Classification



С

Experimental Results: Baseline vs Siamese Network

These tables summarize results obtained with K-NN and SVM classifier for different values of K and different kernels, respectively. As baseline, the combination of all sensors rate results are reported.

	Classification	Anticipation		
K	Baseline	Baseline	Triplet	

Classification

Anticipation

1	63.24%	64.12%	64.65%
3	62.79%	63.61%	64.73%
5	62.79%	63.14%	64.14%
7	61.97%	62.27%	64.55%
9	61.76%	62.15%	64.18%

Classification		Anticipation			
Baseline	e	Baseline		Triplet	
Linear Kernel	RBF	Linear Kernel	RBF	Linear Kernel	RBF
69.55%	73.75%	64.06%	70.29%	57.52%	63.32%

Conclusion

- Our results suggest that multi-modality improves both classification and prediction.
- Considered activities can be anticipated with an accuracy close to the one obtained when the signals are fully observed (i.e., classification task).
- Future works could be devoted to collect bigger labelled multimodal datasets considering different environments and activities, as well as to model attention mechanisms among the different modalities.