

# 11775 Homework 1 Report

Tianyu Xu - [tianyux@andrew.cmu.edu](mailto:tianyux@andrew.cmu.edu)

In this homework, I basically followed the TA's instructions and performed multimedia event detection using provided support code. The key differences between the support code and my revised code are as follows:

1. The cluster number was changed to 1000.
2. I used a modified opensmile configuration file to perform mean normalization on the extracted MFCC features.
3. I used several different kernels to train the SVM model and compared their performances.

At first I set the cluster number to 4000, but the clustering took more than 64 hours and exceeded the time limit, so I changed it to 1000. Also, before the support code is given, I wrote my own training file generator and my own pipelines and did some experiments. These scripts can be found at `/home/tianyux/old_hw1` and `/home/tianyux/very_old_hw1`.

The `run.feature.sh` script first extracts MFCC features of all the videos, and then `select_frames.py` selects a subset of MFCC features to create a dataset for training k-means cluster. After we got the k-means model, we can get new representations of all the videos. Also, I used ASR transcriptions to generate ASR BoW features, and I concatenate ASR features with MFCC features to perform early feature fusion.

The `run.med.sh` script trains three different SVM models - MFCC model, ASR model, fusion model, respectively. For ASR model, it takes feature vectors of dimension of 983, because the vocabulary size is 983. Similarly, fusion model takes feature vectors of dimension of 1983.

The final experiment results are as follows:

Kernel	Feature Type	P001 mAP	P002 mAP	P003 mAP
RBF	MFCC	0.323356	0.46309	0.252863
	ASR	0.208188	0.230023	0.182163
	Fusion	0.313915	0.235122	0.235122
Chi-square	MFCC	0.566816	<b>0.53295</b>	0.354326
	ASR	0.237749	0.177125	0.404604
	Fusion	<b>0.577842</b>	0.465474	<b>0.487691</b>

All the scripts and testing results can be found at: `/home/tianyux/hw1`