

# Data Mining: Learning from Large Data Sets - Fall Semester 2015

mwurm@student.ethz.ch  
merkim@student.ethz.ch  
lwoodtli@student.ethz.ch

October 12, 2015

## Approximate near-duplicate search using Locality Sensitive Hashing

### Problem Description

The goal of the project was to develop a program that finds duplicate videos. The videos were provided as shingles in a text file. There was also a training set for the video data. The program needed to run on a Map Reduce environment and should be developed with a Locality Sensitive Hashing algorithm (LSH).

### Approach of the Team

The team decided that each member will develop the algorithm independently in a first step. Afterwards the implementations were compared. So the details could be discussed and errors were eliminated. At the end the team decided which implementation should be improved for submission.

### Environment

The running environment with a Map Reduce framework was provided. So the Mapper and the Reducer were developed independently and then submitted to the running environment.

### Map

The Mapper uses a Locality Sensitive Hashing algorithm.

In a first step the provided shingles of each video are min-hashed.

```
# each line contains data for a video
for line in sys.stdin:
    # ...
    # iterate over the shingles of a video
    for shingle in shingles:
```

```

# iterate over the hash functions
for i,hash in enumerate(minHashes):
    hashed = (hash[0]*shingle + hash[1]) % 20000
    # keep hashed value only if it is smaller
    if hashed < signature[i]:
        signature[i] = hashed

```

As a second step the Min-hashes are divided in bands by a Locality Sensitive Hash algorithm.

The resulting Min-hashes per band and video are sent to the Reducer together with the original shingles and sorted by band-id and min-hash value. The original shingles were provided in order to allow the Reducer a real comparison of the potential duplicates.

## Reduce

The Reducer iterates over the sorted input and checks for duplicate min-hash values which indicates potential duplicates. To decrease the amount of false positives the Reducer compares in an additional step all shingles of the potential pairs for 90% equality. This real comparison of the videos allows an improvement of the precision. But it also introduces more work to the Reducer. Therefore it is very important to have an optimal Mapper step which makes sure that no duplicates are missed (leads to false negatives) and in the mean time does not provide too many false positives which need to be sorted out by the Reducer (decrease of efficiency).

This shows us that it is really important to choose the right values for the two main parameters of the Map step.

One of them is the number of hash function to use for calculating the Min-hashes.

The other one is the number of rows in a band for the LSH algorithm.

## Improvement

Missed potential duplicates (false negatives) and wrong potential duplicates (false positives) can weaken the result. To improve the result the two main parameters described above can be tuned.

For this project the parameters were tuned by hand. But it would be possible to tune them automatically for some given test data.

## Possible Improvement

To reduce the risk of outputting a video pair more than once there could be a final step that compares the output of all Reducers and eliminates video pairs that occur more than once. This comparison could be performed in an additional Reduce step.

## Conclusion

This project shows that LSH is a nice approach to find near duplicates in a large set of videos. It can be simply developed in a Map Reduce environment.