

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

caifa.zhou@geod.baug.ethz.ch
pungast@student.ethz.ch
llara@student.ethz.ch

December 16, 2015

Explore-Exploit tradeoś in Recommender Systems

Large Scale Image Classification

To classify the images into the categories nature and people, using a Support Vector Machine as classifier and Parallel Stochastic Gradient Descent procedure, a Map and Reduce function was implemented.

Summary: we used a primal version of the Support Vector Machine (implemented in `sklearn.svm.LinearSVC`) and transformed the original features into Random Fourier Features.

The **mapper** is constructed as follows. First, it reads each line, extracts the features and transforms those. Then it adds the transformed features to the batch and calculates the weight vector on that batch, processes the batch and starts a new one until all input is processed. If the batch size is infinite, each mapper treats its whole input as a single batch.

The particular parameters used in the best submission were: infinite batch size and 800 random features.

A batch of examples is processed by first calculating the weight vectors, with $\eta = \frac{1}{\sqrt{t+1}}$, such that if $yw^Tx < 1$ then update the weights, such that $w'_t \leftarrow w + \eta_t yx$ and set $w_{t+1} = \min \left\{ 1, \frac{1/\sqrt{\lambda}}{\|w'_t\|} \right\}$.

As the last step of the processing of one batch, the emit function prints the weight vector from the batch into standard output.

The **reducer** receives the weight vectors from the mapper. A weight vector is thereby represented by one input line, which is parsed into a vector. Then the reducer calculates the average of all the vectors received from mappers and prints the space-separated coefficients of the model.

The groupwork was done as following. First we sat together and discussed the assignment and overall structure of the project. Taivo Pungas created a first draft of the mapper and reducer, Caifa Zhou continuously contributed to the structure of the mapper and both tested different versions. Lara Lingelbach wrote the report.

Extracting Representative Elements

To extract 100 clusters of representative elements from a large dataset, we used k-means++ and k-means in the mapper and reducer, where the mapper outputs 500 points (cluster centers) for each batch and the reducer finds the final 100 based on the output of the mappers.

The **mapper** is constructed as follows. Data is processed in batches of 10000 data points.

The cluster centers of each stream are initialised using k-means++ to find a solution in reasonable time, and then sequential k-means is run. Weights are assigned to all points, whereby points farther away from the existing centers get a higher weight and thus have a higher probability to be selected. Iteratively the subsequent centers are chosen from the remaining data points with probability proportional to its squared distance ($\|x_i - \mu_j\|_2^2$) from the closest existing cluster center μ_j .

Third, the sequential k-means algorithm is implemented to compute representative elements (μ_j) one at a time. The centers (means) are given by the vector μ with μ_1, \dots, μ_k , the algorithm calculates $\partial L / \partial \mu$ with $\partial L(x, \mu) = \min \|x_i - \mu_j\|_2^2$. If μ_i is closest to x , μ_i is replaced by

$$\frac{1}{k_i} \sum_{j=1}^{k_i} x_j$$

where k_i is the number of data points assigned to i^{th} cluster. Last, the emit function prints the centers into stdout.

The **reducer** receives the cluster centers from all mappers. As in the mapper, the reducer finds runs k-means++ in this reduced space to find initial cluster centers. Sequential k-means algorithm is then run on these initial centers to find the final representative elements.

Approximate near-duplicate search using Locality Sensitive Hashing

To detect duplicate videos a locality sensitive hashing solution using MapReduce was implemented. Below, our solution is described in steps.

The **mapper** reads each line of input, extracts the video ID and shingles, deletes recurring elements in the shingles and sorts the shingles. It then calculates the signature matrix column for each video, partitions it into bands, hashes each band and emits the result.

To calculate the signature matrix columns, we create n hash functions of the form $h_{a,b}(x) = ax + b \bmod N$ (where $N = 20001$) by generating parameter vectors a and b of size n containing random nonnegative integers. Using the MinHash approach and each generated $h_{a,b}$, we calculate the column of the signature matrix corresponding to the current video.

To generate candidate pairs, we partition the signature matrix column into b bands and hash each band with a linear hash function (resulting in a bucket ID). The mapper then emits a key-value pair where the key is a string concatenation (band ID + bucket ID) and the value is a tuple (signature matrix column, video ID). In the **reducer**, all videos that have the same key are then compared pairwise to find out whether the videos were actually similar (with bitwise similarity ≥ 0.9 in the signature columns of the two videos).

Each machine uses the same seed when generating random numbers for the hash functions.

Given the constraint $br \leq 1024$, the parameters b and r were chosen to empirically produce the highest F1-score on the training set, resulting in $(b, r) = (20, 50)$.

The general workflow of the mapper and reducer are illustrated in Figure 1.

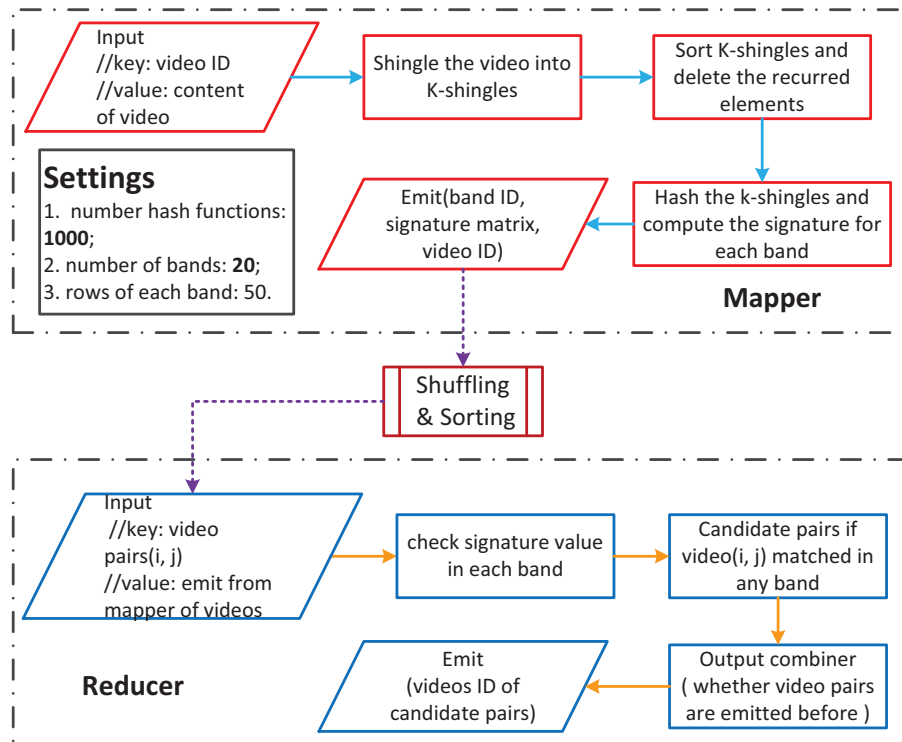


Figure 1: Workflow of the locality sensitive hashing solution using MapReduce.

The groupwork was done according to each member's background. First we sat together and discussed the assignment and overall structure, then Caifa Zhou implemented the mapper, Taivo Pungas created the reducer as well as revised the mapper and Lara Lingelbach generated the report.