

Co-clustering through Optimal Transport

SMAI Project - Team 9

Introduction

- Co-clustering is an unsupervised learning approach that aims to discover homogeneous groups of data instances and features by grouping them simultaneously.
- Many co-clustering methods need number of clusters as input but our method automatically detects the number of co-clusters.
- We use entropy regularized optimal transport between empirical measures defined on data instances and features in order to obtain an estimated joint probability density function represented by the optimal coupling matrix.
- This matrix is further factorized to obtain the induced row and columns partitions using multiscale representations approach.
- The algorithm derived for the proposed method and its kernelized version based on the notion of Gromov-Wasserstein distance are fast, accurate and can determine automatically the number of both row and column clusters.

Optimal Transport

- Optimal transportation theory was first introduced to study the problem of resource allocation.
- Assuming that we have a set of n factories and a set of n mines, the goal of optimal transportation is to move the ore from mines to factories in an optimal way, i.e., by minimizing the overall transport cost.
- If M is the set of mines and F is the set of factories. *Transport plan* is a bijection $T : M \rightarrow F$. In other words, each mine $m \in M$ supplies precisely one factory $T(m) \in F$ and each factory is supplied by precisely one mine. For optimal transport plan, we need to minimize

$$c(T) := \sum_{m \in M} c(m, T(m))$$

Optimal Transport

- More formally, given two empirical probability measures $\hat{\mu}_S = \frac{1}{N_S} \sum_{i=1}^{N_S} \delta_{x_i^S}$ and $\hat{\mu}_T = \frac{1}{N_T} \sum_{i=1}^{N_T} \delta_{x_i^T}$ defined as uniformly weighted sums of Dirac with mass at locations supported on two point sets $X_S = \{\mathbf{x}_i^S \in \mathbb{R}^d\}_{i=1}^{N_S}$ and $X_T = \{\mathbf{x}_i^T \in \mathbb{R}^d\}_{i=1}^{N_T}$, the Monge-Kantorovich problem consists in finding a probabilistic coupling γ matrix defined as a joint probability measures $X_S \times X_T$ with marginals $\hat{\mu}_S$ and $\hat{\mu}_T$ that minimizes the cost of transport w.r.t. some $l : X_S \times X_T \rightarrow \mathbb{R}^+$:

$$\min_{\gamma \in \Pi(\hat{\mu}_S, \hat{\mu}_T)} \langle M, \gamma \rangle_F$$

Where $\langle \cdot, \cdot \rangle_F$ is the Frobenius dot product

Optimal Transport

- $\Pi(\hat{\mu}_S, \hat{\mu}_T) = \{\gamma \in \mathbb{R}_+^{N_S \times N_T} | \gamma \mathbf{1} = \hat{\mu}_S, \gamma^T \mathbf{1} = \hat{\mu}_T\}$ is a set of doubly stochastic matrices and M is a dissimilarity matrix i.e , $M_{ij} = l(\mathbf{x}_i^S, \mathbf{x}_j^T)$ defining the energy needed to move a probability mass from \mathbf{x}_i^S to \mathbf{x}_j^T
- This problem admits a unique solution γ^* and defines a metric on the space of probability measures (called the Wasserstein distance) as follows:

$$W(\hat{\mu}_S, \hat{\mu}_T) = \min_{\gamma \in \Pi(\hat{\mu}_S, \hat{\mu}_T)} \langle M, \gamma \rangle_F$$

- The success of algorithms based on this distance is also due to (Cuturi, 2013) who introduced an entropy regularized version of optimal transport that can be optimized efficiently using matrix scaling algorithm.

Entropic Regularization

- Entropic regularization found its application to the optimal transportation problem through the following objective function:

$$\min_{\gamma \in \Pi(\hat{\mu}_S, \hat{\mu}_T)} \langle M, \gamma \rangle_F - \frac{1}{\lambda} E(\gamma)$$

- Second term $E(\gamma) = -\sum_{i,j}^{N_S, N_T} \gamma_{i,j} \log(\gamma_{i,j})$ is the entropy

Advantages :

- Allows to obtain smoother and more numerically stable solutions compared to the original case and converges to it at the exponential rate.
- It allows to solve optimal transportation problem efficiently using Sinkhorn-Knopp matrix scaling algorithm which can be parallelised.

Why Optimal Transport for Co-Clustering ???

We pose the co-clustering problem as the task of transporting the empirical measure defined on the data instances to the empirical measure defined on the data features. The intuition behind this process is very natural to co-clustering and consists in capturing the associations between instances and features of the data matrix. And the solution of optimal transportation problem can be decomposed into 3 terms as

$$\gamma_{\lambda}^* = \text{diag}(\alpha)\xi_{\lambda}\text{diag}(\beta)$$

where ξ is the approximated joint probability distribution and remaining are diagonal matrices which can be approximated to distributions of data instances and features.

Why Optimal Transport ? - Mathematical View

The optimal co-clustering is taken as the one that minimizes the difference in mutual information between the original random variables and the mutual information between the clustered random variables i.e., $I(X,Y) - I(X',Y')$. The loss in mutual information can be expressed as the distance of $p(X; Y)$ to an approximation $q(X; Y)$ i.e., $KL(p(x,y)/q(x,y))$ where $q(x,y)$ is of the form

$$q(x,y) = p(x/x')p(x',y')p(y/y').$$

This solution is similar to the solution we obtain in Optimal Transport using Sinkhorn Algorithm. If we assume $p(x',y')$ follows Gibbs Distribution then $p(x',y')$ can be approximated to ξ .

Co Clustering - Problem Setup

- Let X and Y be two random variables taking values in the sets $\{x_r\}_{r=1}^n$ and $\{y_c\}_{c=1}^d$, respectively, where subscripts r and c correspond to rows (instances) and columns (features).
- Let the joint probability distribution between X and Y be denoted by $p(X, Y)$ is estimated from the data matrix $A \in \mathbb{R}^{n \times d}$. We further assume that X and Y consist of instances that are distributed w.r.t. probability measures μ_r, μ_c supported on Ω_r, Ω_c where $\Omega_r \subseteq \mathbb{R}^d$ and $\Omega_c \subseteq \mathbb{R}^n$ respectively.
- The problem of co-clustering consists in jointly grouping the set of features and the set of instances into homogeneous blocks by finding two assignment functions C_r and C_c that map as follows: $C_r : \{x_1, \dots, x_n\} \rightarrow \{\hat{x}_1, \dots, \hat{x}_g\}$, $C_c : \{y_1, \dots, y_d\} \rightarrow \{\hat{y}_1, \dots, \hat{y}_m\}$ where g and m denote the number of row and columns clusters, and discrete random variables \hat{X} and \hat{Y} represent the partitions induced by X and Y , i.e., $\hat{X} = C_r(X)$ and $\hat{Y} = C_c(Y)$.

$$\hat{\mu}_r = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i} \text{ and } \hat{\mu}_c = \frac{1}{m} \sum_{i=1}^m \delta_{\mathbf{y}_i}$$

Optimal Transportation - Proposed Approach

- We use optimal transportation to find a probabilistic coupling of the empirical measures defined based on rows and columns of a given data matrix. More formally, for some fixed $\lambda > 0$ we solve the co-clustering problem through the following optimization procedure:

$$\gamma_{\lambda}^* = \operatorname{argmin}_{\gamma \in \Pi(\hat{\mu}_r, \hat{\mu}_c)} \langle M, \gamma \rangle_F - \frac{1}{\lambda} E(\gamma) \quad --(1)$$

- Matrix M is computed using the Euclidean distance, i.e., $M_{ij} = \|x_i - y_j\|^2$. The elements of the resulting matrix γ_{λ}^* provides us with the weights of associations between instances and features: similar instances and features correspond to higher values in γ_{λ}^* . Our intuition is to use these weights to identify the most similar sets of rows and columns that should be grouped together to form co-clusters.

Optimal Transportation - Proposed Approach

- Following (Benamou et al., 2015), this optimization problem can be rewritten as:

$$\min_{\gamma \in \Pi(\hat{\mu}_r, \hat{\mu}_c)} \langle M, \gamma \rangle_F - \frac{1}{\lambda} E(\gamma) = \frac{1}{\lambda} \min_{\gamma \in \Pi(\hat{\mu}_r, \hat{\mu}_c)} \text{KL}(\gamma \| \xi_\lambda)$$

- where $\xi_\lambda = e^{-\lambda M}$ is the Gibbs kernel. Finally, we can rewrite the last expression as follows:

$$\min_{\gamma \in \Pi(\hat{\mu}_r, \hat{\mu}_c)} \text{KL}(\gamma \| \xi_\lambda) = \min_{\gamma \in \mathcal{C}} \text{KL}(\gamma \| \xi_\lambda)$$

- where $\mathcal{C} = \mathcal{C}_1 \cap \mathcal{C}_2$ is the intersection of closed convex subsets given by $\mathcal{C}_1 = \{\gamma \in \mathbb{R}^{d \times d} \mid \gamma \mathbf{1} = \hat{\mu}_r\}$ and $\mathcal{C}_2 = \{\gamma \in \mathbb{R}^{d \times d} \mid \gamma^T \mathbf{1} = \hat{\mu}_c\}$. The solution of the entropy regularized optimal transport can be obtained using Sinkhorn-Knopp algorithm and has the following form

$$\gamma_\lambda^* = \text{diag}(\alpha) \xi_\lambda \text{diag}(\beta)$$

where α and β are the scaling coefficients of the Gibbs kernel ξ_λ

Sinkhorn-Knopp algorithm

Notation :

For two probability vectors r and c

$$U(r, c) := \{P \in \mathbb{R}_+^{d \times d} \mid P\mathbf{1}_d = r, P^T\mathbf{1}_d = c\}.$$

We propose a solution for

$$\text{For } \lambda > 0, d_M^\lambda(r, c) := \langle P^\lambda, M \rangle, \text{ where } P^\lambda = \underset{P \in U(r, c)}{\operatorname{argmin}} \langle P, M \rangle - \frac{1}{\lambda} h(P).$$

Sinkhorn-Knopp algorithm

Algorithm 1 Computation of $\mathbf{d} = [d_M^\lambda(r, c_1), \dots, d_M^\lambda(r, c_N)]$, using Matlab syntax.

Input $M, \lambda, r, C := [c_1, \dots, c_N]$.

$I = (r > 0); r = r(I); M = M(I, :); K = \exp(-\lambda M)$

$u = \text{ones}(\text{length}(r), N) / \text{length}(r);$

$\tilde{K} = \text{bsxfun}(@\text{rdivide}, K, r) \text{ \% equivalent to } \tilde{K} = \mathbf{diag}(1./r)K$

while u changes or any other relevant stopping criterion **do**

$u = 1./(\tilde{K}(C./(K'u)))$

end while

$v = C./(K'u)$

$\mathbf{d} = \text{sum}(u.*((K.*M)v))$

Kernelized version and Gromov-Wasserstein distance

- We'll look at a kernelized version of our method. We first define two similarity matrices $K_r \in \mathbb{R}^{n \times n}$ and $K_c \in \mathbb{R}^{d \times d}$ associated to empirical measures $\hat{\mu}_r, \hat{\mu}_c$.
- Matrices K_r and K_c are defined by calculating the pairwise distances or similarities between rows and columns, respectively, without the restriction of them being positive or calculated based on a proper distance function satisfying the triangle inequality.
- The entropic Gromov-Wasserstein discrepancy is defined as follows:

$$\begin{aligned} \text{GW}(K_r, K_c, \hat{\mu}_r, \hat{\mu}_c) &= \min_{\gamma \in \Pi_{\hat{\mu}_r, \hat{\mu}_c}} \Gamma_{K_r, K_c}(\gamma) - \lambda E(\gamma) \\ &= \min_{T \in \Pi_{\hat{\mu}_r, \hat{\mu}_c}} \sum_{i,j,k,l} L(K_{r_{i,j}}, K_{c_{k,l}}) \gamma_{i,j} \gamma_{k,l} - \lambda E(\gamma). \end{aligned}$$

- where γ is a coupling matrix between two similarity matrices and $L : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+$ is an arbitrary loss-function, usually the quadratic-loss or Kullback-Leibler divergence.

Kernelized version and Gromov-Wasserstein distance

- Based on this definition, one may define the problem of the entropic Gromov-Wasserstein barycenters for similarity or distance matrices K_r and K_c as follows:

$$\min_{K, \gamma_r, \gamma_c} \sum_{i=\{r,c\}} \varepsilon_i \Gamma_{K, K_i}(\gamma_i) - \lambda E(\gamma_i) \quad \text{----(2)}$$

- K is the computed barycenter and $\gamma_r \in \Pi_{\hat{\mu}, \hat{\mu}_r}$, $\gamma_c \in \Pi_{\hat{\mu}, \hat{\mu}_c}$ are the coupling matrices that align it with K_r and K_c , respectively.
- ε_i are the weighting coefficients summing to one, i.e., $\sum_{i=\{r,c\}} \varepsilon_i = 1$ that determine our interest in more accurate alignment between K_r and K or K_c and K .
- In (1) we align rows with columns directly, in (3) our goal is to do it via an intermediate representation given by the barycenter K that is optimally aligned with both K_r and K_c .

Kernelized version and Gromov-Wasserstein distance

- In this case, we obtain the solutions γ_r and γ_c that, similar to previous method and can be decomposed as follows:
- $\gamma_r^* = \text{diag}(\alpha_r)\xi_r \text{diag}(\beta_r)$, $\gamma_c^* = \text{diag}(\alpha_c)\xi_c \text{diag}(\beta_c)$, where $\xi_r = e^{-\lambda M_r}$ and $\xi_c = e^{-\lambda M_c}$ are Gibbs kernels calculated between the barycenter and row and column similarity matrices using any arbitrary loss-function L as explained before.
- Finally, based on the analysis presented above, we further use vectors β_r and β_c to derive row and column partitions.

Computing GW barycenters

Notation :

The simplex of histograms with N bins is $\Sigma_N \stackrel{\text{def.}}{=} \{p \in \mathbb{R}_+^N ; \sum_i p_i = 1\}$.

The entropy of $T \in \mathbb{R}_+^{N \times N}$ is defined as $H(T) \stackrel{\text{def.}}{=} -\sum_{i,j=1}^N T_{i,j}(\log(T_{i,j}) - 1)$

The set of couplings between histograms $p \in \Sigma_{N_1}$ and $q \in \Sigma_{N_2}$ is

$$\mathcal{C}_{p,q} \stackrel{\text{def.}}{=} \{T \in (\mathbb{R}_+)^{N_1 \times N_2} ; T\mathbf{1}_{N_2} = p, T^\top \mathbf{1}_{N_1} = q\}$$

For any tensor $\mathcal{L} = (\mathcal{L}_{i,j,k,l})_{i,j,k,l}$ and matrix $(T_{i,j})_{i,j}$, we define the tensor-matrix multiplication as

$$\mathcal{L} \otimes T \stackrel{\text{def.}}{=} \left(\sum_{k,\ell} \mathcal{L}_{i,j,k,\ell} T_{k,\ell} \right)_{i,j}.$$

Computing GW barycenters

If the loss function L can be written as

$$L(a, b) = f_1(a) + f_2(b) - h_1(a)h_2(b)$$

for functions (f_1, f_2, h_1, h_2) , then, for any $T \in \mathcal{C}_{p,q}$,

$$\mathcal{L}(C, \bar{C}) \otimes T = c_{C, \bar{C}} - h_1(C)Th_2(\bar{C})^\top. \quad (6)$$

where $c_{C, \bar{C}} \stackrel{\text{def.}}{=} f_1(C)p\mathbf{1}_{N_2}^\top + \mathbf{1}_{N_1}q^\top f_2(\bar{C})^\top$ is independent of T .

For square loss $L = L_2$, $f_1(a) = a^2$, $f_2(b) = b^2$, $h_1(a) = a$, $h_2(b) = 2b$ satisfies

For KL loss $L = \text{KL}$, $f_1(a) = a \log(a) - a$, $f_2(b) = b$, $h_1(a) = a$, $h_2(b) = \log(b)$ satisfies

Computing GW barycenters

Algorithm 1 Computation of GW_ε barycenters.

Input: $(C_s, p_s)_{s,p}$

Initialize C .

repeat

// minimize over $(T_s)_s$

for $s = 1$ **to** S **do**

Initialize T_s .

repeat

// compute $c_s = \mathcal{L}(C, C_s) \otimes T_s$ using (6).

$c_s \leftarrow f_1(C) + f_2(C_s)^\top - h_1(C)T_s h_2(C_s)^\top$

// Sinkhorn iterations (3) to compute $\mathcal{T}(c_s, p, q)$

Initialize $a \leftarrow \mathbb{1}$, set $K \leftarrow e^{-c_s/\varepsilon}$.

repeat

$a \leftarrow \frac{p}{Kb}, b \leftarrow \frac{q}{K^\top a}$.

until convergence

Update $T_s \leftarrow \text{diag}(a)K \text{diag}(b)$.

until convergence

end for

// minimize over C using (13).

$C \leftarrow \left(\frac{f'_1}{h'_1} \right)^{-1} \left(\frac{\sum_s \lambda_s T_s^\top h_2(C_s) T_s}{pp^\top} \right)$

until convergence

If $(C_s)_s$ are positive semi-definite matrices

using L^2 loss ,eq(13) becomes

$$C \leftarrow \frac{1}{pp^\top} \sum_s \lambda_s T_s^\top C_s T_s.$$

Using KL loss, eq (13) becomes

$$C \leftarrow \exp \left(\frac{1}{pp^\top} \sum_s \lambda_s T_s^\top \log(C_s) T_s \right)$$

Jump Detection - Detecting number of Clusters

- In order to detect the steps (or jumps) in the approximated marginals, we use a method described in Matei & Meignen, 2012 for multiscale denoising of piecewise smooth signals.
- It determines the significant jumps in the vectors α and β without knowing their number and location. As the proposed procedure deals with non-decreasing functions, we first sort the values of α and β in the ascending order. Since the procedure is identical for both vectors, we only describe it for the vector α .
- We consider that the elements $\{\alpha_i\}_{i=1}^n$ of α are the local averages of a piecewise continuous function $v : [0, 1] \subset \mathbb{R} \rightarrow \mathbb{R}$ on the intervals $I_i^n = [i/n, (i + 1)/n]$ (defined by the uniform subdivision R of step $1/n$ of the interval $[0, 1]$).
- More precisely: $\alpha_i^n = n \int_{I_i} v(t)dt, \quad i = 0, \dots, n - 1.$

Jump Detection

- The detection strategy is based on the cost function: $F(I_i^n) = \sum_{l=i-1}^i |a_{l+1}^n - a_l^n|$ defined for each interval. Therefore, we get the list of the interval suspicious to contain a jump for the subdivision of order n as follows:

$$L^n = \{i^* ; i^* = \operatorname{argmax}_i F(I_i^n)\}$$

- This detection should be refined in order to get only significant jumps in our vector a . To this end we use the multi-scale representation of a as in (Harten, 1989) and we perform this detection on each scale. On the first scale, we get a coarse version of a by averaging:

$$a_{i}^{n/2} = 0.5*(a_{2i}^n + a_{2i+1}^n), \quad i=0, \dots, n/2 - 1.$$

Jump Detection

- Now, by considering the coarse version of α , we obtain a second list $L^{n/2}$ of suspicious intervals as before.
- After that, these two lists merge in the list L jumps as follows: a jump will be considered in the interval I_{2i}^n or I_{2i+1}^n if the interval $I_i^{n/2}$ is also detected as suspicious at the coarse scale. This procedure is iterated $\lceil \log_2 n \rceil$ times and a jump is observed if a chain of detection exists from fine to coarse scales. Finally, the number of clusters is obtained by

$$g = |L_{\text{jumps}}| + 1.$$

Algorithmic implementation

Algorithm 1 Co-clustering through Optimal Transport (CCOT)

Input : \mathcal{A} - data matrix, λ - regularization parameter, n_s - number of sampling

Output: C_r, C_c - partition matrices for rows and columns, g, m - number of row and column clusters

$[n, d] = \text{size}(\mathbf{Z})$

for $i = 1$ **to** n_s **do**

$\mathbf{D}_i = \text{datasample}(\mathbf{Z}, d)$

$\mathbf{M}_i \leftarrow \text{pdist2}(\mathbf{D}_i, \mathbf{D}_i^T)$

$[\boldsymbol{\alpha}_i, \boldsymbol{\beta}_i, \gamma^*] \leftarrow \text{optimal_transport}(\mathbf{M}_i, \lambda)$

$[L_{\text{jumps}}^{\alpha_i}, \mathbf{C}_r^i, g] \leftarrow \text{jump_detection}(\text{sort}(\boldsymbol{\alpha}_i))$

$[L_{\text{jumps}}^{\beta_i}, \mathbf{C}_c^i, m] \leftarrow \text{jump_detection}(\text{sort}(\boldsymbol{\beta}_i))$

$\mathbf{C}_r \leftarrow \text{mode}(\mathbf{C}_r^i)$

$\mathbf{C}_c \leftarrow \text{mode}(\mathbf{C}_c^i)$

$$C_r^i(\mathbf{x}_r) = \begin{cases} 1, & r \leq L_{\text{jumps}}^{\alpha_i}(1) \\ \dots & \\ k, & L_{\text{jumps}}^{\alpha_i}(k-1) < r \leq L_{\text{jumps}}^{\alpha_i}(k) \\ \dots & \\ |L_{\text{jumps}}^{\alpha_i}| + 1, & r > L_{\text{jumps}}^{\alpha_i}(|L_{\text{jumps}}^{\alpha_i}|). \end{cases}$$

Algorithmic Implementation

Algorithm 2 Co-clustering through Optimal Transport with Gromov-Wasserstein barycenters (CCOT-GW)

Input : \mathcal{A} - data matrix, λ - regularization parameter, $\varepsilon_r, \varepsilon_c$ - weights for barycenter calculation

Output: C_r, C_c - partition matrices for rows and columns, g, m - number of row and column clusters

$K_r \leftarrow \text{pdist2}(Z, Z)$

$K_c \leftarrow \text{pdist2}(Z^T, Z^T)$

$[\beta_r, \beta_c, \gamma_r^*, \gamma_c^*] \leftarrow \text{gw_barycenter}(K_r, K_c, \lambda, \varepsilon_r, \varepsilon_c)$

$[L_{\text{jumps}}^{\beta_r}, C_r, g] \leftarrow \text{jump_detection}(\text{sort}(\beta_r))$

$[L_{\text{jumps}}^{\beta_c}, C_c, m] \leftarrow \text{jump_detection}(\text{sort}(\beta_c))$

Experiment

- We used Movie-Lens dataset which is a popular benchmark data set that consists of user-movie ratings, on a scale of one to five, collected from a movie recommendation service gathering 100,000 ratings from 943 users on 1682 movies.
- Our goal is to find homogeneous subgroups of users and films in order to further recommend previously unseen movies that were highly rated by the users from the same group.
- We set the regularisation parameters for the CCOT algorithm as $\lambda = 0.001$ and for CCOT-GW algorithm as $\varepsilon = (0.5, 0.5)$.

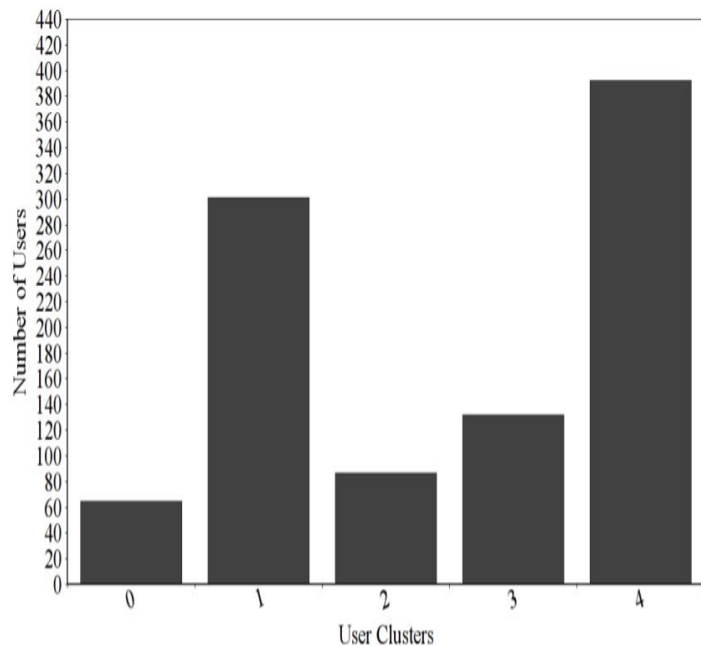
Results

MC5	MC6
Dead man Walking	Diabolique
Golden Eye	All Dogs go to Heaven 2
Usual Suspects	Theodore Rex
CopyCat	Sgt.Bilko

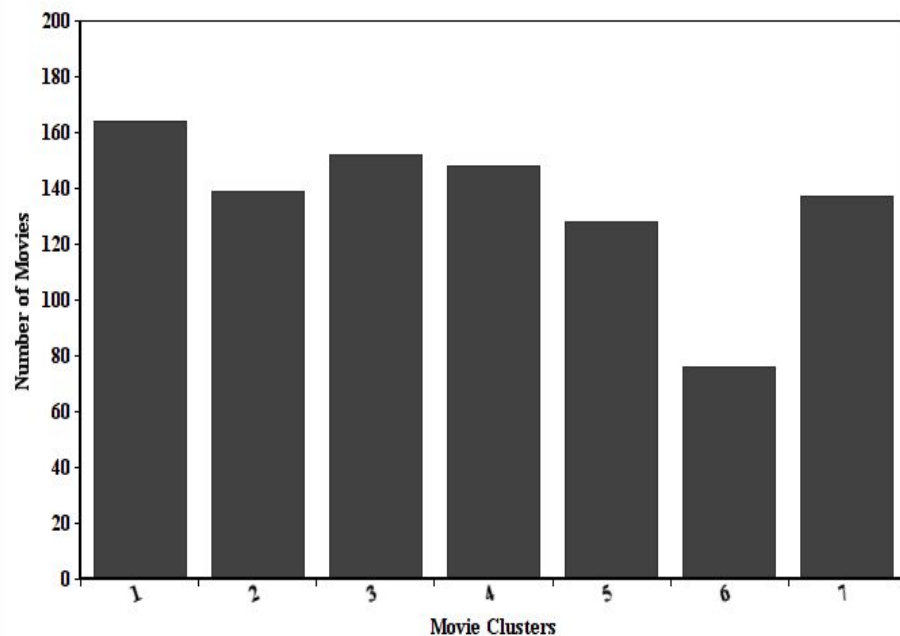
In the results we got the movies in the top 4 rated movies in the 5th and 6th clusters are shown. The movies in 5th cluster are similar in type of genre(all are Thriller Movies) and that of in 6th cluster consists of movies which are less critically acclaimed.

Distribution of Users and Movies in various Clusters

Distribution of Users



Distribution of Movies



Source Code

The Source code is available at <https://github.com/tarungavara/smai-project>

Thank You