



Standard Recognition Computational Work

Academic Year 2022 - 2023

Γ. Tsihrintzis D. Sotiropoulos

Subject: Recommendation Algorithms Using Artificial Neural Networks and Data Clustering Techniques

This work aims to develop algorithms for movie recommendation using artificial neural networks and clustering techniques. The dataset on which you will work can be downloaded from <https://ieee-dataport.org/open-access/imdb-movie-reviews-dataset> and concerns on a set of users $U = \{U_1, U_2, \dots, U_N\}$ who express preferences over a set of movie objects $I = \{I_1, I_2, \dots, I_N\}$. The user $u \in U$ attributes to the movie $i \in I$ the preference degree $R(u, i) \in R_o = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\} \cup \{0\}$, where 0 indicates the absence of an evaluation and not a zero preference score.

The set $X = \{x_1, x_2, \dots, x_T\}$ of available data consists of records of the form $x_k = [u_k, i_k, r_k, t_k]$ where $u_k \in U$ and $i_k \in I$ with $r_k = R(u_k, i_k) \in R_o$. The value t_k corresponds to the date of entry of this record.

Data pre-processing:

1. Find the set of unique users U and the set of unique objects \mathbb{I} .
2. Consider the function $\Phi: U \rightarrow P(I)$ (where $P(I)$ is the dynamo set of \mathbb{I}) which $\forall u \in U$ returns the set $\varphi(u) \subset I$ of objects evaluated by user \mathbb{I} . We can therefore write that $\varphi(u) = \{i \in I: R(u, i) > 0\}$. Restrict the sets of unique users U and unique objects I to the respective sets \hat{U} and \hat{I} such that $\forall u \in \hat{U}, R_{min} \leq |\varphi(u)| \leq R_{max}$ where R_{min} and R_{max} is the minimum required and maximum allowed number of evaluations per user. Consider obviously that $|\hat{U}| = n < N$ and $|\hat{I}| = m < \mathbb{I}$.
3. Create and graphically represent frequency histograms for the number and time range of each user's ratings.
4. Create an alternative representation of the dataset as a set of preference vectors $R = \{R_1, R_2, \dots, R_n\}$ with $R_j = R(u_j) \in R_o^m, \forall j \in [n]$. In particular, we can write that:

$$R_j = \begin{cases} R(u_j, i_1), R(u_j, i_2), \dots, R(u_j, i_m) & \text{if } i_1, i_2, \dots, i_m \in \varphi(u_j) \\ 0, & \text{if } i_1, i_2, \dots, i_m \notin \varphi(u_j) \end{cases}$$

Data Grouping Algorithms

1. Organize the limited set of users \hat{U} into L clusters of the form $\hat{U} = U_{G_1} \cup U_{G_2} \cup \dots \cup U_{G_L}$ such that $U_{G_a} \cap U_{G_b} = \emptyset, \forall a \neq b$ based on the vector representation of their preferences via the set \mathbb{I} . The process of clustering the user preference vectors $R_j = R(u_j) \in R_o^m, \forall j \in [n]$ onto the restricted set of objects \mathbb{I} , can be performed by appropriately parameterizing the **k-means** algorithm by varying the metric that

assesses the distance between two preference vectors for a pair of users u and v as $\text{dist}(R_u, R_v)$. In particular, you can use the following metrics (*):

a. ~~euclidean~~ $\text{dist}(R_u, R_v) = \sqrt{\sum_{k=1}^m |R_u(k) - R_v(k)|^2 \lambda_u(k) \lambda_v(k)}$ (2) such that.

$$\lambda_j(u) = \begin{cases} 1, & \text{if } \lambda_j(u) > 0 \\ 0, & \text{if } \lambda_j(u) = 0 \end{cases} \quad (3)$$

The function $\lambda_j()$ evaluates whether or not user u_j has evaluated object i_k . Obviously, through this function we can if we calculate the value of $|\varphi(u)|$ which evaluates the number of evaluations provided by user u in total as: $|\varphi(u)| = \sum_{k=1}^n \lambda_u(k)$.

b. ~~cosine~~ $\text{dist}(R_u, R_v) = 1 - \frac{\sum_{k=1}^m \lambda_u(k) \lambda_v(k) R_u(k) R_v(k)}{\sqrt{\sum_{k=1}^m \lambda_u(k) R_u(k)^2} \sqrt{\sum_{k=1}^m \lambda_v(k) R_v(k)^2}}$ (4)

- Graphically represent the clusters of users identified by the k-means algorithm for each of the above metrics for various values of the parameter β .
- Comment on the effectiveness of these metrics in assessing the similarity between a pair of user preference vectors R_u and R_v .

* The metrics represented in equations (2) and (4) are computed for each pair of vectors R_u and R_v only for the subset of objects evaluated by both users u and v via the utility function $\lambda()$.

Recommendation Generation Algorithms Using Artificial Neural Networks

- Create an alternative organization of the restricted set of users into L clusters $\hat{U} = U_{G_1} \cup U_{G_2} \cup \dots \cup U_{G_L}$ such that $U_{G_a} \cap U_{G_b} = \emptyset, \forall a, b \in [L]$ with $a \neq b$ using the following metric (**):

$$\text{dist}(u, v) = 1 - \frac{|\varphi(u) \cap \varphi(v)|}{|\varphi(u) \cup \varphi(v)|} \quad (5)$$

** The clustering algorithm for dealing with this query is of your choice. The basic requirement, however, is that it can operate on the square matrix of distances between users described by relation (5).

- Explain what this metric expresses and identify its disadvantages compared to the metrics described in relations (2) and (4) in the light of the particular organisation of users that it may entail.

- b) The metric described by relation (5) can be used in order to determine the set $N(u) = \{u^{(1)}, u^{(2)}, \dots, u^{(k)}\}$ of k nearest neighbors of a user $u_a \in U_G$, $\forall a \in [L]$. Therefore, for each user u_a of each cluster U_{G_a} we can set the vector of personal preferences R_{u_a} as well as the vectors of k nearest neighbors of u_a within cluster U_{G_a} . Objective of this question is to develop a multilayer neural network for each user cluster U_{G_a} , $\forall a \in [L]$ which will approximate the ratings of each user within it by the ratings of its k nearest neighbors via a function of the form:

$$r_{u_a} = f_a(r_{u^{(1)}}, r_{u^{(2)}}, \dots, r_{u^{(k)}}), \forall u_a \in U_G, \forall a \in [L] \quad (6)$$

- c) The set of users in each cluster can be further partitioned into a subset of users for training the neural network U^{train} and a subset of users for testing the performance of the neural network U^{test} such that $U_G = U^{train} \cup U^{test}$ with $U^{train} \cap U^{test} = \emptyset$. The most appropriate configuration of the data for the training of this neural network of the question could be reproduced as follows:

Suppose that the set of users participating in the cluster U_{G_a} is given by $U_{G_a} = \{u_{a,1}, u_{a,2}, \dots, u_{a,n_a}\}$ with $n_a = |U_{G_a}|$. Then the set of feature vectors and the corresponding set of labels could be organized into a table of the form:

$$\begin{bmatrix} Ru(1) & Ru(2) & \dots & Ru(k) \\ a,1 & r_{a,1} & & r_{a,1} \\ Ru(1) & Ru(2) & \dots & Ru(k) \\ a,2 & r_{a,2} & & r_{a,2} \\ \vdots & \vdots & \vdots & \vdots \\ Ru(1) & Ru(2) & \dots & Ru(k) \\ r_{a,n_a} & r_{a,n_a} & & r_{a,n_a} \end{bmatrix} \text{ and } \begin{bmatrix} r_{a,1} \\ r_{a,2} \\ \vdots \\ r_{a,n_a} \end{bmatrix}$$

- d) The approximate accuracy of these neural networks can be measured by the metric of the mean absolute error between the actual and estimated user ratings. Present tables of your results for both training and testing accuracy for each user cluster.

You can work in groups of a maximum of 3 students. The implementation of the project can be done in MATLAB or Python. The code of your project should be accompanied by detailed documentation.

Good luck!